

UC Riverside

UC Riverside Electronic Theses and Dissertations

Title

Computational Analysis of Receptor-Odor Interactions and Prediction of Behavior-Modifying Chemical Space

Permalink

<https://escholarship.org/uc/item/5j67m8w8>

Author

Boyle, Sean Michael

Publication Date

2012

Peer reviewed|Thesis/dissertation

UNIVERSITY OF CALIFORNIA
RIVERSIDE

Computational Analysis of Receptor-Odor Interactions and Prediction of
Behavior-Modifying Chemical Space

A Dissertation submitted in partial satisfaction
of the requirements for the degree of

Doctor of Philosophy

in

Genetics, Genomics, and Bioinformatics

by

Sean Michael Boyle

March 2012

Dissertation Committee:

Dr. Anandasankar Ray, Chairperson
Dr. Anupama Dahanukar
Dr. Thomas Girke

Copyright by
Sean Michael Boyle
2012

The Dissertation of Sean Michel Boyle is approved:

Committee Chairperson

University of California, Riverside

Acknowledgments

Personal Acknowledgements

I would like to acknowledge my parents Laura and Mike Boyle for their constant support and encouragement throughout both my education and life. They have always been a wonderful source of guidance and support while traversing the endeavors of the journey we call life. I would also like to acknowledge the daily guidance, humor, and inspiration provided by my advisor Dr. Anandasankar Ray. He has always been a constant source of wonderful ideas, almost tangible intensity toward science, and at times very needed off topic tangents. I would also like to acknowledge all past and present members of both the Ray and Dahanukar labs for their endless inquisitiveness, humor, and good nature. Without all of you this work would not be so complete and would certainly not have been so enjoyable.

Professional Acknowledgements

- Anandasankar Ray, Anupama Dahanukar, and Thomas Girke for providing research guidance.
- Shane McNally for validating *Drosophila* receptor-odor predictions using electrophysiology.
- Chih-Ying Su, Eliza Kelly-Swift, and John Carlson for validating *Anopheles* receptor-odor predictions using electrophysiology.
- Sana Tharadra for performing *Drosophila* larvae behavior experiments.
- Christine Pham for performing *Drosophila* larvae behavior experiments.
- Dyan MacWilliam for performing CO₂ receptor and repellency validations using electrophysiology.
- Tom Guda for performing Mosquito behavior experiments.
- Will Pitt and Tom Blundell for assisting with structure-based virtual screening.
- Sachin Suardre for validating predicted EthR inhibitors.

ABSTRACT OF THE DISSERTATION

Computational Analysis of Receptor-Odor Interactions and Prediction of Behavior-Modifying Chemical Space

by

Sean Michael Boyle

Doctor of Philosophy, Graduate Program in
Genetics, Genomics, and Bioinformatics
University of California, Riverside, March 2012
Dr. Anandasankar Ray, Chairperson

Coding of information in the peripheral olfactory system and resulting olfactory dependent behavior is thought to depend primarily on two fundamental factors: the interaction of individual odors with different subsets of the odor receptor (Or) repertoire, and the mode of signaling that an individual receptor-odor interaction elicits, activation or inhibition. In order to better understand these processes, we design and implement a structure-based virtual screening approach that identifies common structural features that are highly correlated with odor activity for individual receptors. We then apply these features to rapidly screen for putative ligands *in silico* from a large untested odor space (>240,000 putative volatiles) for the majority of odor receptors in the *Drosophila* antenna, allowing for analysis of odor coding for the majority of receptors for the first time. Functional experiments support a high success rate (~71%) for the screen and we validate numerous new activators and inhibitors for the receptors. Following our initial application in *Drosophila*, we extend our approach to predict activating and inhibiting odors for a large number of important pest and disease vector species including 50 *Anopheles gambiae* Ors (65% validated accuracy), the CO₂ receptors of multiple species (48% validated accuracy), 9 newly identified Citrus Psyllid ORNs, and a large number of functionally distinct mammalian receptors. We next extended our *in silico*

screening approach to identify shared structural features important for a behavioral response to DEET-like repellents for which the molecular target has not yet been identified, identifying ~150 natural compounds as candidate repellents. We select 4 candidates, 3 approved as safe for human food use, and demonstrate that they are strong olfactory and gustatory repellents to both mosquitoes and *Drosophila*. As only a small region of odor space has been explored, there remains potential to uncover previously unidentified patterns of odor coding. Through a combination of *in silico* and electrophysiology screens, we identify odors with ultra-prolonged termination kinetics that are delayed for several minutes, resulting in a memory trace that affects subsequent odor detection. Finally, we successfully perform structure-based virtual screen, identifying potential inhibitors of an important Tuberculosis drug target EthR.

Table of Contents

CHAPTER I: Introduction

- Complexity in the olfactory system p1
- Only a very small portion of the vast odor space has been tested for activity p6
- Or gene identification in insects was largely aided by computational approaches p8
- Applying ligand-based virtual screening to predict receptor-odor interactions p12
- Applying structure-based virtual screening to predict receptor-odor interactions p16
- High-throughput approaches are effective tools for olfactory analysis p18

CHAPTER II: Designing a Ligand-Based Virtual Screening Approach to Decode *Drosophila* Odor Receptor Chemical Space *In Silico*

- Introduction p25
- Results p27
- Discussion p37
- Figures p42
- Tables p60

CHAPTER III: Applying Chemical Informatics to Decode Odor Receptors of Several Important Disease Vector and Pest Insect Species as Well as Mammals

- Introduction p71
- Results p76
- Discussion p95
- Figures p99
- Tables p141

CHAPTER IV: A New Generation of Safe DEET Substitutes that are Strong Olfactory and Gustatory Repellents of Mosquitoes and Flies

- Introduction p201
- Results and Discussion p204
- Figures p214

CHAPTER V: Analyzing Termination Dynamics of Prolonged Activating Odors and their Effects on Receptor-Mediated Olfactory Behavior

- Introduction	p227
- Results	p230
- Discussion	p241
- Figures	p245
- Tables	p265

CHAPTER VI: Applying Structure-Based Virtual Screening to Identify Potent Inhibitors of the Tuberculosis Target EthR

- Introduction	p266
- Results and Discussion	p268
- Figures	p281
- Tables	p301

CHAPTER VII: Methods

- Chemical informatics	p302
- Electrophysiology	p302
- Natural odor compound library	p303
- Pubchem compound library	p303
- eMolecules compound library	p303
- Calculation of 3D conformations	p303
- Calculation of molecular descriptors	p304
- Classification of active compounds	p304
- Calculation of Accumulative Percentage of Actives (APoA)	p305
- Determination of optimized descriptor subsets	p305
- Clustering Ors by most common descriptors	p306
- Clustering compounds by activity of Or	p306
- Clustering Ors by predicted ligand space	p307
- Calculation of Or prediction distribution frequencies	p307
- Or-ligand interaction map	p307
- Computational validation of ligand-based virtual screening (Non-SVM)	p308
- Calculation of LogP and vapor pressure values	p309
- Repellency behavior testing	p309
- Classification of repellent compounds	p311
- Support Vecotor Machine (SVM) predictions	p311
- Computational validation of ligand-based virtual-screening (SVM)	p311
- Vinyl solubility assay	p312
- Olfactory avoidance assay trap assay for <i>Drosophila</i>	p312
- T-maze Assay Methods	p313
- Modified hand in-glove olfactory repellency assay for mosquitoes	p313
- Humidity and warmth attraction assay	p315
- Larval behavior assays	p315
- Assembly of EthR screening library	p316
- Calculating molecular descriptors for EthR analysis	p316

- Applying machine learning with Pipeline Pilot p317
- Performing structure-based virtual screening with Gold p317

REFERENCES p319

List of Figures

CHAPTER I:

- 1.1: Overview schematic describing neuronal wiring of the insect olfactory system p24

CHAPTER II:

- 2.1: Comparing efficacy of various structure analysis methods to analyze odor receptor ligands p42
- 2.2: A receptor-optimized molecular descriptor approach has strong predictive power to find new ligands p44
- 2.3: Accumulated percentage of actives analysis p46
- 2.4: Electrophysiology validates that odorant receptor-optimized molecular descriptors can successfully identify new ligands for *Drosophila* p48
- 2.5: Predicted receptor-odor interactions are highly specific p50
- 2.6: High-ranking odors are more likely to be active than distantly ranked odors p52
- 2.7: Analysis of receptor-odor relationships and breadth of tuning p54
- 2.8: Analysis of receptor-natural odor interactions p56
- 2.9: Predicted odor space and network view of odor coding p58

CHAPTER III:

- 3.1: A molecular descriptor optimized approach is able to explain odor activity for individual *Anopheles* Ors p99
- 3.2: Optimized molecular descriptor sets are able to cluster either active or inhibitory odors p101
- 3.3: Structural relationships between the 75 most diverse aromatic odors p107
- 3.4: Newly optimized molecular descriptor sets are able to cluster odors for aromatically tuned Ors p109
- 3.5: Optimized descriptors for AgOrs 2, 6, and 10 effectively describe training set activity p111
- 3.6: A Support Vector Machine (SVM) integrated approach is highly effective at explaining odor activity for individual Ors p113
- 3.7: SVMs trained using optimized descriptors for AgOrs 2, 6, and 10 effectively describe training set activity p115
- 3.8: A molecular descriptor optimized approach is able to explain odor activity for CO₂ receptors p117
- 3.9: Optimized molecular descriptor sets are able to cluster odors by CO₂ receptor response p119
- 3.10: Active compounds cluster into three distinct structural classes p121
- 3.11: Chemical structures of validated activators and inhibitors p123
- 3.12: A SVM integrated molecular descriptor optimized approach is able to explain odor activity for individual Citrus Psyllid ORNs p125
- 3.13: Optimized molecular descriptor sets are able to cluster ORN activators p127
- 3.14: SVMs trained using descriptor sets that were optimized for individual ORNs effectively describe training set activity p129
- 3.15: A molecular descriptor optimized approach is able to explain odor activity

for individual mammalian ORs	p131
3.16: Mammalian Odorant receptor-optimized molecular descriptors can successfully cluster known ligands	p133
3.17: Optimized descriptors for mammalian ORs effectively describe training set activity	p135
3.18: Analyzing relationships between important features for mammalian ORs	p137
3.19: Analysis of mammalian OR tuning breadth	p139

CHAPTER IV:

4.1: Contribution of olfaction and gustation in DEET avoidance	p214
4.2: Mosquito behavioral assay glove setup	p216
4.3: A chemical informatics method to predict repellency	p218
4.4: Identification of repellents using in-silico screening of a large chemical space	p220
4.5: A new class of mosquito repellents with desirable safety profiles	p222
4.6: Mosquito escape index	p224
4.7: Natural compounds are effective at repelling aedes aegypti in the hand in glove assay	p226

CHAPTER V:

5.1: Activity of a behaviorally important neuron, Ab1A, can be described by an odor receptor neuron-optimized chemical informatics approach	p245
5.2: Electrophysiology validates that ORN-optimized molecular descriptors can successfully identify new ligands for ab1A	p247
5.3: Functional identification of Ultra-Prolonged activators and analysis of their long-term effects on ab1A using electrophysiology	p249
5.4: Ectopic expression confers Or42b UP activation	p251
5.5: Behavioral effects of Ultra-Prolonged activators on odor detection in ab1A	p253
5.6: Behavioral effects of inhibitory odor on odor detection in ab1A	p255
5.7: Functional identification of Ultra-Prolonged activators of additional odor receptors using electrophysiology	p257
5.8: Identification of ultra-prolonged activating odors for additional Ors	p259
5.9: Long-term effects of Ultra-Prolonged activators on odor detection	p261
5.10: Modes of signaling and their behavioral responses	p263

CHAPTER VI:

6.1: Schematics of Ethionamide Activation Pathway	p281
6.2: Solved EthR structures	p283
6.3: The shape of the EthR cavity	p285
6.4: Previously identified EthR ligands bind within the proposed cavity	p287
6.5: Visualizing previously identified EthR ligands	p289
6.6: Single example of successful EthR self-docking	p291
6.7: The results of many EthR site 1 self-docking runs	p293
6.8: The predicted orientation of our strongest identified inhibitors in EthR	p295
6.9: The predicted orientations of 4 modest inhibitors of EthR	p297
6.10: The predicted orientations of the poor inhibitors of EthR	p299

List of Tables

CHAPTER II:

- 2.1: Optimized descriptor sets for each *Drosophila* Or p60
2.2: Top 100 predicted compounds for each *Drosophila* Or p68

CHAPTER III:

- 3.1: Optimized descriptor sets for each *Anopheles* Or p141
3.2: Natural odor library predictions found in the top 500 predictions for each *Anopheles* Or p152
3.3: Predicted odors validated as effective activators for several *Anopheles* Ors p160
3.4: Activity of a large panel of aromatic odors was tested against three aromatically tuned Ors p162
3.5: Optimized descriptor sets for Or2, Or6, and Or10 p164
3.6: Top 75 predicted compounds for each *Drosophila* Or p166
3.7: Dividing training odors into three distinct sets based upon odor structure and receptor response p168
3.8: Optimized descriptors selected for the Aromatic Activator Screen p170
3.9: Optimized descriptors selected for the Broad Activator p172
3.10: Optimized descriptors selected for the Inhibitor Screen p174
3.11: Top predicted natural library compounds for the aromatic activator, broad activator, and inhibitor screens p176
3.12: Predicted odors validated as activators and inhibitors of the CO2 receptor p178
3.13: Optimized descriptor sets for each Citrus Psyllid ORN p180
3.14: Top predicted natural library compounds for Citrus Psyllid ORNs p184
3.15: Optimized descriptor sets for each mammalian OR p188
3.16: Top 100 predicted compounds for each mammalian OR p195

CHAPTER V:

- 5.1: Optimized descriptor sets for ab1A (Or42b) p265

CHAPTER VI:

- 6.1: A breakdown of our structure-based virtual screening accuracy for EthR p301

CHAPTER I: INTRODUCTION

Complexity in the olfactory system

The olfactory system of insects is a very large and intricate neural network, consisting of a peripheral input layer, a middle processing layer, and high level wiring into the brain (Figure 1.1). Odors are detected at the periphery and contextual information about the odors are aggregated and organized in the middle layer. Here additional stages of processing are performed and the information is then provided to the brain for interpretation and decision-making.

The periphery olfactory system in *Drosophila* is housed in two olfactory organs: the antennae and maxillary palps. Both of these organs are covered with many tiny hairs called sensilla, which house between 1 and 4 distinct Olfactory Receptor Neurons (ORNs), however the predominant number of sensilla house 2 (de Bruyne et al., 2001). Interestingly, ORNs generally express only a single odor selective Odor receptor (Or) along with an obligate co-receptor Orco, which is essential for olfaction (Couto et al., 2005; Larsson et al., 2004; Vosshall and Hansson, 2011). *Drosophila* express 60 Or genes, which are translated into 62 proteins, that are housed in ORNs (Robertson et al., 2003). Both Or gene expression by ORNs and housing of ORNs in the antenna and maxillary palps is highly conserved, allowing for sensilla types to be classified according to what ORNs they house and ORNs to be classified by the Or genes they express (Couto et al., 2005; Dobritsa et al., 2003). While the peripheral olfactory system is highly structured and conserved, it also represents a very complex network of gene expression.

It is the response of the Or to environmental odors that determines whether an ORN sends information, in the form of either neuronal activation or inhibition, to the higher centers in the brain (Dobritsa et al., 2003; Hallem et al., 2004). Or responses to odors have been demonstrated to range from total inhibition of the ORN to strong activation (>250 spikes/sec) and each Or is uniquely tuned to odors (Hallem and Carlson, 2006). While some Ors are narrowly tuned and respond to only 1 or a few known ligands, such as DmOr67d to the pheromone cVa, others are more widely tuned and respond to a broad array of both functional groups and molecule sizes (Hallem and Carlson, 2006; van Naters and Carlson, 2007). Interestingly, an individual odor may activate multiple receptors, with each receptor being tuned to a different range of odor concentration (Hallem and Carlson, 2006).

While originally identified in mammals, Ors of insects are unique. Mammalian ORs were first discovered in 1991, resulting in award of the prestigious Nobel prize (Buck and Axel, 1991). Insect Ors were discovered for the important model organism *Drosophila* nearly a decade later in 1999 (Clyne et al., 1999; Vosshall et al., 1999). Being extremely divergent from mammalian ORs, the signaling pathway of insect Ors is currently a hotly contested topic. While mammalian ORs function as true GPCRs, Insect receptors have several unique characteristics, including an inverse orientation in the neuronal membrane and the previously mentioned obligate co-receptor Orco (Belluscio et al., 1998; Benton et al., 2006). Even the phylogenetic similarities between insect Ors are far more divergent than for mammalian ORs (Robertson et al., 2003). It was this dissimilarity between insect Ors and mammalian ORs that was responsible for the 8 additional years that were required for the identification of the first insect Ors, as traditional wet lab techniques attempting to identify genes using mammalian OR

sequence knowledge were highly ineffective and computational approaches designed on the hypothesis of extreme divergence were required (Clyne et al., 1999). Recently two groups have proposed two alternative, albeit slightly overlapping, hypotheses on insect Or function. While one group proposes they function as ligand gated ion channels, another contests that they are both ligand gated and G-protein mediated cyclic-nucleotide-activated cation channels (Kain et al., 2008; Sato et al., 2008; Wicher et al., 2008). It is possible their function is an aggregate of both, as there is compelling evidence supporting each hypothesis.

Chemosensation is not limited to Ors. Identification of the Ors in *Drosophila* paved the way for identification of the Gustatory receptor (Gr) family, which was discovered only one year following Or discovery (Clyne et al., 2000). While gustatory receptors are predominately involved in the chemosensation of chemicals in liquid medium, a pair of Gr receptors (Gr21a and Gr63a) expressed in *Drosophila* is interestingly responsible for detection of CO₂, which is a highly volatile odor (Dahanukar et al., 2001; Kwon et al., 2007). Even more intriguingly, both Gr21a and Gr63a are amongst the rare Grs expressed in the antenna, which almost exclusively houses Or expressing ORNs. As in *Drosophila*, the function of Grs is currently a debated topic including evidence for both GPCR like function and cation channel activity (Ishimoto et al., 2005; Sato et al., 2011; Yao and Carlson, 2010).

Several years after identification of the Or and Gr gene family, a third, more divergent, chemosensory family of Ionotropic Glutamate Receptors (iGluRs) was discovered in *Drosophila* (Benton et al., 2009). While these receptors are related to previously identified known families of glutamate receptors, they appear to be highly

divergent. IRs are expressed in coeloconic sensilla and appear to be tuned to acids, ammonia, and humidity (Abuin et al., 2011; Yao et al., 2005). IRs are believed to function in a heteromeric complex and conduct signals as ion channels (Abuin et al., 2011).

Neuronal wiring of the olfactory system in *Drosophila* is highly complex, containing multiple levels of processing. As discussed above, the peripheral system contains many chemosensory inputs, consisting of three unique classes of receptors, each with their own chemsensory detection spectrum and possibly unique function. The axonal extensions of ORNs project into the antennal lobe, with each ORN class projecting to a unique glomerular location (Fishilevich et al. 2005, Couto et al., 2005). Information processing in the antennal lobe has also been a contested topic, as research groups have debated over the role of local interneurons (LNs). While it was originally claimed that information was faithfully transmitted to higher brain centers without manipulation, it has since been demonstrated that both excitatory and inhibitory local interneurons play an important role in signal processing prior to transmission (Olsen et al., 2007; Olsen and Wilson, 2008; Root et al., 2007; Shang et al., 2007; Wang et al., 2003; Wilson and Laurent, 2005; Wilson et al., 2004). The main reason for the discrepancy of the two differing claims lies in experimental technique. The Wang research group attempted to tease apart the biological question using calcium imaging, which was not as effective in identifying interneuron communication as Rachel Willson's electrophysiological approach. These different experiments provide a wonderful basis for the importance of teasing apart a biological question from multiple angles to obtain insight that may otherwise be missed. After processing in the antennal lobe, olfactory information is sent to higher brain regions by Projection Neurons (PNs), where it is finally

processed in the mushroom body and lateral horn (Jefferis et al., 2007). It is important to note there are both narrowly and broadly wired neuronal systems. While the majority of ORNs appear to provide information to highly interconnected networks of LNs, a few specialized ORNs form direct channels to PNs (Jefferis et al., 2007). These direct line systems have been termed labeled line and are often observed with pheromone detecting ORNs. The entire 11-cis-vaccenyl acetate (cVA) pheromone circuit was beautifully traced from sensory input all of the way to descending output in what is perhaps the best single analysis of olfactory neuronal wiring in *Drosophila* to date, providing an illustrative example for labeled line systems (Ruta et al., 2010). However, these systems appear to be the exception rather than the rule, as processing for the vast majority of ORN input is far too complex for current experimental approaches to have teased apart.

To summarize the system as a whole, the olfactory system is a massive, highly interconnected, and complex neuronal network. It contains inputs from many periphery Ors/ORNs, processes the information in the antennal lobe, and projects processed information into higher brain regions. Luckily for olfactory investigation, the system also has characteristics lending itself to analysis, including the general expression of only a single odor sensitive Or in each ORN, glomerular aggregation of ORN input, and neuronal activity from different ORNs housed in the same sensilla being distinguishable by spike amplitude. While these challenges are considerable, major advances have been made in the field, many of which have been aided by computational approaches.

Only a very small portion of the vast odor space has been tested for activity

It has been estimated that the number of possible carbon based chemical configurations with similar molecular masses to biologically relevant chemicals could exceed 10^{60} unique structures (Bohacek et al., 1996; Dobson, 2004). To put this number into context, it is roughly twice the number of stars than are believed to exist in the universe. Considering that the number of unique chemical structures found in the human body is likely in the hundreds of thousands, this number of potential compounds is staggering (Goto et al., 2002). While volatile chemical space is undoubtedly considerably smaller than 10^{60} , it still likely contains hundreds of thousands of unique compounds, each of which have the potential to activate or inhibit Ors, IRs, or Grs.

Take into consideration that the number of unique Ors expressed by each species has thus far varied from XX on the small end of the scale to over a thousand in mice and each receptor responds independently to a given odor. Clearly, the number of possible Or-odor combinations is staggering. Furthermore, positive classification of a particular odor or blend, such as one produced by a piece of ripe fruit, by the brain requires the integrated response of all expressed Ors to each individual odor in the blend. As one principle goal of olfactory research is to understand how organisms perceive their chemical environment, the aim of understanding the responses of each Or expressed in a particular species to all odors from a large number of behaviorally important source is incredibly challenging, yet still important.

In the 12 years of olfactory research since the discovery of *Drosophila* Ors, 251 unique odors have been experimentally tested for activity in the *Drosophila* olfactory system (de Bruyne et al., 1999; de Bruyne et al., 2001; Dobritsa et al., 2003; Goldman et

al., 2005; Hallem and Carlson, 2006; Hallem et al., 2004; Kreher et al., 2005; Kreher et al., 2008; Kwon et al., 2007; Pelz et al., 2006; Schmucker et al., 2007; Stensmyr et al., 2003; Turner and Ray, 2009; van Naters and Carlson, 2007; Yao et al., 2005).

Responses to these screens have recently been assembled and integrated into a single database (Galizia et al., 2010). While some screens have been narrow and focused on testing a small number of odors against a specific Or, others have involved testing a panel of odors across multiple Ors. For example, In 2001 47 odors were tested for activity against 16 ORN classes (de Bruyne et al., 2001). In the largest screen to date 109 odors were tested for activity against 24 Ors expressed in the antenna of *Drosophila* adults using single unit electrophysiology (Hallem and Carlson, 2006). Additionally, 21 Ors expressed in *Drosophila* larvae were tested for activity against 27 odors in 2008 (Kreher et al., 2008). While these three analyses certainly do not represent all of the analyses performed in *Drosophila*, they do represent the largest bodies of work.

Smaller scale olfactory experimentation has also been performed in other insect and mammalian species. While there is great need for an increased understanding of host or food preference in the many important disease vector and pest species, relatively few experiments have tested a sizeable set of odors for activity against multiple receptors. In the largest screen of any insect species outside of *Drosophila*, 110 odors were recently individually tested for activity against 50 *Anopheles* Ors using single unit electrophysiology (Carey et al., 2010). In the largest single mammalian screen a combination of 62 mouse and human ORs were individually exposed to 63 odors and measured for activity using calcium imaging (Saito et al., 2009). The limited number of screened odors tested on a very narrow set of species makes it difficult to tease apart important olfactory cues and pathways that are responsible for attraction, mating, etc.

The odors tested in many screens are meticulously selected, both to cover a broad range of chemical space and for ecological importance. For example, many odors tested in *Drosophila* were selected due to their previous identification in rotting and ripe fruit, an important food source, and a number of odors tested in the *Anopheles* screen performed by Allison Carey were selected for their previously recognized importance as human host volatiles (Carey et al., 2010; Hallem and Carlson, 2006). Odors that were not selected for ecological significance were chosen to cover a broad area of chemical space. These odors cover an expansive range of carbon chain lengths and contain a broad assortment of functional groups including ketones, aldehydes, carboxylic acids, alcohols, and esters. While the majority of screened odors are aliphatic, a number of aromatic odors are usually included as well. These analyses provide the groundwork from which to further expand upon explored odor space, with the ultimate goal of finding effective and species specific behavior modifying chemicals.

Or gene identification in insects was largely aided by computational approaches

Linda Buck correctly made three very important hypotheses in her identification of the OR gene family (Buck and Axel, 1991). Firstly, ORs were likely to belong to a superfamily of proteins involved in GPCR cellular communication. Secondly, the *OR* gene family should be very large and diverse to recognize so many odors. Thirdly, *OR* expression should be restricted to the olfactory epithelium. Guided by these hypothesis, Linda used degenerate primers of known GPCRs to amplify olfactory epithelium cDNA, resulting in the historic identification of band 13 and subsequently the first *OR* gene family (Buck and Axel, 1991).

As previously stated, the Or gene families of Insects are highly divergent from mammalian Ors. While the same approach was applied to identify the *Drosophila* Or gene family, they were unsuccessful. By 1999 bioinformatics approaches were emerging as successful and important tools for biological analysis. Additionally, the first draft of the *Drosophila* genome was underway, providing a ripe time for a computationally guided approach to Or gene family identification (Adams et al., 2000). As conditions were now suitable, two research groups independently identified the *Drosophila* Or gene family in tandem.

Peter Clyne and Coral Warr from John Carlson's lab teamed up with Junhyong Kim and applied bioinformatics approaches to identify 7-transmembrane proteins from the incomplete 1998 *Drosophila* genome build (Clyne et al., 1999; Kim and Carlson, 2002). Recent evidence had shed light that insect Ors would likely share structural features to GPCRs, yet as previous wet lab attempts had been unsuccessful in identification, it was likely these genes were highly divergent from mammalian ORs. Guided by this insight the team created an algorithm to identify potential members of a large 7TM gene family from the roughly 10% completed *Drosophila* genome build. The algorithm began by identifying Open Reading Frames (ORFs) that were > 300 bases. They then identified important sequence descriptors for 7TM GPCRs, testing 70 parameters for the ability to separate 750 training GPCRs from 1000 non-GPCR proteins by sequence analysis alone. The 5 most important parameters were then applied to screen all predicted ORFs of >300 bases, alternative splicing was considered in order to obtain full length 7TM protein sequences, and the top scoring hits that were not previously known GPCRs were selected. RT-PCR primers were designed for the final list of hits and 2 genes were identified that were only expressed in the antennae. The

well known sequence search tool BLAST was next applied, searching for ORFs with sequence similarity to the two potential Or genes, resulting in identification of 16 Or genes shown to be expressed solely in the antennae and maxillary palps (Altschul et al., 1990). A year later the same approach was applied to identify *Drosophila* Gr genes (Clyne et al., 2000; Kim and Carlson, 2002).

In an independent analysis Leslie Vosshall in Richard Axel's lab identified 10 *Drosophila* Or genes using a combination of wet lab and computational approaches (Vosshall et al., 1999). Vosshall applied difference cloning to identify antennal/maxillary palp specific cDNAs, one of which (dor104) was then applied to scan the *Drosophila* genome for transmembrane proteins using the computational programs Dense Alignment Surface (DAS) and TMAP. 10 of the 12 genes with the highest BLAST determined sequence similarity, albeit very low, to dor104 were demonstrated to be expressed solely in either the antennae or maxillary palps.

Initial insect Or gene discovery has led to an explosion of computationally guided Or identification and analysis. All 60 Or and 68 Gr genes in *Drosophila* had been identified within two years of the initial discovery through iterative application of BLAST on the fully completed *Drosophila* genome (Robertson et al., 2003). Since then, the Or gene families of 12 *Drosophila* species have now been identified, allowing for a comparative analysis identifying both gene conservation prior to subgenera division roughly 40 Mya and providing interesting signs of species specific positive selection (Guo and Kim, 2007). In addition to *Drosophila*, computation has aided in Or family discovery of many other species, including both pest and disease vectors. Examples of species in which Or families have been identified include *Anopheles gambiae*,

Acyrtosiphon pisum (Pea Aphid), *Aedes aegypti*, *Nasonia vitripennis* (Jewel Wasp), *Rhagoletis pomonella* (Apple Maggot), *Bombyx mori* (Silk Moth), *Danaus plexippus* (Monarch Butterfly) and several others (Bohbot et al., 2007; Fox et al., 2001; Robertson et al., 2003; Schwarz et al., 2009; Smadja et al., 2009; Zhan et al., 2011). In a slightly different direction, computational approaches have also allowed for analysis of Or genes across distantly related species, such as was performed for the CO₂ receptors for 12 *Drosophila*, 3 mosquito, the silk moth, and flour beetle species (Robertson and Kent, 2009), providing insight on gene evolution on a very important behavior-regulating receptor.

Or gene family identification for all of these species was made possible by computational advances in genome sequencing, assembly, and search methods able to identify distantly related genes. Without these approaches, identification of the highly divergent Or and Gr families, would certainly have taken far longer to complete, if at all. As sequencing methods continue to improve, the genomes of additional species will begin pouring in, resulting in an increased importance for bioinformatics in Or and Gr gene identification. In turn, the increased volume of Or gene information will further open the door to cross species comparisons. While many of the same tools will likely be applied for Or and Gr gene identification, many new tools will be required for analysis of sequence-activity relationships. Once motifs and protein regions important for odor selectivity are identified, an entirely new field will be available. Traditionally, 3D structures of receptors are required for ligand docking and binding prediction. However, as the Or gene family is highly divergent while likely still retaining core structural features, key motifs or binding regions on Or sequences may be responsible for odor selectivity. As approaches improve it may be possible to computationally predict

receptor-odor interactions from sequence alone. As Or are the largest gene family in several species, this will provide both large challenges and incredible rewards if successful.

Applying ligand-based virtual screening to predict receptor-odor interactions

The field of chemical informatics has been successfully applied to drug discovery by the pharmaceutical industry for years. With the incredible number of chemicals available and the cost required for purchasing and screening, an approach that can computationally identify chemicals with desired features is incredibly useful for initial compound selection (Maldonado et al., 2006). Many distinct chemical informatics approaches have been designed to explain the relationships shared between molecules. Maximum Common Substructure (MCS) determines the largest substructure that is shared between two molecules (Cao et al., 2008). Atom Pairs (AP) identifies atom types with similar distances that are shared between molecules (Carhart et al., 1985). Another approach Molecular Fingerprints (MFs) calculates bit strings, where each 0 or 1 represents the presence or absence of a particular feature, such as a functional group (Nikolova and Jaworska, 2004). Some commercially available programs are available that calculate thousands of individual molecular descriptors for thousands of query compounds, providing quantitative information on many aspects of a compound, such as molecular weight, carbon atom number, predicted vapor pressure, or complex 3-dimensional relationships between atoms of a molecule. Large virtual chemical libraries, which can contain hundreds of thousands or millions of chemicals, can be screened for desired characteristics by applying these approaches.

Interestingly, the relatively newly developing field of chemical informatics has only been minimally applied to olfaction. There are several reasons why chemical informatics approaches analyzing receptor-odor interactions from the odor (ligand) side should be of high priority to the olfactory community. As a large portion of drugs currently on the market target GPCRs and insect receptors are hypothesized to have GPCR like function and structure, it is suggestive that these techniques should apply well to olfactory ligand prediction. This is especially true for mammalian ORs, as they are GPCRs. Additionally, the structures of Ors have not been solved. As Odor receptors are membrane bound and in the case of insect receptors pair with an obligate co-receptor, obtaining x-ray crystal structures suitable for computational docking will be challenging. Finally, we currently have all of the tools required for chemical informatics analysis. Diverse subsets of odors have been tested for activity against several species, resulting in large number of decoded receptors in both mammals and several important insect species.

While chemical informatics approaches are promising, additional challenges unique to olfaction must be overcome. Odors are far smaller than chemicals that are traditionally investigated by the pharmaceutical industry (Wishart et al., 2008). While the reduction in size may make some processes, such as 3D structure optimization, more straightforward, it may also make structural feature differentiation more challenging. Olfactory researchers will need to identify chemical informatics approaches that lend themselves well to the smaller number structural features separating odors than would be typically found between larger pharmaceutical chemicals. Additionally, Odors are volatile while the majority of pharmaceutical targets are not. Volatility could effect how well odors reach their receptors targets and should to be taken into account. The few

attempts at chemical informatics application to olfactory research have been fairly successfully and broadly accepted.

In 2000 Jan Kaluza and Heinz Breer demonstrated that carbon chain length of aliphatic aldehydes was an important feature for the prediction of binding (Kaluza and Breer, 2000). Compounds with similar lengths were more likely to bind to the same OR and just as importantly changing the length of a chain by a single carbon could entirely abolish activity. While this may now appear obvious to the community, prior to large panel screens this feature of odor selectivity had not been documented. It was the first support for Olfactory receptors that similar chemical shape may result in similar activity.

Seven years later Michael Schmucker attempted the first prediction of receptor-odor activity (Schmucker et al., 2007). Schmucker et al. applied Artificial Neural Networks (ANNs) to predict the responses of 21 odors against 7 ORNs. The group calculated ~200 molecular descriptors from MOE (Chemical Computing Group, Montreal) for 47 training odors and 21 test odors. They selected descriptors that were best able to separate the active from inactive molecules and applied them to train the ANN. While the approach was well performed, it unfortunately was only marginally accurate with an accuracy of ~25% (percentage of predicted test odors that activated >50 spikes/sec), which may be due to either molecular descriptor selection, ANN training, or application of ANNs instead of other available approaches.

In 2008 Rafi Haddad from Noam Sobel's laboratory identified molecular descriptors that were correlated among odorants. Haddad et al. calculated 1,664 molecular descriptor values for each odors tested in 7 experimental studies, containing both mammalian and insect receptors (Haddad et al., 2008). The group then selected

which descriptors were the most correlated with activity, resulting in a final set of 32 molecular descriptors that best explained general and non-receptor specific olfactory activity. These 32 descriptors explained activity of the 7 datasets better than carbon atom number or the full set of descriptors. While this approach does not attempt to explain receptor specific selectivity, it was a step in the correct direction.

Most recently in 2010, Sheng Guo and Junhyong Kim applied their own Quantitative Structure Activity Relationship (QSAR) technique to explain the activity of Elissa Hallems analysis where 108 odors were previously tested against 24 *Drosophila* Ors (Guo and Kim, 2010; Hallem and Carlson, 2006). Their model mapped 3D features of each odor and identified which features were most important for receptor activity. Additionally, they performed a sequence analysis for each of the 24 receptors, hypothesizing that the binding pocket of *Drosophila* Ors is on the extracellular halves of its Trans Membrane (TM) domains. The model suggests that the binding pocket is 15 angstroms deep and 6 angstroms wide.

Each of these approaches provides compelling evidence that chemical informatics can be successfully and beneficially applied to olfactory research. However, none of the approaches is successful in explaining receptor specific activity, which should be a major goal of the computational olfaction community. Testing odors using current experimental approaches take a great deal of time and the purchase of a large number of chemicals, many of which may be ineffective. Both of these drawbacks are very expensive. An easy to apply tool that is highly effective at predicting receptor-odor interactions would be immensely beneficial to the olfactory community, allowing for intelligent and quick prioritization of odors for experimental validation. As the number

species for which an initial set of odors has been tested against is currently increasing rapidly, the olfactory field is ripe for such an approach.

Applying structure-based virtual screening to predict receptor-odor interactions

The field of structural-based virtual screening, which involves docking the 3D structure of a chemical into a solved protein structure *in silico*, is regularly used in the field of drug discovery (Kitchen et al., 2004; Nikolova and Jaworska, 2004). Both pharmaceutical companies and academic institutions can have large research groups dedicated to docking libraries of untested chemicals onto the 3D protein structures of potentially drugable targets. The most promising hits are usually selected for experimental validation. The benefits of this system over the ligand-based molecular descriptor prediction techniques are that the resulting dockings can be visually inspected to check that the interaction appears valid, verifying that shape constraints and important polar and non-polar interactions are satisfied. Additionally, manual inspection can reveal important characteristics that can improve structures of current dockings, such as the addition of an atom to fully fill a cavity or addition of an oxygen atom to satisfy a structurally important hydrogen bond. Both ligand-based and structure-based virtual screening can be integrated into a single highly effective pipeline. Molecular descriptors can be applied to select ligands with the most promising characteristics and selected compounds can then be virtually docked into a protein-binding site.

Structure-based virtual screening has been successfully applied for a limited number of mammalian Ors. While structures for insect or mammalian odor receptors have not yet been reported, mammalian receptors are structurally related to rhodopsin,

which has been structurally solved. By applying the structural knowledge from rhodopsin to mammalian Ors research groups have been able to predict a basic structure of mammalian receptors including Rat Or5, Mouse S25, and Rat I7 (Floriano et al., 2000; Singer, 2000; Singer and Shepherd, 1994). The first virtual screening tied to a site directed mutagenesis for functional validation was performed for the mouse mOR-EG in 2005 (Katada et al., 2005). Katada et al. were able to identify the location of an odorant-binding site, as well as how the odor eugenol specifically interacted with it, for the first time, representing a major step forward for the mammalian olfactory community. Since the first successful demonstration of molecular modeling of mammalian ORs, several computational groups have performed similar analyses in additional receptors, including mOR147-9, hOR17-210, and rat OR-I7, without wet lab experimental validation (Khafizov et al., 2007; Kurland et al., 2010; Lai et al., 2008).

In the largest and most complex OR virtual screening attempt to date, 758 compounds from the CAP database were computationally screened for activity against the 5.24 receptor from goldfish, which is activated by all 20 naturally occurring amino acids (Triballeau et al., 2008). Both ligand-based and structure-based virtual screening was applied, reducing the number of potential ligands to 46. Activities of the top 4 predictions were validated using an electro-olfactogram (EOG), each of which produced activation.

Interestingly, structure-based virtual screening has loosely been applied to the human taste receptor TAS2R46 (Brockhoff et al., 2010). Protein chimeras between receptors hTAS2R46 and hTAS2R31 were analyzed to identify important amino acid binding regions, classifying the C-terminus as an important region for ligand selectivity.

Focusing their attention on the extracellular side of trans-membrane spanning alpha helices 6 and 7, they swapped amino acid residues to identify the location of the ligand-binding site. Once a region was identified, a structure-based virtual screen was performed to predict the exact binding site, important interacting residues, and ligand orientation in the membrane bound receptor. This effective blending of wet lab experimental and computational analyses hypothesized the exact orientation of Strychnine on the extracellular side of hTAS2R46.

The field of protein structure determination is rapidly advancing, allowing for determination of structures that were previously untenable. As methodologies continue to advance, our ability to solve membrane bound receptors with heteromeric partners will become possible in the near future, paving the way for structure-based virtual screening of the structurally elusive insect Ors. As the sequences of insect Ors are highly divergent, an additional challenge will be to model the structures of many divergent Ors based upon the few initially solved Insect Ors. As these structures become available, a combined approach involving an initial ligand-based screen of a very large chemical space followed by a focused structure-based virtual screen will likely provide accurate predictions of receptor activations for hundreds of thousands of odors, which would take an unknown number of years to test using wet lab experimentation alone.

High-throughput approaches are effective tools for olfactory analysis

While virtual screening is proving to be an effective tool for identification of receptor-odor interactions, expression analyses are also highly aided by computational approaches. Recent advances in next generation analysis tools, such as microarrays,

RNA-seq, and chip-seq have provided wonderful new ways to understand gene regulation. What would have been impossible just a few years ago is quickly become standard experimental procedure today. Since the initial discovery of the structure and organization of DNA in 1953, a strong charge has been made to understand as much about the genetic makeup of as many species as possible (Watson and Crick, 1953). A major step in this direction was made when Fredrick Sanger and colleagues published the Sanger sequencing method (Sanger and Coulson, 1975). While this method was in no way high throughput by today's standards, it allowed for the sequencing of individual sequences in a straightforward manner that could be performed in a general lab setting. Similar approaches were applied for roughly twenty years until the field-changing introduction of current high throughput techniques. Using today's approaches, such as microarray analysis, RNA-seq, Chip-seq, it is possible to determine the expression of thousands of unique genes of interest or chromatin protein-association in a single experimental protocol. Millions of base pairs can be sequenced. This explosion in sequencing was been greatly aided by computational approaches.

The sheer volume of data produced by a single experiment is astonishing and analysis would simply not be possible without computation. Every aspect from genome assembly to quantification of transcriptional analysis is almost exclusively performed by bioinformatics. For example, transcriptional analysis using Illumina sequencing requires the reading and recording of millions of base pairs. Each read sequence must be checked for quality and aligned to a reference genome. Relative expression profiles must be determined for every gene in the genome prior to promoter sequence, alternative splicing, or SNP analysis is performed. By applying bioinformatics to these approaches a researcher can literally determine SNP placement for every gene in an

individual against a reference genome. As many transcriptional analyses compare gene regulation across experimental conditions, such as time, the entire analysis is performed multiple times and compared, further increasing the complexity. Clearly, computational approaches are advantageous to biological research. High-throughput approaches are increasingly being successfully applied to olfactory research.

Diego Rodriguez-gil from Stuart Firestein's lab recently applied micro-arrays to analyze Or gene expression across mouse development (Rodriguez-Gil et al., 2010). His analysis demonstrated that OR expression is first detectable at embryonic day 9 and that by day 13 expression of nearly all OR genes is detectable. Interestingly, they also noted that gene expression is not constant throughout the course of a mouse life as the "olfactome" expression decreased beyond 3 months of age.

Ewald Grosse-Wilde of Bill Hanson's lab recently identified the Or gene family of *Manduca sexta* using transcriptome analysis by applying 454 sequencing to an antennal cDNA library (Grosse-Wilde et al., 2011). Gene Ontology (GO) annotation using the software Blast2GO was applied to identify olfactory related ontologies. Ors were additionally identified by Blast and HMM profile searches from custom build databases. Predicted amino acid sequences for putative Ors were aligned to previously known *Heliothis* and *Bombyx* sequences, allowing for positive gene identification. Odor Binding Proteins (OBPs) and IR genes were also identified using the same method. In total 54 unique Or gene fragments were identified in a species that contains 73 unique globeruli. If *Manduca* follows the generally applicable one neuron to globerulus map, then roughly 73% of functional Or genes were identified. This approach is considerably significant as it represents application of high throughput transcription analysis in order to identify Or

genes in a species that have not been sequenced by using related species as a reference. As Or identification is traditionally achieved through genome sequencing, this faster and far less expensive transcriptome analysis demonstrates an approach that can be applied to many more species in the future.

R Jason Pitts of Laurence Zwiebel's lab has recently applied transcriptome profiling in *Anopheles gambiae* (Pitts et al., 2011). Both sex specific (male vs female) and olfactory appendage specific (antenna vs whole body) transcriptional regulation was compared using Illumina sequencing (RNA-seq). This approach was able to identify nearly all chemosensory genes from antennal tissue at quantifiable and statistically significant levels. The analysis revealed that expression levels of olfactory genes were present in both sexes, however the levels were much higher in females than males. Additionally, males had a significantly larger number of genes involved in audition. These results were interesting and make sense as it is the female mosquitoes that need to identify hosts for blood feeding and males that need to hear females for mating. This approach demonstrates the sheer volume of information that can be acquired using RNA-seq and its effective application to understand tissue specific regulation in a species and its broader implications.

Each of the previously highlighted analyses demonstrates how high throughput approaches can be successfully and advantageously applied to answer direct questions or identify targeted genes in the field of olfaction. With the continual reduction in cost as well as increased computational analysis tools available for these methods, more information can be gained at a less expensive price. As computational power has steadily increased, analyses that was once impossible can now be easily performed on a

desktop computer. Large numbers of need-driven computational approaches are being designed to assist in answering biological and chemical questions that can now be asked as the technology that was once considered science-fiction is now simply science.

Figure 1.1: Overview schematic describing neuronal wiring of the insect olfactory system

Odors are detected by Olfactor receptor (Or) expressing neurons in the periphery organs (antennae and maxillary palps). The axons of Olfactory Receptor Neurons (ORNs) expressing the same Or proteins, thus being housed in the same sensilla class, extend into unique glomerular locations within the antennal lobe. Local interneurons (LNs) allow for both excitatory and inhibitory information to be shared between specific sets of glomeruli. Projection Neurons (PNs) then send information from the antennal lobe to higher brain centers in the Mushroom Body (MB) and Lateral Horn (LH), where behavioral decisions occur. The figure was taken from (Carey and Carlson, 2011).

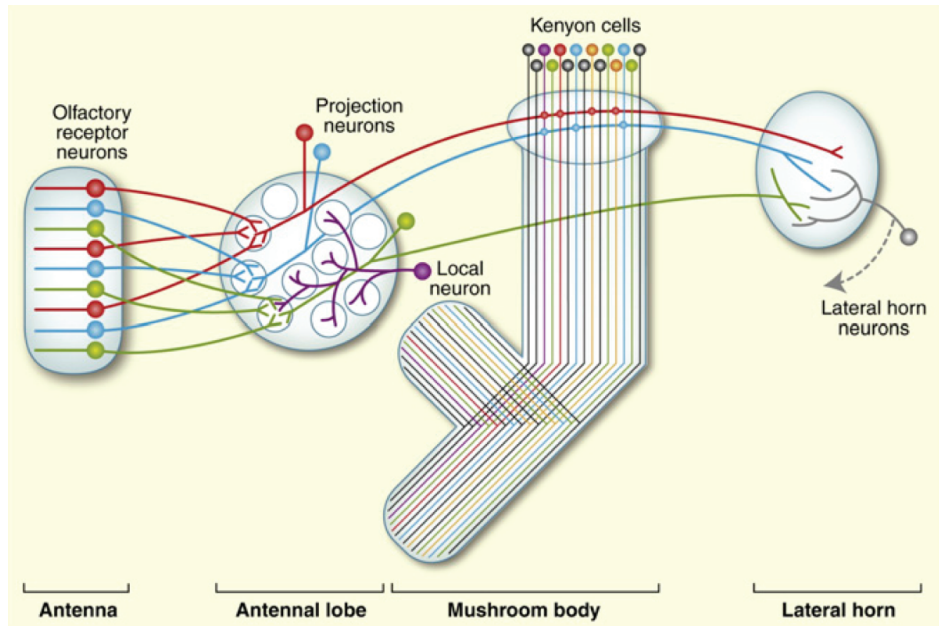


Figure 1.1

CHAPTER II:

Designing a Ligand-Based Virtual Screening Approach to Decode *Drosophila* Odor Receptor Chemical Space *In Silico*

INTRODUCTION

The peripheral olfactory system is unparalleled in its ability to detect and discriminate amongst an extremely large number of volatile compounds in the environment. To detect this wide variety of volatiles, most organisms have evolved large families of receptor genes that typically encode 7-transmembrane proteins expressed in the olfactory neurons (Buck and Axel, 1991; Clyne et al., 1999; Dahanukar et al., 2005; de Bruyne and Baker, 2008; Vosshall et al., 1999). Each volatile chemical in the environment is thought to interact with a specific subset of odor receptors depending upon odor structure and binding sites on the receptor. This precise detection and coding of odors by the peripheral olfactory neurons are subsequently processed, transformed and integrated in the central nervous system to generate specific behavioral responses that are critical for survival such as finding food, finding mates, avoiding predators etc (van der Goes van Naters and Carlson, 2006).

Currently there are two major rate-limiting steps in analysis of peripheral coding in olfaction: a very small proportion of chemical space can be systematically tested for its activity on odor receptors, and a very small fraction of the numerous odor receptors have been tested for responses (Araneda et al., 2000; Hallem and Carlson, 2006; Hallem et al., 2004; Kreher et al., 2008; Pelz et al., 2006; Saito et al., 2009). The challenges for overcoming the rate-limiting steps are enormous. First, volatile chemical space is immense, more than 2,000 odors in the environment have been catalogued from a small

fraction of plant sources alone (Knudsen et al., 2006). Second, the complete 3-D structures of the 7-transmembrane odor receptor proteins have not yet been determined and modeling of protein-odor interactions and sophisticated virtual screening methods are not yet possible except in rare instances (Triballeau et al., 2008). In the decade since the first systematic study of 47 odorants on the *Drosophila* antenna in 2001 (de Bruyne et al., 2001), additional studies have only identified a total of 251 novel active odors (de Bruyne et al., 1999; de Bruyne et al., 2001; Dobritsa et al., 2003; Goldman et al., 2005; Hallem and Carlson, 2006; Hallem et al., 2004; Kreher et al., 2005; Kreher et al., 2008; Kwon et al., 2007; Pelz et al., 2006; Schmuker et al., 2007; Stensmyr et al., 2003; Turner and Ray, 2009; van Naters and Carlson, 2007; Yao et al., 2005), which have been assembled and compared in an online database (Galizia et al., 2010).

Here we overcome this challenge by designing a chemical-informatics platform that is effective and fast. In order to do so we focused our attention on one of the most comprehensive quantitative data sets available, where measurements of responses of 24 *Drosophila* odor receptors to a panel of 109 odorants are known (Hallem and Carlson, 2006). We devised a method to identify molecular structural properties that are shared amongst the activating odorants for each receptor. We then utilize information about these shared molecular features of actives, that are presumably required for binding to a receptor, to perform *in silico* screens on a chemical space of >240,000 chemicals and predict the top 500 hits as new ligands for each of the odorant receptors (Ors). We then use single-unit electrophysiology to validate a subset of predictions *in vivo* and find that our method met with an overall success rate of ~71% in identifying novel ligands, as compared with a low 10% receptor activation rate while using non predicted odors. This approach allows us to create a peripheral coding map of a large

chemical space, which substantially improves our ability to study peripheral olfactory coding and provides a powerful tool for the discovery of novel ligands for Ors.

RESULTS

Analysis of odorant structure

Since the structure of receptor protein complexes is not known, we analyzed receptor-odor interactions by applying the *similarity property principle*, which reasons that structurally similar molecules (e.g. activating odorants) are more likely to have similar properties (Hendrickson, 1991; Martin et al., 2002). Although this general approach has been useful in the area of pharmaceuticals (Keiser et al., 2009; Martin et al., 2002), receptor-odor analysis presents significant additional challenges. Not only are odorant molecules generally smaller in size than pharmaceuticals (average MW of known odors ~3-fold less than FDA approved pharmaceuticals (Wishart et al., 2008)) and therefore offer fewer structural features for differentiation, they are also detected by the receptors with specificity at extremely low concentrations in the volatile phase (Hallem and Carlson, 2006; Kreher et al., 2008). Additionally, odor receptors are differentially tuned and can sometimes appear not to follow distinct structural rules: odors that look structurally different can strongly activate the same receptor, while odors that appear very similar may have very different levels of activity (Hallem and Carlson, 2006)(Figure 2.1A).

General measures of odorant similarity

Similarity in chemical structure can be described and measured quantitatively using multiple approaches, however a single method may not be ideal for every single application (Maldonado et al., 2006). In order to test whether non-optimized approaches would be able to identify similarities in shape of known activators we compared four different approaches: Cerius2 (Accelrys Software Inc), Dragon (Talete), Maximum-Common-Substructure (MCS) (Cao et al., 2008b), and Atom-Pair (AP) (Cao et al., 2008a; Carhart et al., 1985). Cerius2 and Dragon represent collections of 200 and 3,224 molecular descriptors, respectively, that calculate values for a broad range of chemical properties such as molecular weight, functional group counts, and in the case of Dragon, 3 dimensional relationships within molecules. The AP method compares shortest path distances between all atom pairs in a molecule. Lastly, MCS identifies the largest 2 dimensional substructure that exists between two compounds. Using each of these approaches, we computed distances between 109 odors that had previously been tested against 24 Ors from *Drosophila melanogaster* (Hallem and Carlson, 2006) (Figure 2.1B). These represent most of the *Or* genes expressed in the *Drosophila* antenna (Hallem and Carlson, 2006). Upon comparison, we find that none of the four methods were vastly superior and that each method varied in the ability to group known activating odorants “actives” close together in distance as measured for each Or using a method called Accumulative-Percentage-of-Actives (APoA)(Chen and Reynolds, 2002) (Figure 2.3A, See Methods), and value of the Area-Under-the-Curve (AUC). Ultimately, Dragon and Cerius2, which utilize a large number of diverse molecular descriptor values to describe each odor structure, performed better than AP or MCS, suggesting that a more diverse

set of descriptors is better at explaining Or activity than 2D measures alone (Figures 2.2B, 2.3B). Atom-Pair and MCS were subsequently ignored from further development.

Identification of unique subsets of optimized descriptors for each *Drosophila* Or

Individual Ors respond to distinct subsets of ligands with some degree of overlap (Hallem and Carlson, 2006; Kreher et al., 2008). We reasoned that rather than using entire Dragon or Cerius2 descriptor sets, which likely includes a number of measurements for features irrelevant for ligand-binding to an individual Or, judiciously selecting subsets of molecular descriptors suited to cluster actives for an individual receptor may be more effective at defining an Or-specific chemical space. To test this hypothesis, we used a Sequential-Forward-Selection (SFS) method to incrementally create unique optimized descriptor subsets for each Or from an initial combined set of 3,424 descriptors from Dragon and Cerius2 (Whitney, 1971) (See Methods, Figure 2.2A). This optimization-based analysis was performed on the 19 Ors from the dataset with known activating odors, excluding Or82a, since it has but a single known strong active (Hallem and Carlson, 2006).

Not surprisingly, the composition of the optimized descriptor sets varied greatly between Ors (Table 2.1). Molecular descriptors can be categorized from 0 to 3 dimensions. 0-Dimensional (0-D) descriptors define features that can be viewed as not directly being shape dependent, such as molecular weight or vapor pressure. On the other end of the scale, 3-Dimensional (3-D) descriptors define features of molecules in three-dimensional space, such as the distance between two atoms of an odor molecule. We find an overwhelming preference for 3-D and 2-D descriptors compared to 1-D and 0-D descriptors, suggesting that structural shape features are more important for

receptor-odor interactions. The Or-optimized descriptor sets were far superior at grouping together activating odors from the training set (Figure 2.2B) when compared to the non-optimized methods (Dragon, Cerius2, MCS, AP) and a previously identified collection of descriptors that were identified without receptor-specific optimization (Haddad et al., 2008).

Computational validation of optimized descriptor sets

In order to validate the predictive ability of the *Or*-optimized method, we performed 5 independent trials of 5-fold cross-validations followed by a Receiver-Operating-Characteristic (ROC) analysis, an established computational approach (Hastie et al., 2001; Tan et al., 2006) (see methods). Briefly, this involved withholding 20% of the 109 previously tested odors for a receptor. Descriptors were optimized using the remaining 80% odors for training, and ligand-predictions were subsequently performed on the 20% of odors that were withheld. This operation was repeated 5 times for each receptor, each time selecting a different 20% as withheld from the training set. The entire 5-fold operation was repeated 5 times for each receptor and a mean ROC curve representing the prediction accuracy determined. This analysis was possible for 12 *Ors* which had >6 known ligands that activated >100 spikes/sec. The Area-Under-Curve (AUC) value (0.815) is very promising and suggests that the *Or*-optimized descriptor sets are effective at predicting novel ligands (Figure 2.2C).

In addition to performing the 5-fold cross-validation, we also clustered the 109 training odors independently for each *Or* using distances calculated from the previously determined receptor specific descriptor sets (Figure 2.2D). As expected, we find that active odors cluster tightly together for each *Or* (Figure 2.2D). In a few cases, such as

for Or35a and Or98a, not all the highly active compounds are clustered, suggesting the possibility of multiple or flexible binding sites, or imperfect selection of descriptors. Four of the Ors (Or2a, Or23a, Or43a and Or85f) have few known activators, none of which activate the receptors at greater than 150 spikes/sec, however our descriptor optimization approach is still able to cluster each the few weak activators together (Figure 2.2D).

High-throughput *in silico* screening of odor receptors

Since Or-optimized descriptor sets can efficiently group highly active compounds in chemical space, we used them to rank untested compounds according to their distance from known actives for specific Ors. We assembled a natural odor library, which contains 3,197 naturally occurring odors, and a library derived from Pubchem (Bolton et al., 2008), which contains >240,000 compounds with similar molecular weights and atom type compositions to known volatiles (See Methods). We then systematically screened both libraries using the optimized descriptor sets of 19 *D. melanogaster* Ors representing ~5,000,000 receptor-odor interactions *in silico*. We identify the top 500 (0.2%) hits from this vast chemical library for each Or, the top ~100 of which are reported in Table 2.2.

Electrophysiological validation of *in silico* screen and identification of agonists

To validate our *in silico* screen, we obtained several untested odorants (141 total; ~11-23/Or) belonging to the top 500 predicted ligands for 9 different Ors that were available from commercial sources at high purity and reasonable prices. The 9 receptors were selected on the basis of previous functional mapping studies that enable us to

unambiguously identify the antennal olfactory receptor neurons (ORNs) they are housed in (Couto et al., 2005; Hallem et al., 2004). We systematically tested each predicted receptor-odor combination using single-unit electrophysiology to record from the ORNs to which these 9 Ors have been previously mapped (Couto et al., 2005; Hallem et al., 2004). We find that a majority of the predicted ligands evoked responses from the target ORNs; ~71% evoked either activation (>50 spikes/sec above the spontaneous activity) or inhibition (>50% reduction in spontaneous activity) (Fig. 2.6A). These cutoffs were selected based on study from which the training set was obtained and has been used in other studies in the past that use this type of recordings (Hallem and Carlson, 2006; Kreher et al., 2008). Interestingly, the mean vapor pressure of activating odors (11.84 Torr) is 7.5 times higher than of inactive odors (1.58 Torr), raising the possibility that some inactive odors may not be volatilized and delivered at adequate levels to the ORNs. Additionally, we find that ~13% of the predicted compounds we tested showed an inhibitory effect on baseline activity of the respective neuron (Fig. 2.6A). These inhibitors were identified by virtue of structural similarity to known activators suggesting that they may bind to similar sites on the receptor. Thus as an additional benefit our approach may provide a method to identify inhibitors as well. Such inhibitors would not only provide important tools to investigate mechanisms of odor receptor inhibition but could also be used in blocking specific odor-mediated behaviors. Consistent with our observations three of the receptor-odor interactions had been previously identified independently as well, Or22a (Pelz et al., 2006), and Or49b (Hallem et al., 2004). The electrophysiological analysis provides the most important validation of our Or-optimized descriptor-based *in silico* screen.

Odor response spectra of individual Ors

Since we systematically analyzed responses of a large number of new odorants individually, we were able to characterize the odor-response spectra of several antennal ORN classes to these new ligands (Figure 2.4B). New activators are reported for every receptor, and inhibitors are identified for several. Ligand predictions for 2 of the 3 receptors that do not perform as well are Or10a and Or49b that detect aromatic compounds. Their poor performance is explained by the lack of aromatic ligands in the initial training set (13/109) odorants. We find that a >85% of the predicted ligands activate odor receptors Or7a, Or22a, Or59b, Or85a, Or85b, and Or98a (Figure 2.4B).

Specificity of *in silico* predicted ligands

We rigorously examined the rate of false negative predictions for each Or by systematically testing newly identified ligands of each Or against the other non-target receptors using electrophysiology. Of 504 non-target receptor-odor interactions tested, we found that only 10% evoked a response >50 spikes/sec and 3.7% evoked a response >100 spikes/sec (Figure 2.5A). This represents a high degree of specificity, especially considering that the Or-optimized descriptor method did not incorporate any additional computational screening to rule out non-target activators. Additionally, when we plot the percentage of odors that validated as activators when tested using electrophysiology (considering both predicted and non-target receptor-odor interactions), we find that activity is strongly related to predicted odor ranking (Figure 2.5B). Odors which rank closest to known actives for each Or, particularly within the top 500 hits, are far more likely to be activators than odors further away, and there is a drastic drop-off in activating odors present beyond the 1,000 rank. We see the same trend if we plot mean

activity of odors for the same ranking divisions (Figure 2.6). Highly ranked odors have a far higher mean activity than distantly ranked odors.

Relationship between descriptor sets and Or sequence and activity

Since receptor-optimized descriptor sets and the predicted ligand space they define are a function of shared molecular features that a receptor may employ to recognize ligands, we were now in a position to determine how these characteristics correlate with receptor properties such as their known-activity profiles and amino acid sequences. We used hierarchical cluster analysis to create trees that represent the various receptors based on: shared descriptors selected; known activity-based relationships (Hallem and Carlson, 2006); degree of overlap of predicted ligands; and amino acid sequence (Figure 2.7A, See Methods). We found that the maximum overlap in Or relationships is retained between the descriptor and the activity trees, and the descriptor and the cross activity trees with 11 out of 24 Ors present in subgroups that are common in both cases. However, only two subgroups (yellow and grey) are conserved across the 3 trees. The largest shared overlap existing in the descriptor tree suggests that the Or-optimized descriptors link the known and the predicted receptor-odor interactions and that our analysis may expand upon odor receptor activity relationships beyond those previously known from the training data. We also found that the phylogenetic tree has fewer relationships conserved with each of the trees, consistent with previous observations (Hallem et al., 2004) supporting the idea that, while the most conserved amino acid residues in the Ors provide the structure of the tree, they do not correlate strongly with ligand specificity.

Analysis of breadth of predictions for each Or in chemical space

Coding of odors in a large volatile space (>240,000) by a receptor repertoire is virtually impossible to determine experimentally. However, based on the Or-optimized descriptor sets we computationally derived prediction frequency distributions for each of the *Drosophila* Ors in this large chemical space (Figure 2.7B). As expected, we find substantial variation in the distribution frequency of predicted ligands across different receptors. The predicted response profiles support previous observations made with smaller odor panels that the olfactory system can potentially detect thousands of volatile chemicals, many of which the organism may never have encountered in its chemical environment. We computed similar predicted ligand frequency distribution curves to an assembled set of 3,197 known “natural” volatile compounds from plants, humans, and a fragrance collection (Figure 2.8). Plant volatiles constituted a large portion of compounds that are predicted to be ligands for *Drosophila* Ors. To further analyze odor source representation, we classified odors that belong to top 500 prediction lists according to their source, if known, and find that Ors are not specialized for odors from a single source (Figure 2.7C).

Across-receptor activation patterns in *Drosophila*

To study the ensemble activation patterns of odors predicted across all Ors, we analyzed the across-receptor activation patterns of the 3,197 known compounds for 9 receptors (Or7a, 10a, 22a, 47a, 49b, 59b, 85a, 85b, 98a). Surprisingly, we find that only 25% of compounds from the “natural” odor library found in the top 500 predictions for each Or are predicted to activate multiple Ors (Figure 2.7D, lower left panel). If we consider compounds from the Pubchem library in the top 500 predicted actives for each

receptor, we observe further reduction in the proportion of across-receptor activating compounds (Figure 2.7D, upper right). Consistent with this prediction we find that cross-activation by ligands functionally evaluated in this study for 9 receptors is lower than that reported previously using ligands of comparable strength for the same 9 receptors (Hallem and Carlson, 2006) (Figure 2.7D, lower right panel). These data suggest that a number of natural odors may be detected by only a few receptors, particularly at low concentrations.

Producing a systems level view of receptor activity for the *Drosophila* antenna

One of the ultimate challenges of understanding peripheral coding of olfactory information is to be able to map responses of large number of receptors to their ligands. Following the *in silico* analysis, we were able to create such a network-view of peripheral odor coding in the *Drosophila* antenna by mapping all predicted and tested receptor-odor combinations (Figure 2.9A) as has been done previously for mapping drug-target networks (Keiser et al., 2009). Using our chemical informatics pipeline, it becomes possible to infer the network of odor receptors that are activated from complex odor sources without the expensive and time consuming process of purchasing and testing all possible odors.

DISCUSSION

A primary element of the olfactory code is information about odor identity, represented by the characteristic interaction of an odor with the ensemble of olfactory receptors in the nose. Here we report an *in silico* approach to systematically identify ligands from a vast chemical space for a majority of Ors in the antenna of *Drosophila*. We demonstrate that our predictions are accurate using two different validation approaches- computational validations and functional validation using electrophysiology. There is a strong correlation between ranks of predicted ligands to electrophysiological activity.

Obtaining and testing odors using traditional methods is time and cost intensive. Electrophysiology and calcium imaging are consuming processes that require not only a great deal of time to perform, but also the purchase of each odor to be physically tested. Moreover, large plate-based combinatorial chemical libraries, which are commonly implemented in drug discovery in the pharmaceutical industry, are not available for volatile odor libraries for reasonable costs. Since *Drosophila* is a premier model for understanding neurobiology of olfaction, several laboratories over the last 12 years have together screened ~250 odors, activities of which have been and compiled into a valuable database that standardizes across studies (Galizia et al., 2010). In this study we screen >240,000 chemicals and predict >10,000 new ligands which represents a substantial expansion of the known peripheral olfactory code for this important model organism and provides a system-level view of odor detection (Figure 2.9B).

A similar, yet much smaller, analysis applied chemical informatics on *Drosophila* olfactory neuron activities to 47 odorants and screened ligands from 21 untested compounds in *Drosophila* (Schmucker et al., 2007). Although this study had a relatively

modest success rate of ~25% at predicting untested odorants as activators (by applying the same 50 spikes/sec threshold for comparison), it also highlighted that structure-based ligand prediction is a viable method for further development. In another interesting analysis a Quantitative Structure Activity Relationship (QSAR) model was applied to describe odor-activity for *Drosophila* Ors. Important amino acid residues were identified using information from orthologous Or sequences identifying potential odor-binding regions, which was postulated to be 15 angstroms deep and 6 angstroms wide (Guo and Kim, 2010). Our approach is conservative and designed to search for novel odors that share structural features from a previously tested odor panel. Odor molecules are limited in size as well, and may offer a limited scaffold such that novel isofunctional chemotype identification may not be as prevalent as has been seen in other examples of scaffold-hopping (Schneider et al., 2006). However while compounds that share similar values for the optimized descriptors do have structural similarity for selected parts of the molecule, it does not mean that they are not structurally different in other parts of the molecule. In the future, application of machine learning approaches, such as Support Vector Machines (SVMs) to the receptor-optimized molecular descriptor sets, may be useful to further increase the predictive ability. Additionally, we could replace our SFS approach with sequential floating search techniques, which allows for removal, as well as addition, of descriptors in the growing optimized list.

We predict that a number of odorants at low concentrations may be detected by only one or a few receptors. This contrasts a current model of combinatorial coding in which emphasis is placed on the notion that a majority of volatile chemicals, with the exception of pheromones and CO₂, are detected by combinations of several odor receptors. One possible explanation for this disparity could be that our predictions are

fundamentally conservative in nature because we focus only on structurally similar ligands and 7-transmembrane heteromeric receptors may also contain additional unexplored binding sites. Moreover, receptors may respond to compounds ranked beyond the top 500 hits. Another possibility is that previously tested subsets of odors were potentially selected on the basis of strong responses in electroantennograms and behavior assays, which could bias selection of cross-activating odors. In fact it is known that complex fruit odors activate fewer Ors than the number activated by individual odorants at comparable concentrations (Hallem and Carlson, 2006, Semmelhack and Wang, 2009). The architecture of the olfactory code therefore appears to integrate two different models. On the one hand, most odors are detected by a few Ors from the repertoire, which may enhance the specificity of the olfactory system for detection of a large number of odors. On the other hand, 15-20% of odors are predicted to activate several Ors (up to 50%) at the same time, which may serve to aid the olfactory of the system in discriminating between fine concentration changes of important stimuli by having Ors tuned to low and high concentrations such as shown for Or42a and Or42b (Kreher et al., 2008).

By identifying a large number of new ligands for each odor receptor, we can also begin to systematically compare the ligand tuning profiles for each in the endogenous neurons versus the “empty neuron” decoder system. If clear differences were identified, it could enable the identification of underlying reasons such as differences in levels of receptor expression in the neurons, or presence of different Odorant Binding Proteins (OBPs) in the sensillum lymph.

This cheminformatics pipeline can also be applied for system-level analysis of other insects whose receptors and ORNs have been decoded such as mosquitoes

(Carey et al., 2010), and vertebrates such as mice and humans (Saito et al., 2009). The search for novel insect repellents and attractants for species that transmit disease and destroy crops can be greatly assisted by a rational prioritization using such a cheminformatics approach.

Figure 2.1: Comparing efficacy of various structure analysis methods to analyze odor receptor ligands

(A) Comparison of odor structures for known actives and inactive odorants for Or85b and Or98a from (Hallem and Carlson, 2006) (activity shown in parenthesis). **(B)**

Overview of the process by which 4 molecular descriptor methods were compared to determine which one best clusters known actives close together in descriptor space.

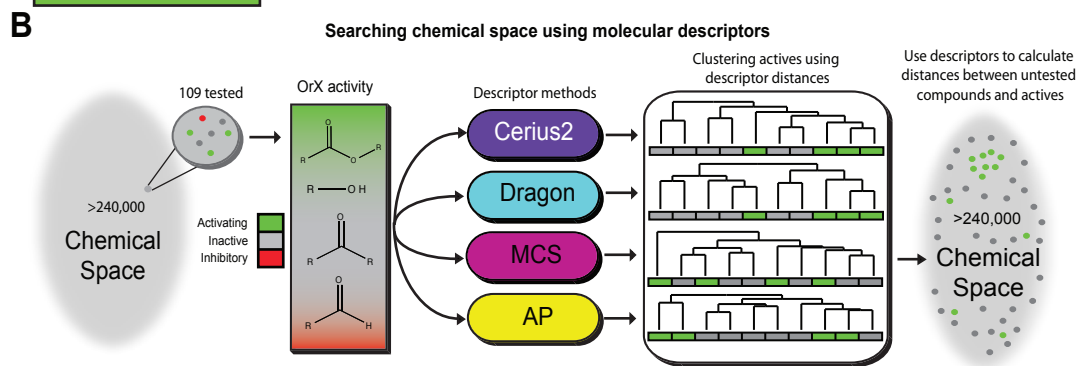
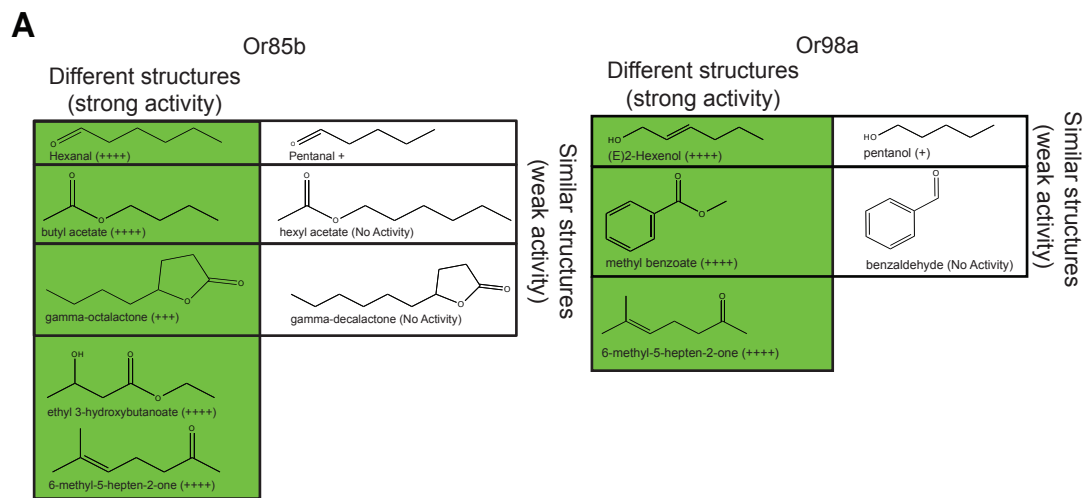


Figure 2.1

Figure 2.2: A receptor-optimized molecular descriptor approach has strong predictive power to find new ligands

(A) Schematic of the cheminformatics pipeline used to identify novel ligands from a larger chemical space. **(B)** Plot of mean APoA values for 19 *Drosophila* Ors calculated using various methods including a previously identified set (Haddad et al., 2008). **(C)** Receiver-operating-characteristic curve (ROC) representing computational validation of ligand predictive ability of the Or-optimization approach. **(D)** Hierarchical cluster analysis of the 109 odorants of the training set (Hallem and Carlson, 2006) using Or-specific optimized descriptor sets to calculate distances in chemical space for odorant receptors with strong activators (green), and odorant receptors with no strong activators (yellow).

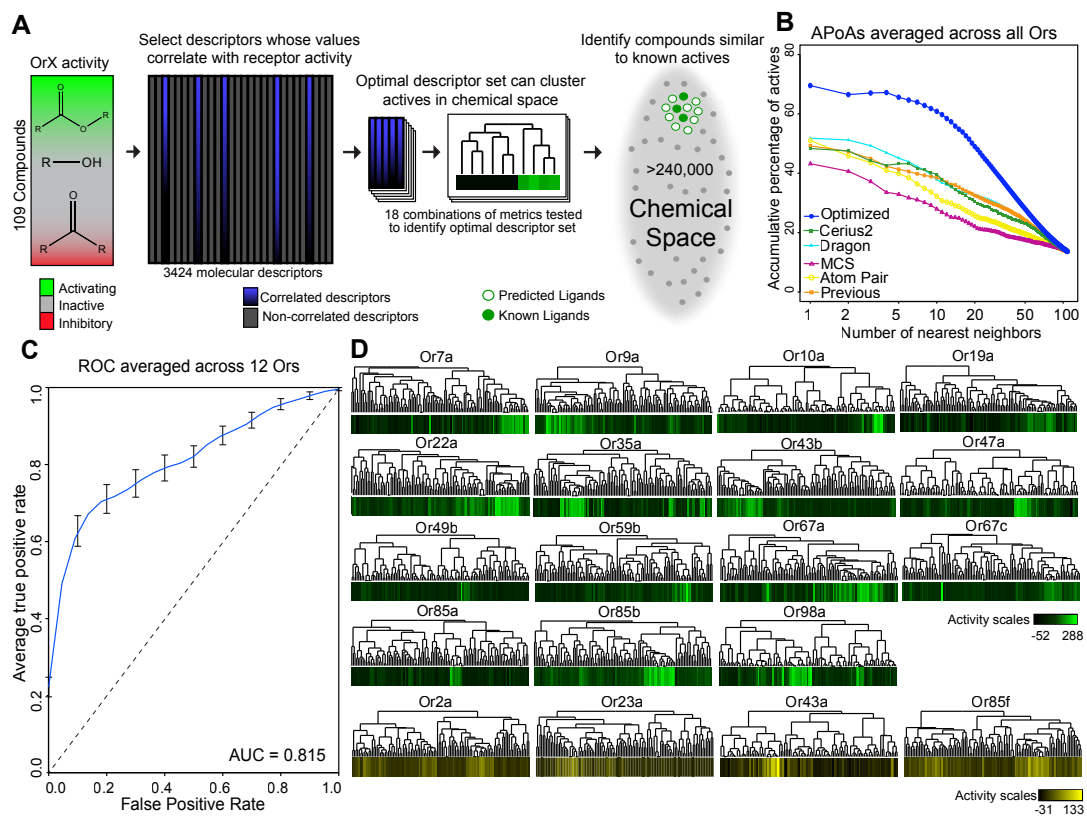


Figure 2.2

Figure 2.3: Accumulated percentage of actives analysis

(A) Representative example for Accumulated Percentage of Actives (APoA) calculation. Green box=active, grey box=inactive. To calculate APoA each active compound was iteratively used as a reference active. Compounds are sorted based upon their increasing descriptor based distance from reference active, and the APoA calculated for each of the other compounds as a ratio of the number of actives over the total number of compounds considered from the reference compound. This process was repeated using each active odorant as a reference active. Reference compound APoAs were averaged to a single mean APoA value. The higher the APoA value while considering a fixed number of nearest neighboring compounds, the greater the proportion of active compounds clustered together. **(B)** Coloured cells mark the method that clusters active ligands best as determined by the highest Area-Under-Curve (AUC) for APoA values.

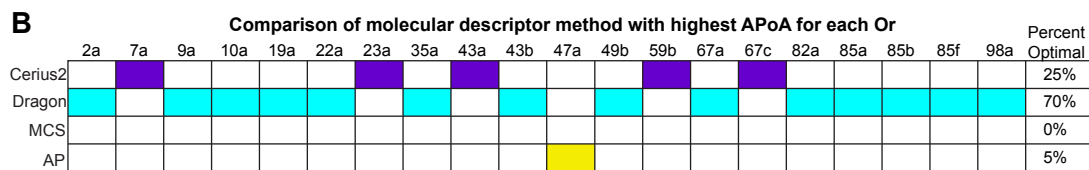
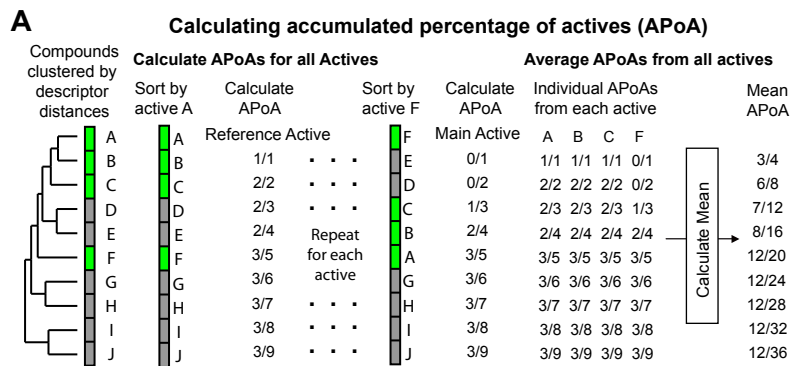


Figure 2.3

Figure 2.4: Electrophysiology validates that odorant receptor-optimized molecular descriptors can successfully identify new ligands for *Drosophila*

(A) Summary of prediction accuracy percentages obtained by electrophysiology validation. Ligands = Agonists (≥ 50 spikes/s) + Inhibitors ($>50\%$ reduction from baseline activity). **(B)** Mean increase in response of neurons to 0.5-sec stimulus of indicated odors (10^{-2} dilution) predicted for each associated Or. Dashed lines indicate the activator threshold (50 spikes/second). ΔH : Or85b (ab3B) = flies lack expression of Or22a in neighboring neuron, thus all observed neuron activation is unambiguously caused by Or85b. N=3, error bars=s.e.m.

A

Classification	Or7a	Or10a	Or22a	Or47a	Or49b	Or59b	Or85a	Or85b	Or98a	Total
Ligands	88%	31%	86%	39%	27%	91%	92%	87%	100%	71%
Agonists (>50 spikes/sec)	63%	31%	81%	33%	18%	64%	69%	70%	92%	58%
Agonists (>100 spikes/sec)	31%	13%	62%	11%	9%	45%	54%	48%	67%	37%
Inhibitors	25%	0%	5%	6%	9%	27%	23%	17%	8%	13%

B

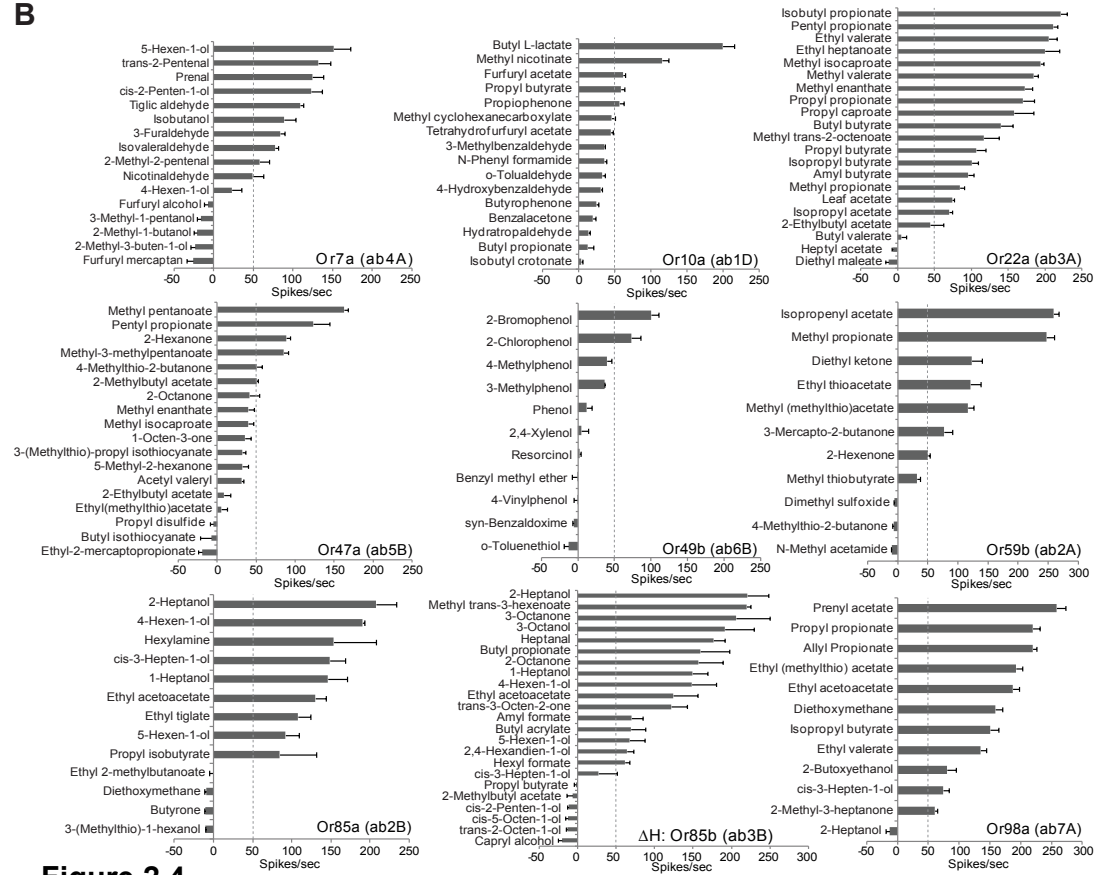


Figure 2.4

Figure 2.5: Predicted receptor-odor interactions are highly specific

(A) Plot of activity (Top) for electrophysiologically tested receptor-odor interactions. (Bottom) Plot indicating locations of predicted receptor-odor combinations (green) and same odorants tested in non-target receptor-odor combinations (gray). **(B)** Plot of percentage of activating odors (>50 spikes/sec) considering all activating or inactive odors (>0 spikes/sec) across ranking bins for all odors tested using electrophysiology.

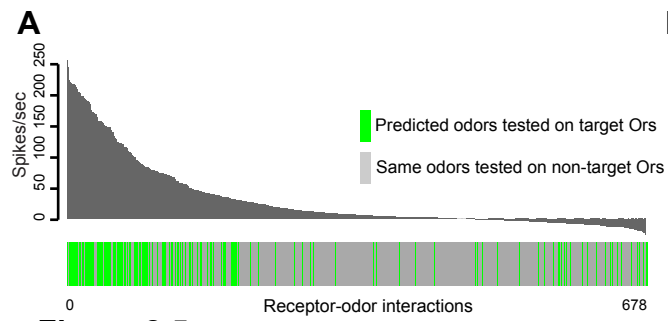


Figure 2.5

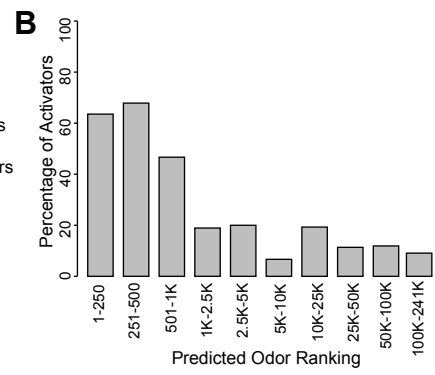


Figure 2.6: High-ranking odors are more likely to be active than distantly ranked odors

Bar plot of the mean electrophysiological activity (in spikes/sec) of all activating or inactive odors (>0 spikes/sec), considering both predicted receptor-odor combinations and same odorants tested in non-target receptor-odor combinations, grouped by predicted rankings.

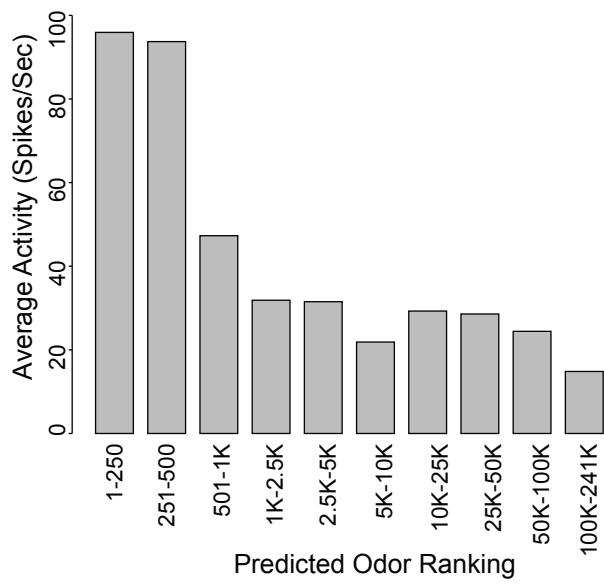


Figure 2.6

Figure 2.7: Analysis of receptor-odor relationships and breadth of tuning

(A) Hierarchical clusters created from Euclidean distance values between *Drosophila* Ors calculated using: (left to right) shared optimized descriptors; known activity to training set odors (Hallem and Carlson, 2006); overlap across top 500 predicted ligands; and Phylogenic tree of receptors (Hallem and Carlson, 2006). Sub clusters shaded with colors or bars. **(B)** Frequency distribution of compounds from the >240K library within the top 15% distance from highest active plotted to generate predicted breadth of tuning curves. Green arrows indicate relative distance of the furthest known activating compound determined by electrophysiology. **(C)** (Left) The numbers of compounds present in the collected volatile library according to source. (Right) The numbers and sources of predicted ligands for the 19 *Drosophila* odor receptors/neurons within the top 500 predicted compounds **(D)** Comparison of plots for percentage of receptors that are: (top left) activated by percentage of known odors from training set (Hallem and Carlson, 2006); (bottom left) predicted to be activated by collected compounds; (top right) predicted to be activated from >240K library; and (bottom right) activated by ligands for 10 shared Ors in this study versus activated by comparable actives previously tested (Hallem and Carlson, 2006).

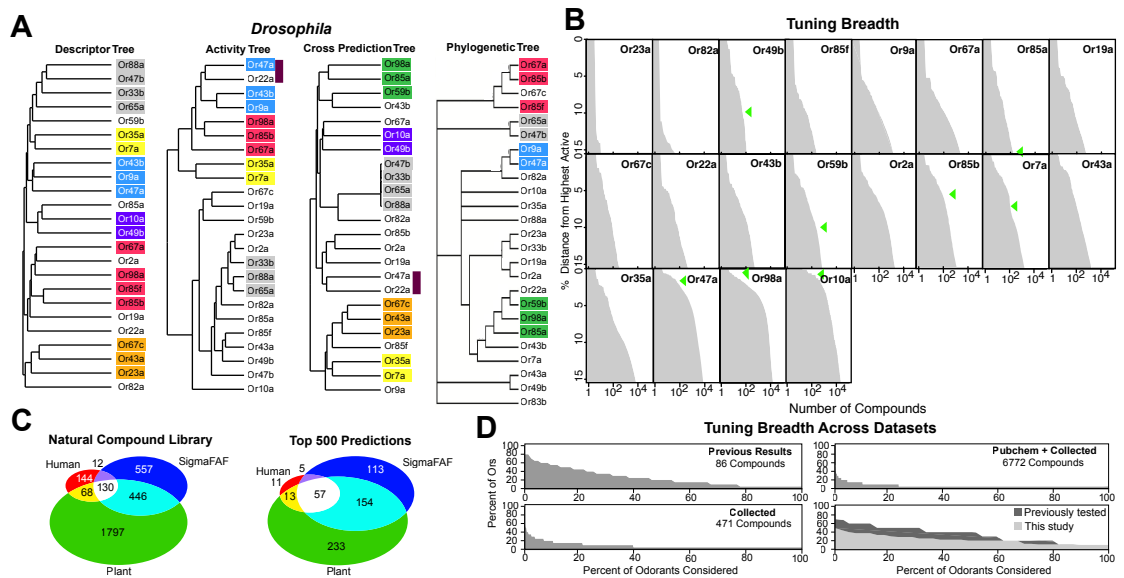


Figure 2.7

Figure 2.8: Analysis of receptor-natural odor interactions

Compounds from the collected compound library that have been cataloged as plant, human and total collected volatiles were ranked according to their relative distance from the compound with highest activity. Frequency distribution of compounds within the top 15% is plotted to generate predicted breadth of tuning curves.

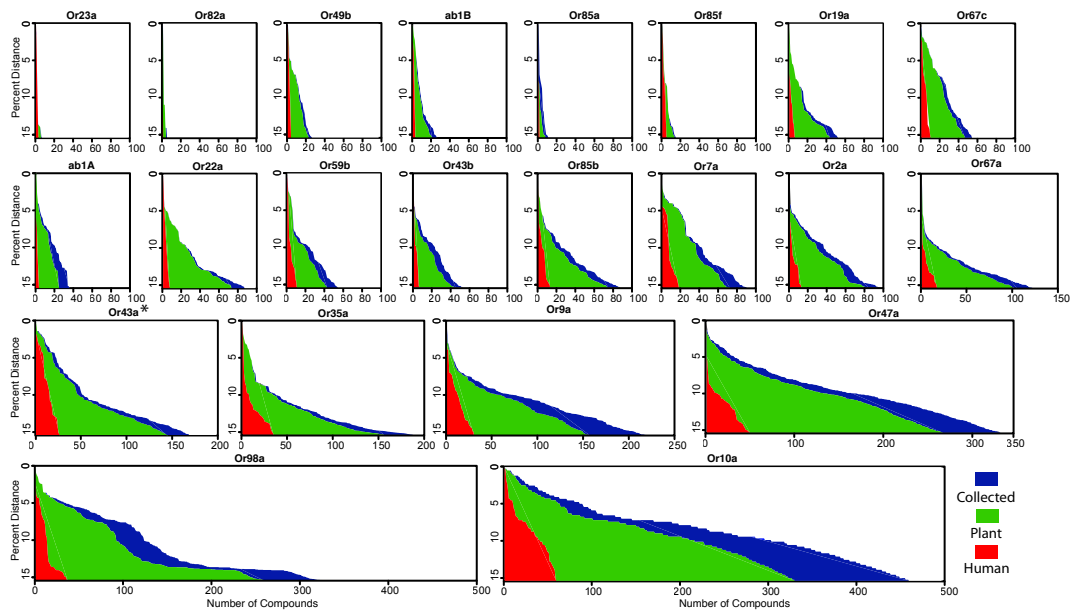


Figure 2.8

Figure 2.9: Predicted odor space and network view of odor coding

(A) *Drosophila* receptor-odor network. Each known interaction (>50 spikes/sec) from this and previous studies (Hallem and Carlson, 2006) is linked by a purple edge. Predicted receptor-odor network (top 500 hits) are linked by light-grey edges. All compounds are represented as small black circles and Ors are represented as large colored circles matching the colors used in (Fig. 4A and SI). **(B)** Expansion of the peripheral olfactory code in this study: large increase in numbers of identified activators and inhibitors. The different sized circles represent the approximate ratio of numbers of previously known ligands (top circles), predicted ligands based on a cutoff of the top 500 predicted compounds per receptor and corrected to the validation success rate (lower, diffuse circles).

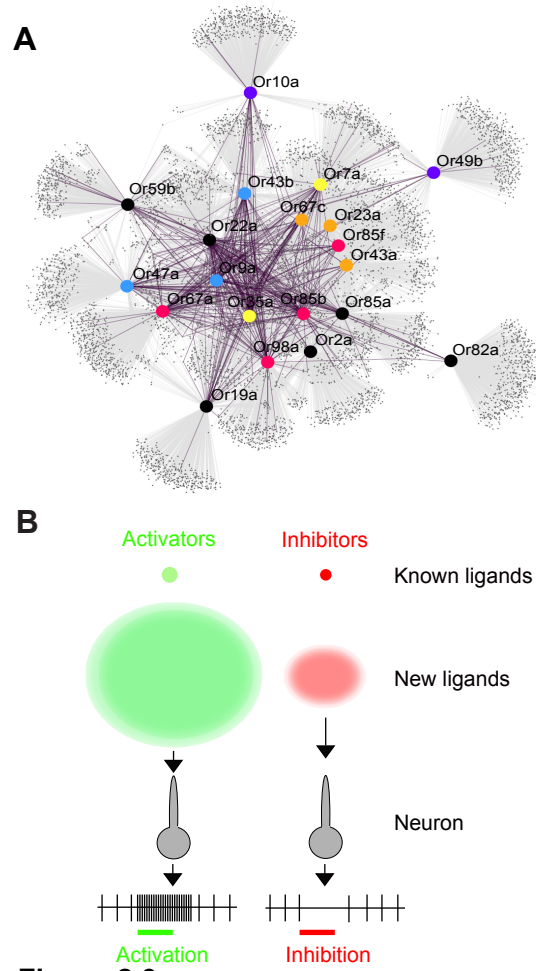


Figure 2.9

Table 2.1: Optimized descriptor sets for each *Drosophila* Or

Optimized descriptors occurrences, symbol, brief description, class, and dimensionality are listed. A summary of the total number of descriptors selected for the receptor repertoire is provided at the beginning. Descriptors are listed in ascending order of when they were selected into the optimized set. Weights indicate the number of times a descriptor was selected in an optimized descriptor set.

Supplementary Table 1

Descriptor #	Weight	Symbol	Description	Class	Dimensionality
Or2a (18 Unique)					
1	Mor18p		3D-MoRSE - signal 18 / weighted by atomic polarizabilities	3D-MoRSE descriptors	3
1	Mor17e		3D-MoRSE - signal 17 / weighted by atomic Sanderson electronegativities	3D-MoRSE descriptors	3
1	Mor28u		3D-MoRSE - signal 28 / unweighted	3D-MoRSE descriptors	3
1	J3D		3D-Balaban index	geometrical descriptors	3
2	O-057		phenol / enol / carboxyl OH	atom-centred fragments	1
1	SICT		structural information content (neighborhood symmetry of 2-order)	information indices	2
1	EEig10x		Eigenvalue 10 from edge adj. matrix weighted by edge degrees	edge adjacency indices	2
1	MATS5e		Moran autocorrelation - lag 5 / weighted by atomic Sanderson electronegativities	2D autocorrelations	2
1	F05[C-O]		frequency of C - O at topological distance 05	2D frequency fingerprints	2
1	HNar		Narumi harmonic topological index	topological descriptors	2
1	MATS8m		Moran autocorrelation - lag 8 / weighted by atomic masses	2D autocorrelations	2
1	G3s		3st component symmetry directional WHIM index / weighted by atomic electrotopological states	WHIM descriptors	3
1	Mor27m		3D-MoRSE - signal 27 / weighted by atomic masses	3D-MoRSE descriptors	3
1	B04[C-O]		presence/absence of C - O at topological distance 04	2D binary fingerprints	2
1	H8v		H autocorrelation of lag 8 / weighted by atomic van der Waals volumes	GETAWAY descriptors	3
1	Mor10v		3D-MoRSE - signal 10 / weighted by atomic van der Waals volumes	3D-MoRSE descriptors	3
1	Mor18v		3D-MoRSE - signal 18 / weighted by atomic van der Waals volumes	3D-MoRSE descriptors	3
2	R8p+		R maximal autocorrelation of lag 8 / weighted by atomic polarizabilities	GETAWAY descriptors	3
Or7a (31 Unique)					
1	MAXDP		maximal electrotopological positive variation	topological descriptors	2
1	MAXDN		maximal electrotopological negative variation	topological descriptors	2
1	B06[C-C]		presence/absence of C - C at topological distance 06	2D binary fingerprints	2
2	HATS1v		leverage-weighted autocorrelation of lag 1 / weighted by atomic van der Waals volumes	GETAWAY descriptors	3
3	Hy		hydrophilic factor	molecular properties	1
1	S_ssO		S_ssO	atomtypes (Cerius2)	1
1	JGT		global topological charge index	topological charge indices	2
2	H-051		H attached to alpha-C	atom-centred fragments	2
2	EEig10d		Eigenvalue 10 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2
1	O-057		phenol / enol / carboxyl OH	atom-centred fragments	1
5	HATS8u		leverage-weighted autocorrelation of lag 8 / unweighted	GETAWAY descriptors	3
1	G2s		2st component symmetry directional WHIM index / weighted by atomic electrotopological states	WHIM descriptors	3
2	Mor16u		3D-MoRSE - signal 16 / unweighted	3D-MoRSE descriptors	3
4	B02[O-O]		presence/absence of O - O at topological distance 02	2D binary fingerprints	2
1	R5p+		R maximal autocorrelation of lag 5 / weighted by atomic polarizabilities	GETAWAY descriptors	3
1	EEig8d		Eigenvalue 08 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2
1	DISPp		d COMMA2 value / weighted by atomic polarizabilities	geometrical descriptors	3
2	C-008		CHR2X	atom-centred fragments	1
1	R4e+		R maximal autocorrelation of lag 4 / weighted by atomic Sanderson electronegativities	GETAWAY descriptors	3
1	EEig9d		Eigenvalue 09 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2
1	nAROH		number of aromatic hydroxyls	functional group counts	1
1	R2m+		R maximal autocorrelation of lag 2 / weighted by atomic masses	GETAWAY descriptors	3
1	nRCOOR		number of esters (aliphatic)	functional group counts	1
1	B02[C-O]		presence/absence of C - O at topological distance 02	2D binary fingerprints	2
1	GATS7m		Geary autocorrelation - lag 7 / weighted by atomic masses	2D autocorrelations	2
1	E2s		2nd component accessibility directional WHIM index / weighted by atomic electrotopological states	WHIM descriptors	3
1	nRCO		number of ketones (aliphatic)	functional group counts	1
1	Mor03m		3D-MoRSE - signal 03 / weighted by atomic masses	3D-MoRSE descriptors	3
1	MATS8m		Moran autocorrelation - lag 8 / weighted by atomic masses	2D autocorrelations	2
1	CIC5		complementary information content (neighborhood symmetry of 5-order)	information indices	2
1	D/Dr06		distance/detour ring index of order 6	topological descriptors	2
Or9a (29 Unique)					
1	BEHp8		highest eigenvalue n. 8 of Burden matrix / weighted by atomic polarizabilities	Burden eigenvalues	2
1	BELv1		lowest eigenvalue n. 1 of Burden matrix / weighted by atomic van der Waals volumes	Burden eigenvalues	2
1	DISPe		d COMMA2 value / weighted by atomic Sanderson electronegativities	geometrical descriptors	3
2	EEig9d		Eigenvalue 09 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2
2	BEHp5		highest eigenvalue n. 5 of Burden matrix / weighted by atomic polarizabilities	Burden eigenvalues	2
1	E2e		2nd component accessibility directional WHIM index / weighted by atomic Sanderson electronegativities	WHIM descriptors	3
1	Mor25m		3D-MoRSE - signal 25 / weighted by atomic masses	3D-MoRSE descriptors	3
1	B03[C-C]		presence/absence of C - C at topological distance 03	2D binary fingerprints	2
3	B07[C-C]		presence/absence of C - C at topological distance 07	2D binary fingerprints	2
1	B01[C-O]		presence/absence of C - O at topological distance 01	2D binary fingerprints	2
1	Atype_H_49		Number of Hydrogen Type 49	atomtypes (Cerius2)	1
1	Infective-80		Ghose-Vishwanadhan-Windolowski antiinfective-like index at 80%	molecular properties	1
2	O-057		phenol / enol / carboxyl OH	atom-centred fragments	1
1	Mor22m		3D-MoRSE - signal 22 / weighted by atomic masses	3D-MoRSE descriptors	3
1	EEig10d		Eigenvalue 10 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2
1	R1u+		R maximal autocorrelation of lag 1 / unweighted	GETAWAY descriptors	3
1	GATS7m		Geary autocorrelation - lag 7 / weighted by atomic masses	2D autocorrelations	2
1	MATS4v		Moran autocorrelation - lag 4 / weighted by atomic van der Waals volumes	2D autocorrelations	2
1	R4e+		R maximal autocorrelation of lag 4 / weighted by atomic Sanderson electronegativities	GETAWAY descriptors	3
1	G3p		3st component symmetry directional WHIM index / weighted by atomic polarizabilities	WHIM descriptors	3
1	Hy		hydrophilic factor	molecular properties	1
1	S_dssC		S_dssC	atomtypes (Cerius2)	1
1	nRCHO		number of aldehydes (aliphatic)	functional group counts	1
1	B08[C-C]		presence/absence of C - C at topological distance 08	2D binary fingerprints	2
1	R2m		R autocorrelation of lag 2 / weighted by atomic masses	GETAWAY descriptors	3
1	HATS5e		leverage-weighted autocorrelation of lag 5 / weighted by atomic Sanderson electronegativities	GETAWAY descriptors	3
1	D/Dr06		distance/detour ring index of order 6	topological descriptors	2
1	RDF030m		Radial Distribution Function - 3.0 / weighted by atomic masses	RDF descriptors	3
2	Jhetv		Balaban-type index from van der Waals weighted distance matrix	topological descriptors	2
Or10a (11 Unique)					
3	S_dO		S_dO	atomtypes (Cerius2)	1
1	BEHm7		highest eigenvalue n. 7 of Burden matrix / weighted by atomic masses	Burden eigenvalues	2
1	E2u		2nd component accessibility directional WHIM index / unweighted	WHIM descriptors	3
1	HATS8m		leverage-weighted autocorrelation of lag 8 / weighted by atomic masses	GETAWAY descriptors	3
1	BELe4		lowest eigenvalue n. 4 of Burden matrix / weighted by atomic Sanderson electronegativities	Burden eigenvalues	2
1	Mor25e		3D-MoRSE - signal 25 / weighted by atomic Sanderson electronegativities	3D-MoRSE descriptors	3
1	B08[C-C]		presence/absence of C - C at topological distance 08	2D binary fingerprints	2
1	JG13		mean topological charge index of order 3	topological charge indices	2
1	ESpm03u		Spectral moment 03 from edge adj. matrix	edge adjacency indices	2
1	nR=Ct		number of aliphatic tertiary C(sp2)	functional group counts	1
2	E2e		2nd component accessibility directional WHIM index / weighted by atomic Sanderson electronegativities	WHIM descriptors	3
Or19a (25 Unique)					
1	Mor31p		3D-MoRSE - signal 31 / weighted by atomic polarizabilities	3D-MoRSE descriptors	3
1	H2m		H autocorrelation of lag 2 / weighted by atomic masses	GETAWAY descriptors	3
1	L1m		1st component size directional WHIM index / weighted by atomic masses	WHIM descriptors	3
1	R1m+		R maximal autocorrelation of lag 1 / weighted by atomic masses	GETAWAY descriptors	3
1	Mor27u		3D-MoRSE - signal 27 / unweighted	3D-MoRSE descriptors	3
1	HATS5u		leverage-weighted autocorrelation of lag 5 / unweighted	GETAWAY descriptors	3
3	GG17		topological charge index of order 7	topological charge indices	2
1	Gs		G total symmetry index / weighted by atomic electrotopological states	WHIM descriptors	3
1	O-057		phenol / enol / carboxyl OH	atom-centred fragments	1
1	H-049		H attached to C3(sp3)/C2(sp2)/C3(sp2)/C3(sp)	atom-centred fragments	1
1	pR08		molecular multiple path count of order 08	walk and path counts	2
2	R7u+		R maximal autocorrelation of lag 7 / unweighted	GETAWAY descriptors	3
2	G3s		3st component symmetry directional WHIM index / weighted by atomic electrotopological states	WHIM descriptors	3
1	R4m+		R maximal autocorrelation of lag 4 / weighted by atomic masses	GETAWAY descriptors	3
1	MATS7p		Moran autocorrelation - lag 7 / weighted by atomic polarizabilities	2D autocorrelations	2
1	R6u+		R maximal autocorrelation of lag 6 / unweighted	GETAWAY descriptors	3

Table 2.1

	1Hy	hydrophilic factor	molecular properties	1
	1ARR	aromatic ratio	constitutional descriptors	0
	1BEHp7	highest eigenvalue n. 7 of Burden matrix / weighted by atomic polarizabilities	Burden eigenvalues	2
	1RDF050v	Radial Distribution Function - 5.0 / weighted by atomic van der Waals volumes	RDF descriptors	3
	1C-005	CH3X	atom-centred fragments	3
	1nRCHO	number of aldehydes (aliphatic)	functional group counts	1
	1nRCOOH	number of carboxylic acids (aliphatic)	functional group counts	1
	1R5m+	R maximal autocorrelation of lag 5 / weighted by atomic masses	GETAWAY descriptors	3
	2C-002	CH2R2	atom-centred fragments	1
Or22a (43 Unique)	1Mor29v	3D-MoRSE - signal 29 / weighted by atomic van der Waals volumes	3D-MoRSE descriptors	3
	1MAXDN	maximal electrotopological negative variation	topological descriptors	2
	1pRCD4	molecular multiple path count of order 04	walk and path counts	2
	1Mor10e	3D-MoRSE - signal 10 / weighted by atomic Sanderson electronegativities	3D-MoRSE descriptors	3
	1Mor27m	3D-MoRSE - signal 27 / weighted by atomic masses	3D-MoRSE descriptors	3
	1R7p+	R maximal autocorrelation of lag 7 / weighted by atomic polarizabilities	GETAWAY descriptors	3
	1S_sCH3	S_sCH3	atomtypes (Cerius2)	1
	2EEg12r	Eigenvalue 12 from edge adj. matrix weighted by resonance integrals	edge adjacency indices	2
	1nRCOOR	number of esters (aliphatic)	functional group counts	1
	4R6u+	R maximal autocorrelation of lag 6 / unweighted	GETAWAY descriptors	3
	1Mor32p	3D-MoRSE - signal 32 / weighted by atomic polarizabilities	3D-MoRSE descriptors	3
	1AlogP98	AlogP98 value	structural (Cerius2)	0
	4O-057	phenol / enol / carboxyl OH	atom-centred fragments	1
	1L3s	3rd component size directional WHIM index / weighted by atomic electrotopological states	WHIM descriptors	3
	1R1v+	R maximal autocorrelation of lag 1 / weighted by atomic van der Waals volumes	GETAWAY descriptors	3
	2nHDon	number of donor atoms for H-bonds (N and O)	functional group counts	1
	2B10[C-C]	presence/absence of C - C at topological distance 10	2D binary fingerprints	2
	1Mor18m	3D-MoRSE - signal 18 / weighted by atomic masses	3D-MoRSE descriptors	3
	1B06[C-O]	presence/absence of C - O at topological distance 04	2D binary fingerprints	2
	2Jhetp	Balaban-type index from polignol distance matrix	topological descriptors	2
	1STN	spanning tree number (log)	topological descriptors	2
	2ESpm15u	Spectral moment 15 from edge adj. matrix	edge adjacency indices	2
	1GATS1v	Geary autocorrelation - lag 1 / weighted by atomic van der Waals volumes	2D autocorrelations	2
	1F03[O-O]	frequency of O - O at topological distance 03	2D frequency fingerprints	2
	1GATS8m	Geary autocorrelation - lag 8 / weighted by atomic masses	2D autocorrelations	2
	2HATS5e	leverage-weighted autocorrelation of lag 5 / weighted by atomic Sanderson electronegativities	GETAWAY descriptors	3
	1DISPv	d COMMA2 value / weighted by atomic van der Waals volumes	geometrical descriptors	3
	1R3v+	R maximal autocorrelation of lag 3 / weighted by atomic van der Waals volumes	GETAWAY descriptors	3
	1E2e	2nd component accessibility directional WHIM index / weighted by atomic Sanderson electronegativities	WHIM descriptors	3
	1Mor32u	3D-MoRSE - signal 32 / unweighted	3D-MoRSE descriptors	3
	2B02[O-O]	presence/absence of O - O at topological distance 02	2D binary fingerprints	2
	1G3e	3st component symmetry directional WHIM index / weighted by atomic Sanderson electronegativities	WHIM descriptors	3
	1HCRs	number of ring secondary C(sp3)	functional group counts	1
	2HOMT	HOMA total	geometrical descriptors	3
	1B05[C-C]	presence/absence of C - C at topological distance 05	2D binary fingerprints	2
	1MATS7m	Moran autocorrelation - lag 7 / weighted by atomic masses	2D autocorrelations	2
	1RDF100m	Radial Distribution Function - 3.0 / weighted by atomic masses	RDF descriptors	3
	1EEg12x	Eigenvalue 12 from edge adj. matrix weighted by edge degrees	edge adjacency indices	2
	1R1m+	R maximal autocorrelation of lag 1 / weighted by atomic masses	GETAWAY descriptors	3
	1MATS4p	Moran autocorrelation - lag 4 / weighted by atomic polarizabilities	2D autocorrelations	2
	1B09[C-O]	presence/absence of C - O at topological distance 09	2D binary fingerprints	2
	1Mor15p	3D-MoRSE - signal 15 / weighted by atomic polarizabilities	3D-MoRSE descriptors	3
	2S_sOH	S_sOH	atomtypes (Cerius2)	1
Or23a (37 Unique)	1ATS3p	Broto-Moreau autocorrelation of a topological structure - lag 3 / weighted by atomic polarizabilities	2D autocorrelations	2
	2O-056	alcohol	atom-centred fragments	1
	1J3D	3D-Balaban index	geometrical descriptors	3
	1BELm5	lowest eigenvalue n. 5 of Burden matrix / weighted by atomic masses	Burden eigenvalues	2
	1TPSA(Tot)	topological polar surface area using N,O,S,P polar contributions	molecular properties	1
	1B08[C-O]	presence/absence of C - O at topological distance 08	2D binary fingerprints	2
	2Mor27v	3D-MoRSE - signal 27 / weighted by atomic van der Waals volumes	3D-MoRSE descriptors	3
	2R6u+	R maximal autocorrelation of lag 6 / unweighted	GETAWAY descriptors	3
	1DISPv	d COMMA2 value / weighted by atomic Sanderson electronegativities	geometrical descriptors	3
	1ESpm12d	Spectral moment 12 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2
	1Mor17m	3D-MoRSE - signal 17 / weighted by atomic masses	3D-MoRSE descriptors	3
	2EEg09d	Eigenvalue 09 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2
	1Hy	hydrophilic factor	molecular properties	1
	2GATS3e	Geary autocorrelation - lag 3 / weighted by atomic Sanderson electronegativities	2D autocorrelations	2
	1GATS8m	Geary autocorrelation - lag 8 / weighted by atomic masses	2D autocorrelations	2
	1R4e+	R maximal autocorrelation of lag 4 / weighted by atomic Sanderson electronegativities	GETAWAY descriptors	3
	1Mor18m	3D-MoRSE - signal 18 / weighted by atomic masses	3D-MoRSE descriptors	3
	2nRCOOH	number of carboxylic acids (aliphatic)	functional group counts	1
	1S_sOH	S_sOH	atomtypes (Cerius2)	1
	1E3m	3rd component accessibility directional WHIM index / weighted by atomic masses	WHIM descriptors	3
	1G3s	3st component symmetry directional WHIM index / weighted by atomic electrotopological states	WHIM descriptors	3
	2BELm6	lowest eigenvalue n. 6 of Burden matrix / weighted by atomic masses	Burden eigenvalues	2
	1GATS1m	Geary autocorrelation - lag 1 / weighted by atomic masses	2D autocorrelations	2
	2EEg08d	Eigenvalue 08 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2
	1F05[C-O]	frequency of C - O at topological distance 05	2D frequency fingerprints	2
	2nHDon	number of donor atoms for H-bonds (N and O)	functional group counts	1
	1EEg10d	Eigenvalue 10 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2
	1R5p+	R maximal autocorrelation of lag 5 / weighted by atomic polarizabilities	GETAWAY descriptors	3
	1BIC	BIC	topological (Cerius2)	1
	2Infective-80	Ghose-Viswanadhan-Wendoloski antiinfective-like index at 80%	molecular properties	1
	1GATS4p	Geary autocorrelation - lag 4 / weighted by atomic polarizabilities	2D autocorrelations	2
	1DISPp	d COMMA2 value / weighted by atomic polarizabilities	geometrical descriptors	3
	1O-057	phenol / enol / carboxyl OH	atom-centred fragments	1
	1Atype_H_49	Number of Hydrogen Type 49	atomtypes (Cerius2)	1
	1GATS5m	Geary autocorrelation - lag 5 / weighted by atomic masses	2D autocorrelations	2
	1B02[O-O]	presence/absence of O - O at topological distance 02	2D binary fingerprints	2
	2GIS	mean topological charge index of order5	topological charge indices	2
Or33b (32 Unique)	6O-057	phenol / enol / carboxyl OH	atom-centred fragments	1
	2EEg08x	Eigenvalue 08 from edge adj. matrix weighted by edge degrees	edge adjacency indices	2
	1DISPv	d COMMA2 value / weighted by atomic van der Waals volumes	geometrical descriptors	3
	1TPSA(NO)	topological polar surface area using N,O polar contributions	molecular properties	1
	5B06[C-C]	presence/absence of C - C at topological distance 06	2D binary fingerprints	2
	4Atype_H_49	Number of Hydrogen Type 49	atomtypes (Cerius2)	1
	2R3v+	R maximal autocorrelation of lag 3 / weighted by atomic van der Waals volumes	GETAWAY descriptors	3
	1G1e	1st component symmetry directional WHIM index / weighted by atomic Sanderson electronegativities	WHIM descriptors	3
	1R2m+	R maximal autocorrelation of lag 2 / weighted by atomic masses	GETAWAY descriptors	3
	4B05[C-O]	presence/absence of C - O at topological distance 05	2D binary fingerprints	2
	1C-006	CH2RX	atom-centred fragments	1
	2TPSA(Tot)	topological polar surface area using N,O,S,P polar contributions	molecular properties	1
	1L/Bw	length-to-breadth ratio by WHIM	geometrical descriptors	3
	1EEg08d	Eigenvalue 08 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2
	3F04[C-O]	frequency of C - O at topological distance 04	2D frequency fingerprints	2
	1BEHvS	highest eigenvalue n. 5 of Burden matrix / weighted by atomic van der Waals volumes	Burden eigenvalues	1
	1Mor30p	3D-MoRSE - signal 30 / weighted by atomic polarizabilities	3D-MoRSE descriptors	3
	1nArCO	number of ketones (aromatic)	functional group counts	1
	1nRCO	number of ketones (aliphatic)	functional group counts	1
	1R1v+	R maximal autocorrelation of lag 1 / weighted by atomic polarizabilities	GETAWAY descriptors	3
	1MATS4p	Moran autocorrelation - lag 4 / weighted by atomic polarizabilities	2D autocorrelations	2
	1nN	number of Nitrogen atoms	constitutional descriptors	0
	1B07[C-C]	presence/absence of C - C at topological distance 07	2D binary fingerprints	2

Table 2.1 Continued

	2]G14	mean topological charge index of order4	topological charge indices	2
	1]rCOOH	number of carboxylic acids (aliphatic)	functional group counts	1
	1]nCconJ	number of non-aromatic conjugated C(sp2)	functional group counts	1
	1]C-005	CH3X	atom-centred fragments	1
	1]G13	mean topological charge index of order3	topological charge indices	2
	1]HATS3p	leverage-weighted autocorrelation of lag 3 / weighted by atomic polarizabilities	GETAWAY descriptors	3
	1]HATS8u	leverage-weighted autocorrelation of lag 8 / unweighted	GETAWAY descriptors	3
	1]E2u	2nd component accessibility directional WHIM index / unweighted	WHIM descriptors	3
	2]H-051	H attached to alpha-C	atom-centred fragments	1
Or35a (51 Unique)	1]ATS4e	Broto-Moreau autocorrelation of a topological structure - lag 4 / weighted by atomic Sanderson electronegativity	2D autocorrelations	2
	2]TPSA(NO)	topological polar surface area using N,O polar contributions	molecular properties	1
	1]Mor27p	3D-MoRSE - signal 27 / weighted by atomic polarizabilities	3D-MoRSE descriptors	3
	8]R6p+	R maximal autocorrelation of lag 6 / weighted by atomic polarizabilities	GETAWAY descriptors	3
	6]rCOOH	number of carboxylic acids (aliphatic)	functional group counts	1
	3]EEig10d	Eigenvalue 10 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2
	2]G5	G total symmetry index / weighted by atomic electrotopological states	WHIM descriptors	3
	9]IG12	mean topological charge index of order2	topological charge indices	2
	3]EEig12r	Eigenvalue 12 from edge adj. matrix weighted by resonance integrals	edge adjacency indices	2
	7]R4e+	R maximal autocorrelation of lag 4 / weighted by atomic Sanderson electronegativities	GETAWAY descriptors	3
	7]Mor28e	3D-MoRSE - signal 28 / weighted by atomic Sanderson electronegativities	3D-MoRSE descriptors	3
	2]MATS7p	Moran autocorrelation - lag 7 / weighted by atomic polarizabilities	2D autocorrelations	2
	2]L3s	3rd component size directional WHIM index / weighted by atomic electrotopological states	WHIM descriptors	3
	6]Mor25v	3D-MoRSE - signal 25 / weighted by atomic van der Waals volumes	3D-MoRSE descriptors	3
	4]Mor30e	3D-MoRSE - signal 30 / weighted by atomic Sanderson electronegativities	3D-MoRSE descriptors	3
	1]HATS8u	leverage-weighted autocorrelation of lag 8 / unweighted	GETAWAY descriptors	3
	7]O-057	phenol / enol / carboxyl OH	atom-centred fragments	1
	3]HATS5m	leverage-weighted autocorrelation of lag 5 / weighted by atomic masses	GETAWAY descriptors	3
	3]Jhetp	Balaban-type index from polarizability weighted distance matrix	topological descriptors	2
	4]IG18	mean topological charge index of order8	topological charge indices	2
	3]Mor04m	3D-MoRSE - signal 04 / weighted by atomic masses	3D-MoRSE descriptors	3
	1]S_dssC	S_dssC	atomtypes (Cerius2)	1
	2]E1m	1st component accessibility directional WHIM index / weighted by atomic masses	WHIM descriptors	3
	2]nHDon	number of donor atoms for H-bonds (N and O)	functional group counts	1
	2]RDF135u	Radial Distribution Function - 13.5 / unweighted	RDF descriptors	3
	2]D/Dr06	distance/detour ring index of order 6	topological descriptors	2
	3]E2s	2nd component accessibility directional WHIM index / weighted by atomic electrotopological states	WHIM descriptors	3
	2]EEig10r	Eigenvalue 10 from edge adj. matrix weighted by resonance integrals	edge adjacency indices	2
	1]G2s	2st component symmetry directional WHIM index / weighted by atomic electrotopological states	WHIM descriptors	3
	3]GATS3p	Geary autocorrelation - lag 3 / weighted by atomic polarizabilities	2D autocorrelations	2
	2]GG1	topological charge index of order 1	topological charge indices	2
	2]Atype_C_18	Number of Carbon Type 18	atomtypes (Cerius2)	1
	1]rRCO	number of ketones (aliphatic)	functional group counts	1
	1]C-005	CH3X	atom-centred fragments	1
	1]Mor27u	3D-MoRSE - signal 27 / unweighted	3D-MoRSE descriptors	3
	2]F08[C-O]	frequency of C - O at topological distance 08	2D frequency fingerprints	2
	3]G3s	3rd component symmetry directional WHIM index / weighted by atomic electrotopological states	WHIM descriptors	3
	3]SIC5	structural information content (neighborhood symmetry of 5-order)	information indices	2
	1]G(N..N)	sum of geometrical distances between N..N	geometrical descriptors	3
	2]nR=Ct	number of aliphatic tertiary C(sp2)	functional group counts	1
	2]E3m	3rd component accessibility directional WHIM index / weighted by atomic masses	WHIM descriptors	3
	1]nArCOOR	number of esters (aromatic)	functional group counts	1
	1]HATS6m	leverage-weighted autocorrelation of lag 6 / weighted by atomic masses	GETAWAY descriptors	3
	1]nArCO	number of ketones (aromatic)	functional group counts	1
	1]Jhete	Balaban-type index from electronegativity weighted distance matrix	topological descriptors	2
	1]G(O..O)	sum of geometrical distances between O..O	geometrical descriptors	3
	1]nCl	number of total tertiary C(sp3)	functional group counts	1
	1]H-051	H attached to alpha-C	atom-centred fragments	1
	1]nN	number of Nitrogen atoms	constitutional descriptors	0
	1]P2s	2nd component shape directional WHIM index / weighted by atomic electrotopological states	WHIM descriptors	3
	1]C-025	R--CR--R	atom-centred fragments	1
Or43a (27 Unique)	2]O-056	alcohol	atom-centred fragments	1
	1]BELm5	lowest eigenvalue n..5 of Burden matrix / weighted by atomic masses	Burden eigenvalues	2
	1]B07[C-O]	presence/absence of C - O at topological distance 07	2D binary fingerprints	2
	1]R5e	R autocorrelation of lag 5 / weighted by atomic Sanderson electronegativities	GETAWAY descriptors	3
	1]TPSA(Tot)	topological polar surface area using N,O,S,P polar contributions	molecular properties	1
	1]R6e+	R maximal autocorrelation of lag 6 / weighted by atomic Sanderson electronegativities	GETAWAY descriptors	3
	2]G17	mean topological charge index of order7	topological charge indices	2
	3]B04[C-C]	presence/absence of C - C at topological distance 04	2D binary fingerprints	2
	1]EEig10d	Eigenvalue 10 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2
	5]B02[O-O]	presence/absence of O - O at topological distance 02	2D binary fingerprints	2
	3]Mor13m	3D-MoRSE - signal 13 / weighted by atomic masses	3D-MoRSE descriptors	3
	3]nHDon	number of donor atoms for H-bonds (N and O)	functional group counts	1
	1]Mor21m	3D-MoRSE - signal 21 / weighted by atomic masses	3D-MoRSE descriptors	3
	1]IX	IX	topological (Cerius2)	2
	1]R1m+	R maximal autocorrelation of lag 1 / weighted by atomic masses	GETAWAY descriptors	3
	2]GATS7m	Geary autocorrelation - lag 7 / weighted by atomic masses	2D autocorrelations	2
	1]BELm6	lowest eigenvalue n..6 of Burden matrix / weighted by atomic masses	Burden eigenvalues	2
	1]E3m	3rd component accessibility directional WHIM index / weighted by atomic masses	WHIM descriptors	3
	2]MATS3e	Moran autocorrelation - lag 3 / weighted by atomic Sanderson electronegativities	2D autocorrelations	2
	1]F04[C-O]	frequency of C - O at topological distance 04	2D frequency fingerprints	2
	1]nRCHO	number of aldehydes (aliphatic)	functional group counts	1
	1]Infective-80	Ghose-Viswanadhan-Wendoloski antiinfective-like index at 80%	molecular properties	1
	1]EEig09x	Eigenvalue 09 from edge adj. matrix weighted by edge degrees	edge adjacency indices	2
	1]GATS1m	Geary autocorrelation - lag 1 / weighted by atomic masses	2D autocorrelations	2
	1]CI2C	complementary information content (neighborhood symmetry of 2-order)	information indices	2
	1]EEig10d	Eigenvalue 10 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2
	1]HATS6u	leverage-weighted autocorrelation of lag 6 / unweighted	GETAWAY descriptors	3
Or43b (29 Unique)	1]EEig04x	Eigenvalue 04 from edge adj. matrix weighted by edge degrees	edge adjacency indices	2
	1]BEHv4	highest eigenvalue n..4 of Burden matrix / weighted by atomic van der Waals volumes	Burden eigenvalues	2
	1]Mor25e	3D-MoRSE - signal 25 / weighted by atomic Sanderson electronegativities	3D-MoRSE descriptors	3
	2]EEig09d	Eigenvalue 09 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2
	1]E1p	1st component accessibility directional WHIM index / weighted by atomic polarizabilities	WHIM descriptors	3
	1]BEHe8	highest eigenvalue n..8 of Burden matrix / weighted by atomic Sanderson electronegativities	Burden eigenvalues	2
	1]R1m+	R maximal autocorrelation of lag 1 / weighted by atomic masses	GETAWAY descriptors	3
	2]B07[C-C]	presence/absence of C - C at topological distance 07	2D binary fingerprints	2
	1]MAXDN	maximal electrotopological negative variation	topological descriptors	2
	1]O-057	phenol / enol / carboxyl OH	atom-centred fragments	1
	1]Infective-80	Ghose-Viswanadhan-Wendoloski antiinfective-like index at 80%	molecular properties	1
	3]B04[C-C]	presence/absence of C - C at topological distance 04	2D binary fingerprints	2
	1]MATS5e	Moran autocorrelation - lag 5 / weighted by atomic Sanderson electronegativities	2D autocorrelations	2
	1]Mor24v	3D-MoRSE - signal 24 / weighted by atomic van der Waals volumes	3D-MoRSE descriptors	3
	1]Mor25v	3D-MoRSE - signal 25 / weighted by atomic van der Waals volumes	3D-MoRSE descriptors	3
	1]BEHp4	highest eigenvalue n..4 of Burden matrix / weighted by atomic polarizabilities	Burden eigenvalues	2
	1]S_sCH3	S_sCH3	atomtypes (Cerius2)	1
	1]HATS3p	leverage-weighted autocorrelation of lag 3 / weighted by atomic polarizabilities	GETAWAY descriptors	3
	1]H7m	H autocorrelation of lag 7 / weighted by atomic masses	GETAWAY descriptors	3
	1]G17	mean topological charge index of order7	topological charge indices	2
	1]STN	spanning tree number (log)	topological descriptors	2
	1]rRCOOH	number of carboxylic acids (aliphatic)	functional group counts	1
	1]MATS6m	Moran autocorrelation - lag 6 / weighted by atomic masses	2D autocorrelations	2
	1]HATS1u	leverage-weighted autocorrelation of lag 1 / unweighted	GETAWAY descriptors	3
	1]EEig10d	Eigenvalue 10 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2

Table 2.1 Continued

Or47a (21 Unique)	1 Atype_H_49	Number of Hydrogen Type 49	atomtypes (Cerius2)	1	
	1 EEig08d	Eigenvalue 08 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2	
	1 nCrS	number of ring secondary C(sp3)	functional group counts	1	
	2 H-047	H attached to C1(sp3)/C0(sp2)	atom-centred fragments	1	
	1 piPC04	molecular multiple path count of order 04	walk and path counts	2	
	2 DISPm	d COMMA2 value / weighted by atomic masses	geometrical descriptors	3	
	1 R7e+	R maximal autocorrelation of lag 7 / weighted by atomic Sanderson electronegativities	GETAWAY descriptors	3	
	1 Mor10p	3D-MoRSE - signal 10 / weighted by atomic polarizabilities	3D-MoRSE descriptors	3	
	1 Mor20u	3D-MoRSE - signal 20 / unweighted	3D-MoRSE descriptors	3	
	1 IC1	information content index (neighborhood symmetry of 1-order)	information indices	2	
	1 nRCOOH	number of carboxylic acids (aliphatic)	functional group counts	1	
	1 EEig01d	Eigenvalue 01 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2	
	2 Infective-80	Ghose-Viswanadhan-Wendolowski antiinfective-like index at 80%	molecular properties	1	
	1 MATS4m	Moran autocorrelation - lag 4 / weighted by atomic masses	2D autocorrelations	2	
	1 GATS5p	Geary autocorrelation - lag 5 / weighted by atomic polarizabilities	2D autocorrelations	2	
	1 PW4	path/walk 4 - Randic shape index	topological descriptors	2	
	1 Mor32p	3D-MoRSE - signal 32 / weighted by atomic polarizabilities	3D-MoRSE descriptors	2	
	1 Mor09e	3D-MoRSE - signal 09 / weighted by atomic Sanderson electronegativities	3D-MoRSE descriptors	3	
	1 TPSA(NO)	topological polar surface area using N,O polar contributions	molecular properties	1	
	1 B04[C-C]	presence/absence of C - C at topological distance 04	2D binary fingerprints	2	
	1 O-057	phenol / enol / carboxyl OH	atom-centred fragments	1	
	Or47b (14 Unique)	1 Atype_H_49	Number of Hydrogen Type 49	atomtypes (Cerius2)	1
		1 ESpm01d	Spectral moment 01 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2
		1 EEig10d	Eigenvalue 10 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2
		1 F2m	2nd component shape directional WHIM index / weighted by atomic masses	WHIM descriptors	2
2 Mor06e		3D-MoRSE - signal 06 / weighted by atomic Sanderson electronegativities	3D-MoRSE descriptors	3	
Or47c (14 Unique)		3 EEig02d	Eigenvalue 02 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2
		5 ESpm03d	Spectral moment 03 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2
		1 nHBonds	number of intramolecular H-bonds (with N,O,F)	functional group counts	1
		4 XSA	average connectivity index chi-5	connectivity indices	2
		1 EEig08x	Eigenvalue 08 from edge adj. matrix weighted by edge degrees	edge adjacency indices	2
		1 C-006	CH2RX	atom-centred fragments	1
		1 nRCHO	number of aldehydes (aliphatic)	functional group counts	1
		2 nRCOOR	number of esters (aliphatic)	functional group counts	1
		1 nRCOOH	number of carboxylic acids (aliphatic)	functional group counts	1
		1 EEig08d	Eigenvalue 08 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2
		1 X4Av	average valence connectivity index chi-4	connectivity indices	2
		1 GATS6m	Geary autocorrelation - lag 6 / weighted by atomic masses	2D autocorrelations	2
		1 EEig07r	Eigenvalue 07 from edge adj. matrix weighted by resonance integrals	edge adjacency indices	2
		1 R2m	R autocorrelation of lag 2 / weighted by atomic masses	GETAWAY descriptors	3
Or49b (37 Unique)		2 nCd-	number of substituted benzene C(sp2)	functional group counts	1
		1 BEHm6	highest eigenvalue n, 6 of Burden matrix / weighted by atomic masses	Burden eigenvalues	2
		2 F04[C-O]	frequency of C - O at topological distance 04	2D frequency fingerprints	2
		1 D/Dr06	distance/detour ring index of order 6	topological descriptors	2
		1 BEHm6	highest eigenvalue n, 6 of Burden matrix / weighted by atomic polarizabilities	Burden eigenvalues	2
		3 H-047	H attached to C1(sp3)/C0(sp2)	atom-centred fragments	2
	1 GATS1m	Geary autocorrelation - lag 1 / weighted by atomic masses	2D autocorrelations	1	
	3 HATS8p	leverage-weighted autocorrelation of lag 8 / weighted by atomic polarizabilities	GETAWAY descriptors	3	
	2 ISH	standardized information content on the leverage equality	GETAWAY descriptors	3	
	1 Mor16e	3D-MoRSE - signal 16 / weighted by atomic Sanderson electronegativities	3D-MoRSE descriptors	3	
	1 JG15	mean topological charge index of order5	topological charge indices	2	
	1 R8e+	R maximal autocorrelation of lag 8 / weighted by atomic Sanderson electronegativities	GETAWAY descriptors	3	
	1 Mor25e	3D-MoRSE - signal 25 / weighted by atomic Sanderson electronegativities	3D-MoRSE descriptors	3	
	2 EEig10d	Eigenvalue 10 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2	
	1 Mor16p	3D-MoRSE - signal 16 / weighted by atomic polarizabilities	3D-MoRSE descriptors	3	
	1 JG14	mean topological charge index of order4	topological charge indices	2	
	1 MATS3p	Moran autocorrelation - lag 3 / weighted by atomic polarizabilities	2D autocorrelations	2	
	3 CIC	CIC	topological (Cerius2)	2	
	1 F2m	2nd component shape directional WHIM index / weighted by atomic masses	WHIM descriptors	2	
	1 nHDon	number of donor atoms for H-bonds (N and O)	functional group counts	1	
	1 Mor03m	3D-MoRSE - signal 03 / weighted by atomic masses	3D-MoRSE descriptors	3	
	2 JG17	mean topological charge index of order7	topological charge indices	2	
	1 Mor23v	3D-MoRSE - signal 23 / weighted by atomic van der Waals volumes	3D-MoRSE descriptors	3	
	1 Mor30e	3D-MoRSE - signal 30 / weighted by atomic Sanderson electronegativities	3D-MoRSE descriptors	3	
	1 IC	IC	topological (Cerius2)	2	
	1 Mor21m	3D-MoRSE - signal 21 / weighted by atomic masses	3D-MoRSE descriptors	3	
	1 Mor13m	3D-MoRSE - signal 13 / weighted by atomic masses	3D-MoRSE descriptors	3	
	1 R7v+	R maximal autocorrelation of lag 7 / weighted by atomic van der Waals volumes	GETAWAY descriptors	3	
	1 piPC07	molecular multiple path count of order 07	walk and path counts	2	
	1 nArOH	number of aromatic hydroxyls	functional group counts	1	
	1 Mor25v	3D-MoRSE - signal 25 / weighted by atomic van der Waals volumes	3D-MoRSE descriptors	3	
	1 Mor08v	3D-MoRSE - signal 08 / weighted by atomic van der Waals volumes	3D-MoRSE descriptors	3	
	1 R6e+	R maximal autocorrelation of lag 6 / weighted by atomic Sanderson electronegativities	GETAWAY descriptors	3	
	1 EEig06x	Eigenvalue 06 from edge adj. matrix weighted by edge degrees	edge adjacency indices	2	
	1 C-001	CH3R / CH4	atom-centred fragments	1	
	1 Mor07m	3D-MoRSE - signal 07 / weighted by atomic masses	3D-MoRSE descriptors	3	
	1 DISPe	d COMMA2 value / weighted by atomic Sanderson electronegativities	geometrical descriptors	3	
1 nR05	number of 5-membered rings	constitutional descriptors	0		
1 Mor07e	3D-MoRSE - signal 07 / weighted by atomic Sanderson electronegativities	3D-MoRSE descriptors	3		
1 EEig09x	Eigenvalue 09 from edge adj. matrix weighted by edge degrees	edge adjacency indices	2		
1 B05[C-O]	presence/absence of C - O at topological distance 05	2D binary fingerprints	2		
1 X5Av	average valence connectivity index chi-5	connectivity indices	2		
1 HATS3p	leverage-weighted autocorrelation of lag 3 / weighted by atomic polarizabilities	GETAWAY descriptors	3		
1 RBu+	R maximal autocorrelation of lag 8 / unweighted	GETAWAY descriptors	3		
1 O-060	Al-O-Ar / Ar-O-Ar / R...O...R / R-O-C=X	atom-centred fragments	1		
2 B04[C-O]	presence/absence of C - O at topological distance 04	2D binary fingerprints	2		
Or59b (23 Unique)	1 piPC06	molecular multiple path count of order 06	walk and path counts	2	
	1 R3u	R autocorrelation of lag 3 / unweighted	GETAWAY descriptors	3	
	1 S_sCH3	S_sCH3	atomtypes (Cerius2)	1	
	4 B06[C-C]	presence/absence of C - C at topological distance 06	2D binary fingerprints	2	
	1 R1e+	R maximal autocorrelation of lag 1 / weighted by atomic Sanderson electronegativities	GETAWAY descriptors	3	
	1 ESpm03u	Spectral moment 03 from edge adj. matrix	edge adjacency indices	2	
	1 EEig10r	Eigenvalue 10 from edge adj. matrix weighted by resonance integrals	edge adjacency indices	2	
	1 EEig08d	Eigenvalue 08 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2	
	1 E1u	1st component accessibility directional WHIM index / unweighted	WHIM descriptors	2	
	1 nCconj	number of non-aromatic conjugated C(sp2)	functional group counts	1	
	1 SP13	shape profile no. 13	Randic molecular profiles	3	
	2 S_00	S_00	atomtypes (Cerius2)	1	
	2 Atype_H_49	Number of Hydrogen Type 49	atomtypes (Cerius2)	1	
	1 EEig10d	Eigenvalue 10 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2	
	1 nHDon	number of donor atoms for H-bonds (N and O)	functional group counts	1	
	1 RBu+	R maximal autocorrelation of lag 8 / unweighted	GETAWAY descriptors	3	
	2 O-057	phenol / enol / carboxyl OH	atom-centred fragments	1	
	1 Mor10v	3D-MoRSE - signal 10 / weighted by atomic van der Waals volumes	3D-MoRSE descriptors	3	
	1 R5m+	R maximal autocorrelation of lag 5 / weighted by atomic masses	GETAWAY descriptors	3	
	1 Mor09e	3D-MoRSE - signal 09 / weighted by atomic Sanderson electronegativities	3D-MoRSE descriptors	3	
	1 nDhp	number of primary alcohols	functional group counts	1	
	1 EEig09d	Eigenvalue 09 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2	
	1 nCrS	number of ring secondary C(sp3)	functional group counts	1	
1 ESpm01d	Spectral moment 01 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2		

Table 2.1 Continued

Or65a (14 Unique)	1 F04(O-O)	frequency of O - O at topological distance 04	2D frequency fingerprints	2	
	2 Mor30m	3D-MoRSE - signal 30 / weighted by atomic masses	3D-MoRSE descriptors	3	
	4 Atype_H_51	Number of Hydrogen Type 51	atomtypes (Cerius2)	1	
	1 EEig08d	Eigenvalue 08 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2	
	2 nArOH	number of aromatic hydroxyls	functional group counts	1	
	2 JGI7	mean topological charge index of order7	topological charge indices	2	
	1 nHBonds	number of intramolecular H-bonds (with N,O,F)	functional group counts	1	
	1 Mor13p	3D-MoRSE - signal 13 / weighted by atomic polarizabilities	3D-MoRSE descriptors	3	
	1 EEig07d	Eigenvalue 07 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2	
	1 B06(C-O)	presence/absence of C - O at topological distance 06	2D binary fingerprints	2	
	1 C-008	CHR2X	atom-centred fragments	1	
	1 EEig08r	Eigenvalue 08 from edge adj. matrix weighted by resonance integrals	edge adjacency indices	2	
	1 B01(C-O)	presence/absence of C - O at topological distance 01	2D binary fingerprints	2	
	2 Mor32e	3D-MoRSE - signal 32 / weighted by atomic Sanderson electronegativities	3D-MoRSE descriptors	3	
	Or67a (37 Unique)	2 AlogP98	AlogP98 value	structural (Cerius2)	0
		8 B04(C-O)	presence/absence of C - O at topological distance 04	2D binary fingerprints	2
		6 F08(C-O)	frequency of C - O at topological distance 08	2D frequency fingerprints	2
		1 GG4	topological charge index of order 4	topological charge indices	2
		3 E2u	2nd component accessibility directional WHIM index / unweighted	WHIM descriptors	3
		2 O-057	phenol / enol / carboxyl OH	atom-centred fragments	1
		1 Mor03v	3D-MoRSE - signal 03 / weighted by atomic van der Waals volumes	3D-MoRSE descriptors	3
		4 X5A	average connectivity index chi-5	connectivity indices	2
		3 Mor10v	3D-MoRSE - signal 10 / weighted by atomic van der Waals volumes	3D-MoRSE descriptors	3
		1 B03(C-O)	presence/absence of C - O at topological distance 03	2D binary fingerprints	2
3 X4A		average connectivity index chi-4	connectivity indices	2	
3 nCt		number of total tertiary C(sp3)	functional group counts	1	
1 C-026		R-CX-R	atom-centred fragments	1	
3 RDF075m		Radial Distribution Function - 7.5 / weighted by atomic masses	RDF descriptors	3	
2 C-008		CHR2X	atom-centred fragments	1	
2 B03(C-C)		presence/absence of C - C at topological distance 03	2D binary fingerprints	2	
1 B01(C-O)		presence/absence of C - O at topological distance 01	2D binary fingerprints	2	
1 nRCHO		number of aldehydes (aliphatic)	functional group counts	1	
1 Jhetv		Balaban-type index from van der Waals weighted distance matrix	topological descriptors	2	
1 L1s		1st component size directional WHIM index / weighted by atomic electrotopological states	WHIM descriptors	3	
1 Hy		hydrophilic factor	molecular properties	1	
2 C-003		CHR3	atom-centred fragments	1	
1 GATS7m		Geary autocorrelation - lag 7 / weighted by atomic masses	2D autocorrelations	2	
1 Mor16e		3D-MoRSE - signal 16 / weighted by atomic Sanderson electronegativities	3D-MoRSE descriptors	3	
1 Mor06u		3D-MoRSE - signal 06 / unweighted	3D-MoRSE descriptors	3	
1 RDF030m		Radial Distribution Function - 3.0 / weighted by atomic masses	RDF descriptors	3	
1 Atype_C_18		Number of Carbon Type 18	atomtypes (Cerius2)	1	
1 F03(O-O)		frequency of O - O at topological distance 03	2D frequency fingerprints	2	
1 nCrs		number of ring secondary C(sp2)	functional group counts	1	
2 nArOH		number of aromatic hydroxyls	functional group counts	1	
1 GATS8m		Geary autocorrelation - lag 8 / weighted by atomic masses	2D autocorrelations	2	
1 Jhete		Balaban-type index from electronegativity weighted distance matrix	topological descriptors	2	
1 EEig13x		Eigenvalue 13 from edge adj. matrix weighted by edge degrees	edge adjacency indices	2	
1 DISPm		d COMMA2 value / weighted by atomic masses	GETAWAY descriptors	3	
1 X3A		average connectivity index chi-3	connectivity indices	2	
1 G(N..N)		sum of geometrical distances between N..N	geometrical descriptors	3	
1 Mor32u		3D-MoRSE - signal 32 / unweighted	3D-MoRSE descriptors	3	
Or67c (24 Unique)		1 BEHe8	highest eigenvalue n..8 of Burden matrix / weighted by atomic Sanderson electronegativities	Burden eigenvalues	2
		1 O-056	alcohol	atom-centred fragments	1
		1 Mor25m	3D-MoRSE - signal 25 / weighted by atomic masses	3D-MoRSE descriptors	3
	1 BELv4	lowest eigenvalue n..4 of Burden matrix / weighted by atomic van der Waals volumes	Burden eigenvalues	2	
	3 B07(C-C)	presence/absence of C - C at topological distance 07	2D binary fingerprints	2	
	1 TPSA(Tot)	topological polar surface area using N,O,S,P polar contributions	molecular properties	1	
	1 DISPm	d COMMA2 value / weighted by atomic masses	geometrical descriptors	3	
	4 HATS6u	leverage-weighted autocorrelation of lag 6 / unweighted	GETAWAY descriptors	3	
	2 EEig08d	Eigenvalue 08 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2	
	2 EEig10d	Eigenvalue 10 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2	
	1 G8	G total symmetry index / weighted by atomic electrotopological states	WHIM descriptors	3	
	3 O-057	phenol / enol / carboxyl OH	atom-centred fragments	1	
	1 B08(C-C)	presence/absence of C - C at topological distance 08	2D binary fingerprints	2	
	1 R1m+	R maximal autocorrelation of lag 1 / weighted by atomic masses	GETAWAY descriptors	3	
	1 BELm5	lowest eigenvalue n..5 of Burden matrix / weighted by atomic masses	Burden eigenvalues	2	
	1 F03(O-O)	frequency of O - O at topological distance 03	2D frequency fingerprints	2	
	1 STN	spanning tree number (log)	topological descriptors	2	
	1 Atype_H_49	Number of Hydrogen Type 49	atomtypes (Cerius2)	1	
	1 H-051	H attached to alpha-C	atom-centred fragments	1	
	1 B01(C-O)	presence/absence of C - O at topological distance 01	2D binary fingerprints	2	
	1 Infective-80	Ghose-Viswanadhan-Wendoloski anti-infective-like index at 80%	molecular properties	1	
	1 Hy	hydrophilic factor	molecular properties	1	
	1 Mor22m	3D-MoRSE - signal 22 / weighted by atomic masses	3D-MoRSE descriptors	3	
	1 JGI7	mean topological charge index of order7	topological charge indices	2	
	Or82a (31 Unique)	1 GG19	topological charge index of order 9	topological charge indices	2
		1 Mor02e	3D-MoRSE - signal 02 / weighted by atomic Sanderson electronegativities	3D-MoRSE descriptors	3
		1 Mor30v	3D-MoRSE - signal 30 / weighted by atomic van der Waals volumes	3D-MoRSE descriptors	3
		1 Mor02v	3D-MoRSE - signal 02 / weighted by atomic van der Waals volumes	3D-MoRSE descriptors	3
		1 Mor30u	3D-MoRSE - signal 30 / unweighted	3D-MoRSE descriptors	3
		2 BLTD48	Verhaar model of Daphnia base-line toxicity from MLOGP (mmol/l)	molecular properties	1
		2 Mor10v	3D-MoRSE - signal 10 / weighted by atomic van der Waals volumes	3D-MoRSE descriptors	3
		2 Atype_H_53	Number of Hydrogen Type 53	atomtypes (Cerius2)	1
		1 O-058	=O	atom-centred fragments	1
1 B02(C-O)		presence/absence of C - O at topological distance 02	2D binary fingerprints	2	
2 R5u+		R maximal autocorrelation of lag 5 / unweighted	GETAWAY descriptors	3	
1 H6e		H autocorrelation of lag 6 / weighted by atomic Sanderson electronegativities	GETAWAY descriptors	3	
1 MATS7p		Moran autocorrelation - lag 7 / weighted by atomic polarizabilities	2D autocorrelations	2	
1 GATS3p		Geary autocorrelation - lag 3 / weighted by atomic polarizabilities	2D autocorrelations	2	
1 Mor18m		3D-MoRSE - signal 18 / weighted by atomic masses	3D-MoRSE descriptors	3	
1 H-051		H attached to alpha-C	atom-centred fragments	1	
2 Mor13p		3D-MoRSE - signal 13 / weighted by atomic polarizabilities	3D-MoRSE descriptors	3	
1 SIC2		structural information content (neighborhood symmetry of 2-order)	information indices	2	
1 Mor32u		3D-MoRSE - signal 32 / unweighted	3D-MoRSE descriptors	3	
1 Mor10m		3D-MoRSE - signal 10 / weighted by atomic masses	3D-MoRSE descriptors	3	
1 nR=Cp		number of terminal primary C(sp2)	functional group counts	1	
1 Mor25p		3D-MoRSE - signal 25 / weighted by atomic polarizabilities	3D-MoRSE descriptors	3	
1 GATS8m		Geary autocorrelation - lag 8 / weighted by atomic masses	2D autocorrelations	2	
1 JGI1		mean topological charge index of order1	topological charge indices	2	
1 E-AD1-mag		E-AD1-mag	topological (cerius2)	2	
1 EEig11x		Eigenvalue 11 from edge adj. matrix weighted by edge degrees	edge adjacency indices	2	
1 B03(O-O)		3D-MoRSE - signal 03 / weighted by atomic Sanderson electronegativities	2D binary fingerprints	2	
1 Mor30e		3D-MoRSE - signal 30 / weighted by atomic Sanderson electronegativities	3D-MoRSE descriptors	3	
1 Rotlbonds		Number of rotatable bonds	structural (Cerius2)	0	
1 EEig09d	Eigenvalue 09 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2		
2 GATS7m	Geary autocorrelation - lag 7 / weighted by atomic masses	2D autocorrelations	2		
Or85a (15 Unique)	1 EEig04r	Eigenvalue 04 from edge adj. matrix weighted by resonance integrals	edge adjacency indices	2	
	2 C-006	CHR2X	atom-centred fragments	1	

Table 2.1 Continued

	3 ATS6e	Brotz-Morreu autocorrelation of a topological structure - lag 6 / weighted by atomic Sanderson electroneg	2D autocorrelations	2
	3 IG15	mean topological charge index of orders5	topological charge indices	2
	2 B07[C-C]	presence/absence of C - C at topological distance 07	2D binary fingerprints	2
	1 nCp	number of terminal primary C(sp3)	functional group counts	1
	2 DISPm	d COMMA2 value / weighted by atomic masses	geometrical descriptors	3
	2 GATS4m	Geary autocorrelation - lag 4 / weighted by atomic masses	2D autocorrelations	3
	1 Mor25p	3D-MoRSE - signal 25 / weighted by atomic polarizabilities	3D-MoRSE descriptors	3
	1 nHDon	number of donor atoms for H-bonds (N and O)	functional group counts	1
	1 EEig09d	Eigenvalue 09 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2
	1 R2m+	R maximal autocorrelation of lag 2 / weighted by atomic masses	GETAWAY descriptors	3
	1 IG14	mean topological charge index of order4	topological charge indices	2
	1 Mor11e	3D-MoRSE - signal 11 / weighted by atomic Sanderson electronegativities	3D-MoRSE descriptors	3
	2 HATS7m	leverage-weighted autocorrelation of lag 7 / weighted by atomic masses	GETAWAY descriptors	3
	Or85b (26 Unique)			
	1 piPCO5	molecular multiple path count of order 05	walk and path counts	2
	1 BLTF96	Verhaar model of Fish base-line toxicity from MLOGP (mmol/l)	molecular properties	1
	2 GATS4p	Geary autocorrelation - lag 4 / weighted by atomic polarizabilities	2D autocorrelations	2
	1 GG17	topological charge index of order7	topological charge indices	2
	3 B05[C-O]	presence/absence of C - O at topological distance 05	2D binary fingerprints	2
	2 O-057	phenol / enol / carboxyl OH	atom-centred fragments	1
	1 Mor27v	3D-MoRSE - signal 27 / weighted by atomic van der Waals volumes	3D-MoRSE descriptors	3
	1 HATS4v	leverage-weighted autocorrelation of lag 4 / weighted by atomic van der Waals volumes	GETAWAY descriptors	3
	1 Gs	G total symmetry index / weighted by atomic electrotopological states	WHIM descriptors	3
	2 Infective-80	Ghose-Viswanadhan-Wendoloski anti-infective-like index at 80%	molecular properties	1
	2 R7u+	R maximal autocorrelation of lag 7 / unweighted	GETAWAY descriptors	3
	2 nC3H	number of unsubstituted benzene C(sp2)	functional group counts	2
	1 B04[C-O]	presence/absence of C - O at topological distance 04	2D binary fingerprints	2
	2 IG17	mean topological charge index of order7	topological charge indices	2
	2 DISPe	d COMMA2 value / weighted by atomic Sanderson electronegativities	geometrical descriptors	3
	1 R4p+	R maximal autocorrelation of lag 4 / weighted by atomic polarizabilities	GETAWAY descriptors	3
	1 EEig12x	Eigenvalue 12 from edge adj. matrix weighted by edge degrees	edge adjacency indices	2
	1 B06[C-O]	presence/absence of C - O at topological distance 06	2D binary fingerprints	2
	1 MATS5e	Moran autocorrelation - lag 5 / weighted by atomic Sanderson electronegativities	2D autocorrelations	2
	1 MATS4m	leverage-weighted autocorrelation of lag 4 / weighted by atomic masses	GETAWAY descriptors	3
	1 MATS6u	leverage-weighted autocorrelation of lag 6 / unweighted	GETAWAY descriptors	3
	1 GATS4m	Geary autocorrelation - lag 4 / weighted by atomic masses	2D autocorrelations	2
	1 F03[O-O]	frequency of O - O at topological distance 03	2D frequency fingerprints	2
	1 HBv	H autocorrelation of lag 8 / weighted by atomic van der Waals volumes	GETAWAY descriptors	3
	1 EEig09d	Eigenvalue 09 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2
	2 Mor16e	3D-MoRSE - signal 16 / weighted by atomic Sanderson electronegativities	3D-MoRSE descriptors	3
	Or85f (53 Unique)			
	1 BEHp8	highest eigenvalue n. 8 of Burden matrix / weighted by atomic polarizabilities	Burden eigenvalues	2
	5 F05[C-O]	frequency of C - O at topological distance 05	2D frequency fingerprints	2
	4 BELM4	lowest eigenvalue n. 4 of Burden matrix / weighted by atomic masses	Burden eigenvalues	2
	1 MATS8m	leverage-weighted autocorrelation of lag 8 / weighted by atomic masses	GETAWAY descriptors	3
	2 B04[C-O]	presence/absence of C - O at topological distance 04	2D binary fingerprints	2
	6 O-057	phenol / enol / carboxyl OH	atom-centred fragments	1
	1 RDF030v	Radial Distribution Function - 3.0 / weighted by atomic van der Waals volumes	RDF descriptors	3
	1 GG17	topological charge index of order 7	topological charge indices	2
	1 Gs	G total symmetry index / weighted by atomic electrotopological states	WHIM descriptors	3
	4 B07[C-C]	presence/absence of C - C at topological distance 07	2D binary fingerprints	2
	1 E2e	2nd component accessibility directional WHIM index / weighted by atomic Sanderson electronegativities	WHIM descriptors	3
	1 MATS2m	Moran autocorrelation - lag 2 / weighted by atomic masses	2D autocorrelations	2
	2 Mor28u	3D-MoRSE - signal 28 / unweighted	3D-MoRSE descriptors	3
	3 BEHp5	highest eigenvalue n. 5 of Burden matrix / weighted by atomic polarizabilities	Burden eigenvalues	2
	2 Infective-80	Ghose-Viswanadhan-Wendoloski anti-infective-like index at 80%	molecular properties	1
	1 MATS4e	leverage-weighted autocorrelation of lag 4 / weighted by atomic Sanderson electronegativities	GETAWAY descriptors	3
	3 IG15	mean topological charge index of orders 4 / weighted by atomic masses	topological charge indices	3
	6 B05[C-O]	presence/absence of C - O at topological distance 05	2D binary fingerprints	2
	2 IG17	mean topological charge index of order7	topological charge indices	2
	2 DISPm	d COMMA2 value / weighted by atomic masses	geometrical descriptors	3
	5 RDF030m	Radial Distribution Function - 3.0 / weighted by atomic masses	RDF descriptors	3
	1 R1e+	R maximal autocorrelation of lag 1 / weighted by atomic Sanderson electronegativities	GETAWAY descriptors	3
	1 MATS8p	leverage-weighted autocorrelation of lag 8 / weighted by atomic polarizabilities	GETAWAY descriptors	3
	1 Atype_H_49	Number of Hydrogen Type 49	atomtypes (Cerius2)	1
	2 Hy	hydrophilic factor	molecular properties	1
	1 Jhetp	Balaban-type index from polarizability weighted distance matrix	topological descriptors	2
	1 HBv	H autocorrelation of lag 8 / weighted by atomic van der Waals volumes	GETAWAY descriptors	3
	2 EEig11d	Eigenvalue 11 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2
	1 MATS8m	Moran autocorrelation - lag 8 / weighted by atomic masses	2D autocorrelations	2
	1 MATS2p	Moran autocorrelation - lag 2 / weighted by atomic polarizabilities	2D autocorrelations	2
	4 B08[C-C]	presence/absence of C - C at topological distance 08	2D binary fingerprints	2
	1 S_sCH3	S_sCH3	atomtypes (Cerius2)	1
	2 MATS1e	leverage-weighted autocorrelation of lag 1 / weighted by atomic Sanderson electronegativities	GETAWAY descriptors	3
	1 nC3m	number of non-aromatic conjugated C(sp2)	functional group counts	1
	1 B04[C-C]	presence/absence of C - C at topological distance 04	2D binary fingerprints	2
	1 S_aasC	S_aasC	atomtypes (cerius2)	1
	1 R8m+	R maximal autocorrelation of lag 8 / weighted by atomic masses	GETAWAY descriptors	3
	1 nRCOOH	number of carboxylic acids (aliphatic)	functional group counts	1
	1 S_sOH	S_sOH	atomtypes (Cerius2)	1
	1 BELe3	lowest eigenvalue n. 3 of Burden matrix / weighted by atomic Sanderson electronegativities	Burden eigenvalues	2
	1 MATS8m	Geary autocorrelation - lag 8 / weighted by atomic masses	2D autocorrelations	2
	1 BEHp4	highest eigenvalue n. 4 of Burden matrix / weighted by atomic polarizabilities	Burden eigenvalues	2
	2 MATS5e	Moran autocorrelation - lag 5 / weighted by atomic Sanderson electronegativities	2D autocorrelations	2
	1 E3s	3rd component accessibility directional WHIM index / weighted by atomic electrotopological states	WHIM descriptors	3
	2 Jhetv	Balaban-type index from van der Waals weighted distance matrix	topological descriptors	2
	1 nR=C	number of aliphatic tertiary C(sp2)	functional group counts	1
	1 nRCHO	number of aldehydes (aliphatic)	functional group counts	1
	1 MATS8v	leverage-weighted autocorrelation of lag 8 / weighted by atomic van der Waals volumes	GETAWAY descriptors	3
	1 Mor28p	3D-MoRSE - signal 28 / weighted by atomic polarizabilities	3D-MoRSE descriptors	3
	1 C-003	CHR3	atom-centred fragments	1
	1 GATS7m	Geary autocorrelation - lag 7 / weighted by atomic masses	2D autocorrelations	3
	1 IG19	mean topological charge index of order9	topological charge indices	2
	1 B03[C-C]	presence/absence of C - C at topological distance 03	2D binary fingerprints	2
	Or88a (19 Unique)			
	3 nHBonds	number of intramolecular H-bonds (with N,O,F)	functional group counts	1
	2 nRCO	number of ketones (aliphatic)	functional group counts	1
	3 GATS6m	Geary autocorrelation - lag 6 / weighted by atomic masses	2D autocorrelations	2
	2 EEig08x	Eigenvalue 08 from edge adj. matrix weighted by edge degrees	edge adjacency indices	2
	1 nFuranes	number of Furanes	functional group counts	1
	1 nArCO	number of ketones (aromatic)	functional group counts	1
	1 ESpm15d	Spectral moment 15 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2
	1 C-005	CHX3	atom-centred fragments	1
	1 O-057	phenol / enol / carboxyl OH	atom-centred fragments	1
	1 L/Bw	length-to-breadth ratio by WHIM	geometrical descriptors	3
	1 nArCOOR	number of esters (aromatic)	functional group counts	1
	1 ESpm15u	Spectral moment 15 from edge adj. matrix	edge adjacency indices	2
	1 E2u	2nd component accessibility directional WHIM index / unweighted	WHIM descriptors	3
	1 EEig08d	Eigenvalue 08 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2
	1 H-051	H attached to alpha-C	atom-centred fragments	1
	1 ESpm14d	Spectral moment 14 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2
	1 GATS7m	Geary autocorrelation - lag 7 / weighted by atomic masses	2D autocorrelations	3
	1 PI13	3D Pettjean shape index	geometrical descriptors	3
	2 X3A	average connectivity index chi-3	connectivity indices	2

Table 2.1 Continued

Or98a (20 Unique)				
1	Lop	Lopping centric index	topological descriptors	2
4	O-057	phenol / enol / carboxyl OH	atom-centred fragments	1
2	B04[C-O]	presence/absence of C - O at topological distance 04	2D binary fingerprints	2
1	GVWAI-80	Ghose-Viswanadhan-Wendoloski drug-like index at 80%	molecular properties	1
1	HATS7p	leverage-weighted autocorrelation of lag 7 / weighted by atomic polarizabilities	GETAWAY descriptors	3
1	HATS5v	leverage-weighted autocorrelation of lag 5 / weighted by atomic van der Waals volumes	GETAWAY descriptors	3
1	MLOGP2	Squared Moriguchi octanol-water partition coeff. (logP ²)	molecular properties	1
2	GATS5e	Gearly autocorrelation - lag 5 / weighted by atomic Sanderson electronegativities	2D autocorrelations	2
1	H-049	H attached to C3(sp3)/C2(sp2)/C3(sp2)/C3(sp)	atom-centred fragments	1
1	WATS8m	Moran autocorrelation - lag 8 / weighted by atomic masses	2D autocorrelations	2
1	nCrS	number of ring secondary C(sp3)	functional group counts	1
3	HATS3p	leverage-weighted autocorrelation of lag 3 / weighted by atomic polarizabilities	GETAWAY descriptors	3
1	G1s	1st component symmetry directional WHIM index / weighted by atomic electrotopological states	WHIM descriptors	3
1	S_aasC	S_aasC	atomtypes (Ceriuss2)	1
1	SP18	shape profile no. 18	Randic molecular profiles	3
1	B05[C-C]	presence/absence of C - C at topological distance 05	2D binary fingerprints	2
1	JGI2	mean topological charge index of order2	topological charge indices	2
1	JGI8	mean topological charge index of order8	topological charge indices	2
1	XAA	average connectivity index chi-4	connectivity indices	2
1	H5e	H autocorrelation of lag 5 / weighted by atomic Sanderson electronegativities	GETAWAY descriptors	3

Table 2.1 Continued

Table 2.2: Top 100 predicted compounds for each *Drosophila* Or

Chemical name or Pubchem compound ID (CIDs), SMILES strings, and distances, of the top ~100 predicted compounds for each Or. All distances represent the minimum distance based on optimized descriptors to the previously known strongest active compound listed in the gray cells for that particular Or.

CHAPTER III:

Applying Chemical Informatics to Decode Odor Receptors of Several Important Disease Vector and Pest Insect Species as Well as Mammals

INTRODUCTION

Olfaction is exceedingly important for a broad array of species living in all ecological regions of this planet (van Naters and Carlson, 2006). While a major effort has been made by the scientific community to decode and analyze Olfactory receptors (Ors) in *Drosophila* over the last decade, there remains a vast number of important pest and disease vector species for which very little is known. Very recently, odor receptors of several high priority pest and disease vector species have been the focus of investigation (Carey et al., 2010; Ditzen et al., 2008; Ghaninia et al., 2007; Hill et al., 2009; Stanczyk et al., 2010; Syed and Leal, 2008; Turner et al., 2011). Olfactory receptors (Ors) or Odor Receptor Neurons (ORNs) in these species are being tested against relatively small panels of odors in an attempt to decode their response profiles. While these panels are limited in size and scope, they do generally consist of multiple functional groups and carbon chain lengths and often a number of odors in the panels were selected for their ecological significance. Humanity could greatly benefit from the identification of highly effective, safe, and receptor-specific odors for these species, ideally a number of which will prove to be behavior modifiers. We demonstrate that our chemical informatics pipeline can be applied to predict active odors for 5 important species.

Anopheles gambiae is the principle disease vector for malaria, which is a devastating disease in many parts of the world, including South America, Africa, India,

and Asia (World Health Organization., 2011). Malaria symptoms include vomiting, retinal damage, anemia, joint pain, convulsions and in some cases coma and death. The World Health Organization (WHO) estimated 216 million episodes of malaria in 2010, resulting in over 655,000 deaths (World Health Organization., 2011). While these numbers do represent a reduction in the incidence of malaria from previous years, growing resistances to current treatments and repellents have been identified. It is estimated that over 2 billion dollars were spent on financing malaria control in 2011, with higher predictions for future years if global control targets are to be achieved. To further complicate matters, current preventative medications are too expensive for widespread use in developing countries. Identification of safer odors with better organoleptic properties that activate or inhibit behaviorally important Ors, such as those used for detection of human hosts, could greatly aid in the global fight against vector-borne diseases, such as malaria. Recently, 110 odors were tested against 50 *Anopheles gambiae* Ors using the “empty neuron” expression system, resulting in newly identified activators and inhibitors (Carey et al., 2010). The 110 compound set consisted of an array of functional groups and sizes, providing a wonderful training set for our computational pipeline.

In addition to *Anopheles* Or responses, we are also interested in predicting odors for the CO₂ receptors of several important species. CO₂ is a primary long-range detection cue that mosquitoes use to identify and navigate towards their blood meal (Gillies, 1980; Grant and Oconnell, 1996). By navigating upstream in CO₂ plumes, mosquitoes are able to identify hosts at long ranges with high effectiveness (Carde and Gibson, 2010; Zwiebel and Takken, 2004). It is believed that other important attractive odors play additional roles in more specific host preference once the mosquito is in close

range (Gillies, 1980; Takken, 1996). As the CO₂ receptor is behaviorally very important, it is not surprising that the genes are conserved across mosquito species including *Anopheles gambiae*, *Aedes aegypti*, and *Culex pipiens*, all of which express a CO₂ receptor in the cpA neuron on the antenna. Interestingly, CO₂ detection is even important in *Drosophila*, which are able to detect CO₂ through gustatory receptors in their ab1C ORN (Turner and Ray, 2009). While no large unified panel of odors has been tested on the CO₂ neuron of the previously mentioned 3 species of mosquito and *Drosophila*, a number of smaller panels have been tested across each and a small panel was tested across all 4 (Turner et al., 2011; Turner and Ray, 2009). As the response profiles for the receptor across species appears to be at least moderately conserved, we integrate the responses of all species into a single unified odor response set, which we use to train our prediction pipeline.

The Asian Citrus Psyllid (ACP), *Diaphorina citri* Kuwayama, is posing a major challenge to the global citrus industry. While originally present only in Asia, this insect has recently invaded North America, including both Florida and Southern California (Halbert and Manjunath, 2004)(<http://www.californiacitrusthreat.com/>). While the insect itself does not cause major damage to citrus trees, the ACP is a vector for the bacterial disease Huanglongbing (HLB), also known as the Citrus Greening Disease, caused by *Candidatus Liberobacter* (Bove, 2006; Dagraca, 1991). If infected, the resulting HLB causes a chronic decline in fruit production. The limited fruit that is produced is often small, misshapen and bitter (Baldwin et al., 2010; Bove, 2006). Unfortunately, there is currently no cure for a tree once it is infected. Several attempts have been made to repel the insect, with limited success, or to super nutrient infected orange trees, which is ultimately not cost effective as a long-term solution (Onagbola et al., 2011). Recently,

our lab has decoded ORNs in the ACP antenna against a panel of 62 odors, several of which are known citrus vegetative volatiles (unpublished data). Each of these odors in the rhinarial plates on each segment of the antenna, identifying ORNs on rhinarial plates 2, 4, 6, and 7. Each rhinarial plate appears to house 3 ORNs, which have been labeled A, B, and C in decreasing order of spike amplitude. While very few ligands for the A neurons have been identified, a number of activators for the B and C neurons in each plate have been classified. If we can identify inexpensive, safe, and effective activators and inhibitors of these ORNs, a number of which will hopefully be attractants and repellents, the citrus industry can apply them as lures for surveillance and repellents to curtail the progressing ACP population, thus stopping the spread of HLB. We use these odors as our training set to predict additional ACP Or activating odors.

Insect and mammalian Ors are functionally distinct from one another. Insect Or proteins are considered to be 7-transmembrane proteins that have a non-traditional inside-out membrane orientation. They are believed to function as ligand-gated ion channels and form a heteromer with an obligate partner Orco, which is required for function (Benton et al., 2006; Sato et al., 2008). Mammalian odor receptors on the other hand are G-protein coupled receptors with a traditional outside-in 7-transmembrane orientation (Zhang and Firestein, 2002). Additionally, Mammals have far larger families of odor receptors (~1000 in mice, ~350 in humans) and thus pose an even greater challenge to examine odor coding (Saito et al., 2009; Saito et al., 2004). These distinctions provide an additional challenge for our ligand discovery pipeline. While our method is highly effective for insect Ors, we now have the opportunity to predict odors for functionally distinct mammalian receptors. 52 mouse and 10 human ORs were recently decoded using panel of 63 odorants. Unlike the previous analyses, these

receptors were decoded using calcium imaging and are reported as concentration required for 50 percent effectiveness (EC50). If successful, predictions for these receptors both demonstrate the breadth of our pipeline and provide a tool that would allow the fragrance industry to identify natural and pre-approved odors that activate specific human odor receptors, providing the possibility of truly tailor made fragrances.

In this chapter we apply our cheminformatics platform to predict activators and inhibitors for *Anopheles* Ors, the mosquito CO₂ receptor, Asian Citrus Psyllid ORNs, Mouse ORs, and Human ORs. Computational validations are performed in all instances, and electrophysiological validations are performed for both the *Anopheles* Ors and CO₂ receptor predictions, demonstrating that our chemical informatics approach can be successfully applied in each of the 5 species. A comparative analysis of predicted receptor-odor relationships between *Drosophila* and mammals is also performed, where interesting properties are uncovered. Most importantly, thousands of new odors are identified for over a hundred odor receptors.

RESULTS

Predictions for *Anopheles gambiae* Odor Receptors

Or tuning of previously identified *Anopheles* receptor-odor interactions

The tuning of the 50 previously tested *Anopheles gambiae* Ors is distinctly different from what was observed in *Drosophila* (Allison 2011, Hallem 2006). Activating odors (>50 spikes/sec) for 21 out of 24 *Drosophila* Ors expressed in the antenna were identified. By contrast, activators for only 35 out of 50 were identified for *Anopheles*. For the 15 Ors that remained, an average of 64% of tested odors were inhibitors. This differs sharply from the other Ors of both *Anopheles* and *Drosophila*, where on average only 30% and 28% of tested odors reduced the activity of Ors below spontaneous respectively. Based on the limited training set there are two hypotheses that can explain this distinction. While it is possible that the true odors for these receptors have not yet been identified, an alternative explanation could be that these receptors are specifically tuned for inhibition. As a result, we have divided the *Anopheles* Ors into 3 different classes, each of which we individually apply the chemical informatics platform to, predicting receptor-odor interactions for the following three classes: Activating (Or57, Or12, Or50, Or38, Or56, Or21, Or20, Or18, Or4, Or75, Or11, Or15, Or39, Or48, Or46, Or30, Or9, Or13, Or1, Or10, Or6, Or2, and Or8), no strong activators (Or32, Or26, Or53, Or16, Or64, Or31, Or27, Or65, Or59, Or63, and Or61), and Inhibitory receptors (Or35, Or25, Or3, Or76, Or44, Or66, Or43, Or54, Or67, Or42, Or41, Or33, and Or45).

Chemical Informatics can successfully explain *Anopheles* receptor-odor activity

It is the molecular structures and properties of odors that determine their activity against an Or. As a result, we have designed a chemical informatics pipeline to identify which structural features are important for the activity of odors for individual Ors and have applied these features to predict the activity of the Or to a large set of untested chemicals, as we have previously done for *Drosophila* (Figure 3.1). A single 3D structure was calculated for each of the odors using the Omega software package (See Methods)(Hawkins et al., 2010). We calculated molecular descriptors, which are mathematical values that describe the structure and features of a molecule, from the 3D conformations to quantitatively explain the odor structures. Dragon (Talete) and Cerius2 (Accelrys), which are commercially available software suites, was applied to calculate 3,424 molecular descriptors for each compound tested against an *Anopheles* Or. We then applied a Sequential Forward Selection (SFS) approach to identify which of the molecular descriptors were most highly correlated with activity, resulting in a single subset of descriptors that describe the activity of the training odors against a single Or (See Methods) (Whitney, 1971). While descriptors for Ors classified as activating or fishing “no strong activators” were selected for their ability to describe activation, descriptors for Ors classified as inhibitory were selected for their ability to describe inhibition.

As observed in *Drosophila*, our approach was able to successfully identify molecular descriptors for each *Anopheles* Or that were highly correlated with activity (Table 3.1). We then applied each Ors optimized descriptors to cluster the training odors in order to visualize how well these descriptors grouped the activating or inhibitory odors (Figure 3.2). We find that our approach is able to successfully cluster together odors

tested against Ors for all three classes. While we had observed similar results in our *Drosophila* predictions, this was the first application of our approach to structurally classify inhibition. Inhibitory odors for the majority of Ors were successfully clustered together using optimized structural features. As inhibition can happen both allosterically and competitively, these results pose two interesting hypotheses. Firstly, the very tight clustering for all inhibitory odors observed for several of the Ors could signify competitive or allosteric inhibition in a single pocket. Secondly, the broader clustering of inhibitory odors observed for several of the Ors may signify two active sites or simply an imperfect selection of descriptors.

Predicting receptor-odor interactions

We next applied our approach to predict the activity of the Ors to a large untested odor panel. Our panel consisted of two odor libraries (See Methods). One panel was a natural odor library containing 3,197 odors emitted by plants, flowers, insects, and mammals. The second was large library of Pubchem odors containing over 240,000 compounds with similar molecular weight (<200) and atom type compositions to those found in nature. We ranked the activity of each compound from both libraries using each of the receptor-optimized descriptor sets, resulting in individual rankings for each Or (See Methods). We list odors from the natural odor library that were ranked in the top 500 predictions for each Or (Table 3.2)

Validating predicted receptor-odor interactions using single unit electrophysiology

Several inexpensive and easily obtainable odors were selected from the natural odor prediction lists of 12 *Anopheles* Ors for experimental validation. Members of Dr. John Carlson's lab at Yale performed the experimental validations using single unit electrophysiology. Each *Anopheles* Or gene to be validated along with the obligate co-receptor was heterologously expressed in the empty neuron system in *Drosophila* (Carey et al., 2010; Dobritsa et al., 2003). As this was the same group that created the initial training set using the same experimental procedure, the validations are expected to be highly accurate. A total of 129 receptor-odor interactions were tested across 12 Ors. On average 65 percent of predicted odors activated the Ors (>50 spikes/sec) (Table 3.3). While this is lower than the 71% we had observed for our *Drosophila* predictions, it is still far better than the average activation rate of 16% that was observed across the Ors from our training set. In *Drosophila*, we had observed that our computational pipeline was more accurate at predicting odors for Ors that respond primarily to aliphatic rather than aromatic compounds. We again see this trend for *Anopheles*, with primarily aliphatic responsive Ors having a high accuracy of 84%, whereas aromatic responsive Ors have a relatively low accuracy of 38%.

Increasing the number of aromatic odors in the training set

In *Drosophila* we had hypothesized that our low aromatic prediction accuracy could be due to a lack of aromatic odors in the training set (Hallem and Carlson, 2006). In order to test this hypothesis in *Anopheles*, we increased the size of our aromatic training set by testing a broad panel of aromatic odors. Ideally, we wanted to test as

broad a set of aromatic compounds as possible, providing the largest amount of information on what aromatic structures are or are not active to our computational pipeline for training. As testing odors for activity is a time intensive and expensive process, we applied chemical informatics to select the most structurally diverse set of aromatic odors for testing. Beginning with all 306 aromatic compounds available in the laboratory of Dr. John Carlson, we applied all 3,224 Dragon descriptors to cluster the 306 odors, producing a cluster derived from non-optimized and broad features. We used the dendrogram to select the 75 most diverse aromatic odors (Figure 3.3). Interestingly, the tree is divided into two main branches with small single ring aromatics being in the left subcluster and large multi ring aromatics residing in the right. The Carlson lab then tested 37 of these diverse odors on the primarily aromatic responsive Or2, Or6, and Or10 using electrophysiology. As expected, only a small percentage (9%) of this structurally broad array of 37 aromatic odors activated the Ors (Table 3.4). Interestingly, this number is similar to the random activation chance observed both for *Drosophila* (10%) and in *Anopheles* (16%).

Predicting receptor-odor interactions using our expanded aromatic training set

We applied our ligand prediction pipeline for Or2, Or6, and Or10, using the original 109 odors, odors validated during the initial round of testing, and the newly tested 37 aromatic odors for training, and optimizing descriptor sets for each of the three receptors. (Table 3.5). Since only 29 of the 109 odors in the original training set were aromatics, the inclusion of the additional 37 odors more than doubled the number of aromatics. This large expansion in size, combined with being a very broad array of structural features, increased the challenge of descriptor selection for our pipeline, as

descriptors that once completely separated active from inactive aromatics now resulted in imperfect classification. Additionally, only 22% of the original molecular descriptors are present in the newly optimized sets and the mean number of descriptors per Or has increased from 21 to 30 for these three receptors, further indicating inclusion of the new aromatics into the training set has increased the complexity of the system.

When we cluster the training odors using the optimized descriptors we find that the descriptor sets are successful in clustering training odors (Figure 3.4). Nearly all activating odors are clustered tightly together for Or2, Or6, and Or10. In order to further validate the predictive ability for these three receptors we performed a 5-fold cross validation as was performed for *Drosophila* (See Methods) (Hastie et al., 2001; Tan et al., 2006). The training sets for each Or were divided into 5 equal partitions. One partition was withheld as a test set while the remaining 4 were used to optimize descriptors. The predictive ability of the system was assessed by performing a Receiver Operating Characteristic (ROC) analysis on the withheld test set. This process was repeated 5 times where each of the 5 partitions was used as the test set once. The whole process was repeated 5 times to improve accuracy. A single mean ROC curve was plotted for each Or and an Area Under the Curve (AUC) was calculated (Figure 3.5). We find that the average AUC value across the three receptors was 0.832, which is very close to the value of 0.815 that we observed in *Drosophila*, suggesting that our approach should show similar effectiveness to what we saw in that species.

We applied the newly determined optimized descriptor sets to rank both the Pubchem and Natural odor libraries (Table 3.6). The three prediction lists are quite different from those previously calculated. The next stage will be to validate these predictions using single-unit electrophysiology.

Predicting receptor-odor interactions using machine learning

In order to improve predictive ability we integrated a Support Vector Machine (SVM), which is a well-known and widely applied machine learning approach (Cortes and Vapnik, 1995), into our prediction pipeline (Figure 3.6). We train the SVM using previously determined optimized descriptor set values for Or2, Or6, and Or10 (See Methods Chapter). The trained SVM is then applied to predict the activity of odors from both the Pubchem and Natural Odor libraries. While we have not yet performed electrophysiological validation of our predictions, we have performed a 5-fold cross validation with the SVM, observing far higher AUC values (Figure 3.7). The mean AUC value of 0.988 is far higher than that of the previous value of 0.832 when SVMs were not applied, suggesting that the approach may be highly successful in predicting activating aromatic compounds for these three receptors.

Identifying safer alternative activators with better organoleptic profiles

One target of our approach is to predict odors from a large library of natural odors, many of which are already approved for use in food and/or fragrances. We compare the organoleptic properties of our newly identified activating odors vs those previously identified and discover that odors from our natural odor library have more favorable profiles for several receptors. For example, we have identified new activators for AgOr8, which is believed to be important for host identification through activation by 1-octen-3-ol, that smell more pleasant to the human nose than those previously known. Previously known activators included 3-octanone, 1-octen-3-ol, and 2-heptanone, the first two of which have been described as smelling either moldy or fungal at high

concentrations (www.thegoodscentcompany.com). 2-heptanone is described as fruity, spicy, or woody, however it is also listed as having a high odor strength resulting in a suggested smelling concentration of 10% or less. Our identified compound 2-Heptanol activates the receptor equally as strongly, is described as smelling like fresh lemon grass, and only has a medium odor strength allowing for a pleasant odor experience even at concentrations as high as 100%. Using our approach it becomes possible to easily explore the organoleptic properties of predicted odors, providing a way to select lures and repellent odors based upon their safety profiles and pleasant aroma.

Predictions for CO₂ receptors in a broad range of species

Integrating response profiles to CO₂ from multiple insect species

The olfactory receptors of insects are highly divergent and very little information can be gained from comparing the full-length sequences of Ors between species. As a result of this divergence, we are very rarely able to compare the response profiles of Ors across species. However, the CO₂ receptor is extremely important in many insect species for host detection or predator avoidance, providing a rare case where receptor response profiles are moderately retained across species. The CO₂ receptor of 3 mosquito species as well as *Drosophila* have recently been identified and exposed to panels of odors (Jones et al., 2007; Kwon et al., 2007; Lu et al., 2007; Turner et al., 2011; Turner and Ray, 2009). While several odors were unique to testing in a single species, a panel of 29 odors was tested across all four species. The response profiles of overlapping odors appear in general similar, suggesting the possibility for pooling the responses of the 4 species into a single metric of CO₂ receptor activity.

We wanted our metric to consider the responses of multiple species for odors where more than one species had been tested, however since responses to some odors were significantly lower in one species than the others, we decided upon selection based on customized criteria. We averaged the two highest activity values for each odor that had been tested on more than 1 species, resulting in a metric that was influenced by multiple species, yet not unfavorably reduced by low activity outliers. The resulting activity values were used as our training set for the remainder of our CO₂ analysis.

The resulting training set contained a large variety of odor structures including both aromatic and aliphatic activators. Visual inspection revealed there were two distinct structural classes of known activating odors (aromatic and aliphatic) and one for inhibiting odors. As we have both recognized the challenges posed in describing aromatic compounds in our analyses of *Drosophila* and *Anopheles* Ors and considered that these two structurally distinct classes of odors may bind in distinct pockets on the receptor, we decided to separate our training odors into three groups with different predictive aims: aromatic ligands, structurally diverse ligands, and aliphatic ligands (Table 3.7). Our aromatic set focused on screening for highly active and aromatic structures and did not consider inhibitors. Our Broad activator screen focused on structurally diverse ligands and tested all previously tested odors including both activators and inhibitors. Finally, we performed an inhibition screen to identify CO₂ receptor inhibitors, where we only considered inhibitory or inactive odors and did not include activators.

Chemical Informatics can successfully explain CO₂ receptor-odor activity

We once again applied our chemical informatics pipeline to predict active odors, in this instance predicting activators and inhibitors for the CO₂ receptor (Figure 3.8). We again calculated a single 3D structure using Omega software (Hawkins et al., 2010) and calculated 3,224 molecular descriptors using Dragon (Talete) for each odor from all three training sets. We then implement our computational pipeline to identify a subset of the 3,224 molecular descriptors that are most important for activity in each of the training sets through application of a Sequential Forward Selection (SFS) method (Whitney, 1971). This method identifies molecular descriptors that are the most highly correlated with odor activity by incrementally building a best descriptor set, which involves adding a single molecular descriptor at a time until the correlation ceases to increase (See Methods Chapter VII). We applied this process independently for each of the training sets, resulting in three uniquely optimized molecular descriptor sets (Tables 3.8, 3.9, 3.10).

While there were many important 2D and 3D descriptors selected for all three of the training sets, several more intuitive 0D and 1D descriptors were selected as well. Chiefly, the number of non-terminal nitrogen atoms and number of aldehydes were important for the aromatic activator training set, the number of aldehydes, presence of a carbon-oxygen bond 6 atoms apart and number of donor atoms for H-bonds were important for the broad activator training set, and the presence of a carbon and oxygen atoms 3 atoms apart, number of aliphatic secondary carbons, presence of carbon and carbon atoms 6 atoms apart, and absence of ring secondary carbons were important for the inhibitor training set. The aromatic activator set was the smallest, containing 11 descriptors, and the broad activator set was the largest with 27.

We clustered each of the training odors by the values of their previously determined optimized descriptor sets in order to visualize how well each descriptor set groups activating or inhibiting odors (Figure 3.9). We observed odors of interest are grouped tightly together for all three training sets, supporting that the optimized descriptors do indeed explain odor activity.

Predicting the response of the CO₂ receptor to a broad panel of untested odors

We next applied the optimized descriptor sets to computationally rank two large libraries of untested odors. We assembled both eMolecules and natural odor libraries containing over 440,000 and 3,197 odors, respectively (See Methods Chapter VII). The eMolecules library contains all chemicals from the eMolecules library of similar molecular weight (<350 MW) and atom type compositions (C, O, N, H, I, Cl, S, F) of known volatile odors. The Natural Odor library contained 3,197 odors of identified floral, vegetative, insect, or mammalian origins. We ranked each odor from both libraries based on their similarity to known activators using the previously determined optimized descriptor sets, resulting in three independent rankings, one for each training set (See Methods Chapter VII). We provide the top 60 predictions for each of the training sets (Table 3.11).

Many predictions validated as either activators or inhibitors of the CO₂ receptor

Electrophysiological tests were conducted by collaborator Dr. Dyan McWilliams for 139 predicted ligands originating from the three prediction lists, that were considered reasonably safe for human use. Each predicted odor was individually tested for activity using single unit electrophysiology on the *Aedes aegypti* CpA neuron, resulting in identification of many new actives for the CO₂ responsive neuron. We observed a broad

range of activity for our validated odors, including both activators and inhibitors (Table 3.12).

In our initial observation we noticed that far fewer (18%) of our predicted CO₂ odors activated the receptor at the threshold (>50 spikes/sec) previously applied for both *Drosophila* (58%) and *Anopheles* (65%) (Table 3.21). As previously described the CO₂ receptor is actually a gustatory receptor, which are believed to function uniquely from Ors. Additionally, the strongest identified activator (143 spikes/sec) for this receptor from either our large screen or any previous screen, which encompasses the responses of over 200 odors, activates significantly lower (~44%) than those observed in either *Drosophila* or *Anopheles* Ors (~250 spikes/sec), leading us to hypothesize that these Gr structured CO₂ receptors may be tuned uniquely from general Ors. As a result we have lowered our activity threshold from 50 spikes/sec to 30 spikes/sec, which loosely relates to the 44% reduction observed from the strongest known activators of other insect species to those observed for the CO₂ receptor. Even with this threshold reduction we observe that only 30% of compounds activated the CO₂ receptor at greater than or equal to 30 spikes/sec.

Interestingly, CO₂ responsive neurons in mosquitoes do not appear to have a spontaneous activity, responding only to the quantity of CO₂ present in the environment, which vary slightly depending on how many people have been in the electrophysiology room that day. As a result, we tested for inhibition by overlaying odors on a 0.15% CO₂ pulse, which on average activates the receptor at 100 spikes/sec, providing a standardized metric by which to compare inhibitory odors. We observed a similar number of tested odors inhibited the CO₂ neuron (18%) as we had observed in *Drosophila*.

In addition to CO₂ receptor structure/function differences, we also must take account of the noise introduced by the multi-species integrated training odor set. As there were no large unified odor panels tested on a single species, we integrated the responses of similarly responding CO₂ receptors for 4 insect species. While this challenge did pose a very interesting intellectual question of whether ligand structural similarity can be compared across species for similarly tuned receptors, it undoubtedly will reduce the accuracy compared to when training is performed on responses of only a single species. Considering these additive challenges, we feel our computational approach was successful in predicting ligands.

Active predicted odors appear to target up to three distinct binding pockets

We next analyzed the structures of our predicted odors that validated as actives. Interestingly, from visual inspection of their structures it appeared the odors could be structurally diverse enough to be acting on distinct binding pockets. In order to identify whether ligands fall into diverse structural classes that could potentially bind to distinct active sites on the receptor, we clustered all activating (>30 spikes/sec) and strongly inhibiting (<-10 spikes/sec) odors. We concatenated all molecular descriptors from each of the three distinct optimized descriptor sets, which we had individually optimized to predict ligands with unique characteristics, creating a single set of optimized descriptors explaining both activation and inhibition of the CO₂ receptor. We then applied hierarchical clustering to organize odor relationships based on the Euclidean distances between odors based on each odors optimized descriptor values (Figures 3.10, 3.11).

The resulting tree had three distinct branches or clusters, each of which contains structurally distinct odor classes. Cluster 1 contains only aliphatic odors that have a

variety of functional groups including aldehydes, ketones, alcohols, and esters. With the exception of the few alcohols, which are H-bond donors, nearly all compounds in this cluster contain an H-bond acceptor group. Cluster 2 contains pyrazine based heterocycles with a variety of small branched side chains. Lastly, cluster 3 contains both small 5 and 6 member non-aromatic cyclic compounds and aromatics.

We hypothesize that each of these three odor classes bind to distinct regions in the CO₂ receptors. As the structures of the CO₂ receptors have not been solved, we can only hypothesize on whether or not they bind to distinct regions, however it is very curious as to whether these distinct structures bind in entirely different protein regions, or to distinct regions in the same binding pocket, or whether the binding pockets of Ors are highly variable, accepting a very wide variety of shapes and sizes.

Interestingly, inhibitory and activating odors are interspersed across sub clusters. While the largest grouping of inhibitors occurs for cluster 2, there are strongly inhibiting odors distinctly residing in both clusters 1 and 3. In the same fashion activating odors exist in all three clusters. If three distinct binding regions do indeed exist within the CO₂ receptors, this would mean that we have identified site selective activators and inhibitors for all three.

It will be interesting to test the proposed three distinct binding region hypothesis by creating odor blends containing odors from each of the clusters. In theory, if these odors do function through distinct binding sites, exposure of a CO₂ receptor to odors from multiple clusters should provide a stronger response than similarly strong odors from within a single cluster, as binding in distinct regions at the same time should further stabilize the active conformation. This can be tested using single unit electrophysiology.

We also plan to test the strongest activating and inhibiting odors from behavioral response. Since we predicted ligands from a large collection of naturally occurring odors, many have safe activity profiles and a few are even approved for use as flavorings, fragrances, or cosmetics. If we can identify safe odors that are highly effective at altering the activity of the CO₂ receptor neuron, we will have made a major step in the crusade for global health.

Citrus Psyllid Odor Receptors

Chemical Informatics can be applied to explain Psyllid receptor-odor activity

Through work performed in our lab, activators for 11 ORNs found on 4 rhinial plates of Citrus Psyllids have been identified. One activator has been found for PR4_A, two have been found for RP6_A, and one has been found for RP7_A, with none of the actives increasing firing rates more than 70 spikes/sec. As our computational approach requires at least two activating odors in order for optimization, we will focus our efforts on 9 AsCP neurons (RP2_B, RP2_C, RP4_A, RP4_C, RP6_A, RP6_B, RP6_C, RP7_B, RP7_C).

We begin our computational pipeline as performed for *Drosophila*, calculating 3D structures with Omega (Hawkins et al., 2010) and 3,224 molecular descriptors using Dragon (Talet) for each of the 61 odors tested against the 9 ORNs (Figure 3.12). Molecular descriptors, which are mathematical values attempting to describe the structure of each compound, are then optimized individually for each ORN using a Sequential Forward Selection (SFS) approach, resulting in a single optimized set of descriptors for each ORN (See Methods Chapter VII) (Whitney, 1971)(Table 3.13). When we cluster the training odors based upon their optimized descriptor values, we find

that each optimized descriptor set is successful in grouping together highly active odors (Figure 3.13).

Predicting the response of the Psyllid ORNs to a broad panel of untested odors

We next used each ORN optimized descriptor set to train a Support Vector Machine (SVM), which is a highly applied machine learning approach (Chang and Lin, 2001; Karatzoglou et al., 2006). The SVM was trained to optimally predict the activity of each odor in the training set using only information provided by each ORNs optimized descriptor set, resulting in one trained SVM for each ORN. We next performed a 5-fold cross validation followed by a Receiver Operating Characteristic (ROC) analysis to validate the predictive ability for each ORN (See Methods Chapter VII). The Area Under the Curve (AUC) was calculated and the ROC analysis was plotted for each ORN (Figure 3.14). The extremely high AUC values, with a mean value of 0.99 across all ORNs, indicates that our odor prediction pipeline was extremely efficient at describing the activity of the training sets for all 9 ORNs.

We next assembled both natural odor and eMolecules libraries, representing large collections of previously untested odors (See Methods Chapter VII). We next applied each trained SVM to predict the responses of all odors from both training sets for each ORN. We provide the top 100 predicted natural odors for each ORN (Table 3.14). Predicted ligands can be validated, identifying strong activators and inhibitors for each ACP ORN. Strong actives can then be tested for their ability to modify ACP behavior. Ideally, we will be able to identify strong attractants for trap lures and strong repellents odors to protect citrus trees.

Mammalian Odor Receptors

Identification of unique subsets of optimal descriptors for mammalian receptors

Insect Or proteins are believed to consist of 7-transmembrane regions that are inversely oriented in the membrane. It is believed that they function as ligand-gated ion channels with a possible cyclic-nucleotide-activated cation channel function along with an obligate co-receptor Orco, forming a heteromer (Benton et al., 2006; Sato et al., 2008; Wicher et al., 2008). Alternatively, mammalian Ors function as G-protein coupled receptors with a traditional outside-in 7-transmembrane orientation and do not require an obligate co-receptor such as Orco (Zhang and Firestein, 2002). In order to test whether the chemical informatics platform would be successful at predicting ligands for the functionally distinct mammalian Ors, we performed a similar analysis on 33 ORs from mouse and 4 ORs from humans. The responses of each of these mammalian receptors to a panel of 62 odorants have been determined by functional expression in heterologous cells and >2 actives have been identified for each (Saito et al., 2009).

We applied our chemical informatics pipeline to optimize descriptors for each of the 37 mammalian ORs (Figure 3.15). 3D conformations for each of the training odors were calculated by Omega (See Methods Chapter VII). Dragon and Cerius2, which are commercially available programs for calculating molecular descriptors, were used to calculate 3,224 and 200 molecular descriptors for each of the training odors (See Methods Chapter VII).

We implemented a Sequential Forward Selection (SFS) approach to select a subset of molecular descriptors that best correlated with the activity of the training set odors (Whitney, 1971) (See Methods Chapter VII) (Figure 3.15). The SFS approach iteratively assembles an optimized subset of descriptors, beginning with a single best

descriptor and then incrementally adding additional descriptors one at a time until the correlation between descriptor values and OR activity fails to increase. The SFS approach is applied for each OR, resulting in 37 unique optimized descriptor sets (Table 3.15).

We applied hierarchical clustering on the Euclidean distances calculated between training odors using optimized descriptor set values for each of the 37 ORs in order to visualize how well our approach brought together activating odors (Figure 3.16). Upon visual inspection, clustering of the 62 training odors by each OR's optimized descriptor set indicated that the receptor-optimized descriptor sets were indeed able to effectively cluster activating odors together.

Computational validation of mammalian OR-optimized descriptor sets and *in-silico* prediction of ligands

Since functional testing of predictions for mammalian receptors are beyond the scope of this study, we performed the well established computational 5-fold cross-validation to determine the predictive ability of the *in silico* approach. ORs with >15 known ligands were selected and for each OR 20% of the compounds (12/60) were excluded as a test set, while the remaining were used as a training set to generate the optimized descriptors. As before we performed a ROC analysis for each of the withheld test-set odors by classifying activity based on distance in the training set derived Or-optimized chemical space. We repeated this operation five times for each receptor, each trial performed by excluding a different subset of odors.

We generated mean ROC curves for each OR and calculated AUC values (Figure 3.17). Using this method we demonstrate that the OR-optimized-descriptor sets

generated using the training odors could accurately identify actives from the test odorants. Both the ROC and ApoA values were comparable, if not better (data not shown), than for the *Drosophila* Ors suggesting that the descriptors are able to efficiently explain the activity of the mammalian Ors.

Predicting mammalian receptor-odor interactions for a large untested odor space

We then applied the OR-optimized descriptors to systematically screen both the Natural Odor and Pubchem libraries *in silico*, comprising over ~8,880,000 receptor-odor interactions and representing 33 mouse ORs, 4 human ORs (See Methods Chapter VII). We identify the top 500 (0.2%) hits from this vast chemical library for each Or/ORN, the top ~100 for each are reported in Table 3.16.

Relationship between descriptor sets and Or sequence and activity

We next analyzed the relationship between ORs based upon their shared molecular descriptor sets, training set activity, predicted odor sets, and phylogenetics, as we had previously performed for *Drosophila* (Figure 3.18). In contrast to what was observed in *Drosophila*, the analysis for the mammalian dataset reveals a greater degree of common relationships across the known-activity, predicted cross-activity and descriptor trees (~77% ORs present in common subgroups). Similarly, the *Drosophila* Or-phylogenetic tree has sparser subgroup relationships conserved with each of the other trees (<45%), as opposed to the mammalian ORs where the majority of subgroups in the phylogenetic tree (>56%) are conserved across the different trees. This difference may reflect the much greater amino-acid similarity across the mammalian receptors (47%) as compared to the divergent *Drosophila* receptors (23%).

Analysis of breadth of predictions for each mammalian OR in chemical space

Determining the breadth of coding in a large volatile chemical space across many receptors is virtually impossible to determine experimentally. However, if we apply the OR-optimized descriptor sets, it becomes possible to computationally predict frequency distributions for each of the mammalian ORs for the odor space of both libraries (>240,000 odors) (Figure 3.19). As expected, we find a great deal of variation across the ORs. While some ORs are predicted to be very broadly tuned, others are decisively narrower. Analyses such as these allow inferences of OR properties that would be very difficult to assess through wet lab experimentation alone.

DISCUSSION

We have applied our chemical informatics pipeline to predict activating and inhibiting odors for a large number of extremely important species, including several pest and disease vector insects. For *Anopheles* Ors we applied our pipeline to predict activators and inhibitors for a large number of receptors. We validated predictions experimentally for 12 Ors, demonstrating an effectiveness similar to that observed for *Drosophila*. For the CO₂ receptor, we considered the responses of multiple species integrating them into a single metric for molecular descriptor optimization. While validations for this receptor were not as successful as seen in either *Drosophila* or *Anopheles* Ors, it still signifies a major progression in the fight for behavior modification on disease vector species. We predicted activating odors for the very first time in ACP, creating a vast expansion upon the previously tested odor space. Finally, we predicted odors for 33 mouse and 4 human Ors.

We have noted challenges involved in predicting aromatic odors in both *Drosophila* and *Anopheles*. In an attempt to improve our accuracy for aromatic responsive Ors we have performed a thorough expansion of aromatic odor space. This is the first time a focused effort has been made to expand upon and explore responses of highly diverse aromatic odors. We applied chemical informatics to select and test 27 aromatic odors that were intelligently selected for their extremely broad composition of molecular features. Further, we integrated SVMs into our prediction pipeline to increase the abilities of our computational approach. Future experimental validation will expose how successful our attempts were.

No computational analyses in olfaction had been previously performed to integrate the responses of multiple species into a single metric for odor prediction. We successfully integrated the responses of 4 divergent species to predict odors for the very important CO₂ receptor. Additionally, while the three mosquito species are attracted to CO₂, *Drosophila* are repelled, further creating an interesting intellectual challenge by integrating odors causing mixed behavioral responses. All of these differences introduce noise into the training set. As a result, the observed lower accuracy (48%) compared to that previously observed in both *Drosophila* (71%) and *Anopheles* (65%) Or predictions was not unexpected. However, due to these multiple challenges we still consider this to be a significant achievement.

It will be interesting to perform experimental validations for both the ACP and mammalian predictions in the future. There are several ways to further improve our computational approach in the future. Sequential Floating Forward Selection (SFFS) allows for the intermittent removal of previously selected descriptors instead of a steady building of an optimized list. Additionally, other predictive approaches beyond SVMs

and other alternative approaches can be considered. Since machine learning has been demonstrated to be highly effective for our pipeline, we plan on applying it for all future analyses.

Figure 3.1: A molecular descriptor optimized approach is able to explain odor activity for individual *Anopheles* Ors

Schematic of our chemical informatics pipeline. Molecular descriptors that are most correlated with activity are selected, resulting in a metric that is able to cluster together highly active odors using important structural features. The optimized descriptor sets can then be applied to predict Or activity against a large panel of odors. This pipeline is applied independently, optimizing descriptors for either activators (Top) or inhibitors (Bottom).

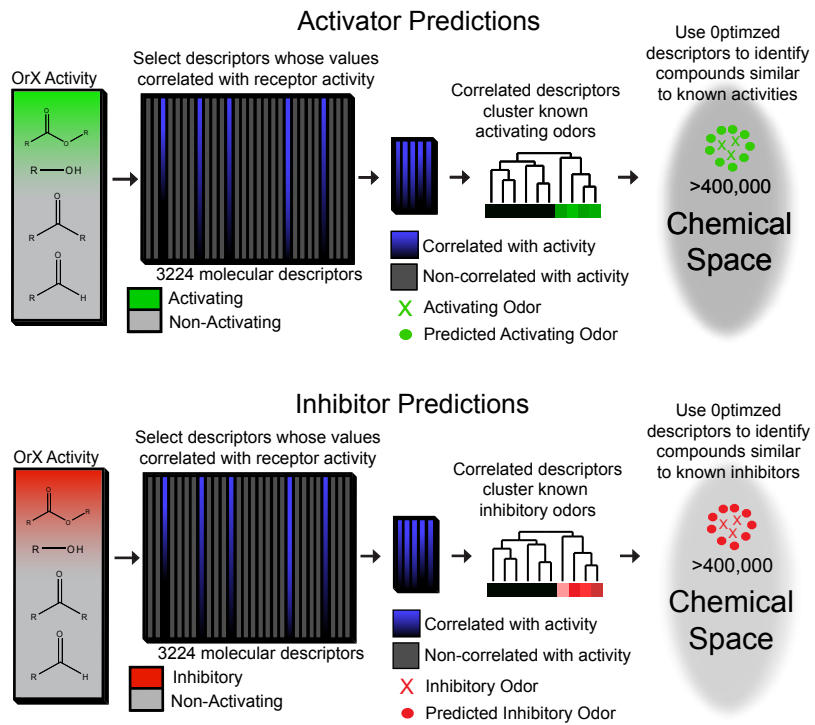


Figure 3.1

Figure 3.2: Optimized molecular descriptor sets are able to cluster either active or inhibitory odors

Optimized molecular descriptors values were applied to cluster training set odors individually for each Or. Clusters are divided into Activators, Fishing Activators, and Inhibitors. Colors for Activators and Fishing activators range from purple to red for lowest activity to highest activity and are individually applied for each Or. Colors for Inhibitors range from blue to bright for highest activity to most inhibitory.

Activators

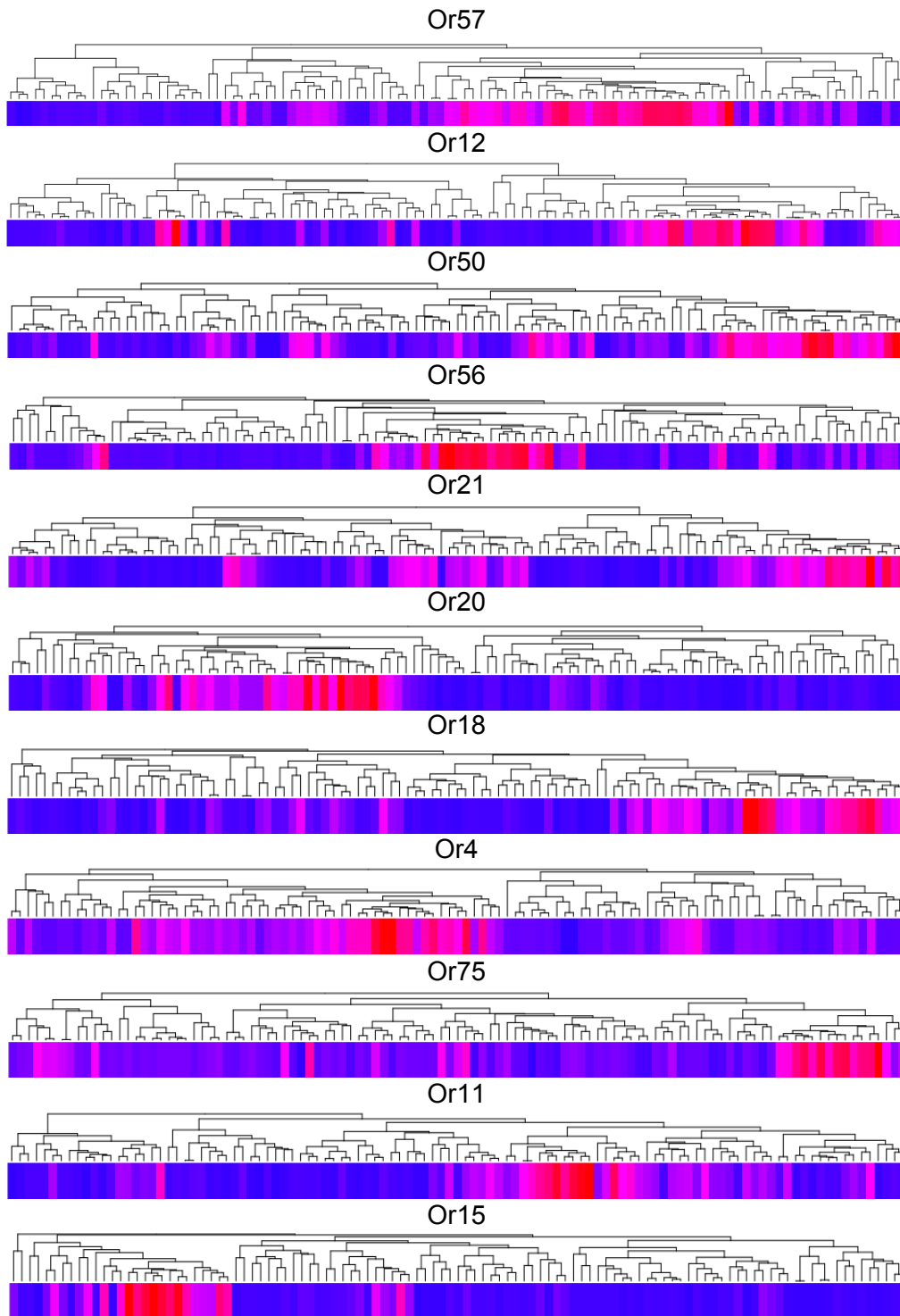


Figure 3.2

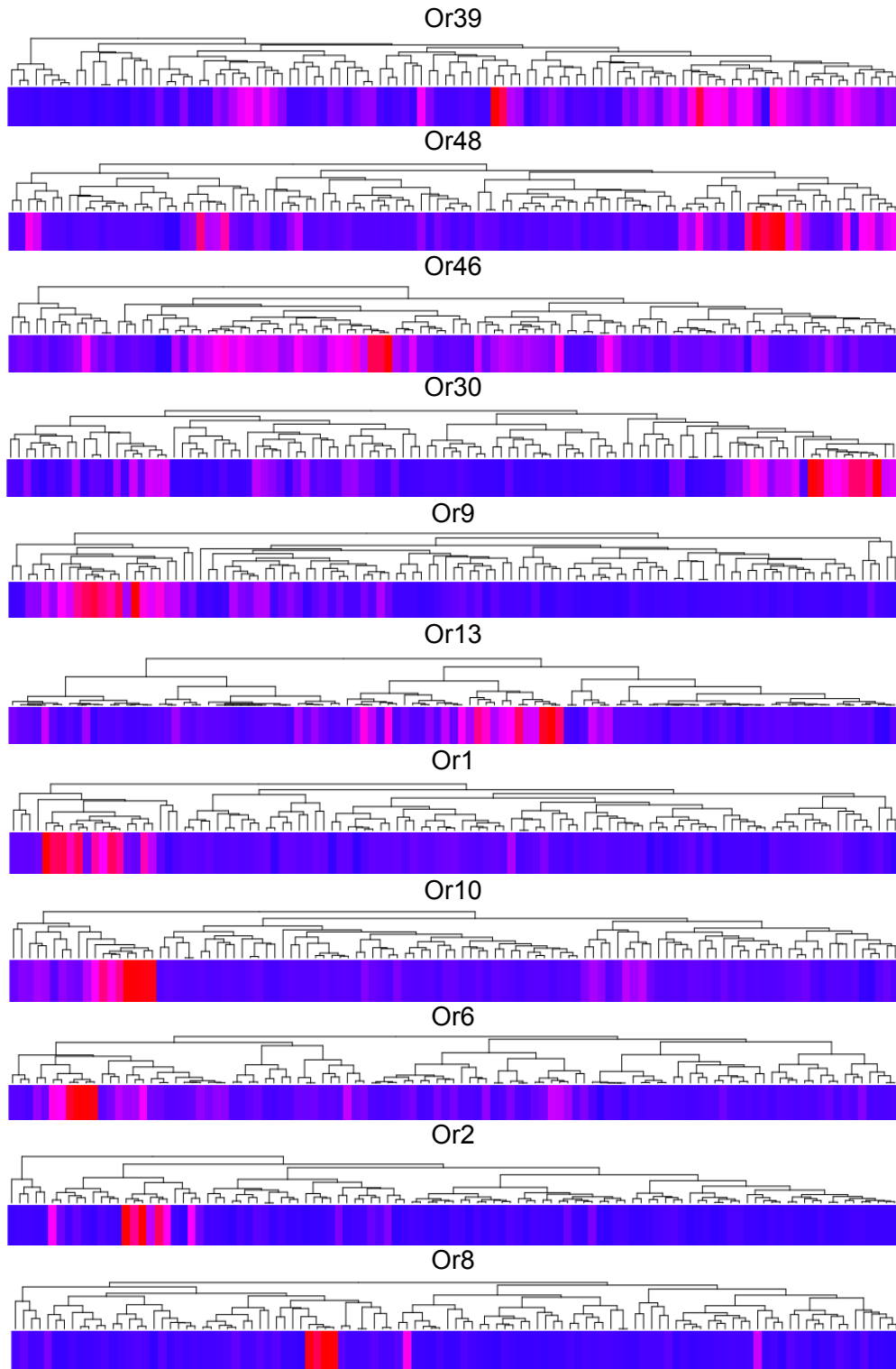


Figure 3.2 Continued

Inhibitors

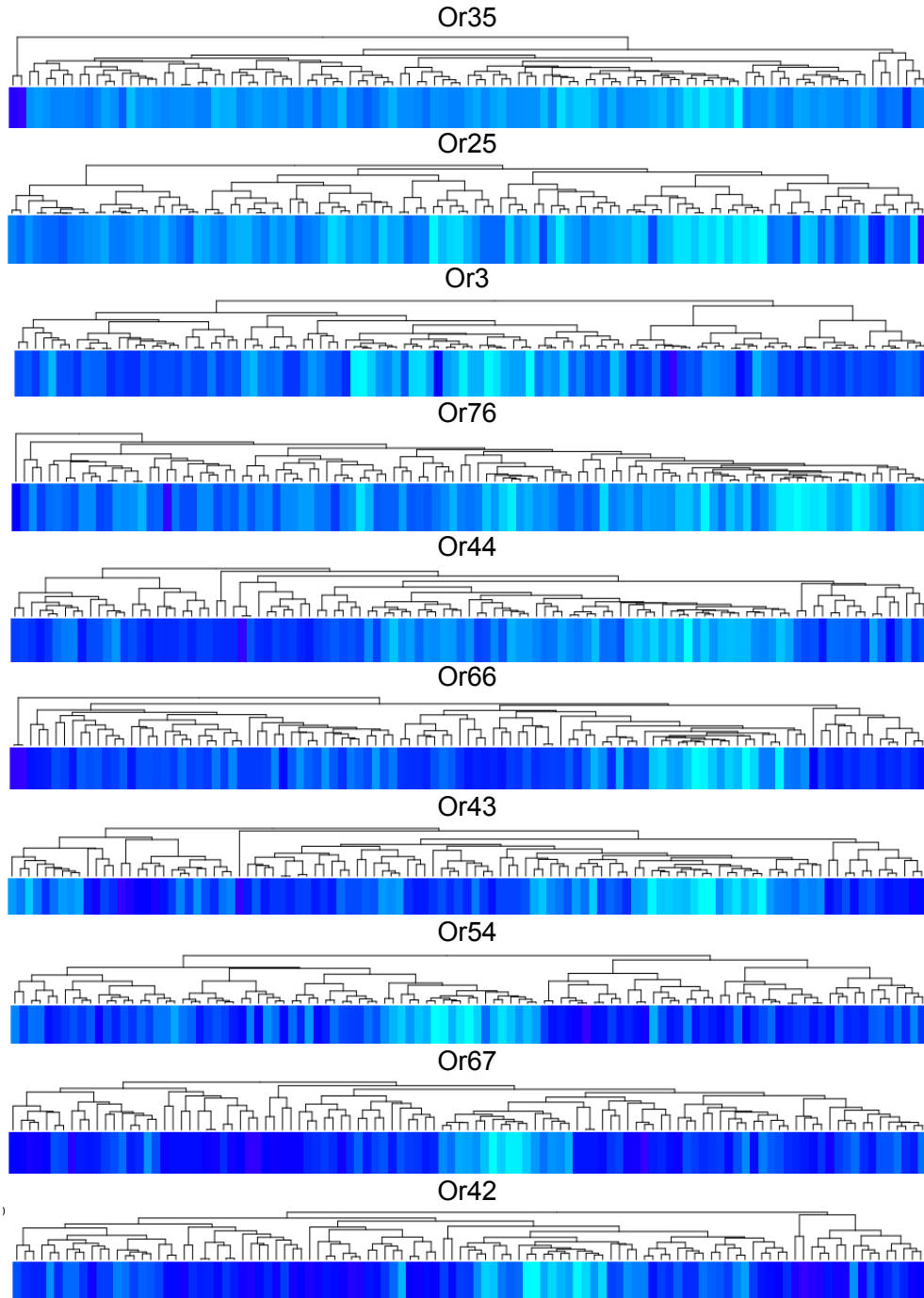


Figure 3.2 Continued

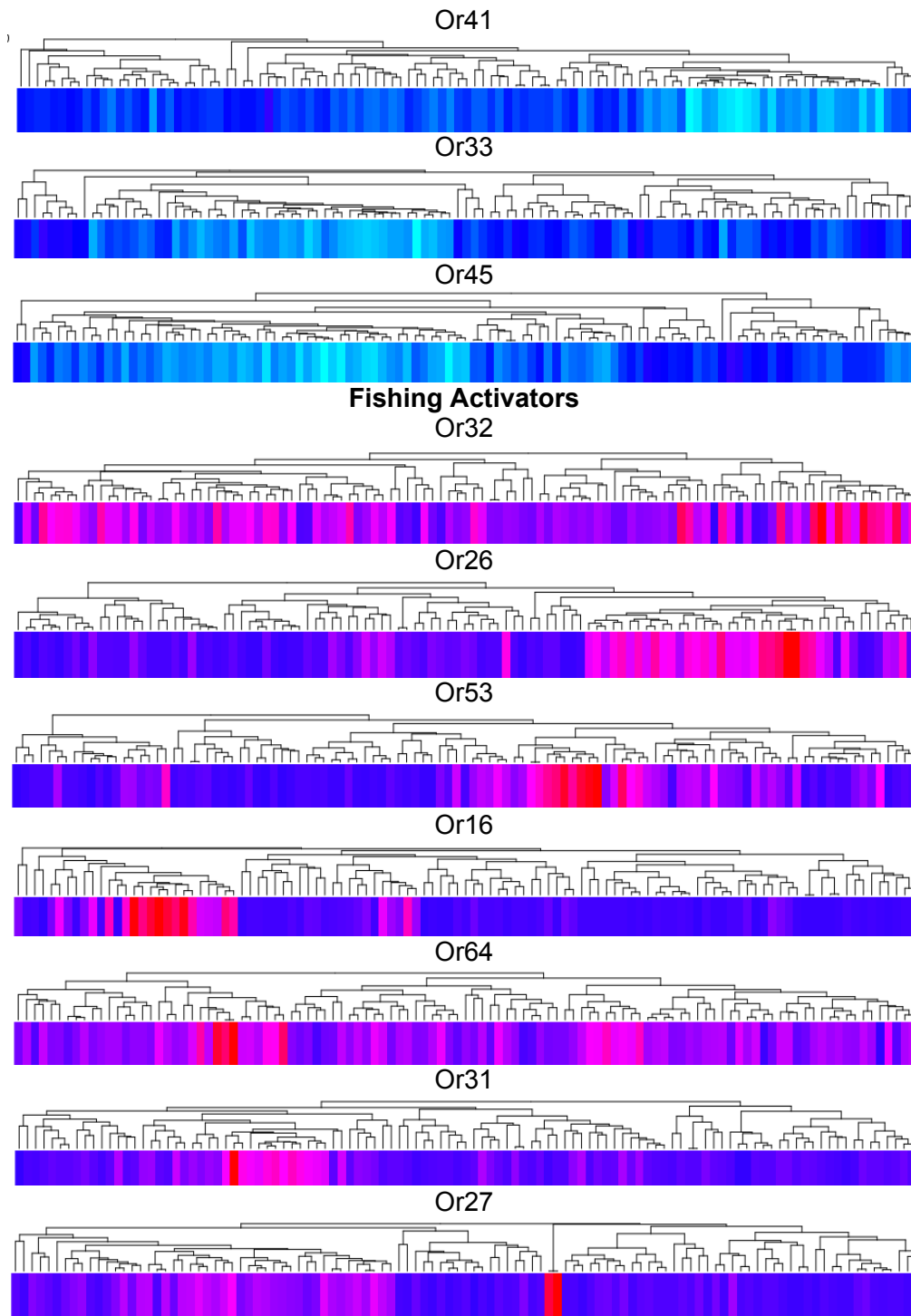


Figure 3.2 Continued

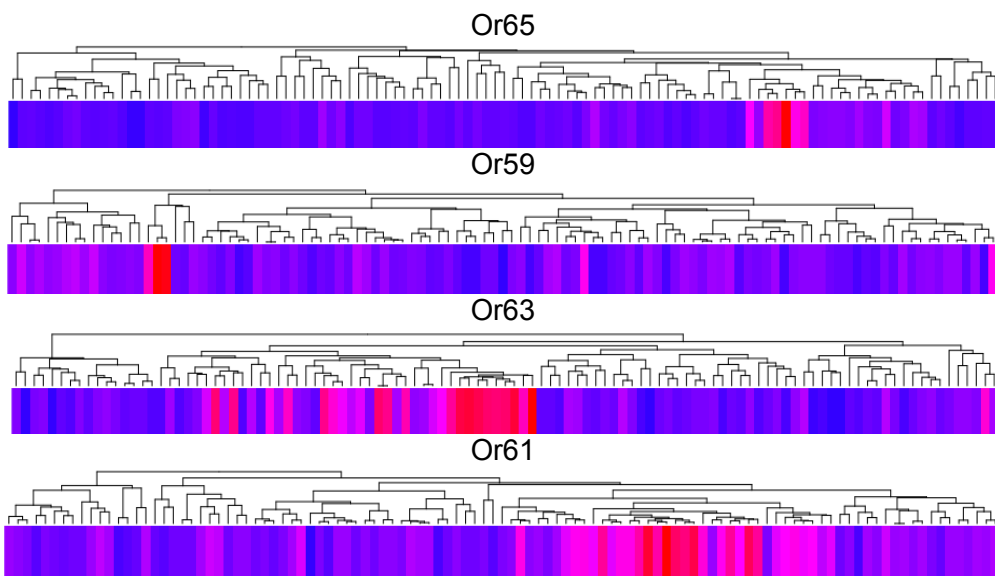


Figure 3.2 Continued

Figure 3.3: Structural relationships between the 75 most diverse aromatic odors

Hierarchical clustering for 75 odors using Euclidean distances between odors calculated using all 3,224 molecular descriptors calculated by Dragon. Compound names are provided for each odor.

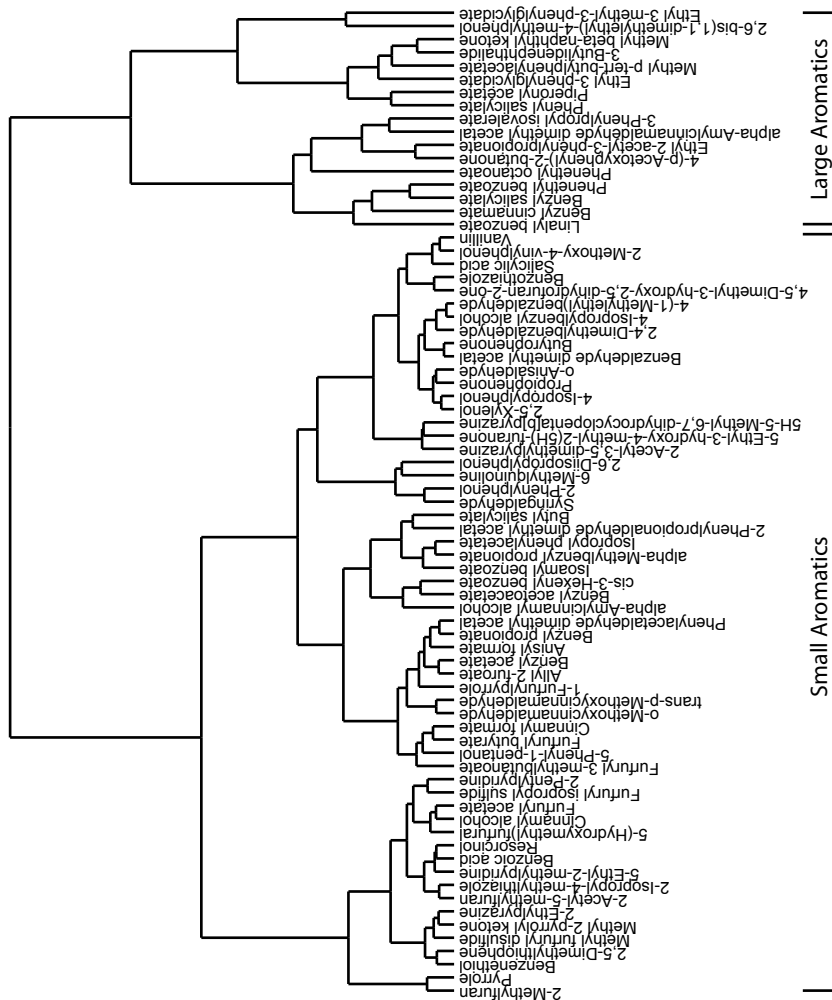


Figure 3.3

Figure 3.4: Newly optimized molecular descriptor sets are able to cluster odors for aromatically tuned Ors

Optimized molecular descriptors values were applied to cluster training set odors individually for each Or. Colors range from gray to black, representing the lowest to highest activity.

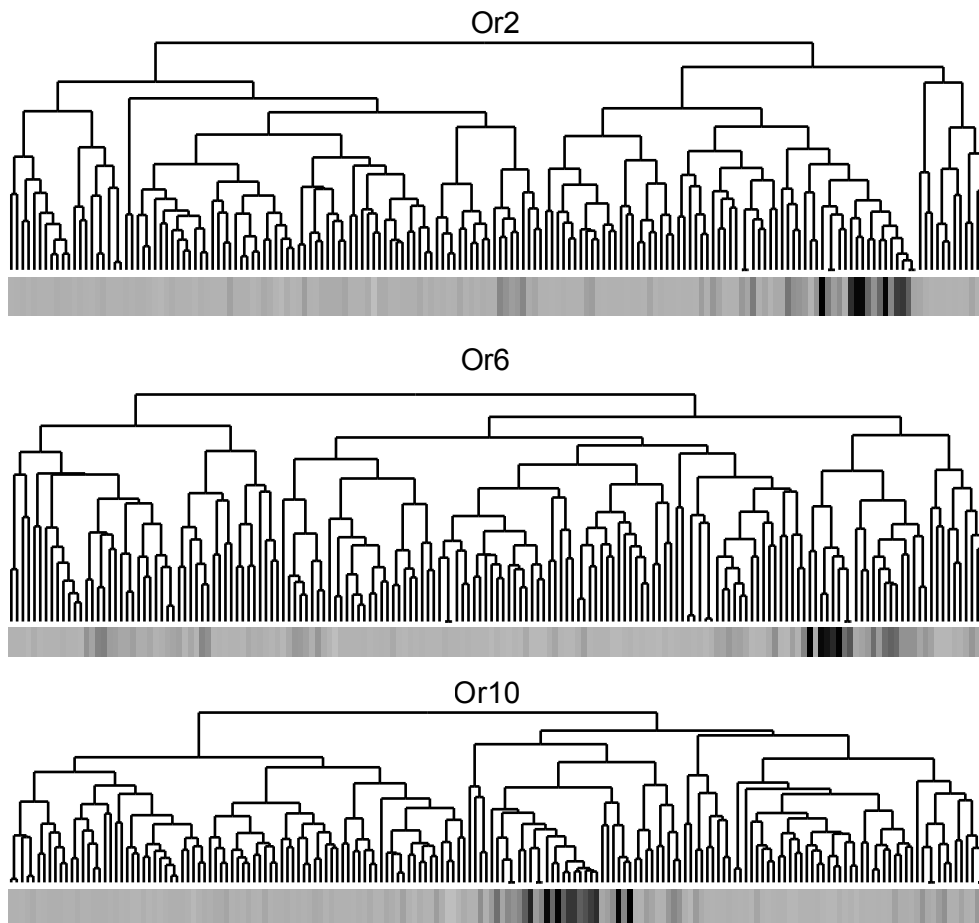


Figure 3.4

Figure 3.5: Optimized descriptors for AgOrs 2, 6, and 10 effectively describe training set activity.

A Receiver-Operating-Characteristic (ROC) curve is plotted for each Or, depicting the computational validation of ligand predictive ability of the Or-optimization approach.

Averaged ROC values for each Or

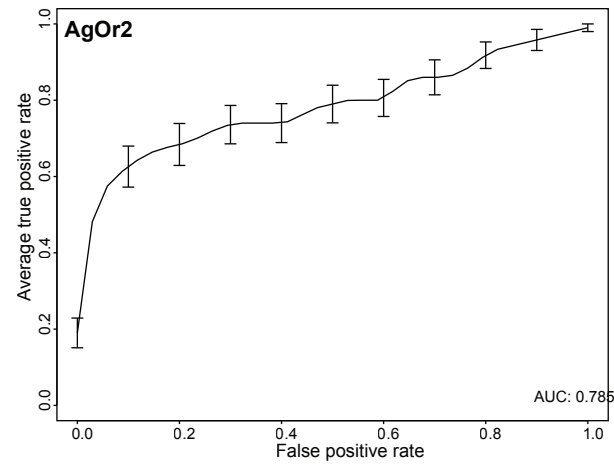
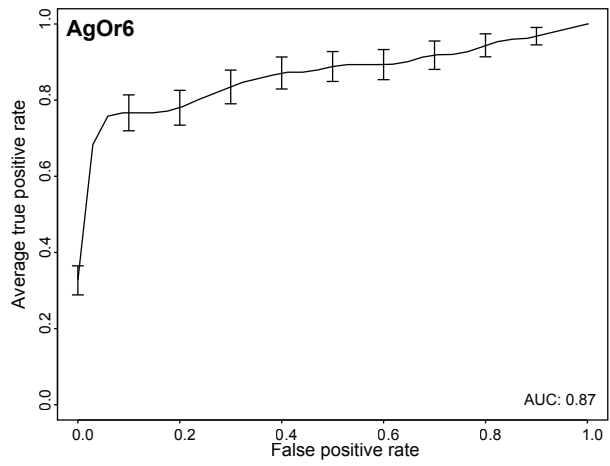
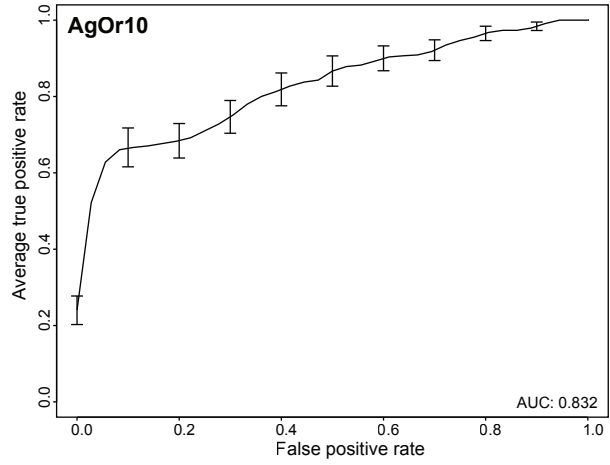


Figure 3.5

Figure 3.6: A Support Vector Machine (SVM) integrated approach is highly effective at explaining odor activity for individual Ors

Schematic of our SVM integrated chemical informatics pipeline. Molecular descriptors that are most correlated with activity are selected, resulting in a metric that is able to cluster together highly active odors using important structural features. The optimized descriptor sets are then be applied to train a SVM to predict Or activity, which is then applied to predict the activity of a large untested odor space.

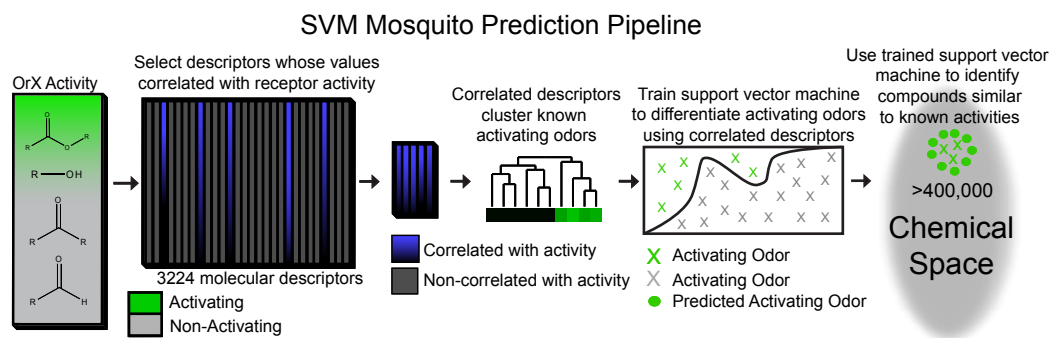


Figure 3.6

Figure 3.7: SVMs trained using optimized descriptors for AgOrs 2, 6, and 10 effectively describe training set activity

A Receiver-Operating-Characteristic (ROC) curve determined from 100 independent iterations of a SVM applied 5-fold cross validation is plotted for each Or, depicting the predictive ability of the Or-optimization approach.

Averaged ROC values for each Or

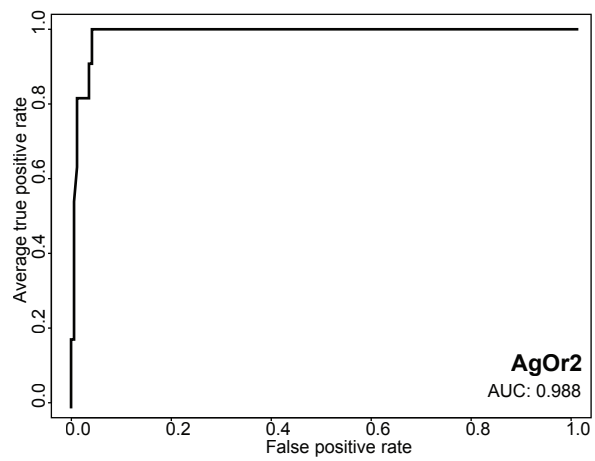
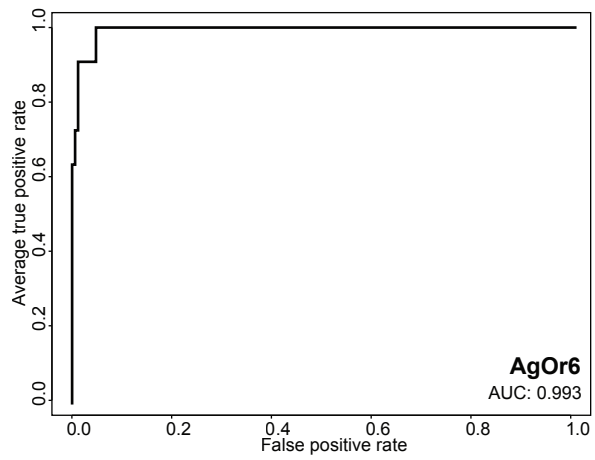
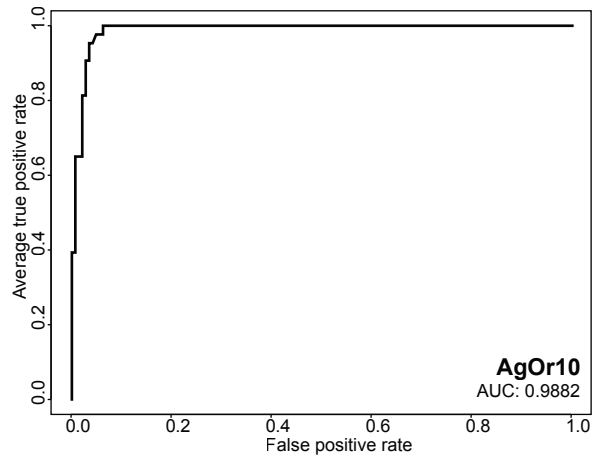


Figure 3.7

Figure 3.8: A molecular descriptor optimized approach is able to explain odor activity for CO₂ receptors

Schematic of our chemical informatics pipeline. Molecular descriptors that are most correlated with activity are selected, resulting in a metric that is able to cluster together highly active odors using important structural features. The optimized descriptor sets can then be applied to predict CO₂ receptor activity against a large panel of odors. This pipeline is applied independently, optimizing descriptors for either activators (Top) or inhibitors (Bottom).

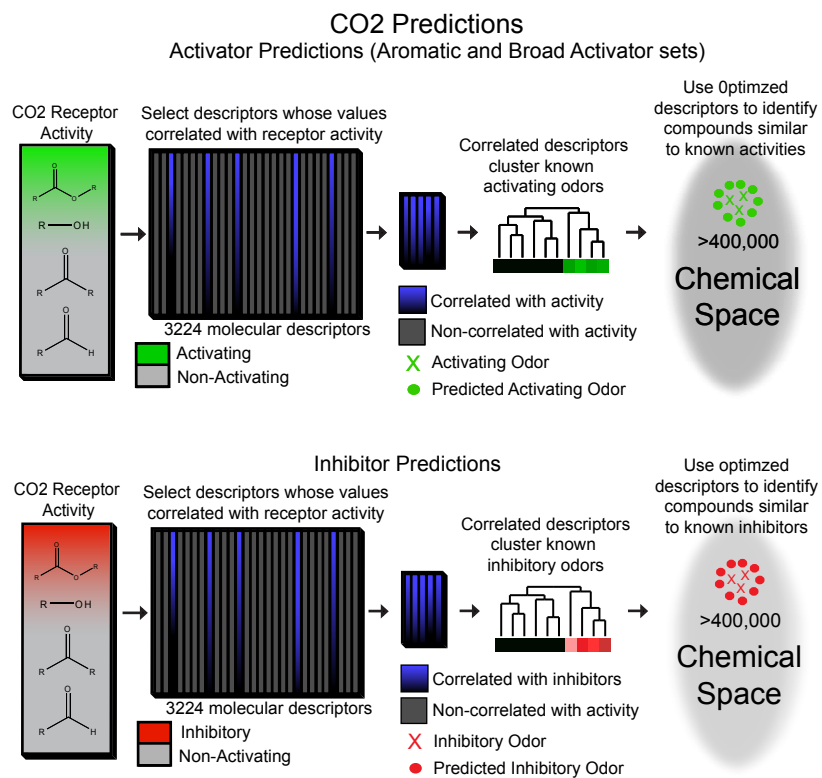
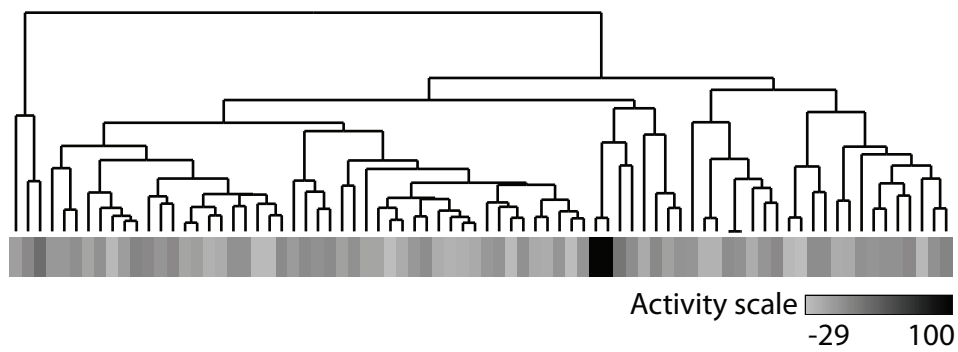


Figure 3.8

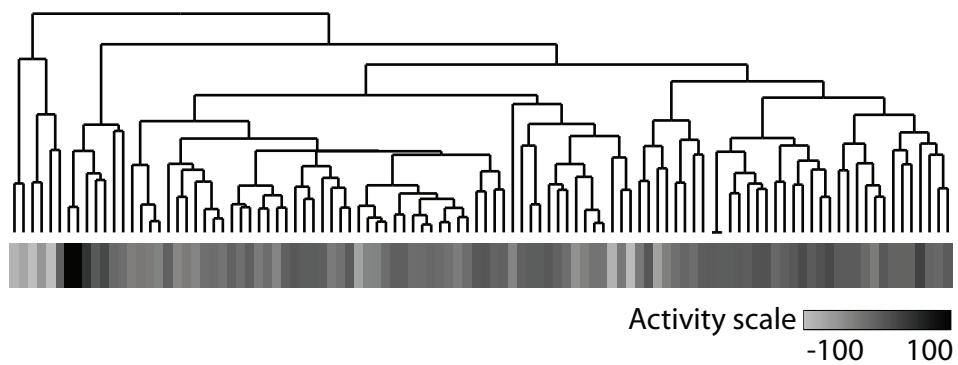
Figure 3.9: Optimized molecular descriptor sets are able to cluster odors by CO₂ receptor response

Optimized molecular descriptors values were applied to cluster training set odors individually for each screening set. Odor activity is represented individually for each screening set, ranging from gray to black in increasing activity.

Aromatic Activator Screen



Broad Activator Screen



Inhibitor Screen

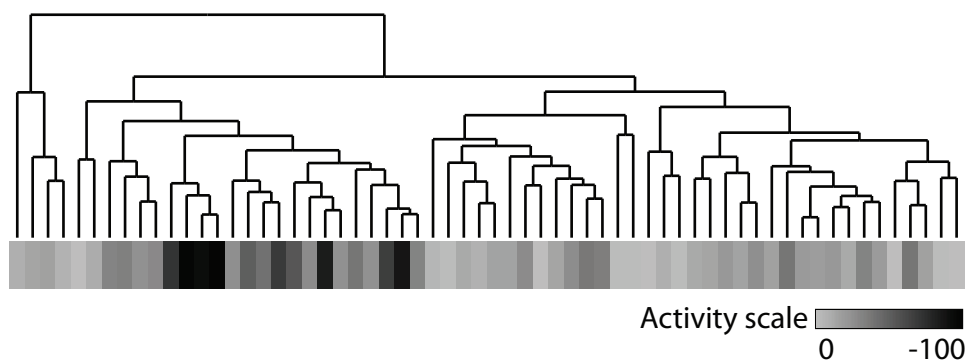


Figure 3.9

Figure 3.10: Active compounds cluster into three distinct structural classes

Active compounds were clustered by the combined 3 sets of optimized descriptors. Hierarchical clustering was performed on compound optimized molecular descriptor values, dividing the chemical structures into three distinct classes that have been outlined by both label and color. CO₂ receptor responses (spikes/sec) have been provided for each compound and green and red boxes label activators or inhibitors, respectively. Representative chemical structures for each class are provided.

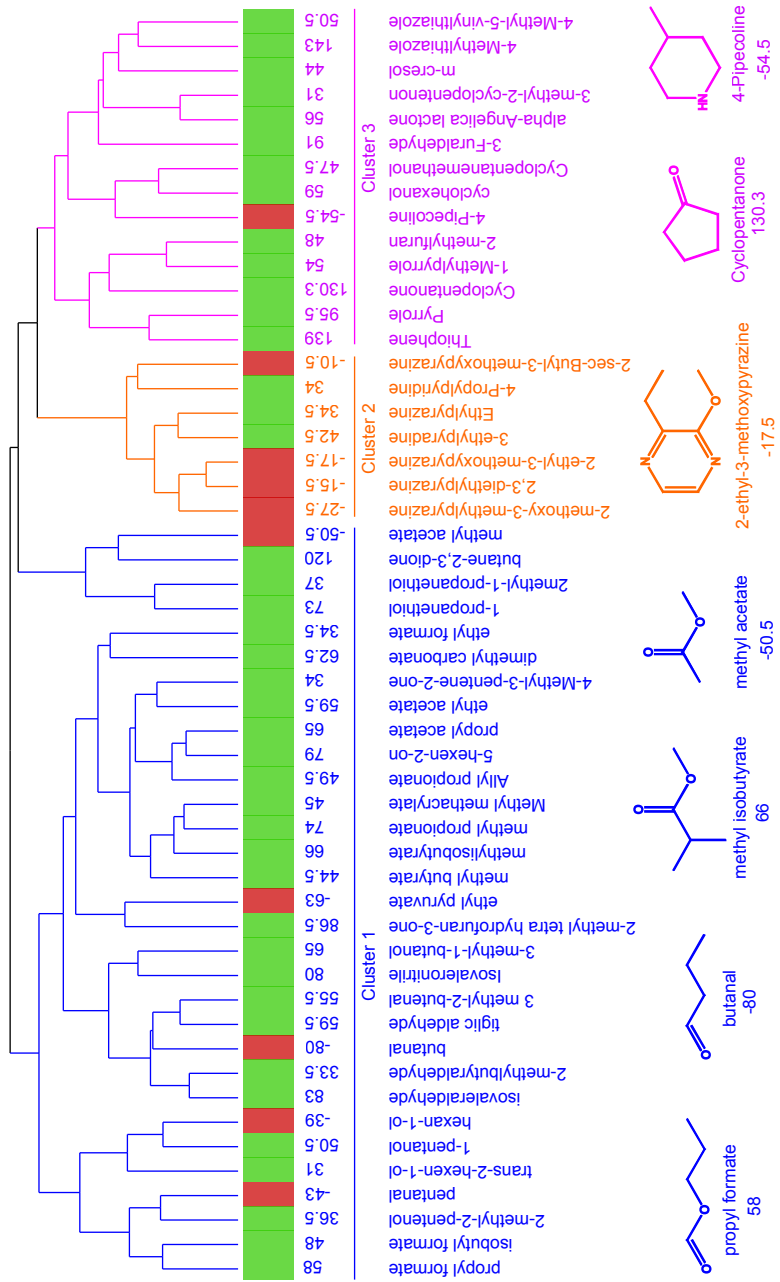
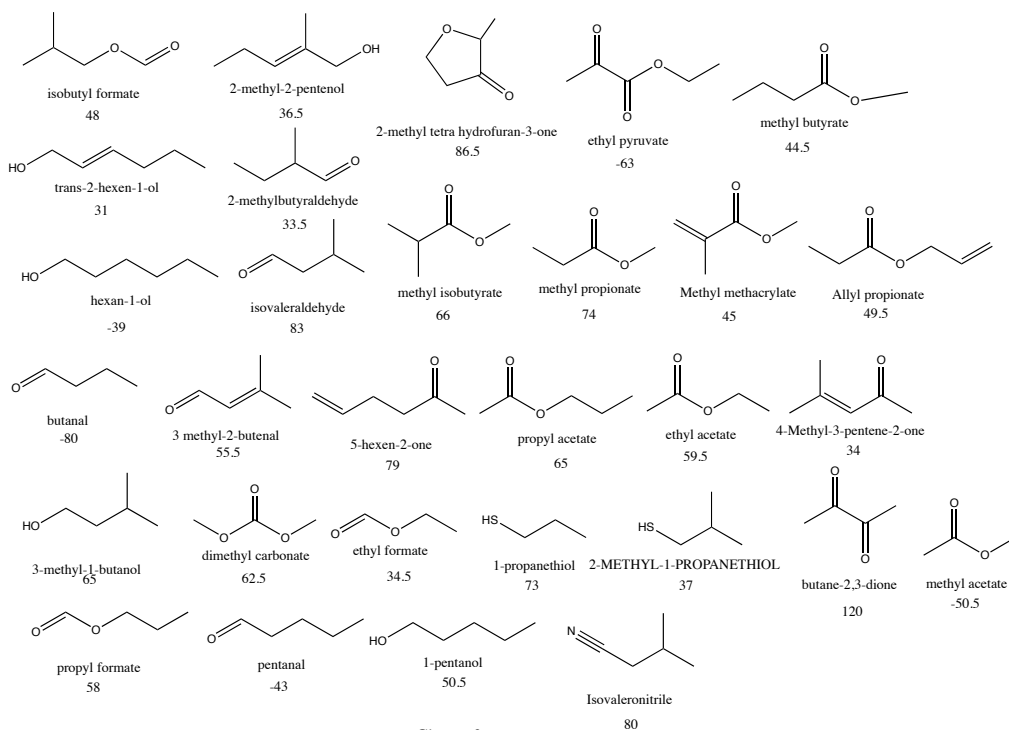


Figure 3.10

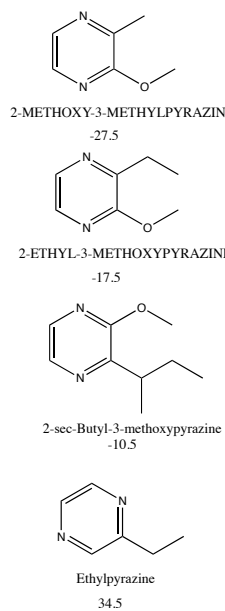
Figure 3.11: Chemical structures of validated activators and inhibitors

The structures for each of the compounds that validated as either activators (>30 spikes/sec) or inhibitors (<0) spikes/sec are provided. Compounds are divided their cluster as described in figure 3.22.

Cluster 1



Cluster 2



Cluster 3

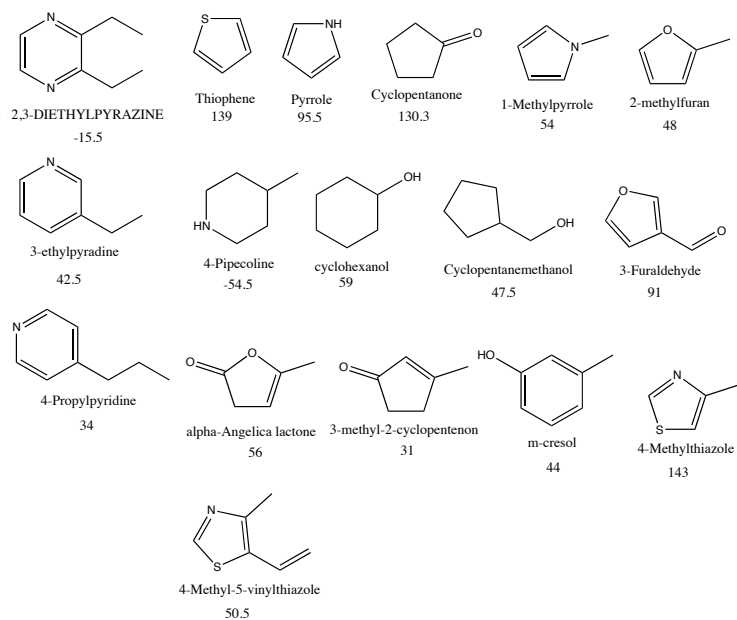


Figure 3.11

Figure 3.12: A SVM integrated molecular descriptor optimized approach is able to explain odor activity for individual Citrus Psyllid ORNs

Schematic of our SVM integrated chemical informatics pipeline. Molecular descriptors that are most correlated with activity are selected, resulting in a metric that is able to cluster together highly active odors using important structural features. The optimized descriptor sets are then be applied to train a SVM to predict Or activity, which is then applied to predict the activity of a large untested odor space.

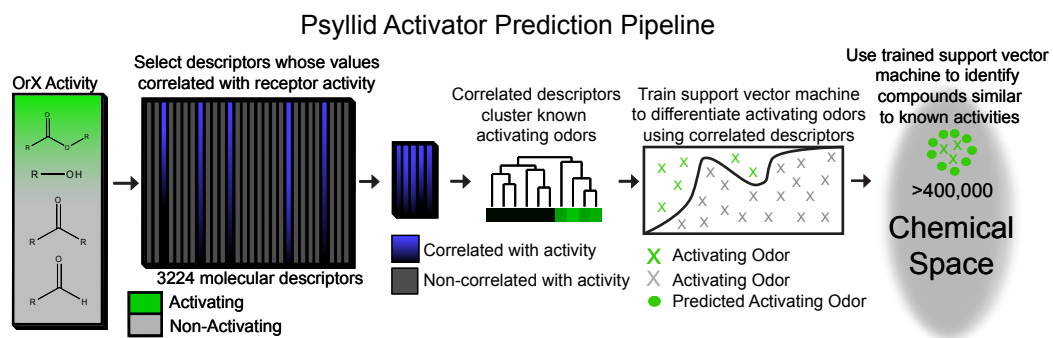


Figure 3.12

Figure 3.13: Optimized molecular descriptor sets are able to cluster ORN activators

Optimized molecular descriptor values were applied to cluster training set odors individually for each ORN. Odor activity is represented individually for each screening set, ranging from gray to black in increasing activity.

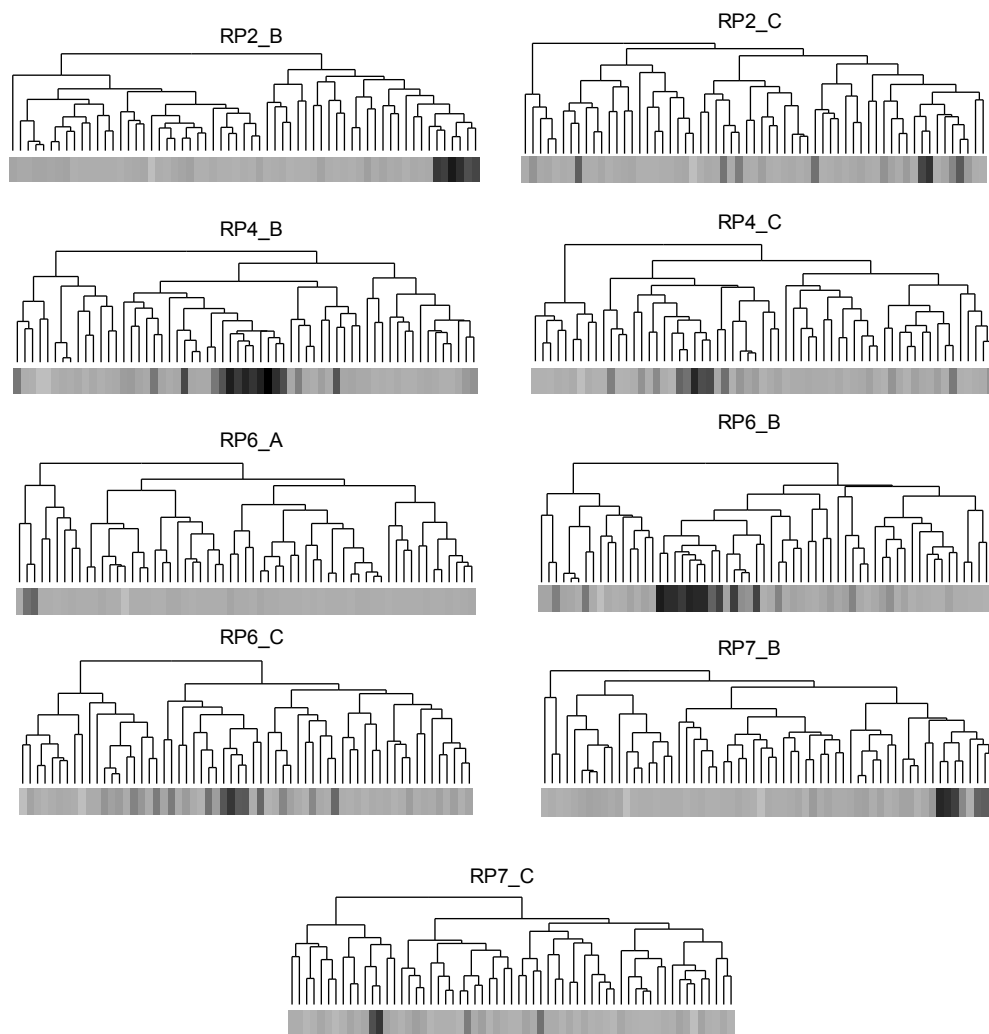


Figure 3.13

Figure 3.14: SVMs trained using descriptor sets that were optimized for individual ORNs effectively describe training set activity

A Receiver-Operating-Characteristic (ROC) curve determined from a SVM applied 5-fold cross validation is plotted for each ORN, depicting the predictive ability of the ORN-optimization approach.

Averaged ROC values for each Or

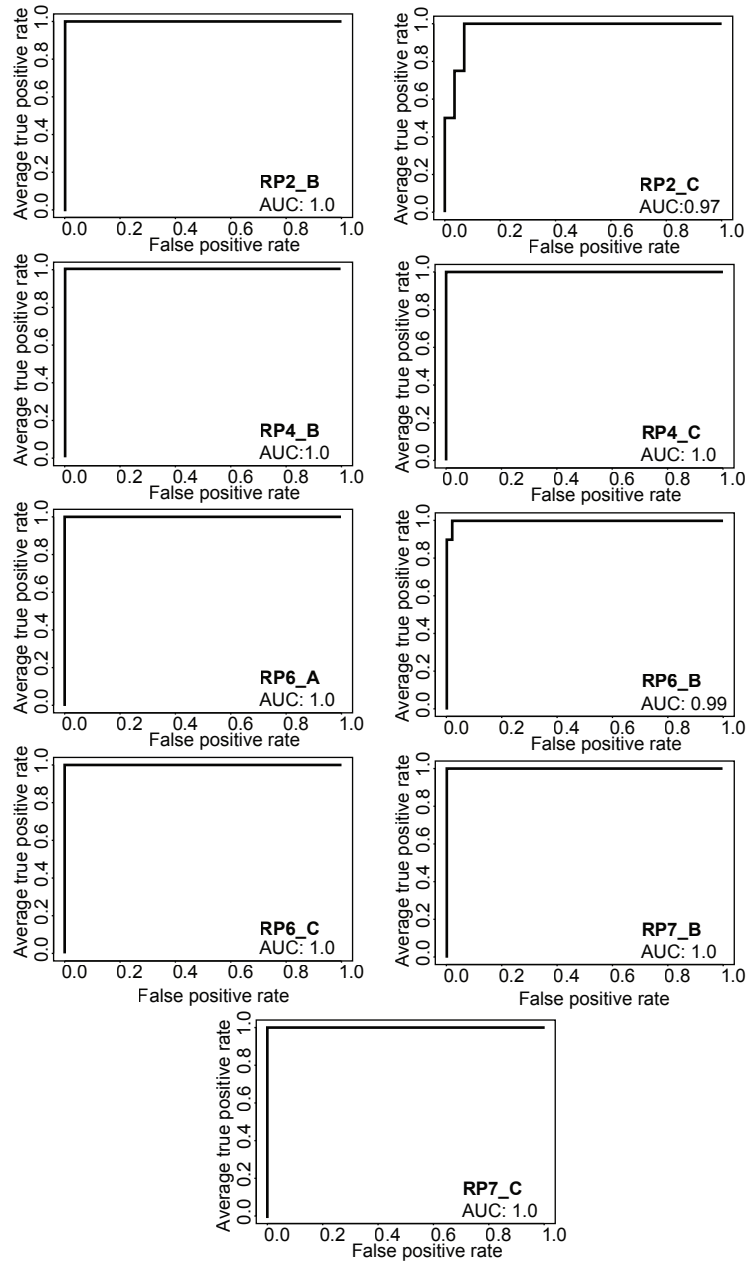


Figure 3.14

Figure 3.15: A molecular descriptor optimized approach is able to explain odor activity for individual mammalian ORs

Schematic of our chemical informatics pipeline. Molecular descriptors that are most correlated with activity are selected, resulting in a metric that is able to cluster together highly active odors using important structural features. The optimized descriptor sets can then be applied to predict Or activity against a large panel of odors.

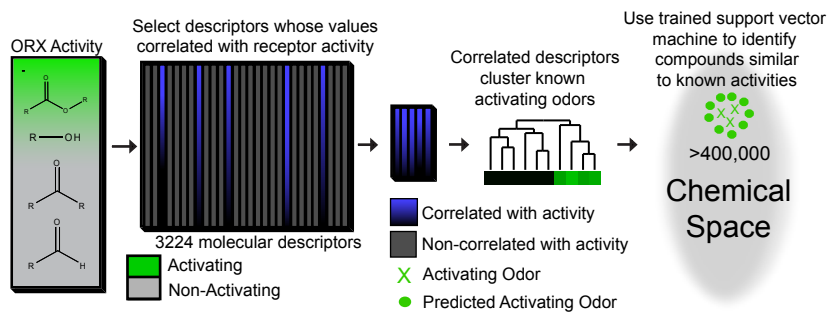


Figure 3.15

Figure 3.16: Mammalian Odorant receptor-optimized molecular descriptors can successfully cluster known ligands

Hierarchical cluster for 33 mouse receptors and 4 human receptors using receptor-optimized descriptor sets using data from (Saito et al., 2009). Known odorant activity scale is indicated using independent color gradient scales.

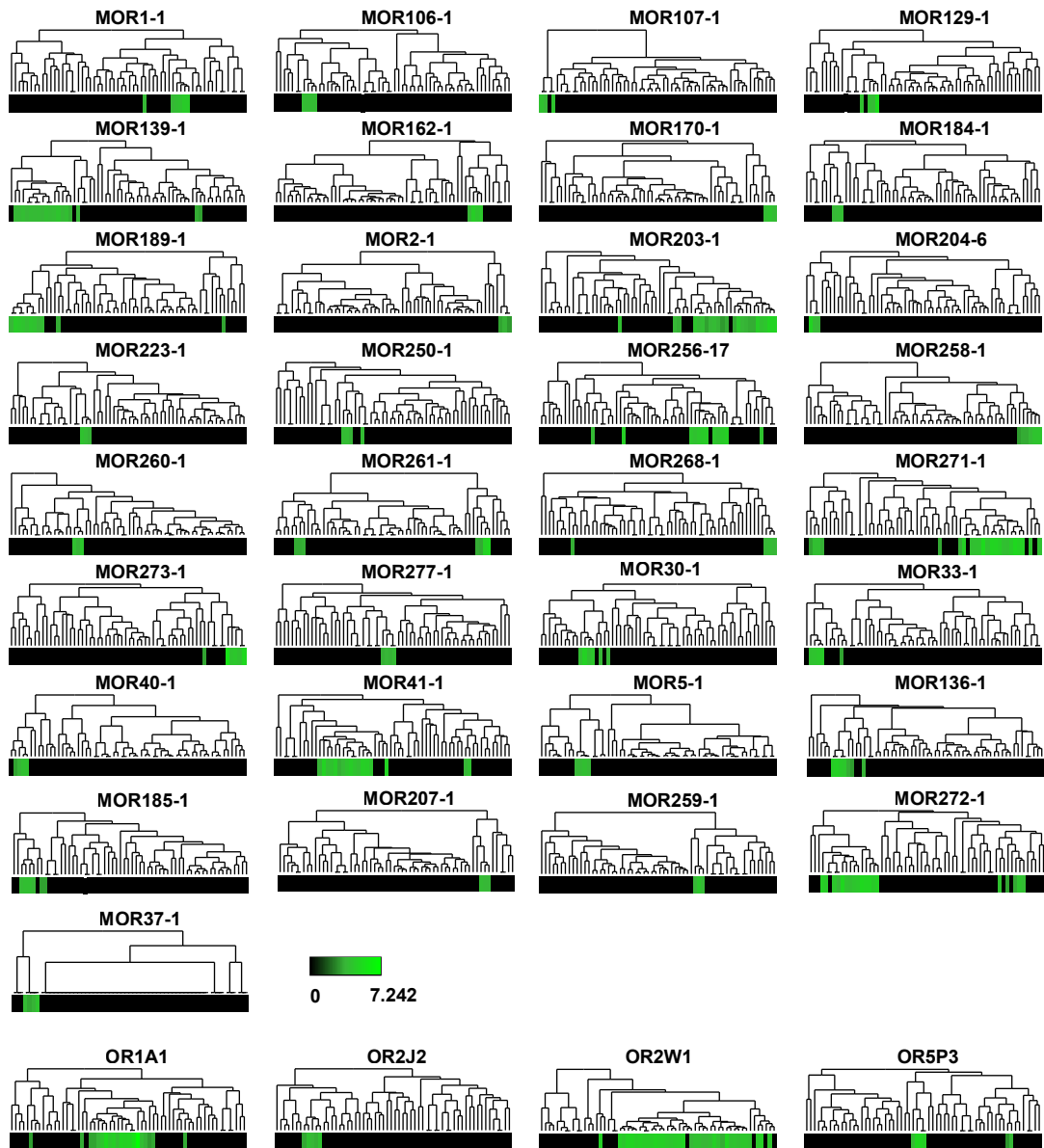


Figure 3.16

Figure 3.17: Optimized descriptors for mammalian ORs effectively describe training set activity

Receiver-operating-characteristic (ROC) curve representing computational validation of ligand predictive ability of the Or-optimization approach. The mean true-positive value from 5 independent 5-fold cross validation runs for 5 mouse and 2 human receptors are plotted.

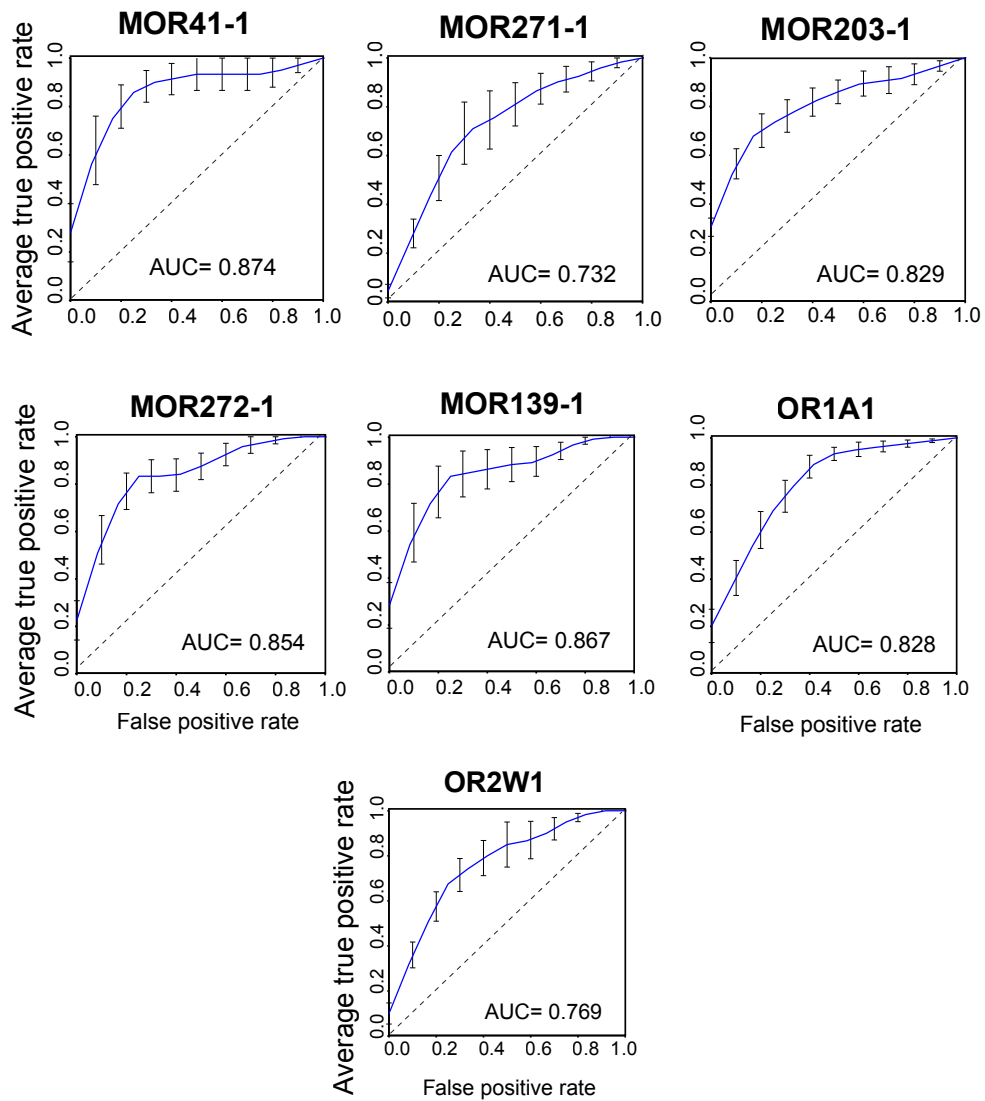


Figure 3.17

Figure 3.18: Analyzing relationships between important features for mammalian ORs

Hierarchical clusters created from Euclidean distance values between mammalian ORs calculated from the following sources: (left to right) shared optimized descriptors; activity to training odor set; similarity across top 500 predicted ligands; and Phylogenetic tree of receptors. Phylogenetic tree of Mammalian ORs calculated from protein sequences. Subclusters are shaded with different colors or bars to enhance ease of comparison.

Figure 3.19: Analysis of mammalian OR tuning breadth

Frequency distribution of compounds from the >240K library within the top 15% distance from highest active plotted to generate predicted breadth of tuning curves.

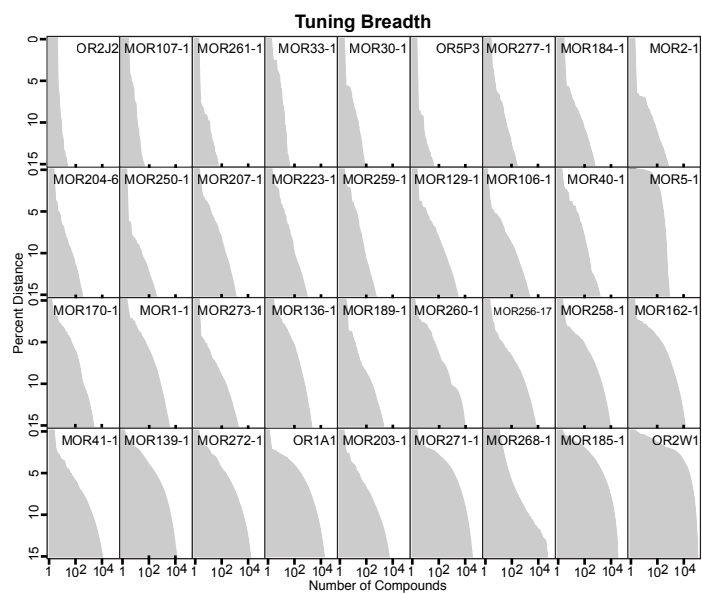


Figure 3.19

Table 3.1: Optimized descriptor sets for each *Anopheles* Or

Optimized descriptor symbols are listed for each Or. Descriptors are listed in the order in which they were included into the list. The same descriptor can be present multiple times for an Or representing the importance of the descriptor for a particular Or. The final correlation between the optimized descriptor set and odor activity is provided.

Or	Desc. Symbol	Or	Desc. Symbol	Or	Desc. Symbol	Or	Desc. Symbol
AgOr57	AlogP98	AgOr12	B06.C.C.	AgOr50	SP05	AgOr38	AlogP98
Correlation	O.057	Correlation	nRCOOH	Correlation	O.057	Correlation	DELs
0.933835	EEig10d	0.8635844	BELe2	0.9031393	piPC08	0.8536688	B06.C.C.
	nRCOOH		EEig10d		H1p		nRCHO
	B07.C.C.		nRCO		R8m.		Atype_C_40.1
	Hy		nRCO		nCconj		R5m.
	Ds		C.040		E2u		PW4
	EEig02d		Hy		Gs		nRCOOH
	Mor27p		B06.C.C..1		JGI6		EEig10d
	B07.C.O.		BEHp4		L2u		MATS2e.1
	Ms		nRCOOH.1		O.057.1		GATS2e
	O.057.1		EEig 11d		B07.C.C.		HATS8e
	HATS8v		HATS8u		piPC05		O.057
	MATS2e		O.057		De		L2u
	EEig10d.1		C.001		JGI7		G.040
	nRCOOH.1		JGI6		H0e		B06.C.C..1
	nRCHO		Ms		MATSev		JGI2
	GV WAI.80		B02.C.C..1		C.040		RDF035m
	O.057.2		O.057.10		X5 Av		B07.C.C.
	F03.C.N.		nR09.1		nR.Ct		Ms
	BELe2		nRCOOH. 11		R Tu.		nR.Ct
	S_dssC		B07.C.C..7		P2u		MATS2e
	O.057.3		C.026.2		piPC08.1		O.057.1
	MATS3v		nRCOOH.12		O.057.2		F03.C.N.
	B07.C.C..1		B07.C.C..8		BEHp4		B06.C.C..2
	O.057.4		nRCOOH.13		B07.C.C..1		JGI3
	nRCO		EEig10d.4		O.060		nRCOOR
	R3u.		JGI1		Mor25e		nRCO
	EEig09r		JGI1.1		X5 Av.1		MATS7v
	nRCOOH.2				C.026		O.057.2
	MATS7v				Hy		E2s
	Hy.1				E2e		GATS2e.1
	C.026				Ds		G2m
	Infective.80				F04.C.O.		Hy
	nRCOOH.3				B07.C.C..2		B07.C.C..1
	B07.C.C..2				HNar		Ms
	MATS2e.1				nCconj.1		nR.Ct.1
	BELe2.1				R7m.		Vindex
	nRCOOH.4				F03.C.N.		B06.C.C..8
	nR09				X4A		Atype_C_40.3
	nRCHO.1				B08.C.C.		MEcc.1
	nR.Ct				B03.C.C.		B07.C.C..5
	O.057.5				RBF		O.057.8
	DISPe				C.026.1		MATS3m.1
	GV WAI.80.1				MATS2e		B06.C.C..9
	nCbH				R1e		E2s
	nRCOOH.5				nR.Ct.1		GATS2e.1
	Hy.2				C.040.1		Hy
	S_dCH2				B07.C.C..3		B07.C.C..1
	MATS4v				B04.C.C.		G.N..O.1
	B07.C.C..3				piPC08.2		O.057.9
	Ms.1				E2e.1		nRCOOR.1
	O.057.6				O.057.3		nR.Ct.2
	EEig10d.2				O.057.4		Mor10u
	G2m						R4e.
	nRCOOH.6						G.N..O.
	B07.C.C..4						B07.C.C..2
	F03.C.N..1						nRNH2
	C.026.1						B07.C.C..6
	Atype_C_18						Yindex
	B02.C.C.						S_aaaC.1
	HATS8e						B06.C.C.. 11
	nRCOOH.7						nRCO.3
	MATS3m						B02.C.C.
	nRCHO.2						B06.C.C..12
	nRCOOH.8						O.057.10
	EEig12x						nCb.
	O.057.7						C.040.2
	nRCO.1						GATS8m
	B07.C.C..5						JGI3.2
							GATS8m
							G1u.1
							nRCOOH.3
							C.026.1
							E2e
							Atype_C_18
							B06.C.C..5
							nRCOOH.1
							R1u.
							EEig10d.2
							JGI3.3
							nRCO.1
							O.057.5
							nRCOOR.3
							nR.Ct.1
							Hy.1
							nR.Ct.2
							Atype_C_40
							B07.C.C..3
							S_aaaC
							JGI1
							B06.C.C..6
							MEcc.2
							O.057.6
							GATS2e.2
							B06.C.C..7
							Yindex.1
							J.5
							J.3
							MATS4v

Table 3.1

Or	Desc. Symbol			Or	Desc. Symbol	Or	Desc. Symbol	Or	Desc. Symbol	Or	Desc. Symbol
AgOr56	Ss	C.024.2	TIC4	AgOr21	SP02	AgOr20	O.058	AgOr18	Mor20m	AgOr4	
Correlation	nHAcc	nRCOOH.3	nRCOOH.8	Correlation	O.057	Correlation	EEig06d	Correlation	BELe2	Correlation	
0.9312913	X4Av	nRCO	S_sNH2.1	0.9052854	piPC08	0.8829924	L2s	0.9201605	O.057	0.896526	
	B07.C.O.	Atype_C_18	EEig09d.3		TPSA.NO.		C.040		EEig10d		
	O.057	nR09.1	O.058.5		B06.C.C.		GATS7m		Gu		
	Hy	nRCOOH.4	MATS2e.7		H1p		X3Av		DISPe		
	EEig09d	B07.C.C..5	O.057.16		Yindex		De		nRCHO		
	MATS2m	C.026.2	GVWAI.80.4		ASP		C.040.1		EEig12r		
	S_dssC	MATS3p	nRCHO.6		O.057.1		EEig10r		MATS4p		
	B07.C.C.	BELe2.2	O.057.17		R8m.		BEHm4		J		
	Ms	Atype_C_40	S_ssO.5		nDB		O.058.1		O.057.1		
	nBnz	nRCHO.3	Mor04m.1		H1e		Hy		JG17		
	nRCHO	nArCHO	MATS2e.8		TPSA.NO..1		JG17		X3A		
	C.040	GVWAI.80.2	nRCOOH.9		EEig10d		E2e		Mor23m		
	BELe2	nRCOOH.5	EEig09d.4		nRCHO		nConj		B06.C.C.		
	O.057.1	nR.Ct.3	EEig09d.5		BEHp4		HNar		nHDon		
	EEig05d	nRNH2			O.057.2		JG16		Mor27p		
	MATS2e	O.057.7			B06.C.C..1		Ku		nRCOOH		
	B04.C.N.	CIC5			B04.C.O.		JG15		GATS3m		
	O.058	S_ssO.1			C.040		B03.C.C.		Mor28e		
	O.057.2	EEig09d.1			X4Av		C.040.2		C.026		
	GVWAI.80	RDF150u			O.057.3		nRCHO		X4A		
	MATS8v	nRCOOH.6			R7m.		MATS4v		S_dssC		
	IC2	nRCO.1			Hy		C.040.3		E2e		
	O.057.3	O.057.8			E2u		nRCO		MATS2v		
	nR.Ct	B07.C.C..6			O.057.4		R8m.		nRCHO.1		
	C.040.1	nPyrroles			Neoplastic.80		IC1		nRCOOH.1		
	C.024	MATS3p.1			EEig10d.1		O.057		Infective.80		
	MATS3v	S_ssO.2			MATS2p		O.060		nR09		
	Xt	O.057.9			O.057.5		MATS6v		HATS3u		
	B07.C.C..1	nRCHO.4			B06.C.C..2		R1e		R8e.		
	nRCOOH	X4A.1			BEHp4.1		L2e		nRCOOH.2		
	F03.C.N.	O.057.10			F03.C.N.		piPC09		JG11		
	O.058.1	B07.C.C..7			Ms		nRCOOH		C.026.1		
	C.026	nR.Ct.4			HATS6m		nR.Ct		RARS		
	MATS2e.1	O.057.11			B07.C.C.		De.1		MATS4v		
	nRCHO.1	C.043.1			O.057.6		PW4		B06.C.C..1		
	C.040.2	S_ssO.3			nRCHO.1		O.060.1		nRCO		
	X4A	MATS7v.1			Yindex.1		E2e.1		nRCOOH.3		
	B07.C.C..2	MATS2e.4			C.040.1		G2m		Gu.1		
	Hy.1	O.057.12			C.026		nRCOOH.1		nRCHO.2		
	O.057.4	Mor04m			B06.C.C..3		nRCHO.1		EEig12x		
	GVWAI.80.1	CIC5.1			ESpm05d		S_aaCH		nRCOOH.4		
	MATS7v	B07.C.C..8			BEHp4.2		C.040.4		X4A.1		
	C.043	O.057.13			nRCOOH		BEHe3		EEig10d.1		
	C.040.3	S_sNH2			EEig10d.2		nRCO.1		MATS5p		
	nR.Ct.1	EEig09d.2			H1e.1		JG16.1		PJ12		
	nR09	O.058.3			H1e.2		Ku.1		PJ12.1		
	Ms.1	Mor24e.1					GATS6v				
	C.026.1	C.026.3					X5Av				
	MATS2e.2	MATS3p.2					C.040.5				
	O.058.2	nR.Ct.5					MATS5v				
	Mor24e	Atype_C_40.1					B07.C.O.				
	G1v	Ms.2					R6e.				
	nRCOOH.1	B07.C.C..9					R6e..1				
	B07.C.C..3	nRCHO.5									
	C.024.1	nRCOOH.7									
	C.040.4	GVWAI.80.3									
	BELe2.1	MATS2e.5									
	EEig10d	S_aaNH									
	nRCOOH.2	RDF130m									
	MATS1v	O.057.14									
	nRCHO.2	S_ssO.4									
	nR.Ct.2	X4A.2									
	O.057.5	O.057.15									
	S_ssO	O.058.4									
	O.057.6	C.026.4									
	Hy.2	MATS2e.6									
	B07.C.C..4	nR.Ct.6									
	MATS2e.3	B07.C.C..10									

Table 3.1 (Continued)

Desc. Symbol	Or	Desc. Symbol	Or	Desc. Symbol	Or	Desc. Symbol	Or	Desc. Symbol	Or	Desc. Symbol
BLTF96	AgOr75	ATS5p	AgOr32	E1s	AgOr26	Mor31v	AgOr11	ATS6e	AgOr15	ARR
F05.C.O.	Correlation	S_sCH3	Correlation	B06.C.O.	Correlation	O.057	Correlation	HNar	Correlation	B06.C.O.
B08.C.O.	0.8265765	BELe5	0.6823186	Mor10p	0.9274835	B08.C.O.	0.889735	R8e.	0.9092773	BELe3
PW4		O.057		Atype_H_47		B03.C.O.		O.057		E2e
nRCOOH		ESpm15u		H.049		Mor27v		E2u		nAB
B05.C.O.		HATS7e		nRCOOH		S_dssC		X5A		Mor10u
Infective.80		H2m		HNar		X4A		HATS7p		HOMA
DISPm		G3s		GATS1v		B07.C.C.		Mor13m		R7m.
BEHp1		DISPe		nCconj		R6e.		EEig10d		SPH
nRCOOH.1		Mor07u		R7v.		Mor25e		Mor25e		E2e.1
HATS5v		ESpm15u.1		GATS7m		Mor27v.1		O.057.1		F04.C.N.
B08.C.O..1		P2s		CIC		nRCOOH		S_ssO		HOMA.1
B05.C.C.		nRCOOH		B06.C.O..1		nDB		E2u.1		nOHs
EEig09d		Infective.80		Mor10u		MATS2v		nDB		GATS8m
DISPm.1		BELe4		E1s.1		HATS8m		X5A.1		GATS3e
nR.Ct		GATS5v		RDF020m		Gs		B06.C.O.		B06.C.C.
Gs		R6u.		R8m.		O.057.1		BELv4		nRCOOH
MATS3e		B09.C.O.		B06.C.O..2		De		BELv4.1		X4A
nRCOOH.2		G3s.1		nRCHO		B07.C.O.				ESpm01d
Mor10u		C.026		DISPm		Mor27v.2				B06.C.C..1
R5e.		ESpm15u.2		R7e		BEHm1				HATS4u
B08.C.O..2		nRCOOH.1		MATS2m		nDB.1				MATS4v
G3s		HATS7p		Mor10p.1		F04.C.O.				B06.C.O..1
EEig10d		C.016		P2s		B08.C.O..1				HATS3u
GATS4p		C.016.1		P2s.1		O.057.2				B06.C.C..2
nRCOOH.3						Mv				BIC
EEig12r						E2u				nRCO
F03.C.N.						O.060				E2e.2
RARS						RDF040u				R2v.
nRCOOH.4						nR.Ct				JG17
R5e..1						MATS3v				nAB.1
R8m.						O.057.3				GATS6v
Mor32v						F02.C.N.				X4A.1
HATS4u						B08.C.O..2				Infective.80
HATS4u.1						piPC05				nRCOOH.1
						C.026				B04.C.N.
						Mor25u				nOHs.1
						R5e.				B06.C.C..3
						B07.C.C..1				X4A.2
						nRCOOH.1				JGT
						Mor27v.3				HOMA.2
						MATS6v				C.026
						R5u.				Mor09u
						B08.C.O..3				G2s
						piPC06				nRCOOH.2
						nCconj				nRCHO
						Mor27v.4				Infective.80.1
						B08.C.O..4				MATS1v
						MATS2v.1				nRCOOH.3
						R1u.				nOHp
						R1u..1				E1p
										R2v
										R6m.
										BIC.1
										BIC.2

Table 3.1 (Continued)

Or	Desc. Symbol	Or	Desc. Symbol	Or	Desc. Symbol	Or	Desc. Symbol	Or	Desc. Symbol	Or	Desc. Symbol
AgOr39	ALOGP	AgOr53	Ss	AgOr16	BELe3	AgOr48	F06.C.O.	AgOr46	Mor25e	AgOr30	
Correlation	RDF035m	Correlation	O.057	Correlation	nHAcc	Correlation	BELv7	Correlation	BEHp3	Correlation	
0.8405153	nRCOOH	0.8962904	EEig09d	0.853112	EEig10r	0.8627068	piPC04	0.8647854	C.006	0.8920303	
	R5e.		nHAcc		ARR		DISPe		E2u		
	RBF		BELe2		R5v.		R8u.		O.057		
	ALOGP.1		B04.C.N.		Mor10e		ESpm03u		EEig10d		
	HATS6e		MATS2v		nRCOOH		F06.C.O..1		HATS7m		
	B09.C.O.		HATS6e		EEig10d		O.057		B04.C.O.		
	MATS8m		O.057.1		R2v.		Mor18m		ESpm15u		
	B05.C.O.		Infective.80		X4A		GATS1v		O.057.1		
	nConj		O.057.2		nRCHO		R8u		S_dssC		
	Mor10e		R7m.		S_dssC		JGT		S_dssC.1		
	F02.C.O.		R1u.		MATS2v		Atype_C_40				
	nRCOOH.1		Mor25p		B04.C.N.		G3s				
	R6u.		BELe2.1		BIC		DISPm				
	Mor27u		MATS3e		JG17		piPC04.1				
	B06.C.O.		O.057.3		B04.C.C.		ESpm05u				
	Ke		nRCHO		EEig11x		S_dssC				
	nHDon		C.040		IC2		AROM				
	B09.C.O..1		EEig08d		nRCOOH.1		nRCOOH				
	R5e..1		Infective.80.1		E1p		B06.C.O.				
	MATS8m.1		F04.C.N.		EEig10d.1		EEig12x				
	MATS3v		E2s		X4A.1		E1s				
	RDF035m.1		Mor32p		Mor10e.1		E1s.1				
	nRCOOH.2		X3Av		TPSA.NO.						
	RBF.1		HATS8u		S_dCH2						
	R6u..1		O.057.4		EEig09d						
	nR.Ct		B06.C.C.		L2u						
	L3s		HATS4e		F04.C.O.						
	Mor23m		O.057.5		E1s						
	nRCHO		Mor10e		E2e						
	S_sOH		E2u		BIC.1						
	HATS6u		RTu.		Mor10e.2						
	EEig09x		EEig09d.1		EEig10d.2						
	R5e..2		MATS5v		C.027						
	C.008		SRW09		nRCOOH.2						
	SPH		C.040.1		Mor32m						
	DISPe		nR.Ct		B04.C.C..1						
	nDB		O.057.6		JG17.1						
	X4A		EEig08d.1		Mor08p						
	B09.C.C.		MATS2e		MATS3m						
	R6u..2		E1p		EEig10d.3						
	RDF035m.2		B04.C.N..1		G1p						
	MATS7m		O.057.7		nRCHO.1						
	MATS7m.1		nArCO		nRCOOH.3						
			nRCHO.1		RDF035m						
			C.040.2		BIC2						
			C.026		JG16						
			BELe2.2		B05.C.O.						
			MATS2v.1		J						
			Infective.80.2		X4A.2						
			nRCOOH		Mor16e						
			G2p		Mor16e.1						
			C.043								
			X4Av								
			Mor27e								
			C.040.3								
			C.026.1								
			BELe2.3								
			MATS6v								
			E2s.1								
			B06.C.C..1								
			Mor10p								
			Mor10p.1								

Table 3.1 (Continued)

Desc. Symbol	Or	Desc. Symbol	Or	Desc. Symbol	Or	Desc. Symbol	Or	Desc. Symbol
Kappa.3.AM	Ks.1	Mor14e	AgOr64	nAB	AgOr13	nAB.7	ARR	ESpm01d
O.057	nRCOOH.8	R3v.	Correlation	ESpm01d	Correlation	MATS3m.2	B05.C.C.	S_sOH
EEig01x	EEig11x	Mor25e	0.7699497	Mv	0.8117099	R7m..2	EEig03r	B06.C.C.
H2m	EEig11x.1	G3s		B07.C.O.		nArCO.1		GATS2p
B06.C.O.		HATS8e		Mor18p		nAB.8		MAXDN
B09.C.O.		MAXDN		MATS3v		E1m		Mor16e
O.058		AROM		B06.C.O.		nAB.9		S_sOH.1
O.057.1		Mor13v		nAB.1		GATS6v.1		C.040
R8p.		Mor13m		H0e		F04.C.N.		IC
ESpm15u		MATS2p		nAB.2		BIC.1		E1s
EEig11d		JG15		EEig11r		Infective.80		B06.C.C..1
G3s		R4e.		Mv.1		X4A		B03.C.N.
nRCOOH		R7m.		R7m.		EEig10r		O.058
Mor18m		G3s.1		HATS3v		nAB.10		GATS2e
MATS3e		Mor13v.1		EEig04x		BELv4.2		E1p
PW5		O.057		nAB.3		MATS2e.1		E1p.1
EEig12r		nRCHO		R3m		nAB.11		
O.057.2		nHDon		F03.C.N.		R7m..3		
G3s.1		Atype_H_51		B06.C.O..1		S_aaaC.1		
ESpm15u.1		C.040		MATS2e		RCI.4		
S_sOH		R2e		nAB.4		nRCOOH.3		
nConj		RTm.		nConj		X5A.3		
B01.C.N.		B06.C.C.		Mv.2		AMW		
B09.C.O..1		nCb.		JG17		G2m		
C.026		nRCOOR		HATS3p		nR.Ct		
X5A		B03.C.C.		Gu		MATS3m.3		
nRCOOH.1		JG15.1		nRCOOH		Mor18m		
Ks		R3v..1		nBM		B06.C.O..2		
EEig12r.1		Mor13v.2		nArCO		X5A.4		
B06.C.O..1		Mor13v.3		E3p		nRCOOH.4		
ESpm15u.2				nAB.5		E2e.1		
nRCOOH.2				MATS3m		O.056.1		
G3s.2				R8p.		R3m.1		
Mor10u				HATSse		R8p..2		
GATS5e				RCI		BIC.2		
nRCOOH.3				nAB.6		BIC.3		
C.026.1				EEig10d				
ESpm15u.3				GATS6v				
Hypnotic.80				X5A				
C.027				BELv4				
nRCOOH.4				X5A.1				
B06.C.O..2				S_dCH2				
ESpm15u.4				nRCOOH.1				
nR.Ct				R7m..1				
PW5.1				MATS4p				
nRCOOH.5				nBM.1				
nCb.				BELv4.1				
B09.C.O..2				nRCHO				
B06.C.O..3				SPH				
S_dCH2				S_aaaC				
ESpm15u.5				S_dO				
HATS6u				RCI.1				
nRCOOH.6				nArCHO				
nArOR				RCI.2				
B03.C.O.				H3m				
L2s				MATS3m.1				
O.057.3				R8p..1				
PW5.2				G1u				
ESpm15u.6				Mp				
R8p..1				E2e				
EEig12x				BIC				
O.057.4				JG17.1				
G3s.3				T.O.O.				
MATS5p				RCI.3				
B09.C.O..3				Mor32u				
PW5.3				EEig12d				
nRCOOH.7				X5A.2				
G3s.4				O.056				
nCb..1				EEig09d				
ESpm15u.7				nRCOOH.2				

Table 3.1 (Continued)

Or	Desc. Symbol	Or	Desc. Symbol	Or	Desc. Symbol	Or	Desc. Symbol	Or	Desc. Symbol	Or	Desc. Symbol
AgOr31	Tl2	AgOr10	nCar	AgOr27	Mor31m	AgOr6	nArCO	AgOr2	C.024	AgOr65	
Correlation	O.057	Correlation	R8e.	Correlation	MEcc	Correlation	RCI	Correlation	HATS6v	Correlation	
0.8778699	HATS8v	0.9310224	C.025	0.8758972	O.057	0.9545822	BELm5	0.8565036	HATS3v	0.863439	
	Mor32e		ZM1V		R8m		O.057		E2e		
	G3s		H.047		MATS2m		ESpm03u		Gs		
	Mor10v		X5A		nBnz		R7m.		Gs.1		
	S_sOH		RTm		Ms		EEig03x				
	B06.C.O.		RCI		MATS8v		nDB				
	ESpm15u		MATS2p		C.006		Mor23m				
	R6u.		R8u.		B09.C.O.		E2e				
	nRCOOH		nBnz		R2e.		Kp				
	E3p		HATS5m		O.057.1		nCconj				
	nHDon		Kp		nCb.		S_sOH				
	nCrs		JGI1		B08.C.C.		ESpm01d				
	R5u.		O.060		X3A		nAB				
	GATS6v		nBnz.1		S_aasC		F04.C.O.				
	nRCOOH.1		GATS6m		C.040		B06.C.O.				
	B06.C.O..1		JGI6		nCq		GATS3e				
	Mor23m		B04.C.C.		MATS8v.1		GATS3e.1				
	Mor32e.1		Mor25p		MATS5p						
	HATS6u		MATS3v		C.024						
	IVDE		O.056		nRCOOH						
	O.057.1		nRCOOR		Neoplastic.80						
	nOHp		Mor25v		C.040.1						
	D.Dr05		H.051		X3A.1						
	B06.C.O..2		E3m		nCb..1						
	nRCHO		P2v		nRCHO						
	O.057.2		nOHs		MATS3e						
	B06.C.O..3		C.006		B07.C.C.						
	C.008		EEig04r		R6u.						
	nRCOOH.2		R4v.		C.040.2						
	R6u..1		Mor11p		MATS6v						
	C.026		GATS6v		C.024.1						
	E1s		B04.C.C..1		REIG						
	MATS5v		JGI6.1		R8m.						
	Mor32e.2		nBnz.2		MEcc.1						
	G3s.1		MATS3e		O.057.2						
	nCconj		HOMA		C.034						
	nOHp.1		HATS5m.1		O.057.3						
	EEig12x		RBF		MATS1v						
	JGI5		Mor07u		nOHp						
	nRCOOH.3		C.027		C.040.3						
	G3s.2		JGI7		HATS8v						
	JGI6		AROM		O.057.4						
	Mor10p		JGI1.1		EEig08d						
	R5e.		R8u..1		EEig08d.1						
	B09.C.O.		RCI.1								
	HATS6u.1		RCI.2								
	GATS6v.1										
	nRCOOH.4										
	GATS5v										
	nArCO										
	B05.C.O.										
	RDF125u										
	C.026.1										
	C.026.2										

Table 3.1 (Continued)

Desc. Symbol	Or	Desc. Symbol	Or	Desc. Symbol	Or	Desc. Symbol	Or	Desc. Symbol	Or	Desc. Symbol
nArOH	R6u.	R6m.	AgOr59	S_dCH2	AgOr63	nH	AgOr61	nHAcc	AgOr35	Hypnotic.80
TPSA.NO.	Ke.1	GATS2p	Correlation	nOHs	Correlation	Dv	Correlation	EEig09d	Correlation	GGI9
REIG	E3s	F01.C.O.	0.7551107	R8u.	0.8400413	O.057	0.8538082	Mor17e	0.8460269	R1u
HATS4m	nAB.2	J3D		nROH		B04.C.O.		O.057		UI
B07.C.O.	R1m.	B07.C.O.		JGI4		H2m		MATS2e		S_dsCH
E1u	MATS4p.1	B06.C.C.		Hbond.acceptor		H.049		EEig10d		R8v.
SPAM	nRCHO.2	nDB		JGI7		R5e.		Gs		S_sOH
MATS4p	nRCOOH.3	JGI1		JGI5		O.057.1		PW4		GGI9.1
S_dO	nArCO.2	Atype_C_40		R1v.		Mor13p		B08.C.O.		nDB
Gu	nAB.3	ESpm05u		Mor17m		G3s		O.057.1		EEig09d
Mor16v	nArOR	EEig09d		ESpm01d		B08.C.O.		O.060		GATS4e
GATS7m	G2u.1	GATS4e		JGI5.1		GATS5v		PW4.1		JGI9
Mor32e	nAB.4	GATS4e.1		nR.Cs		Vindex		nHDon		DISPv
C.040	JGI5			nR.Cs.1		nRCOOH		O.057.2		S_dssC
B04.C.C.	Mor04m.1					Ui		Infective.80		JGI7
nRCHO	nAB.5					B09.C.O.		MATS2v		O.056
EEig11r	R8m.					B06.C.C.		MorZ3m		Mor16p
AlogP	Mor16p.3					MATS7v		MorZ3m.1		E1u
nArOH.1	R6u..1					MATS5e				DISPm
G2u	Atype_C_40					G3s.1				GGI9.2
nRCOOH	F03.C.N.					nCconj				Hypnotic.80.1
Mor10e	MATS4p.2					B08.C.O..1				JGI3
nArCO	nCconj.1					PW5				S_dsCH.1
MATS2v	HATS1p					B06.C.C.1				Mor09u
R8u.	E1u.3					EEig10d				R8v..1
nAB	E1u.4					B09.C.O..1				JGI4
F04.C.O.						G3s.2				JGI4.1
H0e						MATS4e				
Ke						O.057.2				
E1u.1						R8m.				
JGI7						GATS6v				
Mor16p						HATS6u				
Hy						Mor17m				
SPAM.1						Mor17m.1				
nRCOOH.1										
R4m.										
nAB.1										
H.049										
MATS4v										
nR.Cl										
C.040.1										
nR09										
HATS3p										
MATS7v										
E1u.2										
E2e										
Mor16p.1										
nCconj										
nOHs										
HATS4p										
nHDon										
MATS2e										
nArCO.1										
C.040.2										
B04.C.C..1										
nRCHO.1										
EEig11d										
JGI7.1										
B04.C.C..2										
Mor16p.2										
nArCOOR										
J3D										
HATSp										
nBM										
HATSe										
nRCOOH.2										
C.006										
Mor04m										
MATS2e.1										
Infective.80										

Table 3.1 (Continued)

Or	Desc. Symbol	Or	Desc. Symbol	Or	Desc. Symbol	Or	Desc. Symbol	Or	Desc. Symbol	Or	Desc. Symbol	
AgOr25	ESpm11u	AgOr3	BLTF96	AgOr76	EEig10d	X4A.3	AgOr44	ATS8m	AgOr66	EEig10x	MATS4p.1	AgOr43
Correlation	nRCOOH	Correlation	nRCOOH	Correlation	nRCOOH	Neoplastic.80	Correlation	EEig12x	Correlation	T12	BEHp5.1	Correlation
0.7611209	nOHp	0.8205165	Infective.80	0.7581273	R1v	nRCOOH.6	0.8356094	nRCOOH	0.866496	DISPe	EEig12d.1	0.8854268
	EEig12x		B04.C.O.		MATS8m	EEig10d.13		EEig10r		SPH	JGI4.3	
	HATS3p		JGI7		nOHT	EEig09d.1		B03.C.O.		B05.C.O.	nRCOOH.3	
	B09.C.O.		nCconj		EEig10d.1	G2m.1		B07.C.O.		B08.C.O.	EEig10d.6	
	nN		G2v		X4A	MATS7p.1		Mor08m		Mor08p	B05.C.O..6	
	nCconj		nRCOOH.1		JGI7	nRCOOH.7		nRCOOH.1		TPSA.NO.	B09.C.O..4	
	nOHp.1		F02.C.N.		MATS4v	nR.Ct		EEig12r		EEig10d	SPH.2	
	nRCOOH.1		B09.C.O.		R1e.	nOHT.4		JGI7		JGI4	C.026.1	
	nThiophenes		PW5		B07.C.O.	B02.C.C..2		nRCHO		F04.C.N.	E1u	
	MATS5v		ATS1p		EEig10d.2	EEig10d.14		MATS3m		EEig12d	E1u.1	
	B09.C.O..1		E2s		MATS6v	MATS8m.5		SPH		HATS4u		
	ESpm14d		Infective.80.1		S_dscC	nRCOOH.8		DISPe		EEig09d		
	E1m		nRCOOH.2		nRCOOH.1	EEig10d.15		Gs		DISPe.1		
	HATS8m		EEig04x		EEig10d.3	MATS4v.1		O.057		GATS4v		
	GATS5e		Hy		B02.C.C.	X4A.4		Infective.80		B09.C.O.		
	nOHp.2		EEig10d		G2m	B06.C.O.		Mor16e		MEcc		
	O.057		nRCOOH.3		MATS8m.1	EEig10d.16		Mor16e.1		PW5		
	F03.C.N.		B02.C.C.		E2s	E2s.2				EEig10d.1		
	nCrs		B08.C.C.		EEig10d.4	E2s.3				nRCOOH		
	E1s		nArCO		MATS7p					B05.C.O..1		
	nRCOOH.2		B04.C.O..1		nOHT.1					Mor10u		
	nOHp.3		nRCHO		EEig10d.5					G2s		
	EEig12d		H.052		MATS4p					MATS8m		
	ESpm15u		nRCOOH.4		RDF085e					Mor13p		
	E1u		PW5.1		EEig09d					DISPe.2		
	EEig09d		Infective.80.2		X4A.1					B05.C.O..2		
	S_dsCH		RDF140m		nRCOOH.2					B09.C.O..1		
	X5A		B04.C.O..2		EEig10d.6					Infective.80		
	DISPe		nRCOOH.5		MATS8m.2					nHBonds		
	Ke		JGI7.1		nRCO					C.026		
	Ke.1		PW5.2		GVWAI.80					MATS5v		
			F02.C.N..1		S_aaaC					SPH.1		
			E2s.1		nRCOOH.3					B09.C.O..2		
			Mor10e		CIC4					B05.C.O..3		
			EEig10d.1		nOHp					EEig10d.2		
			nRCOOH.6		MATS8m.3					JGI4.1		
			G2v.1		MATS4e					nRCOOH.1		
			nCconj.1		EEig10d.7					MEcc.1		
			X5A		nOHT.2					EEig09d.1		
			Infective.80.3		nRCOOH.4					F04.C.N..1		
			nRCOOH.7		JGI7.1					Mor32e		
			C.043		EEig10d.8					C.008		
			B08.C.C..1		nCconj					BEHp5		
			PW5.3		B04.C.N.					EEig10d.3		
			Hy.1		R7m.					MATS4p		
			EEig09d		MATS2e					nPyridines		
			nRCOOH.8		X4A.2					B09.C.O..3		
			CIC		E2e					B05.C.O..4		
			MATS8m		EEig13r					MATS8m.1		
			GATS5e		EEig10d.9					HATS4p		
			MATS2m		PJ13					nRCOOH.2		
			EEig04x.1		MATS8v					Infective.80.1		
			Infective.80.4		E2s.1					JGI4.2		
			X5A.1		EEig10d.10					EEig10d.4		
			X5A.2		MATS4p.1					B05.C.O..5		
					B02.C.C..1					Mor10e		
					nRCHO					DISPm		
					nRCOOH.5					B05.C.C.		
					EEig10d.11					B10.C.O.		
					nOHT.3					EEig10d.5		
					JGI7.2					nThiophenes		
					nRCO.1					MEcc.2		
					JGI6					G2s.1		
					E1u					O.057		
					MATS6v.1					F04.C.N..2		
					EEig10d.12					Mor32p		
					MATS4p.2					Mor13p.1		
					MATS8m.4					GATS8v		

Table 3.1 (Continued)

Desc. Symbol	Or	Desc. Symbol	Or	Desc. Symbol	Or	Desc. Symbol	Or	Desc. Symbol	Or	Desc. Symbol	Or	Desc. Symbol
SNar	AgOr54	T11	HATS4p	AgOr67	VRv2	AgOr42	BLTF96	AgOr41	Mor29v	AgOr33	EEig06x	
nRCOOH	Correlation	EEig10d	MATS5m	Correlation	EEig12d	Correlation	DISPe	Correlation	O.057	Correlation	nRCOOH	
EEig10d	0.8997646	nRCOOH	nOHp	0.8871336	S_sOH	0.84541	PW5	0.8698281	EEig10d	0.890219	EEig10d	
Hy		BAC	EEig11x.5		UI		Infective.80		RBF		Hy	
GNar		BLTD48	nRCOOH.6		EEig09d		TPSA.NO.		DISPe		Mor32v	
O.057		JG14	EEig09d.4		nDB		HNar		H2m		Infective.80	
Infective.80		B08.C.O.	JG14.5		H2m		DISPm		Hy		O.057	
PHI		B03.C.O.	MATS8v.1		MATS5e		E2s		Dv		Xindex	
MATS3p		EEig11x	G3s		G3v		B08.C.O.		O.057.1		EEig09d	
MATS8v		R2m.	MATS2p		S_sOH.1		Mor07u		Infective.80		MATS4v	
BEHp1		EEig09d	H2m		JG17		B05.C.O.		EEig09d		B09.C.O.	
HATS3p		B05.C.C.	H2m.1		BEHm1		B09.C.O.		JG14		nRCOOH.1	
nRCOOH.1		RDF140e			E2u		Hy		JG110		Mor10u	
EEig09d		GATS8v			B09.C.O.		nRCOOH		E1e		R8m	
G2u		MATS5e			O.057		EEig10d		MATS4e		RARS	
nHBonds		HATS5m			ESpm01d		JG14		JG17		EEig10d.1	
B09.C.O.		Mor30m			EEig12d.1		L3m		O.057.2		F03.O.O.	
S_sNH2		nRCOOH.1			IC1		X5Av		BEHm1		Mor30m	
BEHp1.1		Infective.80			JG14		HATS8m		DISPm		MATS5e	
EEig12d		GATS4e			EEig09d.1		E2s.1		HATS5e		RDF150e	
E1u		B09.C.O.			ESpm15d		JG110		B09.C.O.		Infective.80.1	
EEig10d.1		R4e			E1u		MATS3e		C.026		nRCOOH.2	
E1m		Ui			HATS8m		Mor10e		SPAM		G1u	
nRCOOH.2		nRCOOH.2			G3v.1		DISPm.1		nRCOOH		EEig09d.1	
nRCHO		EEig10d.1			EEig11d		PW5.1		EEig10x		EEig09d.2	
PW5		JG14.1			nROH		EEig09d		nN			
BEHp1.2		nN			R5p.		MATS6p		PW3			
Mor32v		MATS8v			DISPm		nRCOOH.1		S_dssC			
HATS3p.1		R3u.			Infective.80		Infective.80.1		O.057.3			
Infective.80.1		MATS3p			O.057.1		GATS5e		nCrs			
O.057.1		B09.C.O..1			E2e		GATS5e.1		DISPm.1			
Jhetp		EEig11x.1			nN				DISPm.2			
EEig09d.1		Mor28u			B08.C.O.							
E1u.1		nHBonds			Jhetp							
E1u.2		JG14.2			EEig09d.2							
		EEig10d.2			O.056							
		E1u			RDF125m							
		Ts			nR05							
		EEig09d.1			C.026							
		B05.C.C..1			G3v.2							
		nArOR			nRCOOH							
		B09.C.O..2			nDB.1							
		MATS8m			HATS8m.1							
		B05.C.C..2			HATS8m.2							
		EEig11x.2										
		nRCOOH.3										
		MATS6v										
		R5e.										
		Mor10e										
		JG14.3										
		EEig09d.2										
		B02.C.C.										
		RDF150v										
		EEig11x.3										
		GATS4e.1										
		nRCOOH.4										
		C.034										
		MATS5p										
		GATS8v.1										
		EEig11x.4										
		B09.C.O..3										
		R5m.										
		JG14.4										
		EEig09d.3										
		R2m..1										
		EEig04r										
		nRCOOH.5										
		EEig10d.3										
		GATS4e.2										
		B09.C.O..4										

Table 3.1 (Continued)

Or	Desc. Symbol
AgOr45	nRCOOH
Correlation	ATS1m
0.8828472	R2m.
	Hy
	B09.C.O.
	nRCOOH.1
	Mor13m
	nRCOOH.2
	nDB
	Mor10p
	B03.C.C.
	GATS5e
	EEig10d
	nRCOOH.3
	nN
	EEig12x
	MATS4p
	nOHp
	nRCOOH.4
	GGI8
	R2m..1
	nRCOOH.5
	B02.C.C.
	RDF140e
	MATS6p
	Infective.80
	R2m..2
	F03.C.N.
	nRCOOH.6
	JGI4
	nRCOOH.7
	JGI8
	MATS3m
	C.006
	R2m..3
	S_dssC
	nRCOOH.8
	nCrs
	B03.C.C..1
	RDF150p
	Hy.1
	HATS4p
	nRCOOH.9
	EEig09d
	GATS5p
	nHBonds
	B09.C.O..1
	nOHp.1
	nRCOOH.10
	R2m..4
	R2m..5

Table 3.1 Continued

Table 3.2: Natural odor library predictions found in the top 500 predictions for each *Anopheles* Or

The prediction type (Activators, Fishing Activators, or Inhibitors), chemical structure, and predicted distance from the known training activators is listed for each odor.

Compounds marked in grey boxes represent odors from the training library.

Table 3.3: Predicted odors validated as effective activators for several Anopheles Ors

Subsets of purchasable odors from prediction lists for several Ors were validated using single unit electrophysiology. The Or names, number of activators over the number of odors tested, accuracy percentage, and primary class (aliphatic or aromatic) of responding odors is listed for each validated Or. The total accuracy, as well as accuracies for aliphatic and aromatic responsive Ors, is also provided.

Ors	Active/Tested	Accuracy	Responding Class
AgOr1	5/9	56%	Aliphatic
AgOr2	1/6	17%	Aromatic
AgOr4	12/13	92%	Aliphatic
AgOr6	2/5	40%	Aromatic
AgOr8	7/9	78%	Aliphatic
AgOr10	8/14	57%	Aromatic
AgOr12	8/8	100%	Mix
AgOr20	2/15	13%	Mix
AgOr30	16/16	100%	Aliphatic
AgOr48	7/8	88%	Aliphatic
AgOr53	5/12	42%	Mix
AgOr56	13/14	93%	Aliphatic
Total		65%	

Aliphatic Only	84%
Aromatic Only	38%

Table 3.3

Table 3.4: Activity of a large panel of aromatic odors was tested against three aromaticly tuned Ors

Responses of the previously selected 37 aromatic odors are reported for Or2, Or6, and Or10. The number of odors tested, number activating at greater than 50 spikes/sec, and percentage are listed for each Or. An average activity is provided below.

General Aromatic Assay			
Odor Receptor	Num Tested	Num >50 Spikes/Sec	Percentage
AgOr2	37	4	11%
AgOr6	37	3	8%
AgOr10	37	3	8%
Average			9%

Table 3.4

Table 3.5: Optimized descriptor sets for Or2, Or6, and Or10

Optimized descriptor symbols and final descriptor-activity correlations are listed for each *Anopheles* Or. Descriptors are listed in ascending order of when they were selected into the optimized set. The number of times a descriptor was selected in an optimized descriptor set represent weights.

Or	Descriptors	Or	Descriptors	Or	Descriptors
AgOr6	ARR	AgOr2	ATS6m	AgOr10	Dv
Correlation	R8e.	Correlation	nCconj	Correlation	ARR
0.7075761	Mor07u	0.7844095	H1v	0.7818453	Infective.80
	BEHm7		ASP		Mor18p
	GATS2v		EEig11r		JG16
	C.040		X4A		nBnz
	Mor16e		nArCO		HATS6m
	C.024		E3p		X5A
	Infective.80		O.057		B07.C.C.
	nFuranes		Gs		N.075
	ARR.1		piPC09		HOMA
	H.049		GATS3e		Infective.80.1
	Mor30m		B06.C.O.		nROH
	R8e..1		ESpm01d		HATS5m
	R3u		Mor23m		C.006
	Mor07e		nROH		Mor28u
	DISPv		Infective.80		RCI
	X4A		E2e		B07.C.C..1
	B06.C.O.		B04.C.N.		MATS3e
	C.001		B06.C.O..1		JG16.1
	TPSA.NO.		piPC05		piPC05
			EEig11d		Infective.80.2
			O.057.1		RCI.1
			H1m		nFuranes
			nRCHO		EEig14r
			R3e.		RCI.2
			nRCOOR		GATS6v
			MATS3e		N.075.1
			Gm		EEig11d
			Infective.80.1		ARR.1
			nROH.1		HATSv
			Mor28e		JG17
			JG14		PCR
			MATS7v		Infective.80.3
			nR.Cp		Ds
			PJ13		C.043
			C.006		JG16.2
			O.057.2		Mv
			piPC09.1		EEig13d
			Mor08p		HOMA.1
			nPyridines		Infective.80.4
			R6e.		X5A.1
			B06.C.O..2		G3e
			ESpm01d.1		PCD
			EEig11d.1		B07.C.C..2
			O.057.3		Mor32m
			R3e..1		nFuranes.1
			G3s		RCI.3
			ESpm05u		nPyridines

Table 3.5

Table 3.6: Top 75 predicted compounds for each *Drosophila* Or

The chemical structure and predicted distance from the known training activators is listed for each odor. All distances represent the minimum distance based on optimized descriptors to the previously known strongest active compound listed in the gray cells for that particular Or.

Or2		Or6		Or10	
Structure	Distance	Structure	Distance	Structure	Distance
C(=O)(C)c1ccccc1		0 C(=O)(C)c1ccccc1	0	0 Cc1c(ccccc1)O	0
Cc1c(ccccc1)O		0 C(=O)(C)c1sccc1	0	0 Cc1cc(ccc1)O	0
Cc1cc(ccc1)O		0 C(=O)(C)c1sccn1	0	0 Cc1ccc(cc1)O	0
Cc1ccc(cc1)O		0 CCC(C1=CC=CC=C1)=O	0	0 C(C)c1c(cccc1)O	0
C(C)c1c(cccc1)O		0 CC1=CC=C(C(C)=O)O1	0	0 CC1=CNC2=C1C=CC=C2	0
c12c([nH]cc1)ccc2		0 CCC(C1=CC=CC=C1)=O	0	0 COC1=CC=C(C)C=C1	0
CC1=CNC2=C1C=CC=C2		0 CC(=O)C1=CC=CO1	2.814136	0 COC1=CC=CC=C1C	0
N1C=CC=C1		0 CC(=O)C1=CC=CN1	3.024555	0 COC1=CC=C(C)C=C1C	0
CC1=CC(O)=CC=C1		0 CC(=O)C1=CC=CC=C1N	3.589921	0 COC1=CC=CC=C1	0
SC1=CC=CC=C1		0 CC1=CC(C(O1)C)C(=O)C	4.542537	Cn1c2c(cc1)cccc2	0.9605277
CCC1=CC(=CC=C1)O	0.51988	CC(=O)C1=CN=CC=C1	4.603079	CCC1=CC=CC=C1S	1.08787
c1(ccccc1)C(C)N	0.7157186	C1=CC=C(C=C1)C(=O)C#N	4.853174	CC1=CC=CC=C1S	1.182823
Cn1c2c(cc1)cccc2	1.037082	CC(=O)C1=CC=CC=C1O	5.11137	CC1=C(C(=CC=C1)O)C	1.216329
CCC1=CC=CC=C1S	1.090404	CC(=O)C1=NC=CN=C1	5.311086	CC1=CC(=C(C=C1)C)O	1.239482
C1=CC=C(C(=C1)O)Cl	1.202171	CC1=NC=CN=C1C(=O)C	5.32461	CC1=CC=C(C(=C1)O)C	1.239482
CC1=CC=CC=C1OC	1.261832	C1C(O1)C2=CC=CC=C2	5.368613	CCC1=CC(=CC=C1)O	1.324363
CC(=O)C1=CC=CC=C1O	1.263328	CC1=CC=C(O1)C=O	5.706266	c1(ccccc1)C(C)N	1.382239
CC1=CC(=CC=C1)O)O	1.27613	O=C(C1=CC=CO1)OC	5.861467	CC1=C(C(=C(C=C1)C)O)C	1.445959
CC1=C(C(=CC=C1)C)O	1.279905	CC1=C2C(=O)CCN2C=C1	5.918714	CC1=C(C(=CC=C1)O)O	1.446803
CC1=C(C(=CC=C1)O)C	1.28218	C1=CC=C(C=C1)CS	5.951914	COc1c(ccccc1)N	1.465265
NC1=CC=CC=C1S	1.305137	CCOC(O)C1=CC=CO1	5.958834	CNc1ccccc1	1.483803
NC1C=CC=C1	1.318499	CC1=CC=C(C=C1)OC	6.049608	CC1=C(C(=C(C=O)C)C)C	1.627195
CC1=C(C=C(C=C1)O)C	1.346765	COc1=CC=CC=C1	6.09279	CCC1=CC(=CC(=C1)C)O	1.716798
CNc1ccccc1	1.360506	CC1=CC=CC=C1OC	6.180274	C1=CC=C(C(=C1)N)O	1.807295
CC(=O)C1=CC=CC=C1N	1.372227	COc(O)C1=CN=CC=C1	6.185672	CC(=O)C1=CC=CC=C1N	1.828949
CC1=C(C(=CC=C1)O)O	1.407024	C1=COC(=C1)CS	6.205407	CCC1=C(C(=CC(=C1)O)O	1.857206
CC1=CC(=C(C=C1)C)O	1.426872	CC1=CC=C(C=C1)C(=O)C	6.300831	CC1=CC=CC(=C1)C=O	1.860195
CC1=CC(=C(C=C1)O)C	1.426872	O=CC1=CC=CC=C1NC	6.304842	CC(=O)C1=CC=CC=C1O	1.916893
CC1=CC(=CC=C1)O)C	1.431107	CC1=CC=CC(=C1)C=O	6.348972	CC(C)C1=CC=CC=C1O	1.917121
COc1=CC=CC=C1	1.441395	CC1=C(S(C(=N1)C)C(=O)C	6.383765	C1C(O1)C2=CC=CC=C2	2.006981
COc1c(ccccc1)N	1.526856	n1cc(ccc1)C=O	6.396123	CC1=CC(=CC(=C1)O)O	2.040846
C1=CC(=CC=C1)O)O	1.547724	CNc1ccccc1	6.452874	CC1=CC(=C(C=C1)O)OC	2.170293
COc1=CC=CC=C1O	1.549906	Nc1c(ccccc1)C=O	6.516774	CC1=C(C=C(C=C1)O)C	2.424889
C1=CC=C(C(=C1)CO)O	1.562781	C1=CC=C(C=C1)N	6.516774	CC1=CC=CC=C1C=O	2.456347
C1=CC(=CC=C1)O	1.612072	CC1=CN=C(C(=N1)C)C	6.529618	Nc1c(ccccc1)C=O	2.472857
CSC1=CC=CC=C1	1.635403	CCOC(O)C1=CC=CC=C1	6.54385	C1=CC=C(C=C1)N	2.478257
c1(c(ccccc1)O)O	1.638952	COc1c(ccccc1)C=O	6.641815	CC1=CC(O)=C(O)C=C1O	2.545385
CC1=CC(=C(C=C1)Cl)O)Cl	1.665124	COc(O)C1=CC=CN1	6.65084	CC1=C(C(=CC=C1)C)O	2.54936
CC(C)C1=CC=CC=C1O	1.747473	CCC1=CC=CO1	6.67399	O=CC1=CC=CC=C1NC	2.561711
C1=CSC(=C1)S	1.760378	C(C)c1c(ccccc1)C=O	6.678888	COc1=CC=CC=C1O	2.582564
C1=CC=C(C(=C1)N)O	1.773475	CC1=NC=CN=C1	6.69373	NC1=CC=CC=C1S	2.663355
c1(ccccc1)C(C)C)O	1.77454	OC(C=O)c1ccccc1	6.706399	CC1=C(C(=CC=C1)C)S	2.664136
CC(=O)C1=CC=CN1	1.784645	O=C1(C(C)=O)OCCC1	6.764351	c1(ccccc1)CO	2.713267
CC1=C(C(=C(C=C1)C)S	1.865646	CCC1=CC=CC=C1S	6.799576	C1=CC=C(C(=C1)CO)O	2.729375
C1=CC(=C(C=C1)Cl)Cl	1.879134	Cc1c(ccc1)C)O	6.824733	c1(ccccc1)C(C)C)O	2.737525
CC1=CC=NC2=CC=CC=C12	1.881836	C1=CC=C(C=C1)C(C#N)O	6.85846	CSC1=CC=CC=C1	2.753418
C1=CC=C2C(C1)C=CC=N2	1.925142	CC1C(=CCC1)C(=O)C	6.886481	CC1=CC(=CC(=C1)O)C	2.79423
C1=CC(=C(C=C1)Cl)O	1.927863	CSC1=CC=CC=C1	6.892301	C1=CC=C(C=C1)S	2.800785
CC1=CC(=C(C=C1)O)OC	1.96065	COc1=CC=C(C=C1)C=O	6.905708	CC1=CNC(=C1)C(=O)OC	2.835645
C1=CC=C(C(=C1)O)Br	2.024861	CC1=CCCC(=O)C1	6.941003	COc1=CC(=CC=C1)OC	2.932035
[C-]#[N+]CC1=CC=CC=C1	2.06791	CCC1C=CC(=O)O1	6.988239	COc1=CC=CC=C1OC	2.942067
c1(ccccc1)CO	2.099141	CC(=O)C(=O)C1=CC=CC=C1	6.995212	C(C)c1c(cccc1)C=O	3.022379
CCC1=C(C(=CC(=C1)O)O	2.147843	c1(ccccc1)COc	7.014262	CC(=O)C1=CC=CN1	3.047612
C1=CC(=CC(=C1)O)C=O	2.169262	c1(ccccc1)CC(=O)C	7.018756	Cc1c(ccc1)C)O	3.071884
C1=CC=C(C=C1)C(C#N)O	2.171679	CC1=CC(=O)CCC1	7.043093	CC1=C(C(=CC=C1)C)O	3.142583
s1ccccc1	2.200401	CC1=NC=CN=C1OC	7.0771	COc(O)C1=CC=CN1	3.199062
C1=CC(=CC=C1C=O)O	2.20377	CC1=CNC(=C1)C(=O)OC	7.116798	CC1=C(S(C=N1)C)C=C	3.232885
O=C1=CC=C(C=C1)C(=O)C#N	2.228871	Oc1cc(cc1)C=O	7.117114	c1(c(ccccc1)O)O	3.240345
COc1=CC=CC=C1OC	2.27902	c1(ccccc1)CO	7.1293	C1=CC(=CC=C1)O	3.263079
O=CC1=CC=C(O)C=C1C	2.285233	OCC1cc(cc1)C=O	7.197465	C1=CC=C(C(=C1)C=O)O	3.287156
C1=CC=C(C(=C1)C)O)O	2.308407	o1c(ccc1)C=O	7.219843	Oc1c(ccc1)C)C=O	3.296966
c1(ccccc1)C#N	2.315736	CCC1CC(=O)O1	7.222036	CC1=NC=NC2=CC=CC=C12	3.301889
Oc1c(c(ccc1)C)C=O	2.348356	CC(C)C=CC(=O)C	7.232699	C1=CC=C(C(=C1)O)Cl	3.336167
O=CC1=CC=CC=C1NC	2.357368	CC1CCCC(=O)C1=O	7.235711	C1=CC(=CC(=C1)O)C=O	3.365463
C1=CC=C(C=C1)CS	2.373006	CC(=O)C1=NCCS1	7.239078	O=CC1=CC=C(O)C=C1C	3.3787
CC1=NC=NC2=CC=CC=C12	2.393298	c1(ccccc1)CC#N	7.29472	C1=CC(=CC=C1)CO)O	3.382749
[nH]1ccccc1	2.398815	CC1=CC(=C(S1)C)C(=O)C	7.298992	C1=CC(=CC(=C1)O)O	3.391496
CCC1=CC(=CC(=C1)C)O	2.429883	Oc1ccc(cc1)C(=O)C	7.307233	C1=CC(=CC=C1)O)O	3.392312
c1(ccccc1)C(=O)O	2.440054	CC1C(CCC1)C(=O)C	7.312247	C1=CC=C(C=C1)CS	3.409539
c1(ccccc1)CC#N	2.450592	CC1=CC=CC=C1C(=O)O	7.32185	CC1=CC=C(S1)C=O	3.438549
C1C(O1)C2=CC=CC=C2	2.481837	COc1c(ccccc1)N	7.336948	n1c[nH]cc1CCN	3.481199
CC1=CC=CC=C1C(=O)O	2.485752	CCC1=NC(=CS1)C	7.352985	C1CC2=CC=CC=C2NC1	3.50207
Nc1c(ccccc1)C=O	2.490524	NC1=CC=CC=C1S	7.35831	CC1=C2C(=CC=C1)N=CC=N2	3.505628
C1=CC=C(C=C1)N	2.490524	CCNC1=CN=C1	7.365535	CC1=CC=NC2=CC=CC=C12	3.53354
CC1=CC=CC=C1C=O	2.533502	C1=C(NC=N1)CCN	7.365535	CC1=C(C=CO1)SC	3.541116

Table 3.6

Table 3.7: Dividing training odors into three distinct sets based upon odor structure and receptor response

Odors were divided into three different screening categories (Aromatic Activator Screen, Broad Activator Screen, and Inhibitor Screen). For each odor the name, final activity value, and inclusion in each screening set are listed. Odor activities were normalized to range from -100 (maximum observed inhibition) to 100 (maximum observed activation).

Odor Name	Final Activity	Aromatic Activator Screen	Broad Activator Screen	Inhibition Screen
butanal	-85	No	Yes	Yes
pentanal	-51	No	Yes	Yes
hexanal	-32	No	Yes	Yes
heptanal	-21	Yes	Yes	Yes
octanal	-20	Yes	Yes	Yes
butanol	-25	Yes	Yes	Yes
pentanol	-41	No	Yes	Yes
hexanol	-70	No	Yes	Yes
heptanol	-38	No	Yes	Yes
octanol	-35	No	Yes	Yes
butanone	-25	Yes	Yes	Yes
pentanone	-28	Yes	Yes	Yes
hexanone	-19	Yes	Yes	Yes
heptanone	-12	Yes	Yes	Yes
octanone	-18	Yes	Yes	Yes
butyl acetate	-28	Yes	Yes	Yes
pentyl acetate	-15	Yes	Yes	Yes
hexyl acetate	-12	Yes	Yes	Yes
heptyl acetate	-9	Yes	Yes	Yes
octyl acetate	-15	Yes	Yes	Yes
butyric acid	-94	No	Yes	Yes
pentanoic acid	-26	Yes	Yes	Yes
hexanoic acid	-16	Yes	Yes	Yes
heptanoic acid	-14	Yes	Yes	Yes
octanoic acid	-21	Yes	Yes	Yes
pentane	-31	No	Yes	Yes
hexane	-25	Yes	Yes	Yes
heptane	-29	Yes	Yes	Yes
octane	-34	No	Yes	Yes
2,3-butanedione	-99	No	Yes	Yes
1-octen-3-ol	-27	Yes	Yes	Yes
Ethanol	-16	Yes	Yes	Yes
3-octanol	-14	Yes	Yes	Yes
Methanol	-14	Yes	Yes	Yes
Nonanol	-12	Yes	Yes	Yes
Eugenol Methyl Ether	-9	Yes	Yes	Yes
Acetic Acid	-7	Yes	Yes	Yes
g-valerolactone	-5	Yes	Yes	Yes
Fenchone	-2	Yes	Yes	Yes
Isoamyl Acetate	-2	Yes	Yes	Yes
Limonene	-2	Yes	Yes	Yes
Menthol	-2	Yes	Yes	Yes
E2-hexenal	0	Yes	Yes	Yes
Geranyl Acetate	0	Yes	Yes	Yes
Methional	0	Yes	Yes	Yes
Eugenol	1	Yes	Yes	No
4-methylphenol	3	Yes	Yes	No
Isopropyl Alcohol	3	Yes	Yes	No
Carvone	4	Yes	Yes	No
Phenylethanone	5	Yes	Yes	No
Anisole	6	Yes	Yes	No
Benzaldehyde	6	Yes	Yes	No
Benzophenone	8	Yes	Yes	No
Citronellal	8	Yes	Yes	No
Geraniol	8	Yes	Yes	No
Ethyl Acetate	8	Yes	Yes	No
Methylsalicylate	13	No	Yes	No
Thymol	15	No	Yes	No
Cyclohexanone	48	No	Yes	No
Indole	21	Yes	Yes	No
2-methylphenol	24	No	Yes	No
methyl pyruvate	-100	No	Yes	Yes
propionyl bromide	-88	No	Yes	Yes
propionyl chloride	-73	No	Yes	Yes
propionaldehyde	-68	No	Yes	Yes
2,3-pentanedione	-55	No	Yes	Yes
2-heptanol	-39	No	Yes	Yes
2-(propylamino)-ethanol	-39	No	Yes	Yes
butyl chloride	-39	No	Yes	Yes
propionic acid	-32	No	Yes	Yes
2-methyl-3-heptanone	-26	Yes	Yes	Yes
3-heptanol	-16	Yes	Yes	Yes
4-(methylthio)-1-butanol	-15	Yes	Yes	Yes
4-hydroxy-2-butanone	-11	Yes	Yes	Yes
2,5-dimethylthiophene	-9	Yes	Yes	Yes
6-methyl-5-hepten-2-ol	0	Yes	Yes	Yes
1,5-pentanediol	0	Yes	Yes	Yes
1-hepten-3-ol	0	Yes	Yes	Yes
3-decanone	1	Yes	Yes	No
pyruvic acid	2	Yes	Yes	No
3-nonanone	2	Yes	Yes	No
4-heptanone	2	Yes	Yes	No
2-hexanol	2	Yes	Yes	No
1-bromohexane	3	Yes	Yes	No
1-hexanethiol	3	Yes	Yes	No
hexylsilane	3	Yes	Yes	No
phenylacetaldehyde	3	Yes	Yes	No
1-iodohexane	3	Yes	Yes	No
2,4,5-trimethylthiazole	5	Yes	Yes	No
ethyl valerate	5	Yes	Yes	No
cis-2-hexene	5	Yes	Yes	No
3-methyl-2-pentene	5	Yes	Yes	No
methoxyacetone	6	Yes	Yes	No
1-chlorohexane	8	Yes	Yes	No
cis-3-hexen-1-ol	10	Yes	Yes	No
fluoroacetone	10	Yes	Yes	No
acetophenone	15	No	Yes	No
2-acetylthiophene	31	No	Yes	No
pyridine	99	Yes	Yes	No
thiazole	100	Yes	Yes	No
2-ethyl-3,5(6)-dimethylpyrazine	8	Yes	Yes	No
2,5-dimethylpyrazine	26	Yes	Yes	No
pyrazine	-8	Yes	Yes	Yes
naphthalene	-14	Yes	Yes	Yes

Table 3.7

Table 3.8: Optimized descriptors selected for the Aromatic Activator Screen

Optimized descriptor symbols, brief descriptions, classes, dimensionality, and occurrences are listed for the aromatic activator screen. Descriptors are listed in ascending order of when they were selected into the optimized set. Weights indicate the number of times a descriptor was selected in an optimized descriptor set.

symbol	brief description	class	dimensionality	occurrence
N1.075	R--N--R / R--N--X	atom-centred fragments	2	1
R3v.	R maximal autocorrelation of lag 3 / weighted by atomic van der Waals volumes	GETAWAY descriptors	3	1
H.049	H attached to C3(sp3)/C2(sp2)/C3(sp2)/C3(sp)	atom-centred fragments	2	1
nRCHO	number of aldehydes (aliphatic)	functional group counts	1	1
nN	number of Nitrogen atoms	constitutional descriptors	1	1
ISH	standardized information content on the leverage equality	GETAWAY descriptors	3	1
EEig07d	Eigenvalue 07 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2	1
piPC04	molecular multiple path count of order 04	walk and path counts	2	1
MATS4e	Moran autocorrelation - lag 4 / weighted by atomic Sanderson electronegativities	2D autocorrelations	2	1
ESpm14d	Spectral moment 14 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2	1
Mor12m	3D-MoRSE - signal 12 / weighted by atomic masses	3D-MoRSE descriptors	3	1

Table 3.8

Table 3.9: Optimized descriptors selected for the Broad Activator

Optimized descriptor symbols, brief descriptions, classes, dimensionality, and occurrences are listed for the broad activator screen. Descriptors are listed in ascending order of when they were selected into the optimized set. Weights indicate the number of times a descriptor was selected in an optimized descriptor set.

symbol	brief description	class	dimensionality	occurrence
HNar	Narumi harmonic topological index	topological descriptors	2	1
R3v+	R maximal autocorrelation of lag 3 / weighted by atomic van der Waals volumes	GETAWAY descriptors	3	4
HATS3m	leverage-weighted autocorrelation of lag 3 / weighted by atomic masses	GETAWAY descriptors	3	1
Mor13p	3D-MoRSE - signal 13 / weighted by atomic polarizabilities	3D-MoRSE descriptors	3	1
ISH	standardized information content on the leverage equality	GETAWAY descriptors	3	2
P1s	1st component shape directional WHIM index / weighted by atomic electrotopological states	WHIM descriptors	3	1
R4e+	R maximal autocorrelation of lag 4 / weighted by atomic Sanderson electronegativities	GETAWAY descriptors	3	1
nRCHO	number of aldehydes (aliphatic)	functional group counts	1	2
JGI2	mean topological charge index of order2	topological charge indices	2	2
E1u	1st component accessibility directional WHIM index / unweighted	WHIM descriptors	3	2
MATS5m	Moran autocorrelation - lag 5 / weighted by atomic masses	2D autocorrelations	2	1
STN	spanning tree number (log)	topological descriptors	2	2
DISPe	d COMMA2 value / weighted by atomic Sanderson electronegativities	geometrical descriptors	3	1
B06.C.O.	presence/absence of C - O at topological distance 06	2D binary fingerprints	2	1
X4A	average connectivity index chi-4	connectivity indices	2	4
JGI3	mean topological charge index of order3	topological charge indices	2	1
De	D total accessibility index / weighted by atomic Sanderson electronegativities	WHIM descriptors	3	2
Mor25u	3D-MoRSE - signal 25 / unweighted	3D-MoRSE descriptors	3	1
nRCOX	number of acyl halogenides (aliphatic)	functional group counts	1	1
B03.O.O.	presence/absence of O - O at topological distance 03	2D binary fingerprints	2	1
nHDon	number of donor atoms for H-bonds (N and O)	functional group counts	1	1
MATS3e	Moran autocorrelation - lag 3 / weighted by atomic Sanderson electronegativities	2D autocorrelations	2	1
RBF	rotatable bond fraction	constitutional descriptors	1	1
GATS5m	Geary autocorrelation - lag 5 / weighted by atomic masses	2D autocorrelations	2	1
C.008	CHR2X	atom-centred fragments	2	1
Mor13v	3D-MoRSE - signal 13 / weighted by atomic van der Waals volumes	3D-MoRSE descriptors	3	1
R6u.	R maximal autocorrelation of lag 6 / unweighted	GETAWAY descriptors	3	1

Table 3.9

Table 3.10: Optimized descriptors selected for the Inhibitor Screen

Optimized descriptor symbols, brief descriptions, classes, dimensionality, and occurrences are listed for the inhibitor screen. Descriptors are listed in ascending order of when they were selected into the optimized set. Weights indicate the number of times a descriptor was selected in an optimized descriptor set.

symbol	brief description	class	dimensionality	occurrence
ATS1p	Broto-Moreau autocorrelation of lag 1 (log function) weighted by polarizability	2D autocorrelations	2	1
Gu	total symmetry index / unweighted	WHIM descriptors	3	6
PW5	path/walk 5 - Randic shape index	Topological indices	2	1
H-052	H attached to C0(sp3) with 1X attached to next C	Atom-centred fragments	2	4
H4m	H autocorrelation of lag 4 / weighted by mass	GETAWAY descriptors	3	3
Rtu+	R total index / unweighted	GETAWAY descriptors	3	3
HATS6m	leverage-weighted autocorrelation of lag 6 / weighted by mass	GETAWAY descriptors	3	3
B03[C-O]	Presence/absence of C - O at topological distance 3	2D Atom Pairs	2	1
nR=Cs	number of aliphatic secondary C(sp2)	Functional group counts	1	2
EEig09r	Eigenvalue 09 from edge adj. matrix weighted by resonance integrals	edge adjacency indices	2	1
Mor16m	3D-MoRSE - signal 16 / weighted by atomic masses	3D-MoRSE descriptors	3	3
X5A	average connectivity index chi-5	connectivity indices	2	1
EEig02t	Eigenvalue 02 from edge adj. matrix weighted by resonance integrals	edge adjacency indices	2	2
RDF055m	Radial Distribution Function - 5.5 / weighted by atomic masses	RDF descriptors	3	1
EEig04d	Eigenvalue 04 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2	2
B06[C-C]	presence/absence of C - C at topological distance 06	2D binary fingerprints	2	1
JG14	mean topological charge index of order4	topological charge indices	2	1
RDF085m	Radial Distribution Function - 8.5 / weighted by atomic masses	RDF descriptors	3	1
Mor08u	3D-MoRSE - signal 08 / unweighted	3D-MoRSE descriptors	3	1
MATS5e	Moran autocorrelation - lag 5 / weighted by atomic Sanderson electronegativities	2D autocorrelations	2	2
B02[C-C]	presence/absence of C - C at topological distance 02	2D binary fingerprints	2	1
nCrS	number of ring secondary C(sp3)	functional group counts	1	1
X4AV	average valence connectivity index chi-4	connectivity indices	1	1
R7e+	R maximal autocorrelation of lag 7 / weighted by atomic Sanderson electronegativities	GETAWAY descriptors	2	1
EEig08r	Eigenvalue 08 from edge adj. matrix weighted by resonance integrals	edge adjacency indices	3	1
E3s	3rd component accessibility directional WHIM index / weighted by atomic electrotopological states	WHIM descriptors	3	1

Table 3.10

Table 3.11: Top predicted natural library compounds for the aromatic activator, broad activator, and inhibitor screens

The chemical structure and predicted distance from the known training activators is listed for each predicted odor. All distances represent the minimum optimized descriptor distance to a training set compound, which are listed in gray.

Broad Activator		Aromatic Activator		Inhibitors	
Structure	Distance	Structure	Distance	Structure	Distance
O=C1CCCCC1	0	C1=NC=CC=C1	0	CCCC(O)=O	0
C1=NC=CC=C1	0	S1C=NC=C1	0	CC(C)(C)=O	0
S1C=NC=C1	0	C1=NC=CN1	2.297704	O=C(C)C(OC)=O	0
O=C1CCCCC1	2.037369	N1N=CC=C1	2.768945	O=C(C)(C)=O	1.193469
O=C1NCCCC1	2.16724	CC1=CSC=N1	3.028959	NCCC(O)=O	2.028579
C1=CC=C1	2.997682	CCCCC1=NC=CC=C1	3.534064	CCC(C)(O)=O	2.579103
O=C1C=CC=C1	3.977348	CN1C=CC=C1	3.721854	CCC(N)(C)=O	2.752901
CC1CCNCC1	4.39789	CCC1=NC(=CN=C1)C	3.758666	CC(N)(C)=O	3.088135
C1=NC=CN1	4.619738	C=C/C=C/C/C/C/C/C/C	3.778543	CC(C)(OC)=O	3.128039
O=C1C=CC=CO1	4.919501	CC1=CN=CC(O)=N1	3.815402	CC(OCC)=O	3.133545
N1N=CC=C1	4.943772	C1=CC=CS1	3.817309	CC(C)/C=N/O	3.16974
CN1CCCC1	5.167015	CC1=NC=CN=C1CC	3.82998	CC(C)(C)(O)=O	3.179965
CC(O)=CC1CC(C)=CCO1	5.609835	CCC1=NC=CN=C1CC	3.853768	CC(C)C(O)=O	3.318447
CC(O)=C[C@H]1C[C@@H](C)CCO1	5.751913	C12=CC=CC=C1N=CC=N2	3.889099	C/C(C)=O=C/C	3.442493
O=C1CCCC1	5.94834	CCC1=NC=CN=C1	3.914359	C/C=C(C)(C)[O-]=O	3.452961
CC(O)=CC1C=C(C)CCO1	5.971966	CC1=NC=CN=C1OC	4.038081	C=C(C)CCO	3.736961
O=C1CC(C)(O)CCO1	5.996197	COCC1=NC=CN=C1CC	4.03881	O=C(O)COC	3.778532
CC(O)=CC1OCCC(C1)C	6.075926	C=C=C1=CN=CC(C)=N1	4.048223	CC(C)(OC)=O	3.80734
OC1CCCCC1	6.143672	CC1=COC=C1	4.074861	CC(C)=CCO	3.874121
S1SCCC1	6.220288	CC1=NC=C(C)C=C1	4.124326	CC(COC)=O	3.874164
C1C=CCS1	6.313776	C/C=C/C=C/C/C/C/C/C	4.131876	C/C(C)=O=C/C	4.012943
OC1(C)=CCCC1	6.407874	CC1=CC2=C(C=C1)N=CC=C2	4.139062	CCC(C)=O	4.027003
OC1(C)CCCC1	6.471678	CC(C)=C/C=C/C=C/C	4.187466	C=C(C)C=O	4.210918
O=C1C=C(C)CCCCC1	6.523937	CC(C)=C/C=C/C=C/C	4.235145	C=CC(C)=O	4.220098
CC(C)C(C)=CC(C)=CC(C)C	6.54101	CC(C)=C/C=C/C=C/C	4.260037	CC(C)=O	4.256449
C1=CC=CC=C1	6.558737	CC(C)=C/C=C/C=C/C	4.26035	CC(C)C#N	4.2874
CC1=NC=C(C)C=C1	6.577361	CC(C)=C/C=C/C=C/C	4.262767	O=C(O)(N)CO	4.307723
N1C=CC=C1	6.588811	CC1=NC=CN=C1CC(C)C	4.265956	CC(C)C	4.368743
CN1C=CC=C1	6.60897	COCC1=NC=CN=C1C(C)C	4.290186	NC(OCC)=O	4.378132
O=C(NCCC(C)C)C1=CC=CC=C1	6.612483	CC(C)=C/C=C/C=C/C	4.295488	CC(OC)OC	4.394431
CC1SCCN1	6.61741	CC(C)=C/C=C/C=C/C	4.295748	C/C=C(C)CO	4.410392
O=C(O)(C)C1=CC=CC=C1	6.638106	CC(C)=C/C=C/C=C/C	4.312156	C=C(C)C(C)=C	4.446982
CCCC(C1=CC=CC=C1)=O	6.651088	CC(C)=C/C=C/C=C/C	4.317797	CC(C)(C)=O	4.523044
O=C(N1CCCCC1)C2=CCCCC2	6.675009	CC(C)=C/C=C/C=C/C	4.328344	CC(C)C#N	4.540592
CC(O)CC1=NC=CS1	6.701117	CC(C)=C/C=C/C=C/C	4.347131	O=C(O)(C)C	4.545043
CCCCCCC1CCCC(O1)=O	6.717317	O=C(C)C=CC=C1	4.35801	CC(CCC)=O	4.576413
CCCCCCC1CCCC(O1)=O	6.743187	CC(C)=C/C=C/C=C/C	4.371786	CC(C)C(C)=O	4.608949
CC1CCCC(C)N1	6.751586	CC(C)=C/C=C/C=C/C	4.375425	O=C(C)C(O)=O	4.696431
O=CC1=CC=C1	6.853976	CC(C)=C/C=C/C=C/C	4.38817	N[C@@H](C(C)C)C(O)=O	4.786178
O=C(OCC=C(C)C)C1=CC=CC=C1	6.86057	CC(C)=C/C=C/C=C/C	4.401625	CC(C)N+([O-])=O	4.809623
CC1=CC=C1	6.867993	CC(C)=C/C=C/C=C/C	4.413842	CC(C)O=C	4.818327
O=C(O)CC(OCC=C(C)C)=O	6.899062	O=C(C)C=C/C1=CC=CC=C1	4.423901	CC(C)C(O)=O	4.835928
OC1=C=CC(O)=C1	6.902425	CC1=NC=C(C)C(N)=C1C	4.441315	CC(C)=CC(O)=O	4.870955
O=C(N(C)CC)CC1=CC=CC=C1	6.905865	N1=CCCC1	4.45126	CC(OC)(C)=O	4.88651
CC(C)N(C)CCCC1=O	6.908113	C/C=C/C=C/C=C/C	4.455448	CC(C)(C)C(O)=O	4.932188
C/C(O)C(OCC=C(C)C)=O	6.915907	C1(C)C=CC=CC=CC=CC	4.455697	CC(C)C(O)=O	4.94746
CC(O)C(C)C=CC=C1=O	6.94	O=CC1=CC=C(C)C=C1	4.48037	CC(O)COC	4.960194
O=C1C/C=C/C/C=C/C/C=C1	6.96282	CC(C)=C/C=C/C=C/C	4.484765	CC(O)C(C)=O	4.969207
CC(C)C=C(C)C=CC(O1)=O	6.971452	CC(C)=C/C=C/C=C/C	4.4849	O=CN(C)C	4.971724
CCCCC1=NC=CC=C1	6.99577	CC(C)=C/C=C/C=C/C	4.495566	CC(O)C=O	4.979604
CCCCC1=C=C(C)C=C1	7.00054	CC1=CN=C(C)C(C)C(N)=N1	4.495643	O=C(C)C=C/C	5.003101
CCCCC1=C=C(C)C=C1	7.008492	CC1CCCC1	4.507732	CC(C)=CC#N	5.003739
CCCCC1=C=C(C)C=C1	7.01097	CC1=C(C)N=CS1	4.524113	CC(C)C(O)=O	5.00945
C/C=C/C(O)C=C=C(C)=O	7.011722	CCCC(C1=NC(C)=CN=C1)C	4.534104	CC(C)C(O)=O	5.013806
OCCCCC1=CC=CC=C1	7.012331	C1(C2NCCC2)=CC=CN=C1	4.551649	CC(C)(O)C(O)=O	5.020397
CCCC(O)C=C=CC=C1=O	7.016306	C1(C)C@H2NCCC2)=CC=CN=C1	4.551649	C[C@@H](O)CC(O)=O	5.047488
CC(O)C(C)C=CC=C1=O	7.018781	CCCCC1=CC=CC=C1	4.559646	CC(C)C(O)=O	5.057977
CC(O)C(C)C=CC=C1=O	7.0234	CCCCC1=CC=CC=C1	4.559792	CC(C)C=O	5.059612
CC(O)C(C)C=CC=C1=O	7.030379	CC(C)=C/C=C/C=C/C	4.564464	CC(C)C=O	5.151031
OC1=CC=CC=C1	7.03727	CC(C)=C/C=C/C=C/C	4.57072	CC(OCC)=S	5.164894
OC1=CC=C(C)C=C1	7.045313	CC(C)=C/C=C/C=C/C	4.573949	SC(C)C(O)=O	5.213845
CCCC1OCC(O1)C	7.057154	CC(C)=C/C=C/C=C/C			
CCCCCCC1NCCC(C)CC1=O	7.066965	CC(C)=C/C=C/C=C/C			

Table 3.11

Table 3.12: Predicted odors validated as activators and inhibitors of the CO₂ receptor

The overall accuracy for CO₂ predictions considering all three predicted sets. The accuracy and percentage accuracy are listed for each class.

Activity	Accuracy	Percentage
Activators (>50 spikes/sec)	25/139	18%
Activators (>30 spikes/sec)	42/139	30%
Inhibitors (<-5 spikes/sec)	25/139	18%
Total		48%

Table 3.12

Table 3.13: Optimized descriptor sets for each Citrus Psyllid ORN

Optimized descriptor symbols and final descriptor-activity correlations are listed for each Citrus Psyllid ORN. Descriptors are listed in ascending order of when they were selected into the optimized set. The number of times a descriptor was selected in an optimized descriptor set represent weights.

Or	Desc.Symbol	Or	Desc.Symbol	Or	Desc.Symbol	Or	Desc.Symbol
RP2_B	R7m	RP2_C	ESpm10d	RP2_A	ATS5v	RP4_C	SP02
Correlation	QYYe	Correlation	L1m	Correlation	ATS4e	Correlation	ATS6m
0.8289435	Infective.80	0.746558	Mor30p	0.8315387	DISPv	0.807409	R5p.
	RDF065m		nHDon		R7p.		ATS5m
	JGI6		ESpm10d.1		H5v		SP12
	RDF080v		Mor27p		nHDon		Mor29e
	nHDon		E3m		As		R7m.
	Mor27e		H.051		B05.C.C.		R8u.
	H8v		Mor24p		nDB		RDF020e
	E1u		Mor27p.1		E2s		Mor23p
	R2v.		HATS8m		Mor26v		H.052
	Mor03m		Mor16u		Infective.80		R8p.
	R7m.1		R1u.		B05.C.C..1		Mor16p
	Mor32u		ESpm06d		L.Bw		BEHm7
	H6m		R4v.		piPC08		SP12.1
	H8m		Mor13p		G3p		R6m.
	Mor16u		B04.C.C.		RDF080e		EEig08x
	O.056		L3s		nR.Cp		HATS3p
	Mor27e.1		EEig08x		BELm1		G3u
	JGI6.1		nHDon.1		nHDon.1		nCb.
	EEig12x		Mor27e		PJI3		B07.C.O.
	QYYe.1		Mor30m		H8u		R6m..1
	R8p.		nR.Cs		Mor24m		H7u
			Mor30e		DISPv.1		B04.C.C.
			Mor10u		nCconj		EEig11x
			C.006		B05.C.C..2		MATS8m
			R5v.		L2s		RTu.
			Mor04m		nHDon.2		BEHm7.1
					H4m		R7m..1
					G3u		R8e.
					nR04		C.025
					Infective.80.1		H4p
					nRCO		CIC3
					RDF075e		Mor30u
					DISPv.2		Mor16m
					O.057		R6m..2
					H5e		B07.C.O..1
							BEHm7.2
							EEig11x.1
							BEHm7.3
							EEig08x.1
							R8e..1
							R7m..2
							R5p..1
							SP12.2
							BEHm7.4
							Mor16p.1
							nRCO
							H7u.1

Table 3.13

Or	Desc.Symbol	Or	Desc.Symbol	Or	Desc.Symbol	Or
RP6_A	nRCHO	RP6_B	ATS5v	RP6_C	SP01	RP7_B
Correlation	nCconj	Correlation	Xu	Correlation	Mor13u	Correlation
0.8951214	EEig10d	0.8758008	nHDon	0.8019783	SP11	0.8016844
	EEig03d		nDB		RDF020u	
	ESpm14d		B05.C.C.		H7u	
	H5m		Mor05e		R6m.	
	H8u		E2s		Mor16p	
	Mor27e		Infective.80		R7m.	
	nROH		PJI2		RTu.	
	F01.C.O.		piPC08		Mor22e	
	H.051		nCconj		R8e.	
	R2p		B05.C.C..1		H5m	
			L.Bw		RDF050v	
			Infective.80.1		JGI5	
			G3p		RDF020e	
			Mor15v		R8v.	
			R7p.		R6m..1	
			R6u		EEig08x	
			H8u		Mor30e	
			B05.C.C..2		HATS4u	
			DISPv		Mor16p.1	
			nHDon.1		RTu..1	
			H5e		R7m..1	
			RDF080e		CIC3	
			E1s		R6m..2	
			nR04		H7e	
			E2s.1		C.025	
			nDB.1		B04.C.C.	
			Infective.80.2		B07.C.O.	
			Mor03v		nOht	
			nOHs		nR03	
			PJI3		RDF085p	
			C.026		R8u.	
			L.Bw.1		H5m.1	
			R6v.		RDF020u.1	
			nDB.2		RTu..2	
			B05.C.C..3		Mor16p.2	
			Infective.80.3		R6m..3	
			nR10		H7u.1	
			DISPv.1		B07.C.O..1	
			B05.C.C..4		R8e..1	
			JGI7			
			Mor28v			
			E2s.2			
			nCconj.1			
			Infective.80.4			
			B05.C.C..5			
			G3u			

Table 3.13 Continued

Desc.Symbol	Or	Desc.Symbol
Mor02p	RP7_C	nRCHO
Mor27e	Correlation	EEig10x
R7p.	0.8422402	nCconj
DISPm		EEig03d
Infective.80		RDF020e
JGI5		DISPv
B07.C.C.		Mor30u
HATS4e		PJI3
Mor23p		G3s
Mor32m		GATS2p
CIC2		RDF080e
R7m.		nROH
JGI4		R3p
Mor10m		RDF020e.1
X5A		EEig10x.1
EEig10r		ESpm01d
B07.C.C..1		nRCOOR
H3m		H.051
Mor03p		GATS5e
F05.C.O.		H8u
R6m.		H.051.1
JGI5.1		Mor30e
B07.C.C..2		nROH.1
Mor27e.1		EEig09x
HATS4e.1		Jhetp
E2e		Mor30u.1
B07.C.C..3		RDF020e.2
Mor03m		
B04.C.C.		
EEig10r.1		
CIC2.1		
B07.C.C..4		
JGI4.1		

Table 3.13 Continued

Table 3.14: Top predicted natural library compounds for Citrus Psyllid ORNs

The chemical structure and predicted activity are listed for the top predicted activators for each ORN. Compounds in gray are training set odors.

SMILES	RP2_B	Pred Activity	SMILES	RP2_C	Pred Activity	SMILES	RP4_B	Pred Activity
C/C1=C(C)C(C)C(C)C=C/C(C)C=C/C1		74.35797	CCC=C/C=C/O		45.26686	CC1=CC=C(C(C)C)C=C1		163.3939
CC(O)CC(C)C=O		71.94693	CC(C)=CC(C)C=C=C		44.17876	CC(C1=CC=CC=C1)=O		139.1738
O=C1C(C)C(C)C(C)C(C)C@H]2[H]C(C)C(C)C@H]12[H]		71.84607	CCCC(O)C=O		42.36729	CC12CCCC(C1C2)=C/C		135.1755
C=C(C)C=C(C)C=C(C)C		63.22675	CCCC(O)C=O		30.8423	CC1(C)C@H]2[H]C(C)C=C/C(C)C@H]12[H]		128.784
CCCCC(C)C=O		58.49042	O=CC1=CC=CC=C1		26.24606	CCCCC(C)C=O		117.5805
CC(C)C(C)C=O		48.96719	C=CC(O)CCCC		23.14035	CC(C)C=C/C=C/O		116.1112
CC(C)C(C)C=O		36.11657	CC1=CC(C)C(C)C=C1		21.80823	CC(O)C(C)C=O		112.9869
C/C1=C(C)C(C)C=C/C(C)C=C/C1		74.35797	C/C=C/C(C)C=O		51.43602	CC(C1=NC=CS1)C		210.9587
CC(O)CC(C)C=O		71.94693	CCCC=CC=O		45.26686	CCCC1=CC=CO1		201.5245
O=C(C)C(C)C(C)C=O		71.28565	CC(C)=C/C(C)C=O		44.65782	CC(C)CC1=NC=CS1		195.0182
C/C1=C(C)C(C)C=C/C(C)C(C)C@H]1C)=C(C)C		71.24567	CC(C)C=O		44.38818	CC(C1=CC=CC=C1)C		194.1631
CC(C)C(C)C=O		68.60787	CC1=CC(C)C(C)C=C1		44.17876	CC1=C(C)C=CC=C1		193.5604
CC(C)C(C)C=O		68.13362	CCCC(O)C=O		42.65992	COCC1=CC=CC=C1		191.4757
CCCC(O)C(C)C=O		68.07507	CCCC(O)C=O		42.36729	CC1=C(C)C=CC=C1		191.0865
CC(C)C(C)C=O		66.2354	C/C=C/C(C)C=O		42.26823	COCC1=CC=CC=C1		190.5536
CC(O)CC(C)C=O		65.7931	CCC(C)C=C/O		41.84369	CCCC1=CC=CC=C1		188.6928
CC1=CC=C(O)C(C)C=O		63.96134	CC(C)C(C)C=O		40.55074	CCO1=CC=CC=C1		188.517
O=C/C=C(C)C=O		63.22675	CCC(C)C=C/O		40.52789	CC1=C(C)C=CC=C1		187.4761
CCCCC(C)C=O		63.10904	CC1CC2=C(C)C=CC=C12		40.00807	CC1=C(C)C(C)C=C1		184.7897
CC(C)C(C)C=O		62.47671	CCCCC(C)C=O		39.88164	CO(C)C1=CC=CC=C1		178.6651
CC(O)C(C)C=O		62.39316	CC(C)C(C)C=O		39.71849	CC1=C(C)C=CC=C1		178.4423
CC(C)C(C)C=O		62.34478	CC(C)C=C/O		39.4017	CC1=C(C)C=CC=C1		177.9384
C=C(C)C=C(C)C=C(C)C(C)C1		61.23283	C=CC(O)C=C		38.99514	C1(C2=CC=CC=C2)C=C1		175.7152
C=C(C)C(C)C=O		60.88536	C=C(C)C=C/C(C)C=C		38.95081	CC1=C(C)C=CC=C1		174.1878
C=C(C)C=C/C(C)C(C)C1		59.00969	CCCCC(O)C=O		38.01496	CCO1=CC=CC=C1		174.0245
CC(C)C(C)C=O		58.49042	O=C(O)C(C)C=O		37.8518	CCO1=CC=CC=C1		169.4623
CC(C)C(C)C=O		56.55607	C=C(C)C=C/C(C)C=O		37.74838	CCCCC1=CC=CO1		168.622
CCCCC(C)C=O		56.5295	C/C(C)C=O		37.48334	NC(C)C1=CC=CC=C1		167.7722
CC(C)C(C)C=O		56.44963	C=C(C)C=C/C(C)C=O		37.47729	CCCCC1=CC=CO1		166.5448
CC1=C(C)C(C)C=C(C)C@H]2[C]C@H]1C)=C(C)C		56.43485	CC(C)C(C)C=O		37.34352	CC1=C(C)C=CC=C1		163.3939
CC(C)C(C)C=O		55.76329	O=C(C)C=C/C		37.33052	COCC1=CC=CC=C1		163.2257
C=C(O)C1(C)C=CC=C(C)C		54.60038	C/C=C/C(C)C=O		37.21464	CC(C)C(C)C=O		161.5202
O=C(O)C1=CC=CC=C1		54.0418	COCC1=CC=CC=C1		37.11306	SOCC1=CC=CC=C1		161.2887
O=C/C=C(C)C=O		53.745	CCCCC(C)C=C/O		36.75797	CC1=CC=CC=C1		160.0841
COCC(C)C=C/O		53.73228	CCCCC(C)C=O		36.4936	SCC1=CC=CC=C1		159.0215
C=C(O)C(C)C=C(O)C		53.71671	C=C(C)C=O		36.43435	CCO1C2CCOC(O2)C=C1		158.5089
CC(C)C(C)C=O		53.64558	C=C(C)C=C(C)C=C		36.26165	CC1=NC=CC=C1		158.0411
CC(C)C(C)C=O		53.12195	O=C1(C)C=C(C)C=CC1		36.06916	CCCCC1=CC=CO1		156.9442
O=C(N)C(C)C=C(C)C=C1		52.23282	O=C(C)C=C/C		36.01778	CS1=CC=CC=C1		156.4988
CC(C)C(C)C=O		52.18524	CCCCC(C)C=C/O		36.00254	[C-]N(C)C1=CC=CC=C1		154.3357
CCCCC(C)C=O		52.16527	CCC(O)C=C/C=C/O		35.75418	CCCC1CC(O)C=C		154.2606
CC(C)C(C)C=O		52.08181	CC(C)C(C)C=O		35.71133	O=CC1=CC=CC=C1		153.2292
CCCCC(C)C=O		51.98421	O=C/C=C(C)C=CC		35.50662	CC1CC(O)C(C)C=C1		152.8799
CC(C)C(C)C=C(C)C=C(C)C2		51.6773	CC(C)C=C/O		35.3759	CCCC(O)C=S		152.8392
CC(C)C(C)C=C(O)C1		51.29109	CC(C)C(C)C=O		35.34795	SC1=CC=CC=C1		152.6492
O=C(C)C(C)C=C(O)C		51.17537	CC(C)C=C/O		35.05646	CCO1=CC=CC=C1		152.1006
CC(C)C(C)C=O		51.00948	C=C(C)C=O		35.03119	CC1=C(C)C=CC=C1		150.7356
CC(C)C(C)C=O		51.10806	O=CC(C)C		35.02802	CS(C)C1=CC=CC=C1		148.9938
CC(O)C(C)C=C(C)C=C(O)C		51.00987	CC(C)C=C/C=C/O		35.00791	CC1=NC=CC(C)C=C1		144.4554
C=CC(C)C(C)C=O		50.96592	CC(C)C=C/C=C/O		34.89028	CCO1=CC=CC=C1		144.0637
O=CC1=CC=C(C)C=C1		50.59055	O=C(O)C(C)C=O		34.68742	CC1OC2(C)C(C)C=C1		142.4315
CO(C)C(C)C=C(C)C		50.49973	COCC1=CC=CC=C1		34.60386	CC1(C)C2CC(C)C=C1		142.1961
CC1(C)C2C(C)C=C(C)C=C1		50.36754	CCCC		34.57437	CCCCC#N		141.7605
O=C(C)C(C)C=O		50.33509	CCCC(O)C=C		34.19629	O=CC1=CC=CC=C1		141.4747
CC(C)C(C)C=C(O)C		50.29364	O=C1=C(C)C=CC=C1		34.17189	CC1C=CC2(C)C(C)C=C1		140.7053
CC(C)C(C)C=O		50.09946	O=C(C)C=C		33.9544	CC(O)C1=CC=CC=C1		140.5338
C(C)C(C)C=C(N)O		49.90491	C=C(C)C(C)C=O		33.84005	CC1=NC=CC=C1		139.2927
CC(C)C=C1=CC=CC=C1		49.89894	CC(C)C=C/O		33.43157	COCC1=CC=CC=C1		139.2339
CCCCC(C)C=O		48.96719	CC(O)C(C)C=C(O)C		33.40156	CC(C)C=C/C=C/O		139.1738
CC(C)C(C)C=O		48.87246	CC(C)C(C)C=O		33.32586	CC(O)C1=CC=CC=C1		137.8998
CC(O)C(C)C=C(O)C		48.54704	CC(O)C		33.32402	COCC1=CC=CC=C1		137.7904
CC1CC=C2C1(C3C(C)C)CC2		48.32759	CCCC(O)C=C/C=C/O		33.27204	O=C(C)C(C)C=C1		137.15
CCCCC(C)C=O		48.08597	CC(C)C=C(O)C		33.02956	O=C(O)C1=CC=CC=C1		137.0934
O=C(C)C(C)C=C(C)C		48.03822	O=C(C)C=C		32.95811	CC12CCCC(O2)C(C)C=C1		136.7338
CC(C)C(C)C=C(C)C		47.67373	O=C(N)C(C)C=C1		32.91734	NC1=CC=CC=C1		136.328
O=C(C)C=C(C)C=C(O)C		47.46852	CCCCC(C)C=C/O		32.77198	CC(C)C1=CC=CC=C1		135.7478
O=C(C)C=C(C)C=C(O)C		47.18779	O=C(C)C=C/C=C/O		32.71905	NC1=CC=CC=C1		135.7411
O=C(C)C=C(C)C=C(O)C		47.10583	C/C=C/C(C)C=C/O		32.70471	C=C1CC2(C)C(C)C=C1		135.1755
CC(C)C=C(C)C=C(O)C		46.94897	CCCCC(C)C=C/O		32.552	CC(C)C#N		134.4416
O=C(N)C(C)C=C(C)C=C1		46.86099	CCO1=CC=CC=C1		32.41418	COCC1=CC=CC=C1		134.3196
CCCCC(C)C=O		46.73148	CCCCC(C)C=C/O		32.03659	CCCC1CC(O)C=C1		133.6848
C=C(C)C(C)C=C(O)C		46.5108	CC(O)C(C)C=C(O)C		31.94634	CC(C)C(C)C=O		132.8239
CCCCC(C)C=C(O)C		46.33825	CC(O)C(C)C=O		31.59891	O=CC1=CC=CC=C1		132.788
CC(O)C(C)C=C(O)C		46.16463	CC(O)C(C)C=C(O)C		31.5376	CC1=C(C)C=CC=C1		131.9495
CCCCC(C)C=C(O)C		46.15281	CC(C)C=C/C(C)C=C		31.51514	CC(O)C1=CC=CC=C1		131.9484
C=C(C)C(C)C=C(O)C		45.83685	CCCCC(C)C=C/O		31.44217	O=C(O)C1=CC=CC=C1		131.7989
CC(C)C(C)C=C(O)C		45.48793	CC(C)C=C(C)C=C(O)C		31.42432	COCC1=CC=CC=C1		131.7978
C=CCCCC		45.46762	COCC1=CC=CC=C1		31.38461	CC(O)C1=CC=CC=C1		131.719
SSCC1=CC=CO1		45.46385	CC(C)C=C/O		31.37606	CC(C)C(C)C=O		131.6372
CC(C)C(C)C=C(O)C		45.44599	CC=C(C)C=C(O)C		31.37038	CC(C)C1=CC=CC=C1		131.534
O=C(C)C(C)C=C(O)C		44.91413	CC1=CC=CC=C1		31.31489	CC1=C(C)C=C1		130.8618
O=C(C)C(C)C=C(O)C		44.851	CC(O)C(C)C=C(O)C		31.30609	O=CC(C)C(C)C		130.8068
CC(C)C(C)C=C(O)C		44.74891	CC(C)C=C(C)C=C(O)C		31.17691	CC1=C(C)C=CC=C1		129.885
CCCCC(C)C=C(O)C		44.65088	NC1CCCC1		31.17583	CC(C)C#N		129.1444
O=C(C)C=C(C)C=C(O)C		44.49086	CC1=CC=C(C)C=C1		31.11419	O=CC(C)C=C(C)C=C1		128.9196
CC(O)C(C)C=C(O)C		44.25152	C=C(C)C=C		31.06791	CCCC1CC(O)C=C1		128.6957
CC(C)C=C(C)C=C(O)C		44.13589	CC(O)C(O)C=O		31.04751	CCCCC(C)C(C)C		128.5922
CCCCC(C)C=C(O)C		44.07246	CC(C)C=C/C=C/O		30.89184	CCCCC(C)C=C		128.5503
C=C(C)C(C)C=C(O)C		43.77837	CCCCC(C)C=C/O		30.8423	CC(=O)C1=CC=CC=C1		128.3256
CCCCC(C)C=C(O)C		43.65871	CC(C)C(C)C=C(O)C		30.83327	CC1=CC2(C)C(C)C=C1		128.1545
CC(C)C=C(C)C=C(O)C		43.61321	C=C(C)C=C/O		30.81392	CC(O)C(C)C=C(O)C		128.0658
CC1=NC=CC=C1		43.43477	CC(C)C=C/O		30.78685	CC1(C)C(C)C=C1		127.756
CCCCC(C)C=C(O)C		43.1834	CCC(C)C(C)C=C(O)C		30.61718	CCO1=CC=CC=C1		127.5711
CC(C)C(C)C=C(O)C		43.09237	CCCCC(C)C=C(O)C		30.53002	CC(C)C(C)C=C		127.5105
CCCCC(C)C=C(O)C		43.08604	C(C)C(C)C=C(O)C		30.51805	COCC(C)C		127.4569
C=CCCCC		43.04082	CC(C)C=C(O)C=C(O)C		30.41643	COCC1=CC=CC=C1		127.3973
O=C(O)C(C)C=C(O)C		42.88664	CC(C)C=C(O)C		30.40109	CC(C)C(C)C=C(O)C		127.2072
CC(C)C=C(C)C=C(O)C		42.71533	CC(O)C=C		30.34567	CC(O)C(C)C=C		126.9944
CCCCC(C)C=C(O)C		42.6727	O=CCCC		30.31383	CC(C)C=C=CC=C1		126.5322
CC1=CC(C)C=C(C)C=C1		42.66021	CC(C)C=C(C)C=C(O)C		30.29812	CC(C)C=C=CC=C1		126.26

Table 3.14

RP4_C			RP6_A			RP6_B		
SMILES	Pred	Activity	SMILES	Pred	Activity	SMILES	Pred	Activity
CCCCCO	88.03097	88.03097	CC(C)=O	55.52031	55.52031	CC(C)C(C)C(C)C		134.0782
CC(C)=CCO	81.96087	81.96087	CC(C)C(C)C	50.76842	50.76842	CC(C)=CC(C)C(C)C		130.4358
CC(C)C(C)C	70.86368	70.86368	CC(C)=O	55.52031	55.52031	CC(C)=CC(C)C(C)C		125.9246
C=C(C)C	56.16688	56.16688	CC(C)C(C)C	55.23	115	CC(C)C(C)C		123.529
C=C(C)C	53.44199	53.44199	CC(C)C(C)C	51.5647	51.5647	CC(C)C(C)C		120.9403
O=C(C)C	41.21922	41.21922	O=C(C)C	50.94559	50.94559	CC(C)C(C)C		119.8744
CC(C)C(C)C	40.38888	40.38888	CC(C)C(C)C	50.76842	50.76842	CC(C)C(C)C		106.2103
CC(C)C(C)C	88.57896	88.57896	CC(C)C(C)C	49.23892	49.23892	CC(C)C(C)C		97.64069
CCCCCO	88.03097	88.03097	CC(C)C(C)C	46.17592	46.17592	CC(C)=CC(C)C		154.2551
CCCCCO	82.40683	82.40683	CC(C)C(C)C	44.74221	44.74221	CC(C)C(C)C		153.9955
CC(C)C(C)C	81.96087	81.96087	CC(C)C(C)C	44.28983	44.28983	CC(C)C(C)C		152.9753
CCCCCO	77.84128	77.84128	CC(C)C(C)C	44.2	1182	CC(C)C(C)C		149.1351
CC(C)C(C)C	77.38877	77.38877	CC(C)C(C)C	44.03917	44.03917	CC(C)C(C)C		148.7203
C=C(C)C	76.88151	76.88151	CC(C)C(C)C	43.72791	43.72791	CC(C)C(C)C		143.9979
CC(C)C(C)C	70.86368	70.86368	CC(C)C(C)C	42.81228	42.81228	CC(C)C(C)C		143.4794
CC(C)C(C)C	70.27421	70.27421	CC(C)C(C)C	41.48486	41.48486	CC(C)C(C)C		141.5474
CC(C)C(C)C	69.68066	69.68066	CC(C)C(C)C	41.48082	41.48082	CC(C)C(C)C		139.4834
CC(C)C(C)C	68.15898	68.15898	CC(C)C(C)C	41.43222	41.43222	CC(C)C(C)C		138.4367
CC(C)C(C)C	67.27872	67.27872	CC(C)C(C)C	41.34868	41.34868	CC(C)C(C)C		137.8212
CC(C)C(C)C	66.08893	66.08893	CC(C)C(C)C	40.79962	40.79962	CC(C)C(C)C		136.8536
CCCCCO	65.7308	65.7308	CC(C)C(C)C	40.45953	40.45953	CC(C)C(C)C		135.7359
CC(C)C(C)C	65.19794	65.19794	CC(C)C(C)C	40.2463	40.2463	CC(C)C(C)C		135.3627
CC(C)C(C)C	64.4434	64.4434	CC(C)C(C)C	40.08869	40.08869	CC(C)C(C)C		134.9339
CC(C)C(C)C	64.29792	64.29792	CC(C)C(C)C	39.8543	39.8543	CC(C)C(C)C		134.1139
CC(C)C(C)C	63.89223	63.89223	CC(C)C(C)C	39.68227	39.68227	CC(C)C(C)C		134.0782
CC(C)C(C)C	60.78221	60.78221	CC(C)C(C)C	38.46738	38.46738	CC(C)C(C)C		133.5617
C=C(C)C	59.32255	59.32255	CC(C)C(C)C	37.37	114	CC(C)C(C)C		132.9786
CC(C)C(C)C	58.79435	58.79435	CC(C)C(C)C	37.3392	37.3392	CC(C)C(C)C		132.7151
CC(C)C(C)C	58.25072	58.25072	CC(C)C(C)C	37.00024	37.00024	CC(C)C(C)C		131.8252
CC(C)C(C)C	58.21219	58.21219	CC(C)C(C)C	36.90013	36.90013	CC(C)C(C)C		131.1524
CC(C)C(C)C	58.07201	58.07201	CC(C)C(C)C	36.8585	36.8585	CC(C)C(C)C		131.9053
CC(C)C(C)C	58.02782	58.02782	CC(C)C(C)C	36.69429	36.69429	CC(C)C(C)C		130.8624
CC(C)C(C)C	57.7253	57.7253	CC(C)C(C)C	35.66337	35.66337	CC(C)C(C)C		130.6334
O=C(C)C	56.69925	56.69925	CC(C)C(C)C	35.30428	35.30428	CC(C)C(C)C		130.4358
C=C(C)C	56.53	56.53	CC(C)C(C)C	34.93723	34.93723	CC(C)C(C)C		129.5921
CC(C)C(C)C	56.16688	56.16688	CC(C)C(C)C	34.80002	34.80002	CC(C)C(C)C		129.4238
CC(C)C(C)C	55.37078	55.37078	CC(C)C(C)C	34.45637	34.45637	CC(C)C(C)C		128.56
CC(C)C(C)C	55.27365	55.27365	CC(C)C(C)C	34.2548	34.2548	CC(C)C(C)C		126.9605
O=C(C)C	55.19322	55.19322	CC(C)C(C)C	33.70733	33.70733	CC(C)C(C)C		126.9013
O=C(C)C	54.64541	54.64541	CC(C)C(C)C	33.21461	33.21461	CC(C)C(C)C		125.9246
CC(C)C(C)C	54.07587	54.07587	CC(C)C(C)C	31.90149	31.90149	CC(C)C(C)C		124.8349
C=C(C)C	53.94009	53.94009	CC(C)C(C)C	31.59909	31.59909	CC(C)C(C)C		124.8332
CC(C)C(C)C	53.44199	53.44199	CC(C)C(C)C	31.14121	31.14121	CC(C)C(C)C		122.8646
C=C(C)C	52.41753	52.41753	CC(C)C(C)C	31.08018	31.08018	CC(C)C(C)C		123.6242
O=C(C)C	52.16408	52.16408	CC(C)C(C)C	31.03749	31.03749	CC(C)C(C)C		123.591
C=C(C)C	52.36687	52.36687	CC(C)C(C)C	29.83913	29.83913	CC(C)C(C)C		122.529
C=C(C)C	52.15408	52.15408	CC(C)C(C)C	29.71734	29.71734	CC(C)C(C)C		123.9075
C=C(C)C	51.85418	51.85418	CC(C)C(C)C	29.61341	29.61341	CC(C)C(C)C		122.6161
C=C(C)C	51.58962	51.58962	CC(C)C(C)C	29.36574	29.36574	CC(C)C(C)C		122.4707
C=C(C)C	50.88131	50.88131	CC(C)C(C)C	29.09899	29.09899	CC(C)C(C)C		122.2686
C=C(C)C	50.85841	50.85841	CC(C)C(C)C	28.85484	28.85484	CC(C)C(C)C		122.2946
O=C(C)C	50.86649	50.86649	CC(C)C(C)C	28.75354	28.75354	CC(C)C(C)C		121.2901
C=C(C)C	50.35582	50.35582	CC(C)C(C)C	28.32705	28.32705	CC(C)C(C)C		120.9415
C=C(C)C	50.149	50.149	CC(C)C(C)C	28.20997	28.20997	CC(C)C(C)C		120.8793
C=C(C)C	49.62951	49.62951	CC(C)C(C)C	27.82103	27.82103	CC(C)C(C)C		120.9403
CC(C)C(C)C	49.24534	49.24534	CC(C)C(C)C	27.55956	27.55956	CC(C)C(C)C		119.9674
CC(C)C(C)C	49.20329	49.20329	CC(C)C(C)C	27.54716	27.54716	CC(C)C(C)C		119.8744
CC(C)C(C)C	49.01636	49.01636	CC(C)C(C)C	27.33422	27.33422	CC(C)C(C)C		119.8238
CC(C)C(C)C	48.94532	48.94532	CC(C)C(C)C	26.61405	26.61405	CC(C)C(C)C		119.3759
CC(C)C(C)C	48.47395	48.47395	CC(C)C(C)C	25.84082	25.84082	CC(C)C(C)C		119.3573
C=C(C)C	48.31301	48.31301	CC(C)C(C)C	25.62622	25.62622	CC(C)C(C)C		119.3014
CC(C)C(C)C	48.17758	48.17758	CC(C)C(C)C	25.47008	25.47008	CC(C)C(C)C		119.2656
C=C(C)C	48.14932	48.14932	CC(C)C(C)C	25.15138	25.15138	CC(C)C(C)C		118.642
O=C(C)C	48.08631	48.08631	CC(C)C(C)C	25.14498	25.14498	CC(C)C(C)C		118.2089
CC(C)C(C)C	48.05477	48.05477	CC(C)C(C)C	24.59309	24.59309	CC(C)C(C)C		118.1076
O=C(C)C	47.82076	47.82076	CC(C)C(C)C	23.99376	23.99376	CC(C)C(C)C		117.0701
CC(C)C(C)C	47.81427	47.81427	CC(C)C(C)C	23.89949	23.89949	CC(C)C(C)C		116.8756
CC(C)C(C)C	47.64478	47.64478	CC(C)C(C)C	23.58684	23.58684	CC(C)C(C)C		116.8452
CC(C)C(C)C	47.06617	47.06617	CC(C)C(C)C	22.7958	22.7958	CC(C)C(C)C		116.1719
CC(C)C(C)C	46.97055	46.97055	CC(C)C(C)C	22.71871	22.71871	CC(C)C(C)C		116.0629
CC(C)C(C)C	46.92427	46.92427	CC(C)C(C)C	22.56	111	CC(C)C(C)C		115.2144
CC(C)C(C)C	46.84	116	CC(C)C(C)C	21.89512	21.89512	CC(C)C(C)C		114.4769
CC(C)C(C)C	46.45414	46.45414	CC(C)C(C)C	21.57717	21.57717	CC(C)C(C)C		114.2716
O=C(C)C	45.75258	45.75258	CC(C)C(C)C	21.27533	21.27533	CC(C)C(C)C		114.1827
CC(C)C(C)C	45.67465	45.67465	CC(C)C(C)C	20.49882	20.49882	CC(C)C(C)C		113.986
CC(C)C(C)C	45.54189	45.54189	CC(C)C(C)C	20.17205	20.17205	CC(C)C(C)C		112.9504
O=C(C)C	45.3955	45.3955	CC(C)C(C)C	19.93297	19.93297	CC(C)C(C)C		112.6495
CC(C)C(C)C	45.29639	45.29639	CC(C)C(C)C	19.87415	19.87415	CC(C)C(C)C		112.624
O=C(C)C	45.09746	45.09746	CC(C)C(C)C	19.81655	19.81655	CC(C)C(C)C		112.1986
CC(C)C(C)C	44.3957	44.3957	CC(C)C(C)C	19.52323	19.52323	CC(C)C(C)C		111.6846
CC(C)C(C)C	44.29157	44.29157	CC(C)C(C)C	19.21828	19.21828	CC(C)C(C)C		110.9251
O=C(C)C	44.27086	44.27086	CC(C)C(C)C	19.18392	19.18392	CC(C)C(C)C		110.8575
CC(C)C(C)C	44.07604	44.07604	CC(C)C(C)C	19.0844	19.0844	CC(C)C(C)C		110.2391
CC(C)C(C)C	44.07124	44.07124	CC(C)C(C)C	18.8576	18.8576	CC(C)C(C)C		109.7006
C=C(C)C	43.82648	43.82648	CC(C)C(C)C	18.75997	18.75997	CC(C)C(C)C		109.6941
CC(C)C(C)C	43.76267	43.76267	CC(C)C(C)C	17.82433	17.82433	CC(C)C(C)C		108.6
CC(C)C(C)C	43.67166	43.67166	CC(C)C(C)C	17.35709	17.35709	CC(C)C(C)C		107.924
O=C(C)C	43.43836	43.43836	CC(C)C(C)C	16.94432	16.94432	CC(C)C(C)C		107.4248
CC(C)C(C)C	43.3631	43.3631	CC(C)C(C)C	16.89494	16.89494	CC(C)C(C)C		107.4013
O=C(C)C	43.12843	43.12843	CC(C)C(C)C	16.89059	16.89059	CC(C)C(C)C		107.3825
CC(C)C(C)C	43.0524	43.0524	CC(C)C(C)C	16.85496	16.85496	CC(C)C(C)C		107.1016
CC(C)C(C)C	42.9705	42.9705	CC(C)C(C)C	16.82522	16.82522	CC(C)C(C)C		107.0201
CC(C)C(C)C	42.48007	42.48007	CC(C)C(C)C	16.71336	16.71336	CC(C)C(C)C		106.9824
C=C(C)C	42.31857	42.31857	CC(C)C(C)C	16.69362	16.69362	CC(C)C(C)C		106.9291
CC(C)C(C)C	42.26556	42.26556	CC(C)C(C)C	16.54525	16.54525	CC(C)C(C)C		106.7174
CC(C)C(C)C	42.0531	42.0531	CC(C)C(C)C	16.28343	16.28343	CC(C)C(C)C		106.7335
O=C(C)C	41.98418	41.98418	CC(C)C(C)C	15.97413	15.97413	CC(C)C(C)C		106.2103
CC(C)C(C)C	41.96767	41.96767	CC(C)C(C)C	15.79651	15.79651	CC(C)C(C)C		105.9842
C=C(C)C	41.92689	41.92689	CC(C)C(C)C	15.71584	15.71584	CC(C)C(C)C		105.7342
CC(C)C(C)C	41.67139	41.67139	CC(C)C(C)C	15.27856	15.27856	CC(C)C(C)C		105.613

Table 3.14 Continued

RP6_C			RP7_B			RP7_C		
SMILES	Pred	Activity	SMILES	Pred	Activity	SMILES	Pred	Activity
CCCC(O)C=O	86.13873		C1=C(O)C(O)C(O)C=O	126.4566		CC(C)C(O)C=O	36.52692	
OC1=CC(O)C(O)C=O	79.27182		C=C1C(O)C(O)C(O)C=O	119.7731		CCCC(O)C=O	30.50133	
OC(O)C=O	87.20345		OC1=C(O)C(O)C=O	95.38008		CC(O)C(O)C=O	16.94263	
OC(C1=CC(O)C=O)C=O	65.13451		CC(C)C(O)C=O	76.21054		CC(C)C(O)C=O	36.52692	
OC1=CC(O)C(O)C=O	58.16259		CCCC(O)C=O	67.77761		CCCC(O)C=O	33.78351	
OC(O)C(O)C=O	54.28428		OC1=C(O)C(O)C(O)C=O	126.4566		O=CC(O)C=O	33.33575	
OC(O)C(O)C=O	86.13873		OC1=C(O)C(O)C(O)C=O	95.38008		CC(O)C(O)C=O	31.71206	
OC(O)C(O)C=O	83.38456		CC(C)C(O)C=O	92.59738		CC(C)C(O)C=O	31.70961	
OC1=CC(O)C(O)C=O	79.27182		CC(C)C(O)C(O)C=O	91.15.114		CC(O)C(O)C=O	31.0054	
O=C(O)C(O)C=O	79.13522		CC(C)C(O)C(O)C=O	89.91984		CCCC(O)C=O	30.50133	
C=CC(O)C=O	77.20856		CC(O)C(O)C(O)C=O	87.99201		CCCC(O)C=O	29.72137	
C=CC(O)C(O)C=O	76.38981		CC(C)C(O)C(O)C=O	87.32837		C=CC(O)C(O)C=O	29.40163	
CCCC(O)C=O	76.16308		CC(C)C(O)C(O)C=O	87.14336		O=CCCC(O)C=O	28.98876	
OC(O)C(O)C=O	74.89678		CC(O)C(O)C(O)C=O	83.85426		CC(C)C(O)C=O	28.26886	
O=C(O)C(O)C(O)C=O	74.12276		CC(C)C(O)C(O)C=O	85.36532		CC(O)C(O)C(O)C=O	26.77087	
O=C1C(O)C(O)C=O	73.47661		CC(O)C(O)C(O)C=O	82.94325		CCCC(O)C(O)C=O	26.07484	
C=O(O)C(O)C=O	73.1929		O=C(O)C(O)C(O)C=O	81.97053		CCCC(O)C=O	25.96475	
OC(O)C(O)C=O	71.946.11		OC1=C(O)C(O)C(O)C=O	80.79637		CC(O)C(O)C=O	25.59708	
OC(O)C(O)C=O	71.24633		O=C(O)C(O)C(O)C=O	80.65999		CC(O)C(O)C(O)C=O	25.02162	
O=C(O)C(O)C=O	71.14553		CC(C)C(O)C(O)C=O	80.58242		CCCC(O)C(O)C=O	24.15837	
OC(O)C(O)C(O)C=O	71.05457		CC(O)C(O)C(O)C=O	80.41422		CCCC(O)C(O)C=O	23.89254	
C=CC(O)C=O	70.38179		CC(O)C(O)C(O)C=O	80.22339		CC(O)C(O)C=O	23.88196	
C=C(O)C(O)C(O)C=O	70.27329		OC1=C(O)C(O)C(O)C=O	79.99125		CCCC(O)C=O	23.52.119	
NCCCC(O)C=O	69.99448		CC(N)C(O)C(O)C=O	76.5694		CC(O)C(O)C=O	21.55237	
C=C(O)C(O)C=O	69.2848		OC1=C(O)C(O)C(O)C=O	76.23738		CC(O)C(O)C(O)C(O)C=O	21.16626	
CC(O)C(O)C(O)C=O	69.24917		CCCC(O)C(O)C=O	76.21054		O=C(O)C(O)C(O)C(O)C=O	18.74665	
C=C(O)C(O)C(O)C(O)C=O	68.68112		NCCCC(O)C(O)C=O	75.85999		CC(O)C(O)C(O)C(O)C=O	18.6151	
O=C(O)C(O)C(O)C(O)C=O	68.14354		CC(O)C(O)C(O)C=O	75.52013		CCCC(O)C(O)C(O)C=O	18.52774	
C=C(O)C(O)C(O)C(O)C=O	67.73656		CC(O)C(O)C(O)C=O	75.43823		SCC1=CC(O)C=O	18.04152	
OC(O)C(O)C=O	67.64093		CC(O)C(O)C(O)C=O	74.71706		O=C(O)C(O)C(O)C(O)C=O	17.93662	
OC(O)C(O)C=O	67.29345		OC1=C(O)C(O)C(O)C=O	74.65248		SCC(O)C(O)C=O	17.85101	
C=C(O)C(O)C(O)C=O	66.99099		O=C(O)C(O)C(O)C(O)C=O	74.49208		O=CC(O)C(O)C(O)C=O	17.37472	
C=C(O)C(O)C(O)C=O	66.84646		OC(O)C(O)C(O)C(O)C=O	74.0.1133		O=CC(O)C(O)C(O)C=O	17.32718	
C=C(O)C(O)C(O)C=O	66.64964		OC(O)C(O)C(O)C(O)C=O	73.71874		CC(O)C(O)C(O)C(O)C=O	17.25056	
OC(O)C(O)C(O)C=O	65.98814		OC(O)C(O)C(O)C(O)C=O	73.5887		CC(O)C(O)C(O)C(O)C=O	17.15448	
O=C(O)C(O)C(O)C(O)C=O	65.85213		CC(O)C(O)C(O)C=O	73.46298		CC(O)C(O)C(O)C=O	17.12756	
OC(O)C(O)C(O)C(O)C=O	65.65882		CCCC(O)C(O)C=O	73.3584		CC(O)C(O)C(O)C(O)C=O	16.86213	
OC(O)C(O)C(O)C(O)C=O	65.13451		O=C(O)C(O)C(O)C(O)C=O	73.31343		CC(O)C(O)C(O)C(O)C=O	16.84203	
OC(O)C(O)C(O)C(O)C=O	63.74982		OC(O)C(O)C(O)C(O)C=O	72.78.11		SCC(O)C(O)C=O	16.78584	
CCCC(O)C(O)C=O	63.63551		NCCCC(O)C(O)C=O	72.23238		O=C(O)C(O)C(O)C=O	16.65321	
C=CC(O)C=O	63.27646		C=C(O)C(O)C(O)C(O)C=O	71.73195		O=C(O)C(O)C(O)C(O)C=O	16.58147	
CCCC(O)C(O)C(O)C=O	62.21891		CCCC(O)C(O)C(O)C=O	71.16331		CC(O)C(O)C(O)C(O)C=O	16.53789	
C=CC(O)C=O	62.04977		O=C(O)C(O)C(O)C(O)C=O	70.22063		OC(O)C(O)C(O)C(O)C=O	16.53871	
O=C(O)C(O)C=O	61.81868		CCCC(O)C(O)C=O	70.09967		O=CCCC(O)C=O	16.48318	
OC1=N2C=CC(O)C=O	61.46244		CC(O)C(O)C(O)C=O	69.75241		CC(O)C(O)C(O)C=O	16.30734	
OC(O)C(O)C(O)C=O	61.36557		CC(O)C(O)C(O)C=O	69.04661		SCC(O)C(O)C=O	16.24069	
O=C(O)C(O)C(O)C(O)C=O	60.99415		C=C(O)C(O)C(O)C(O)C=O	68.91515		CCCC(O)C(O)C(O)C=O	16.19149	
C=C(O)C(O)C(O)C(O)C=O	60.81307		CC(O)C(O)C(O)C=O	68.86344		CC(O)C(O)C(O)C=O	16.10514	
C=CC(O)C=O	60.4865		CC(O)C(O)C(O)C(O)C=O	68.28992		CCC(O)C=O	16.09149	
C=C(O)C(O)C(O)C=O	60.14782		CCCC(O)C(O)C=O	67.77761		O=C(O)C(O)C(O)C(O)C=O	16.08724	
OC1=N2C=CC(O)C=O	59.98227		CC(O)C(O)C(O)C=O	67.68886		OC1=C(O)C(O)C(O)C(O)C=O	16.08161	
C=C(O)C(O)C=O	59.95997		CC(O)C(O)C(O)C=O	67.50108		CC(N)C(O)C(O)C=O	16.05204	
OC1=CC(O)C(O)C(O)C=O	59.7741		OC1=C(O)C(O)C(O)C(O)C=O	67.2853		O=C(O)C(O)C(O)C=O	16.0168	
OC(O)C(O)C=O	59.55189		O=C(O)C(O)C(O)C(O)C=O	67.07072		OC1=C(O)C(O)C(O)C(O)C=O	15.94205	
O=C(O)C(O)C(O)C(O)C=O	59.29277		CC(O)C(O)C(O)C=O	66.89404		C(O)C(O)C(O)C(O)C(O)C=O	15.74377	
CC(O)C(O)C(O)C=O	59.19607		[C-]N+CC1=CC(O)C=O	66.39249				
CC(O)C(O)C(O)C=O	59.08062		SCC(O)C(O)C(O)C=O	66.19192				
OC(O)C(O)C(O)C(O)C=O	58.72099		OC1=C2C=CC(O)C(O)C=O	66.18288				
O=C(O)C(O)C(O)C(O)C=O	58.68.119		CC(O)C(O)C(O)C=O	66.96761				
OC1=CC(O)C(O)C(O)C=O	58.16259		CC(O)C(O)C(O)C(O)C=O	65.222				
OC1=CC(O)C(O)C(O)C(O)C=O	58.04722		CC(O)C(O)C(O)C(O)C=O	65.18269				
OC(O)C(O)C(O)C(O)C=O	57.98341		OC1=C(O)C(O)C(O)C=O	64.98162				
OC(O)C(O)C(O)C=O	57.90524		CC(O)C(O)C(O)C(O)C=O	64.94732				
C=CC(O)C=O	57.6607		C=CC(O)C(O)C=O	64.92842				
O=C(O)C(O)C(O)C=O	57.63992		C=CC(O)C(O)C(O)C(O)C=O	64.86062				
OC(O)C(O)C(O)C(O)C=O	56.97463		OC(O)C(O)C(O)C(O)C=O	64.52239				
OC(O)C(O)C(O)C=O	56.07965		CC(O)C(O)C(O)C(O)C=O	64.39563				
CC(O)C(O)C(O)C(O)C=O	55.89354		O=C(O)C(O)C(O)C(O)C=O	64.20792				
CC(O)C(O)C(O)C=O	55.65697		CCCC(O)C(O)C=O	64.0447				
CC(O)C(O)C(O)C=O	55.32886		O=C(O)C(O)C(O)C=O	63.7268				
CC(O)C(O)C(O)C(O)C=O	54.53139		OC1=C2C=CC(O)C(O)C(O)C=O	63.94465				
C12=CC(O)C(O)C(O)C=O	54.33279		CC(O)C(O)C(O)C(O)C(O)C=O	62.80598				
OC1=C(O)C(O)C(O)C(O)C=O	54.28428		O=C(O)C(O)C(O)C=O	62.68223				
OC(O)C(O)C(O)C(O)C(O)C=O	54.23887		CC(O)C(O)C(O)C(O)C=O	62.66521				
O=C(O)C(O)C(O)C(O)C(O)C=O	54.08839		CC(O)C(O)C(O)C(O)C(O)C=O	62.63064				
O=C(O)C(O)C(O)C(O)C(O)C(O)C=O	53.91696		CC(O)C(O)C(O)C(O)C(O)C=O	62.58525				
OC(O)C(O)C(O)C(O)C(O)C=O	53.80537		C=CC(O)C(O)C(O)C=O	62.19846				
OC(O)C(O)C(O)C(O)C(O)C=O	53.80223		O=C(O)C(O)C(O)C(O)C(O)C=O	61.90239				
C12=CC(O)C(O)C(O)C(O)C=O	53.84288		OC(O)C(O)C(O)C(O)C(O)C=O	61.85532				
O=C(O)C(O)C(O)C(O)C(O)C=O	53.42226		C=C(O)C(O)C(O)C(O)C=O	61.75396				
O=C(O)C(O)C(O)C=O	53.40861		CC(O)C(O)C(O)C(O)C=O	61.65226				
O=C(O)C(O)C(O)C(O)C(O)C=O	53.16238		C=C(O)C(O)C(O)C(O)C(O)C(O)C=O	61.64239				
O=C(O)C(O)C(O)C(O)C(O)C(O)C=O	52.99559		OC1=C(O)C(O)C(O)C(O)C=O	61.638.11				
O=C(O)C(O)C(O)C(O)C(O)C(O)C=O	52.29142		CC(O)C(O)C(O)C(O)C=O	61.58498				
OC(O)C(O)C(O)C(O)C=O	52.16263		CC(O)C(O)C(O)C(O)C(O)C=O	61.33144				
OC(O)C(O)C(O)C(O)C(O)C(O)C=O	51.91039		O=CCCC(O)C=O	61.20754				
OC(O)C(O)C(O)C(O)C(O)C(O)C=O	51.85885		OC1=C(O)C(O)C(O)C(O)C(O)C=O	61.05316				
O=C(O)C(O)C(O)C(O)C(O)C(O)C=O	51.52354		C=C(O)C(O)C(O)C(O)C(O)C=O	61.01872				
OC(O)C(O)C(O)C(O)C(O)C(O)C=O	51.519.11		CC(O)C(O)C(O)C(O)C=O	60.98861				
OC(O)C(O)C(O)C(O)C(O)C(O)C=O	51.09141		OC(O)C(O)C(O)C(O)C=O	60.7091				
OC(O)C(O)C(O)C(O)C(O)C(O)C=O	50.99655		CC(O)C(O)C(O)C(O)C(O)C=O	60.625.11				
OC(O)C(O)C(O)C(O)C(O)C(O)C=O	49.9.1146		OC1=C2C=CC(O)C(O)C(O)C=O	60.32029				
OC(O)C(O)C(O)C(O)C(O)C(O)C=O	49.79867		CCCC(O)C(O)C=O	60.2745				
OC(O)C(O)C(O)C(O)C(O)C(O)C=O	49.77356		O=C(O)C(O)C(O)C(O)C(O)C=O	60.11139				
O=C(O)C(O)C(O)C(O)C(O)C(O)C=O	49.71818		CC(O)C(O)C(O)C(O)C=O	59.55436				
O=C(O)C(O)C(O)C(O)C(O)C(O)C=O	49.69014		C=CC(O)C(O)C(O)C(O)C=O	59.54975				
OC(O)C(O)C(O)C=O	49.61974		CC(O)C(O)C(O)C(O)C=O	59.50504				
C=CC(O)C=O	49.44057		C=C(O)C(O)C(O)C(O)C(O)C=O	59.12154				
O=C1SC2=CC(O)C=O	48.96835		CC(O)C(O)C(O)C(O)C=O	59.03413				

Table 3.14 Continued

Table 3.15: Optimized descriptor sets for each mammalian OR

Optimized descriptors occurrences, symbol, brief description, class, and dimensionality are listed. A summary of the total number of descriptors selected for the receptor repertoire is provided at the beginning. Descriptors are listed in ascending order of when they were selected into the optimized set. Weights indicate the number of times a descriptor was selected in an optimized descriptor set.

Descriptor Class	Type	Counts for all Ors
GETAWAY descriptors		109
atom-centred fragments		49
2D autocorrelations		48
RDF descriptors		48
3D-MoRSE descriptors		46
WHIM descriptors		43
functional group counts		33
2D binary fingerprints		26
Burden eigenvalues		23
edge adjacency indices		21
geometrical descriptors		21
topological descriptors		14
2D frequency fingerprints		13
topological charge indices		12
atomtypes (Cerius2)		11
molecular properties		11
walk and path counts		7
constitutional descriptors		6
Randic molecular profiles		5
topological (Cerius2)		4
information indices		4
connectivity indices		3
structural (Cerius2)		1
eigenvalue-based indices		1
charge descriptors		0

Dimensionality Counts (Weights Included)	
Num zero dimensional descriptors:	7
Num one dimensional descriptors:	104
Num two dimensional descriptors:	176
Num three dimensional descriptors:	272

Origin (Weights Included)	
Num Dragon descriptors:	546
Num Cerius2 descriptors:	13

Dimensionality Counts (Weights Excluded)	
Num zero dimensional descriptors:	7
Num one dimensional descriptors:	37
Num two dimensional descriptors:	93
Num three dimensional descriptors:	155

Origin (Weights Excluded)	
Num unique Dragon descriptors:	284
Num unique Cerius2 descriptors:	8

Odor Receptor Name	Weight	Symbol	Description	Class	Dimensionality	
MOR1.1	2	Mor17m	3D-MoRSE - signal 17 / weighted by atomic masses	3D-MoRSE descriptors	3	
	8	H-051	H attached to alpha-C	atom-centred fragments	1	
	2	R6p+	R maximal autocorrelation of lag 6 / weighted by atomic polarizabilities	GETAWAY descriptors	3	
	4	Mor23e	3D-MoRSE - signal 23 / weighted by atomic Sanderson electronegativities	3D-MoRSE descriptors	3	
	5	Mor30m	3D-MoRSE - signal 30 / weighted by atomic masses	3D-MoRSE descriptors	3	
	4	R5v+	R maximal autocorrelation of lag 5 / weighted by atomic van der Waals volumes	GETAWAY descriptors	3	
	1	Mor32e	3D-MoRSE - signal 32 / weighted by atomic Sanderson electronegativities	3D-MoRSE descriptors	3	
	1	IGI7	mean topological charge index of order 7	topological charge indices	2	
	2	E1s	1st component accessibility directional WHIM index / weighted by atomic electrotopological states	WHIM descriptors	3	
	1	nCt	number of total tertiary C(sp3)	functional group counts	1	
	1	nArCO	number of ketones (aromatic)	functional group counts	1	
	1	O-058	presence/absence of C - O at topological distance 07	atom-centred fragments	1	
	1	B07[C-O]	Geary autocorrelation - lag 2 / weighted by atomic masses	2D binary fingerprints	2	
	2	GAT52m	Radial Distribution Function - 11.0 / weighted by atomic Sanderson electronegativities	2D autocorrelations	2	
	1	RDF110e	number of non-aromatic conjugated C(sp2)	RDF descriptors	3	
	1	nConj	CH2RX	functional group counts	1	
	1	C-006	number of ring secondary C(sp3)	atom-centred fragments	1	
	1	nCs	Geary autocorrelation - lag 7 / weighted by atomic masses	functional group counts	1	
	1	GAT57m	CHR3	atom-centred fragments	1	
	1	C-003	Moran autocorrelation - lag 8 / weighted by atomic masses	atom-centred fragments	1	
	1	MAT58m		2D autocorrelations	2	
	MOR106.1	1	Mor25p	3D-MoRSE - signal 25 / weighted by atomic polarizabilities	3D-MoRSE descriptors	3
		1	BEHe6	highest eigenvalue n. 6 of Burden matrix / weighted by atomic Sanderson electronegativities	Burden eigenvalues	2
1		IC1	information content index (neighborhood symmetry of 1-order)	information indices	2	
2		C-006	CH2RX	atom-centred fragments	1	
2		Rte+	R maximal index / weighted by atomic Sanderson electronegativities	GETAWAY descriptors	3	
1		piPC07	molecular multiple path count of order 07	walk and path counts	2	
1		RDF045e	Radial Distribution Function - 4.5 / weighted by atomic Sanderson electronegativities	RDF descriptors	3	
1		nRCOOH	number of carboxylic acids (aliphatic)	functional group counts	1	
1		R7v+	R maximal autocorrelation of lag 7 / weighted by atomic van der Waals volumes	GETAWAY descriptors	3	
1		HOMT	HOMA total	geometrical descriptors	3	
1		RDF035m	Radial Distribution Function - 3.5 / weighted by atomic masses	RDF descriptors	3	
1		H-049	H attached to C3(sp3)/C2(sp2)/C3(sp2)/C3(sp)	atom-centred fragments	1	
1		SHP2	average shape profile index of order 2	Randic molecular profiles	3	
MOR107.1	15	nCq	number of total quaternary C(sp3)	functional group counts	1	
	1	Mor07v	3D-MoRSE - signal 07 / weighted by atomic van der Waals volumes	3D-MoRSE descriptors	3	
	1	BEL3	lowest eigenvalue n. 3 of Burden matrix / weighted by atomic van der Waals volumes	Burden eigenvalues	2	
	1	R4p+	R maximal autocorrelation of lag 4 / weighted by atomic polarizabilities	GETAWAY descriptors	3	
	1	E1u	1st component accessibility directional WHIM index / unweighted	WHIM descriptors	3	
	1	JG14	mean topological charge index of order 4	topological charge indices	2	
	1	DISPV	d COMMAZ value / weighted by atomic van der Waals volumes	geometrical descriptors	3	
	1	nR06	number of 6-membered rings	constitutional descriptors	0	
	1	RDF040m	Radial Distribution Function - 4.0 / weighted by atomic masses	RDF descriptors	3	
	1	R8e+	R maximal autocorrelation of lag 8 / weighted by atomic Sanderson electronegativities	GETAWAY descriptors	3	
	1	piPC07	molecular multiple path count of order 07	walk and path counts	2	
	1	L2s	2nd component size directional WHIM index / weighted by atomic electrotopological states	WHIM descriptors	3	
	1	piPC08	molecular multiple path count of order 08	walk and path counts	2	
	1	nDB	number of double bonds	constitutional descriptors	0	
	1	H-051	H attached to alpha-C	atom-centred fragments	1	
	1	E1m	1st component accessibility directional WHIM index / weighted by atomic masses	WHIM descriptors	3	
	1	B05[C-O]	presence/absence of C - O at topological distance 05	2D binary fingerprints	2	
	1	BELp3	lowest eigenvalue n. 3 of Burden matrix / weighted by atomic polarizabilities	Burden eigenvalues	2	
	MOR129.1	1	Dv	D total accessibility index / weighted by atomic van der Waals volumes	WHIM descriptors	3
		2	HATS7e	leverage-weighted autocorrelation of lag 7 / weighted by atomic Sanderson electronegativities	GETAWAY descriptors	3
		1	F01[C-C]	frequency of C - C at topological distance 01	2D frequency fingerprints	2
		1	Mor10u	3D-MoRSE - signal 10 / unweighted	3D-MoRSE descriptors	3
		1	HATS5m	leverage-weighted autocorrelation of lag 5 / weighted by atomic masses	GETAWAY descriptors	3
2		H1v	H autocorrelation of lag 1 / weighted by atomic van der Waals volumes	GETAWAY descriptors	3	
1		Mor11v	3D-MoRSE - signal 11 / weighted by atomic van der Waals volumes	3D-MoRSE descriptors	3	
1		B05[C-C]	presence/absence of C - C at topological distance 05	2D binary fingerprints	2	
1		RDF085e	Radial Distribution Function - 8.5 / weighted by atomic Sanderson electronegativities	RDF descriptors	3	
1		Dm	D total accessibility index / weighted by atomic masses	WHIM descriptors	3	
1		H0m	H autocorrelation of lag 0 / weighted by atomic masses	GETAWAY descriptors	3	
1		Mor10e	3D-MoRSE - signal 10 / weighted by atomic Sanderson electronegativities	3D-MoRSE descriptors	3	
1		D/D06	distance/detour ring index of order 6	topological descriptors	2	
1		HATS6m	leverage-weighted autocorrelation of lag 6 / weighted by atomic masses	GETAWAY descriptors	3	
1		MATS5p	Moran autocorrelation - lag 5 / weighted by atomic polarizabilities	2D autocorrelations	2	
1		RDF035m	Radial Distribution Function - 3.5 / weighted by atomic masses	RDF descriptors	3	
MOR136.1		1	H7e	H autocorrelation of lag 7 / weighted by atomic Sanderson electronegativities	GETAWAY descriptors	3
	1	H0v	H autocorrelation of lag 0 / weighted by atomic van der Waals volumes	GETAWAY descriptors	3	
	1	R5m	R autocorrelation of lag 5 / weighted by atomic masses	GETAWAY descriptors	3	
	1	RDF030m	Radial Distribution Function - 3.0 / weighted by atomic masses	RDF descriptors	3	
	1	S_dssC	atomtypes (cerius2)	atomtypes (cerius2)	1	
MOR139.1	1	Mor05m	3D-MoRSE - signal 05 / weighted by atomic masses	3D-MoRSE descriptors	3	
	1	BELe1	lowest eigenvalue n. 1 of Burden matrix / weighted by atomic Sanderson electronegativities	Burden eigenvalues	2	
	1	R5u+	R maximal autocorrelation of lag 5 / unweighted	GETAWAY descriptors	3	
	1	SHP2	average shape profile index of order 2	Randic molecular profiles	3	

Table 3.15

	2H1v	H autocorrelation of lag 1 / weighted by atomic van der Waals volumes	GETAWAY descriptors	3
	1HATS1u	leverage-weighted autocorrelation of lag 1 / unweighted	GETAWAY descriptors	3
	1TPSA(NO)	topological polar surface area using N,O polar contributions	molecular properties	1
	1piPC06	molecular multiple path count of order 06	walk and path counts	2
	1HBm	H autocorrelation of lag 8 / weighted by atomic masses	GETAWAY descriptors	3
	1GVWAI-80	Ghose-Viswanadhan-Wendoloski drug-like index at 80%	molecular properties	1
	1R1v	R autocorrelation of lag 1 / weighted by atomic van der Waals volumes	GETAWAY descriptors	3
	1GATS2m	Geary autocorrelation - lag 2 / weighted by atomic masses	2D autocorrelations	2
	1Mor24e	3D-MoRSE - signal 24 / weighted by atomic Sanderson electronegativities	3D-MoRSE descriptors	3
	1EEig09d	Eigenvalue 09 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2
MOR162.1	1HOMA	Harmonic Oscillator Model of Aromaticity index	geometrical descriptors	3
	1HATS5m	leverage-weighted autocorrelation of lag 5 / weighted by atomic masses	GETAWAY descriptors	3
	1E2e	2nd component accessibility directional WHIM index / weighted by atomic Sanderson electronegativities	WHIM descriptors	3
	1H2e	H autocorrelation of lag 2 / weighted by atomic Sanderson electronegativities	GETAWAY descriptors	3
	1R6v+	R maximal autocorrelation of lag 6 / weighted by atomic van der Waals volumes	GETAWAY descriptors	3
	1P2e	2nd component shape directional WHIM index / weighted by atomic Sanderson electronegativities	WHIM descriptors	3
	1MAT55p	Moran autocorrelation - lag 5 / weighted by atomic polarizabilities	2D autocorrelations	2
	1RC1	Jug RC index	geometrical descriptors	3
	1HATS6m	leverage-weighted autocorrelation of lag 6 / weighted by atomic masses	GETAWAY descriptors	3
	1RDF035m	Radial Distribution Function - 3.5 / weighted by atomic masses	RDF descriptors	3
	1H1e	H autocorrelation of lag 1 / weighted by atomic Sanderson electronegativities	GETAWAY descriptors	3
MOR170.1	3C-02S	R-CR-R	atom-centred fragments	1
	1F0S[C-O]	frequency of C - O at topological distance 05	2D frequency fingerprints	2
	1R6v+	R maximal autocorrelation of lag 6 / weighted by atomic van der Waals volumes	GETAWAY descriptors	3
	2E2e	2nd component accessibility directional WHIM index / weighted by atomic Sanderson electronegativities	WHIM descriptors	3
	1S_aaCH	S_aaCH	atomtypes (cerius2)	1
	2Rt+	R maximal index / weighted by atomic Sanderson electronegativities	GETAWAY descriptors	3
	1R8u+	R maximal autocorrelation of lag 8 / unweighted	GETAWAY descriptors	3
	1ATS2p	Broto-Moreau autocorrelation of a topological structure - lag 2 / weighted by atomic polarizabilities	2D autocorrelations	2
	1MAT55p	Moran autocorrelation - lag 5 / weighted by atomic polarizabilities	2D autocorrelations	2
	1RDF045v	Radial Distribution Function - 4.5 / weighted by atomic van der Waals volumes	RDF descriptors	3
	1P2p	2nd component shape directional WHIM index / weighted by atomic polarizabilities	WHIM descriptors	3
	1R7p+	R maximal autocorrelation of lag 7 / weighted by atomic polarizabilities	GETAWAY descriptors	3
	1R5m	R autocorrelation of lag 5 / weighted by atomic masses	GETAWAY descriptors	3
	1HOMA	Harmonic Oscillator Model of Aromaticity index	geometrical descriptors	3
MOR184.1	1S_dCH2	S_dCH2	atomtypes (cerius2)	1
	1GATS1m	Geary autocorrelation - lag 1 / weighted by atomic masses	2D autocorrelations	2
	1H-047	H attached to C1(sp3)/C0(sp2)	atom-centred fragments	1
	1Qindex	Quadratic index	topological descriptors	2
	1DISPv	d COMMA2 value / weighted by atomic van der Waals volumes	geometrical descriptors	3
	1L2v	2nd component size directional WHIM index / weighted by atomic van der Waals volumes	WHIM descriptors	3
	1nCIC	number of rings	constitutional descriptors	0
	1HATS2v	leverage-weighted autocorrelation of lag 2 / weighted by atomic van der Waals volumes	GETAWAY descriptors	3
	1Mor07u	3D-MoRSE - signal 07 / unweighted	3D-MoRSE descriptors	3
	1R7e	R autocorrelation of lag 7 / weighted by atomic Sanderson electronegativities	GETAWAY descriptors	3
	1MPC06	molecular path count of order 06	walk and path counts	2
	1EEig09r	Eigenvalue 09 from edge adj. matrix weighted by resonance integrals	edge adjacency indices	2
	1nCIR	number of circuits	constitutional descriptors	0
	1RDF045m	Radial Distribution Function - 4.5 / weighted by atomic masses	RDF descriptors	3
	1BELp3	lowest eigenvalue n. 3 of Burden matrix / weighted by atomic polarizabilities	Burden eigenvalues	3
	1EEig09d	Eigenvalue 09 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2
	1L2p	2nd component size directional WHIM index / weighted by atomic polarizabilities	WHIM descriptors	3
	1HATS7p	leverage-weighted autocorrelation of lag 7 / weighted by atomic polarizabilities	GETAWAY descriptors	3
	1R4p+	R maximal autocorrelation of lag 4 / weighted by atomic polarizabilities	GETAWAY descriptors	3
	1Mor10m	3D-MoRSE - signal 10 / weighted by atomic masses	3D-MoRSE descriptors	3
MOR185.1	1BAC	Balaban centric index	topological descriptors	2
	1HATS2m	leverage-weighted autocorrelation of lag 2 / weighted by atomic masses	GETAWAY descriptors	3
	1R8u+	R maximal autocorrelation of lag 8 / unweighted	GETAWAY descriptors	3
	1X5A	average connectivity index chi-5	connectivity indices	2
	1E1e	1st component accessibility directional WHIM index / weighted by atomic Sanderson electronegativities	WHIM descriptors	3
	1Rt+	R maximal index / weighted by atomic Sanderson electronegativities	GETAWAY descriptors	3
	1E2e	2nd component accessibility directional WHIM index / weighted by atomic Sanderson electronegativities	WHIM descriptors	3
	1GATS2m	Geary autocorrelation - lag 2 / weighted by atomic masses	2D autocorrelations	2
	1H1p	H autocorrelation of lag 1 / weighted by atomic polarizabilities	GETAWAY descriptors	3
	1HATS7p	leverage-weighted autocorrelation of lag 7 / weighted by atomic polarizabilities	GETAWAY descriptors	3
	1RDF035m	Radial Distribution Function - 3.5 / weighted by atomic masses	RDF descriptors	3
	1SHP2	average shape profile index of order 2	1D molecular profiles	3
	1R2p+	R maximal autocorrelation of lag 2 / weighted by atomic polarizabilities	GETAWAY descriptors	3
MOR189.1	1nCrs	number of ring secondary C(sp3)	functional group counts	1
	1V-DIST-mag	V-DIST-mag	topological (cerius2)	2
	1J3D	3D-Balaban index	geometrical descriptors	3
	2Atype_C_40	Number of Carbon Type 40	atomtypes (Cerius2)	1
	1EEig11r	Eigenvalue 11 from edge adj. matrix weighted by resonance integrals	edge adjacency indices	2
	1R4m	R autocorrelation of lag 4 / weighted by atomic masses	GETAWAY descriptors	3
	1Mor07p	3D-MoRSE - signal 07 / weighted by atomic polarizabilities	3D-MoRSE descriptors	3
	1GVWAI-80	Ghose-Viswanadhan-Wendoloski drug-like index at 80%	molecular properties	1
	1RDF025m	Radial Distribution Function - 2.5 / weighted by atomic masses	RDF descriptors	3
	1B07[C-C]	presence/absence of C - C at topological distance 07	2D binary fingerprints	2
	1R4u+	R maximal autocorrelation of lag 4 / unweighted	GETAWAY descriptors	3
	1Mor22e	3D-MoRSE - signal 22 / weighted by atomic Sanderson electronegativities	3D-MoRSE descriptors	3
	1O-058	#NAME?	atom-centred fragments	1
	1EEig07d	Eigenvalue 07 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2
	1EEig11d	Eigenvalue 11 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2
MOR2.1	1RDF045m	Radial Distribution Function - 4.5 / weighted by atomic masses	RDF descriptors	3
	2GVWAI-80	Ghose-Viswanadhan-Wendoloski drug-like index at 80%	molecular properties	1
	1BELe7	lowest eigenvalue n. 7 of Burden matrix / weighted by atomic Sanderson electronegativities	Burden eigenvalues	2
	1RDF050m	Radial Distribution Function - 5.0 / weighted by atomic masses	RDF descriptors	3
	1Mor14e	3D-MoRSE - signal 14 / weighted by atomic Sanderson electronegativities	3D-MoRSE descriptors	3
	1EEig12x	Eigenvalue 12 from edge adj. matrix weighted by edge degrees	edge adjacency indices	2
	1Mor16m	3D-MoRSE - signal 16 / weighted by atomic masses	3D-MoRSE descriptors	3
	2EEig12r	Eigenvalue 12 from edge adj. matrix weighted by resonance integrals	edge adjacency indices	2
	1Mor26p	3D-MoRSE - signal 26 / weighted by atomic polarizabilities	3D-MoRSE descriptors	3
	1Mor09u	3D-MoRSE - signal 09 / unweighted	3D-MoRSE descriptors	3
	1nConj	number of non-aromatic conjugated C(sp2)	functional group counts	1
	1Mor18m	3D-MoRSE - signal 18 / weighted by atomic masses	3D-MoRSE descriptors	3
	1R3u+	R maximal autocorrelation of lag 3 / unweighted	GETAWAY descriptors	3
	1B07[C-O]	presence/absence of C - O at topological distance 07	2D binary fingerprints	2
	1F02[C-C]	frequency of C - C at topological distance 02	2D frequency fingerprints	2
MOR203.1	6B06[C-C]	presence/absence of C - C at topological distance 06	2D binary fingerprints	2
	1H-049	H attached to C3(sp3)/C2(sp2)/C1(sp2)/C0(sp)	atom-centred fragments	1
	1Mor08p	3D-MoRSE - signal 08 / weighted by atomic polarizabilities	3D-MoRSE descriptors	3
	1Hbond acceptor	Number of Hydrogen bond acceptors	structural (Cerius2)	0
	1RDF075u	Radial Distribution Function - 7.5 / unweighted	RDF descriptors	3
	1R6v	R autocorrelation of lag 6 / weighted by atomic van der Waals volumes	GETAWAY descriptors	3

Table 3.15 Continued

	3nRCHO	number of aldehydes (aliphatic)	functional group counts	1
	2nRCOOH	number of carboxylic acids (aliphatic)	functional group counts	1
	1MATS4m	Moran autocorrelation - lag 4 / weighted by atomic masses	2D autocorrelations	2
	1G3s	3st component symmetry directional WHIM index / weighted by atomic electrotopological states	WHIM descriptors	3
	1D/Dr05	distance/detour ring index of order 5	topological descriptors	2
	1nACD	number of ketones (aromatic)	functional group counts	1
	1B10[C-C]	presence/absence of C - C at topological distance 10	2D binary fingerprints	2
	1MATS6m	Moran autocorrelation - lag 6 / weighted by atomic masses	2D autocorrelations	2
	1BEHm7	highest eigenvalue n. 7 of Burden matrix / weighted by atomic masses	Burden eigenvalues	2
	1MAXDN	maximal electrotopological negative variation	topological descriptors	2
	1L2s	2nd component size directional WHIM index / weighted by atomic electrotopological states	WHIM descriptors	3
	1BELE3	lowest eigenvalue n. 3 of Burden matrix / weighted by atomic Sanderson electronegativities	Burden eigenvalues	2
	1RDF095u	Radial Distribution Function - 9.5 / unweighted	RDF descriptors	3
MOR204.6				
	1nCconj	number of non-aromatic conjugated C(sp2)	functional group counts	1
	1S_dsCH	S_dsCH	atomtypes (cerius2)	1
	2HATS3m	leverage-weighted autocorrelation of lag 3 / weighted by atomic masses	GETAWAY descriptors	3
	1DISPv	d COMMAZ value / weighted by atomic van der Waals volumes	geometrical descriptors	3
	1R4u+	R maximal autocorrelation of lag 4 / unweighted	GETAWAY descriptors	2
	2D/Dr06	distance/detour ring index of order 6	topological descriptors	2
	1R6m	R autocorrelation of lag 6 / weighted by atomic masses	GETAWAY descriptors	3
	1Mor10m	3D-MORSE - signal 10 / weighted by atomic masses	3D-MORSE descriptors	3
	1EEig09r	Eigenvalue 09 from edge adj. matrix weighted by resonance integrals	edge adjacency indices	2
	1R7e	R autocorrelation of lag 7 / weighted by atomic Sanderson electronegativities	GETAWAY descriptors	3
	1R3u+	R maximal autocorrelation of lag 3 / unweighted	GETAWAY descriptors	3
	1R2v	R autocorrelation of lag 2 / weighted by atomic van der Waals volumes	GETAWAY descriptors	3
	1F05[C-O]	frequency of C - O at topological distance 05	2D frequency fingerprints	2
	1HATS7e	leverage-weighted autocorrelation of lag 7 / weighted by atomic Sanderson electronegativities	GETAWAY descriptors	3
MOR207.1				
	3C-025	R-CR-R	atom-centred fragments	1
	1F05[C-O]	frequency of C - O at topological distance 05	2D frequency fingerprints	2
	1R6v+	R maximal autocorrelation of lag 6 / weighted by atomic van der Waals volumes	GETAWAY descriptors	3
	2Ez+	2nd component accessibility directional WHIM index / weighted by atomic Sanderson electronegativities	WHIM descriptors	3
	1S_sacH	S_sacH	atomtypes (cerius2)	1
	2Rte+	R maximal index / weighted by atomic Sanderson electronegativities	GETAWAY descriptors	3
	1R8u+	R maximal autocorrelation of lag 8 / unweighted	GETAWAY descriptors	3
	1ATS2p	Broto-Moreau autocorrelation of a topological structure - lag 2 / weighted by atomic polarizabilities	2D autocorrelations	2
	1MATS5p	Moran autocorrelation - lag 5 / weighted by atomic polarizabilities	2D autocorrelations	2
	1RDF045v	Radial Distribution Function - 4.5 / weighted by atomic van der Waals volumes	RDF descriptors	3
	1P2p	2nd component shape directional WHIM index / weighted by atomic polarizabilities	WHIM descriptors	3
	1R7p+	R maximal autocorrelation of lag 7 / weighted by atomic polarizabilities	GETAWAY descriptors	3
	1R5m	R autocorrelation of lag 5 / weighted by atomic masses	GETAWAY descriptors	3
	1HOMA	Harmonic Oscillator Model of Aromaticity index	geometrical descriptors	3
	1R7v+	R maximal autocorrelation of lag 7 / weighted by atomic van der Waals volumes	GETAWAY descriptors	3
	1H1p	H autocorrelation of lag 1 / weighted by atomic polarizabilities	GETAWAY descriptors	3
MOR273.1				
	1P2s	2nd component shape directional WHIM index / weighted by atomic electrotopological states	WHIM descriptors	3
	2Jhetv	Balaban-type index from van der Waals weighted distance matrix	topological descriptors	2
	1GATS5m	Geary autocorrelation - lag 5 / weighted by atomic masses	2D autocorrelations	2
	1HATS1u	leverage-weighted autocorrelation of lag 1 / unweighted	GETAWAY descriptors	3
	1Mor20u	3D-MORSE - signal 20 / unweighted	3D-MORSE descriptors	3
	1JG16	mean topological charge index of order6	topological charge indices	2
	1R1e	R autocorrelation of lag 1 / weighted by atomic Sanderson electronegativities	WHIM descriptors	3
	1Mor32u	3D-MORSE - signal 32 / unweighted	3D-MORSE descriptors	3
	1EEig08r	Eigenvalue 08 from edge adj. matrix weighted by resonance integrals	edge adjacency indices	2
	1BEHm3	highest eigenvalue n. 3 of Burden matrix / weighted by atomic masses	Burden eigenvalues	2
MOR250.1				
	1Dv	D total accessibility index / weighted by atomic van der Waals volumes	WHIM descriptors	3
	6H-047	H attached to C1(sp3)/C0(sp2)	atom-centred fragments	1
	1BEHm5p	highest eigenvalue n. 5 of Burden matrix / weighted by atomic masses	Burden eigenvalues	2
	3C-006	CHRX	atom-centred fragments	1
	2MATS5p	Moran autocorrelation - lag 5 / weighted by atomic polarizabilities	2D autocorrelations	2
	1O-057	phenol / enol / carboxy OH	atom-centred fragments	1
	2D/Dr06	distance/detour ring index of order 6	topological descriptors	2
	2Mor24p	3D-MORSE - signal 24 / weighted by atomic polarizabilities	3D-MORSE descriptors	3
	3P1s	1st component shape directional WHIM index / weighted by atomic electrotopological states	WHIM descriptors	3
	1EEig07d	Eigenvalue 07 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2
	1R7p+	R maximal autocorrelation of lag 7 / weighted by atomic polarizabilities	GETAWAY descriptors	3
	3E1v	1st component accessibility directional WHIM index / weighted by atomic van der Waals volumes	WHIM descriptors	3
	2L2s	2nd component size directional WHIM index / weighted by atomic electrotopological states	WHIM descriptors	3
	2Mor10e	3D-MORSE - signal 10 / weighted by atomic Sanderson electronegativities	3D-MORSE descriptors	3
	1H1v	H autocorrelation of lag 1 / weighted by atomic van der Waals volumes	GETAWAY descriptors	3
	4R4m	R autocorrelation of lag 4 / weighted by atomic masses	GETAWAY descriptors	3
	1RDF035m	Radial Distribution Function - 3.5 / weighted by atomic masses	RDF descriptors	3
	1BELm1	lowest eigenvalue n. 1 of Burden matrix / weighted by atomic masses	Burden eigenvalues	2
	2HATS7m	leverage-weighted autocorrelation of lag 7 / weighted by atomic masses	GETAWAY descriptors	3
	1F04[C-O]	frequency of C - O at topological distance 04	2D frequency fingerprints	2
	1B04[O-O]	presence/absence of O - O at topological distance 04	2D binary fingerprints	2
	1B05[C-C]	presence/absence of C - C at topological distance 05	2D binary fingerprints	2
	2C-003	CHR3	atom-centred fragments	1
	3H1p	H autocorrelation of lag 1 / weighted by atomic polarizabilities	GETAWAY descriptors	3
	2Mor11m	3D-MORSE - signal 11 / weighted by atomic masses	3D-MORSE descriptors	3
	1EEig07x	Eigenvalue 07 from edge adj. matrix weighted by edge degrees	edge adjacency indices	2
	3GVWAI-80	Ghose-Viswanadhan-Wendoloski drug-like index at 80%	molecular properties	1
	2JG17	mean topological charge index of order7	topological charge indices	2
	3pPC07	molecular multiple path count of order 07	walk and path counts	2
	1nR09	number of 9-membered rings	constitutional descriptors	0
	1RDF060p	Radial Distribution Function - 6.0 / weighted by atomic polarizabilities	RDF descriptors	3
	1nRCOOH	number of carboxylic acids (aliphatic)	functional group counts	1
	1Mor10m	3D-MORSE - signal 10 / weighted by atomic masses	3D-MORSE descriptors	3
	1GATS2m	Geary autocorrelation - lag 2 / weighted by atomic masses	2D autocorrelations	2
	2BELE1	lowest eigenvalue n. 1 of Burden matrix / weighted by atomic Sanderson electronegativities	Burden eigenvalues	2
	1ICI	information content index (neighborhood symmetry of 1-order)	information indices	2
	1MATS4m	Moran autocorrelation - lag 4 / weighted by atomic masses	2D autocorrelations	2
	1nAR	number of ethers (aromatic)	functional group counts	1
	1D/Dr05	distance/detour ring index of order 5	topological descriptors	2
	1RDF060u	Radial Distribution Function - 6.0 / unweighted	RDF descriptors	3
MOR256.17				
	7BIC	BIC	topological (cerius2)	2
	1EEig01x	Eigenvalue 01 from edge adj. matrix weighted by edge degrees	edge adjacency indices	2
	6HATS6m	leverage-weighted autocorrelation of lag 6 / weighted by atomic masses	GETAWAY descriptors	3
	5nOHp	number of primary alcohols	functional group counts	1
	3S_sscCH	S_sscCH	atomtypes (cerius2)	1
	3nARCO	number of ketones (aromatic)	functional group counts	1
	3H-049	H attached to C3(sp3)/C2(sp2)/C3(sp2)/C3(sp)	atom-centred fragments	1
	3nR=Ct	number of aliphatic tertiary C(sp2)	functional group counts	1
	3HTm	H total index / weighted by atomic masses	GETAWAY descriptors	3
	1Mor30v	3D-MORSE - signal 30 / weighted by atomic van der Waals volumes	3D-MORSE descriptors	3
	2GATS6m	Geary autocorrelation - lag 6 / weighted by atomic masses	2D autocorrelations	2
	1RDF130m	Radial Distribution Function - 13.0 / weighted by atomic masses	RDF descriptors	3
	2RDF040e	Radial Distribution Function - 4.0 / weighted by atomic Sanderson electronegativities	RDF descriptors	3
	1nR10	number of 10-membered rings	constitutional descriptors	0
	1B09[C-S]	presence/absence of C - S at topological distance 09	2D binary fingerprints	2

Table 3.15 Continued

	1 C-008	CHR2X	atom-centred fragments	1
	1 XSA	average connectivity index chi-5	connectivity indices	2
	1 MATS7m	Moran autocorrelation - lag 7 / weighted by atomic masses	2D autocorrelations	2
	1 Mor24v	3D-MoRSE - signal 24 / weighted by atomic van der Waals volumes	3D-MoRSE descriptors	3
	1 Mor27e	3D-MoRSE - signal 27 / weighted by atomic Sanderson electronegativities	3D-MoRSE descriptors	3
	1 EEig09d	Eigenvalue 09 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2
MOR258.1				
	3 R1v	R autocorrelation of lag 1 / weighted by atomic van der Waals volumes	GETAWAY descriptors	3
	1 BELm2	lowest eigenvalue n. 2 of Burden matrix / weighted by atomic masses	Burden eigenvalues	2
	2 H1p	H autocorrelation of lag 1 / weighted by atomic polarizabilities	GETAWAY descriptors	3
	1 E1u	1st component accessibility directional WHIM index / unweighted	WHIM descriptors	3
	1 HATS2m	leverage-weighted autocorrelation of lag 2 / weighted by atomic masses	GETAWAY descriptors	3
	1 JG12	mean topological charge index of order2	topological charge indices	2
	1 R1p	R autocorrelation of lag 1 / weighted by atomic polarizabilities	GETAWAY descriptors	3
	1 R06[C-C]	presence/absence of C - C at topological distance 06	2D binary fingerprints	2
	1 L2s	2nd component size directional WHIM index / weighted by atomic electrotopological states	WHIM descriptors	3
	1 DISPm	d COMMA2 value / weighted by atomic masses	geometrical descriptors	3
	1 SHP2	average shape profile index of order 2	Randic molecular profiles	3
	1 H1u	H autocorrelation of lag 1 / unweighted	GETAWAY descriptors	3
	1 P2s	2nd component shape directional WHIM index / weighted by atomic electrotopological states	WHIM descriptors	3
MOR259.1				
	1 nCb-	number of substituted benzene C(sp2)	functional group counts	1
	1 G3s	3st component symmetry directional WHIM index / weighted by atomic electrotopological states	WHIM descriptors	3
	1 R4m	R autocorrelation of lag 4 / weighted by atomic masses	GETAWAY descriptors	3
	1 MPC06	molecular path count of order 06	walk and path counts	2
	1 DISPV	d COMMA2 value / weighted by atomic van der Waals volumes	geometrical descriptors	3
	1 R1p	R autocorrelation of lag 1 / weighted by atomic polarizabilities	GETAWAY descriptors	3
	1 F05[C-O]	frequency of C - O at topological distance 05	2D frequency fingerprints	2
	1 R2v+	R maximal autocorrelation of lag 2 / weighted by atomic van der Waals volumes	GETAWAY descriptors	3
	1 E1u	1st component accessibility directional WHIM index / unweighted	WHIM descriptors	3
	1 GATS2m	Geary autocorrelation - lag 2 / weighted by atomic masses	2D autocorrelations	2
	1 E2e	2nd component accessibility directional WHIM index / weighted by atomic Sanderson electronegativities	WHIM descriptors	3
	1 H1u	H autocorrelation of lag 1 / unweighted	GETAWAY descriptors	3
MOR260.1				
	1 RDF110v	Radial Distribution Function - 11.0 / weighted by atomic van der Waals volumes	RDF descriptors	3
	2 R2m+	R maximal autocorrelation of lag 2 / weighted by atomic masses	GETAWAY descriptors	3
	1 JG18	mean topological charge index of order8	topological charge indices	2
	1 Hy	hydrophilic factor	molecular properties	1
	1 O-056	alcohol	atom-centred fragments	1
	1 H-050	H attached to heteroatom	atom-centred fragments	1
	1 MATS8m	Moran autocorrelation - lag 8 / weighted by atomic masses	2D autocorrelations	2
	1 RDF080e	Radial Distribution Function - 8.0 / weighted by atomic Sanderson electronegativities	RDF descriptors	3
	1 HATS6m	leverage-weighted autocorrelation of lag 6 / weighted by atomic masses	GETAWAY descriptors	3
MOR261.1				
	1 RDF110e	Radial Distribution Function - 11.0 / weighted by atomic Sanderson electronegativities	RDF descriptors	3
	2 C-006	CHR2X	atom-centred fragments	1
	1 RDF075v	Radial Distribution Function - 7.5 / weighted by atomic van der Waals volumes	RDF descriptors	3
	1 RDF095u	Radial Distribution Function - 9.5 / unweighted	RDF descriptors	3
	1 DISPm	d COMMA2 value / weighted by atomic masses	geometrical descriptors	3
	1 C-001	CH3R / CH4	atom-centred fragments	1
	1 BEHm8	highest eigenvalue n. 8 of Burden matrix / weighted by atomic masses	Burden eigenvalues	2
	1 RDF110v	Radial Distribution Function - 11.0 / weighted by atomic van der Waals volumes	RDF descriptors	3
MOR268.1				
	2 ATS8e	Broto-Moreau autocorrelation of a topological structure - lag 8 / weighted by atomic Sanderson electronegati	2D autocorrelations	2
	1 C-006	CHR2X	atom-centred fragments	1
	1 TIE	E-state topological parameter	topological descriptors	2
	1 RDF050m	Radial Distribution Function - 5.0 / weighted by atomic masses	RDF descriptors	3
	3 B09[C-O]	presence/absence of C - O at topological distance 09	2D binary fingerprints	2
	5 B07[C-C]	presence/absence of C - C at topological distance 07	2D binary fingerprints	2
	4 H-051	H attached to alpha-C	atom-centred fragments	1
	3 S_ sssCH	S_ sssCH	atomtypes (cerius2)	1
	1 MATS6m	Moran autocorrelation - lag 6 / weighted by atomic masses	2D autocorrelations	2
	2 G2p	2st component symmetry directional WHIM index / weighted by atomic polarizabilities	WHIM descriptors	3
	1 RDF115e	Radial Distribution Function - 11.5 / weighted by atomic Sanderson electronegativities	RDF descriptors	3
	1 H-049	H attached to C3(sp3)/C2(sp2)/C3(sp2)/C3(sp)	atom-centred fragments	1
	1 O-057	phenol / enol / carboxyl OH	atom-centred fragments	1
	1 RDF130u	Radial Distribution Function - 13.0 / unweighted	RDF descriptors	3
	1 G3u	3st component symmetry directional WHIM index / unweighted	WHIM descriptors	3
	1 nRCHO	number of aldehydes (aliphatic)	functional group counts	1
	1 nRKO	number of ketones (aromatic)	functional group counts	1
	1 nRCOOH	number of carboxylic acids (aliphatic)	functional group counts	1
	1 GATS1m	Geary autocorrelation - lag 1 / weighted by atomic masses	2D autocorrelations	2
	1 nCconj	number of non-aromatic conjugated C(sp2)	functional group counts	1
	1 MATS5p	Moran autocorrelation - lag 5 / weighted by atomic polarizabilities	2D autocorrelations	2
MOR271.1				
	1 H-050	H attached to heteroatom	atom-centred fragments	1
	1 PHI	Kier flexibility index	topological descriptors	2
	3 JG16	mean topological charge index of order6	topological charge indices	2
	2 RDF075m	Radial Distribution Function - 7.5 / weighted by atomic masses	RDF descriptors	3
	2 H-049	H attached to C3(sp3)/C2(sp2)/C3(sp2)/C3(sp)	atom-centred fragments	1
	1 SPH	sphericity	geometrical descriptors	3
	2 Hy	hydrophilic factor	molecular properties	1
	1 EEig09x	Eigenvalue 09 from edge adj. matrix weighted by edge degrees	edge adjacency indices	2
	1 RTU+	R maximal index / unweighted	GETAWAY descriptors	3
	1 GATS8m	Geary autocorrelation - lag 8 / weighted by atomic masses	2D autocorrelations	2
	1 nRCOOH	number of carboxylic acids (aliphatic)	functional group counts	1
	1 RDF110p	Radial Distribution Function - 11.0 / weighted by atomic polarizabilities	RDF descriptors	3
	1 GATS4m	Geary autocorrelation - lag 4 / weighted by atomic masses	2D autocorrelations	2
	1 MATS5p	Moran autocorrelation - lag 5 / weighted by atomic polarizabilities	2D autocorrelations	2
	1 R8m+	R maximal autocorrelation of lag 8 / weighted by atomic masses	GETAWAY descriptors	3
MOR272.1				
	1 BLTF96	Verhaar model of Fish base-line toxicity from MLOGP (mmol/l)	molecular properties	1
	1 DISP8e	d COMMA2 value / weighted by atomic Sanderson electronegativities	geometrical descriptors	3
	1 B08[C-O]	presence/absence of C - O at topological distance 08	2D binary fingerprints	2
	1 HATS6p	leverage-weighted autocorrelation of lag 6 / weighted by atomic polarizabilities	GETAWAY descriptors	3
	1 MATS6m	Moran autocorrelation - lag 6 / weighted by atomic masses	2D autocorrelations	2
	1 RDF110v	Radial Distribution Function - 11.0 / weighted by atomic van der Waals volumes	RDF descriptors	3
	1 G3s	3st component symmetry directional WHIM index / weighted by atomic electrotopological states	WHIM descriptors	3
	1 H-049	H attached to C3(sp3)/C2(sp2)/C3(sp2)/C3(sp)	atom-centred fragments	1
	1 H-050	H attached to heteroatom	atom-centred fragments	1
	1 R2m+	R maximal autocorrelation of lag 2 / weighted by atomic masses	GETAWAY descriptors	3
	1 GATS8m	Geary autocorrelation - lag 8 / weighted by atomic masses	2D autocorrelations	2
	1 HATS5e	leverage-weighted autocorrelation of lag 5 / weighted by atomic Sanderson electronegativities	GETAWAY descriptors	3
	1 H3m	H autocorrelation of lag 3 / weighted by atomic masses	GETAWAY descriptors	3
	1 BEHm3	highest eigenvalue n. 3 of Burden matrix / weighted by atomic masses	Burden eigenvalues	2
MOR273.1				
	1 P2s	2nd component shape directional WHIM index / weighted by atomic electrotopological states	WHIM descriptors	3
	2 Jhetv	Balaban-type index from van der Waals weighted distance matrix	topological descriptors	2
	1 GATS5m	Geary autocorrelation - lag 5 / weighted by atomic masses	2D autocorrelations	2
	1 HATS1u	leverage-weighted autocorrelation of lag 1 / unweighted	GETAWAY descriptors	3

Table 3.15 Continued

	1 Mor2Du	3D-MORSE - signal 20 / unweighted	3D-MORSE descriptors	3
	1 JG16	mean topological charge index of order6	topological charge indices	2
	1 R1e	R autocorrelation of lag 1 / weighted by atomic Sanderson electronegativities	GETAWAY descriptors	3
	1 Mor32u	3D-MORSE - signal 32 / unweighted	3D-MORSE descriptors	3
	1 EEig08r	Eigenvalue 08 from edge adj. matrix weighted by resonance integrals	edge adjacency indices	2
	1 BEHm3	highest eigenvalue n. 3 of Burden matrix / weighted by atomic masses	Burden eigenvalues	2
MOR277.1				
	1 HATS3v	leverage-weighted autocorrelation of lag 3 / weighted by atomic van der Waals volumes	GETAWAY descriptors	3
	4 E1e	1st component accessibility directional WHIM index / weighted by atomic Sanderson electronegativities	WHIM descriptors	3
	1 HATS5m	Moran autocorrelation - lag 5 / weighted by atomic masses	2D autocorrelations	2
	1 E2p	2nd component accessibility directional WHIM index / weighted by atomic polarizabilities	WHIM descriptors	3
	2 RDF035m	Radial Distribution Function - 3.5 / weighted by atomic masses	RDF descriptors	3
	2 R2m+	R maximal autocorrelation of lag 2 / weighted by atomic masses	GETAWAY descriptors	3
	1 R8e+	R maximal autocorrelation of lag 8 / weighted by atomic Sanderson electronegativities	GETAWAY descriptors	3
	3 SHP2	average shape profile index of order 2	Randic molecular profiles	3
	1 GATS6m	Geary autocorrelation - lag 6 / weighted by atomic masses	2D autocorrelations	2
	1 BELp3	lowest eigenvalue n. 3 of Burden matrix / weighted by atomic polarizabilities	Burden eigenvalues	2
	1 Jhetp	Balaban-type index from polarizability weighted distance matrix	topological descriptors	2
	1 H-051	H attached to alpha-C	atom-centred fragments	1
	1 EEig09r	Eigenvalue 09 from edge adj. matrix weighted by resonance integrals	edge adjacency indices	2
	1 BELe3	lowest eigenvalue n. 3 of Burden matrix / weighted by atomic Sanderson electronegativities	Burden eigenvalues	2
	1 nRCOOR	number of esters (aliphatic)	functional group counts	1
	1 E2u	2nd component accessibility directional WHIM index / unweighted	WHIM descriptors	3
	1 JG15	mean topological charge index of order5	topological charge indices	2
	1 GATS3m	Geary autocorrelation - lag 3 / weighted by atomic masses	2D autocorrelations	2
	1 L2u	2nd component size directional WHIM index / unweighted	WHIM descriptors	3
	1 S_dssC	S_dssC	atomtypes (cerius2)	1
	1 G3e	3st component symmetry directional WHIM index / weighted by atomic Sanderson electronegativities	WHIM descriptors	3
	1 B06[C-C]	presence/absence of C - C at topological distance 06	2D binary fingerprints	2
	1 L2p	2nd component size directional WHIM index / weighted by atomic polarizabilities	WHIM descriptors	3
MOR30.1				
	1 B09[C-O]	presence/absence of C - O at topological distance 09	2D binary fingerprints	2
	1 O-056	alcohol	atom-centred fragments	1
	2 B07[C-O]	presence/absence of C - O at topological distance 07	2D binary fingerprints	2
	1 RDF085v	Radial Distribution Function - 8.5 / weighted by atomic van der Waals volumes	RDF descriptors	3
	5 H-051	H attached to alpha-C	atom-centred fragments	1
	1 RDF075m	Radial Distribution Function - 7.5 / weighted by atomic masses	RDF descriptors	3
	3 C-006	CH2RX	atom-centred fragments	1
	1 Lbw	length-to-breadth ratio by WHIM	geometrical descriptors	3
	1 EEig07d	Eigenvalue 07 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2
	1 RDF095v	Radial Distribution Function - 9.5 / weighted by atomic van der Waals volumes	RDF descriptors	3
	1 E2p	2nd component accessibility directional WHIM index / weighted by atomic polarizabilities	WHIM descriptors	3
	2 GATS7m	Geary autocorrelation - lag 7 / weighted by atomic masses	2D autocorrelations	2
	1 Mor30m	3D-MORSE - signal 30 / weighted by atomic masses	3D-MORSE descriptors	3
	1 Mor09u	3D-MORSE - signal 09 / unweighted	3D-MORSE descriptors	3
	1 EEig06d	Eigenvalue 06 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2
	1 GATS8m	Geary autocorrelation - lag 8 / weighted by atomic masses	2D autocorrelations	2
	1 Infective-80	Ghose-Viswanadhan-Wendolowski antiinfective-like index at 80%	molecular properties	1
	1 R6e+	R maximal autocorrelation of lag 6 / weighted by atomic Sanderson electronegativities	GETAWAY descriptors	3
	1 GATS2m	Geary autocorrelation - lag 2 / weighted by atomic masses	2D autocorrelations	2
	1 RDF045e	Radial Distribution Function - 4.5 / weighted by atomic Sanderson electronegativities	RDF descriptors	3
MOR33.1				
	1 F07[C-O]	frequency of C - O at topological distance 07	2D frequency fingerprints	2
	1 nRCOOH	number of carboxylic acids (aliphatic)	functional group counts	1
	1 DISPv	d COMMA2 value / weighted by atomic van der Waals volumes	geometrical descriptors	3
	1 F04[C-O]	frequency of C - O at topological distance 04	2D frequency fingerprints	2
	2 Mor17e	3D-MORSE - signal 17 / weighted by atomic Sanderson electronegativities	3D-MORSE descriptors	3
	1 H-051	H attached to alpha-C	atom-centred fragments	1
	1 C-006	CH2RX	atom-centred fragments	1
	1 GATS7m	Geary autocorrelation - lag 7 / weighted by atomic masses	2D autocorrelations	2
	1 H-050	H attached to heteroatom	atom-centred fragments	1
MOR37.1				
	1 B09[C-O]	presence/absence of C - O at topological distance 09	2D binary fingerprints	2
	1 O-056	alcohol	atom-centred fragments	1
	1 B08[C-O]	presence/absence of C - O at topological distance 08	2D binary fingerprints	2
MOR40.1				
	2 RDF110v	Radial Distribution Function - 11.0 / weighted by atomic van der Waals volumes	RDF descriptors	3
	1 H-051	H attached to alpha-C	atom-centred fragments	1
	1 Mor17e	3D-MORSE - signal 17 / weighted by atomic Sanderson electronegativities	3D-MORSE descriptors	3
	1 C-006	CH2RX	atom-centred fragments	1
	1 GATS7m	Geary autocorrelation - lag 7 / weighted by atomic masses	2D autocorrelations	2
MOR41.1				
	1 HVcpx	graph vertex complexity index	information indices	2
	1 Kappa-3	Kappa-3	topological (cerius2)	2
	4 GATS3m	Geary autocorrelation - lag 3 / weighted by atomic masses	2D autocorrelations	2
	8 nRCOOH	number of carboxylic acids (aliphatic)	functional group counts	1
	3 H-049	H attached to C3(sp3)/C2(sp2)/C3(sp2)/C3(sp)	atom-centred fragments	1
	2 Mor10v	3D-MORSE - signal 10 / weighted by atomic van der Waals volumes	3D-MORSE descriptors	3
	2 RDF085p	Radial Distribution Function - 8.5 / weighted by atomic polarizabilities	RDF descriptors	3
	3 EEig08r	Eigenvalue 08 from edge adj. matrix weighted by resonance integrals	edge adjacency indices	2
	4 BIC	BIC	topological (cerius2)	2
	5 GATS8m	Geary autocorrelation - lag 8 / weighted by atomic masses	2D autocorrelations	2
	3 H6m	H autocorrelation of lag 6 / weighted by atomic masses	GETAWAY descriptors	3
	1 BEHv3	highest eigenvalue n. 3 of Burden matrix / weighted by atomic polarizabilities	Burden eigenvalues	2
	1 JG13	mean topological charge index of order3	topological charge indices	2
	4 RDF085u	Radial Distribution Function - 8.5 / unweighted	RDF descriptors	3
	1 GATS2m	Geary autocorrelation - lag 2 / weighted by atomic masses	2D autocorrelations	2
	4 Mor22e	3D-MORSE - signal 22 / weighted by atomic Sanderson electronegativities	3D-MORSE descriptors	3
	3 nC4	number of total quaternary C(sp3)	functional group counts	1
	1 E1p	1st component accessibility directional WHIM index / weighted by atomic polarizabilities	WHIM descriptors	3
	1 RDF050v	Radial Distribution Function - 5.0 / weighted by atomic van der Waals volumes	RDF descriptors	3
	2 PW3	path/walk 3 - Randic shape index	topological descriptors	2
	1 GVWAI-80	Ghose-Viswanadhan-Wendolowski drug-like index at 80%	molecular properties	1
	1 C-008	CH2RX	atom-centred fragments	1
	1 H6v	H autocorrelation of lag 6 / weighted by atomic van der Waals volumes	GETAWAY descriptors	3
	1 MATS3m	Moran autocorrelation - lag 3 / weighted by atomic masses	2D autocorrelations	2
	1 Atype_C_18	Number of Carbon Type 18	atomtypes (Cerius2)	1
	1 BEHv3	highest eigenvalue n. 3 of Burden matrix / weighted by atomic van der Waals volumes	Burden eigenvalues	2
	1 RDF130u	Radial Distribution Function - 13.0 / unweighted	RDF descriptors	3
	1 nRCHO	number of aldehydes (aliphatic)	functional group counts	1
	1 B05[C-O]	presence/absence of C - O at topological distance 05	2D binary fingerprints	2
	1 MEcc	molecular eccentricity	geometrical descriptors	3
	1 XZA	average connectivity index chi-2	connectivity indices	2
MOR5.1				
	1 nRCOOH	number of carboxylic acids (aliphatic)	functional group counts	1
	1 F07[C-O]	frequency of C - O at topological distance 07	2D frequency fingerprints	2
	1 F04[C-O]	frequency of C - O at topological distance 04	2D frequency fingerprints	2
	1 O-057	phenol / enol / carboxyl OH	atom-centred fragments	1
	1 Mor17e	3D-MORSE - signal 17 / weighted by atomic Sanderson electronegativities	3D-MORSE descriptors	3
OR1A1				

Table 3.15 Continued

	2 HATS8u	leverage-weighted autocorrelation of lag 8 / unweighted	GETAWAY descriptors	3
	1 E1s	1st component accessibility directional WHIM index / weighted by atomic electrotopological states	WHIM descriptors	3
	2 RSV+	R maximal autocorrelation of lag 5 / weighted by atomic van der Waals volumes	GETAWAY descriptors	3
	1 Mor02p	3D-MorSE - signal 02 / weighted by atomic polarizabilities	3D-MorSE descriptors	3
	1 DISPe	d COMMA2 value / weighted by atomic Sanderson electronegativities	geometrical descriptors	3
	3 B06[C-C]	presence/absence of C - C at topological distance 06	2D binary fingerprints	2
	2 nRCHO	number of aldehydes (aliphatic)	functional group counts	1
	1 Mor21p	3D-MorSE - signal 21 / weighted by atomic polarizabilities	3D-MorSE descriptors	3
	1 JG13	mean topological charge index of order3	topological charge indices	2
	1 O-057	phenol / enol / carboxyl OH	atom-centred fragments	1
	1 RDF100m	Radial Distribution Function - 10.0 / weighted by atomic masses	RDF descriptors	3
	1 MATS5m	Moran autocorrelation - lag 5 / weighted by atomic masses	2D autocorrelations	2
	1 R6m+	R maximal autocorrelation of lag 6 / weighted by atomic masses	GETAWAY descriptors	3
	1 RDF095u	Radial Distribution Function - 9.5 / unweighted	RDF descriptors	3
	1 nRCOOH	number of carboxylic acids (aliphatic)	functional group counts	1
	1 ESpm06d	Spectral moment 06 from edge adj. matrix weighted by dipole moments	edge adjacency indices	2
	1 BELm5	lowest eigenvalue n. 5 of Burden matrix / weighted by atomic masses	Burden eigenvalues	2
OR2J2				
	3 E1s	1st component accessibility directional WHIM index / weighted by atomic electrotopological states	WHIM descriptors	3
	1 F06[C-O]	frequency of C - O at topological distance 06	2D frequency fingerprints	2
	1 DISPv	d COMMA2 value / weighted by atomic van der Waals volumes	geometrical descriptors	3
	1 JG15	mean topological charge index of order5	topological charge indices	2
	3 H-051	H attached to alpha-C	atom-centred fragments	1
	1 H1u	H autocorrelation of lag 1 / unweighted	GETAWAY descriptors	3
	4 DISPm	d COMMA2 value / weighted by atomic masses	geometrical descriptors	3
	1 BEHv8	highest eigenvalue n. 8 of Burden matrix / weighted by atomic van der Waals volumes	Burden eigenvalues	2
	2 H-049	H attached to C3(sp3)/C2(sp2)/C3(sp2)/C3(sp)	atom-centred fragments	1
	1 B06[C-O]	presence/absence of C - O at topological distance 06	2D binary fingerprints	2
	1 O-057	phenol / enol / carboxyl OH	atom-centred fragments	1
	1 Mor10u	3D-MorSE - signal 10 / unweighted	3D-MorSE descriptors	3
	1 H8m	H autocorrelation of lag 8 / weighted by atomic masses	GETAWAY descriptors	3
	2 F01[C-C]	frequency of C - C at topological distance 01	2D frequency fingerprints	2
	2 HATS5p	leverage-weighted autocorrelation of lag 5 / weighted by atomic polarizabilities	GETAWAY descriptors	3
	1 MATS5p	Moran autocorrelation - lag 5 / weighted by atomic polarizabilities	2D autocorrelations	2
	1 H6m	H autocorrelation of lag 6 / weighted by atomic masses	GETAWAY descriptors	3
	1 S_dssc	S_dssc	atomtypes (cerius2)	1
	1 HATS0e	leverage-weighted autocorrelation of lag 0 / weighted by atomic Sanderson electronegativities	GETAWAY descriptors	3
	1 R7p+	R maximal autocorrelation of lag 7 / weighted by atomic polarizabilities	GETAWAY descriptors	3
	1 VEA2	average eigenvector coefficient sum from adjacency matrix	eigenvector-based indices	2
OR2W1				
	2 B05[C-C]	presence/absence of C - C at topological distance 05	2D binary fingerprints	2
	2 HATS5p	leverage-weighted autocorrelation of lag 5 / weighted by atomic polarizabilities	GETAWAY descriptors	3
	1 RDF110m	Radial Distribution Function - 11.0 / weighted by atomic masses	RDF descriptors	3
	1 R2m+	R maximal autocorrelation of lag 2 / weighted by atomic masses	GETAWAY descriptors	3
	1 MATS5p	Moran autocorrelation - lag 5 / weighted by atomic polarizabilities	2D autocorrelations	2
OR5P3				
	2 nCconj	number of non-aromatic conjugated C(sp2)	functional group counts	1
	1 HATS3m	leverage-weighted autocorrelation of lag 3 / weighted by atomic masses	GETAWAY descriptors	3
	3 R7e	R autocorrelation of lag 7 / weighted by atomic Sanderson electronegativities	GETAWAY descriptors	3
	1 Yindex	Balaban Y index	information indices	2
	1 R7e+	R maximal autocorrelation of lag 7 / weighted by atomic Sanderson electronegativities	GETAWAY descriptors	3
	2 BLTD48	Verhaar model of Daphnia base-line toxicity from MLOGP (mmol/l)	molecular properties	1
	1 R6m	R autocorrelation of lag 6 / weighted by atomic masses	GETAWAY descriptors	3
	1 H-047	H attached to C1(sp3)/C0(sp2)	atom-centred fragments	1
	1 R4m	R autocorrelation of lag 4 / weighted by atomic masses	GETAWAY descriptors	3
	3 GATS2m	Geary autocorrelation - lag 2 / weighted by atomic masses	2D autocorrelations	2
	1 DISPm	d COMMA2 value / weighted by atomic masses	geometrical descriptors	3
	2 Mor10u	3D-MorSE - signal 10 / unweighted	3D-MorSE descriptors	3
	1 RDF050p	Radial Distribution Function - 5.0 / weighted by atomic polarizabilities	RDF descriptors	3
	1 B04[O-O]	presence/absence of O - O at topological distance 04	2D binary fingerprints	2
	1 BELm1	lowest eigenvalue n. 1 of Burden matrix / weighted by atomic masses	Burden eigenvalues	2
	1 Mor11v	3D-MorSE - signal 11 / weighted by atomic van der Waals volumes	3D-MorSE descriptors	3
	2 HATS7e	leverage-weighted autocorrelation of lag 7 / weighted by atomic Sanderson electronegativities	GETAWAY descriptors	3
	1 B05[C-C]	presence/absence of C - C at topological distance 05	2D binary fingerprints	2
	1 H-049	H attached to C3(sp3)/C2(sp2)/C3(sp2)/C3(sp)	atom-centred fragments	1
	2 B07[C-C]	presence/absence of C - C at topological distance 07	2D binary fingerprints	2
	1 nRCOOH	number of carboxylic acids (aliphatic)	functional group counts	1
	1 Mor18p	3D-MorSE - signal 18 / weighted by atomic polarizabilities	3D-MorSE descriptors	3
	1 HATS2v	leverage-weighted autocorrelation of lag 2 / weighted by atomic van der Waals volumes	GETAWAY descriptors	3

Table 3.15 Continued

Table 3.16: Top 100 predicted compounds for each mammalian OR

Chemical name or Pubchem compound ID (CIDs), SMILES strings, and distances, of the top ~100 predicted compounds for each Or. All distances represent the minimum distance based on optimized descriptors to the previously known strongest active compound listed in the gray cells for that particular Or.

CHAPTER IV:

A New Generation of Safe DEET Substitutes that are Strong Olfactory and Gustatory Repellents of Mosquitoes and Flies

INTRODUCTION

Blood-feeding insects transmit deadly diseases such as malaria, Dengue, Filariasis, West Nile Fever, Yellow fever, Sleeping sickness and Leishmaniasis to hundreds of millions of people, causing untold suffering and more than a million deaths every year. Insect repellents can be very effective in reducing vectorial capacity by blocking the contact between blood-seeking insects and humans, however they are seldom used in disease-prone areas of Africa and Asia due to high costs and need for continuous application on skin.

N,N-Diethyl-*m*-toluamide (DEET) has remained the primary insect repellent used for more than 60 years in the developed world. DEET is a solvent, melting several forms of plastics, synthetic fabrics, painted and varnished surfaces (Krajick, 2006). Additionally, DEET has been shown to inhibit mammalian cation channels and human acetylcholinesterase, which is also inhibited by carbamate insecticides commonly used in disease endemic areas (Corbel et al., 2009). These concerns are enhanced by the requirement of direct and continuous application of DEET to every part of exposed skin in concentrations that can be as high as 30-100 percent. Several instances of increased resistance to DEET have also been reported in flies (Reeder et al., 2001), *Anopheles albimanus* (Klun et al., 2004), and *Aedes aegypti* (Stanczyk et al., 2010). Moreover, mosquito strains with resistance to pyrethroid insecticides, the main line of defense against mosquitoes in developing countries, are spreading (Butler, 2011). There is

therefore an urgent need to develop safer and improved alternatives to DEET. The other major barrier in developing new repellents is the time and cost of development. It has been suggested that >\$30M and several years may be required for identification (Gupta, 2007) and adequate human-safety analysis of new repellent chemistries.

A major limitation to finding effective substitutes is that the molecular targets in adult mosquitoes through which DEET causes repellence are unknown. Recent studies have put forward a few competing theories about mechanism of action, but demonstration of a causal relationship is lacking. Some reports show that pure DEET causes inhibition (Ditzen et al., 2008) or mild electrophysiological modulations of neural responses to weakly activating odors in *Drosophila* antennal olfactory neurons (Pellegrino et al., 2011), but whether such effects cause repellency in either *Drosophila* or mosquitoes are unknown. On the other hand a DEET-responsive receptor, Or42a, has been identified in the maxillary palps of *Drosophila*, however Or42a is also known to be a broad-spectrum detector for several food-associated odorants that elicit attraction (Syed et al., 2011). Other studies in *Drosophila* have shown that DEET is also detected through bitter taste neurons causing contact-avoidance (Lee et al., 2010; Weiss et al., 2011). Although the three proposed mechanisms confound appropriate target choice for high-throughput screening, *Drosophila* might provide a tractable model for discovery of new DEET substitutes via direct behavioral screening of compounds.

Here, we describe identification of novel repellent odors that are safe, inexpensive, and effective at repelling *D. melanogaster* and *A. aegypti*. We first demonstrate that *D. melanogaster* and *A. aegypti* may utilize their gustatory system and olfactory system to different extents to avoid DEET, with mosquitoes relying more heavily on olfactory pathways. We then use a novel *chemical informatics* approach that

uses supervised training from known repellents to identify important structural features that are responsible for avoidance behavior. Using these features we predict novel repellents from a very large untested odor space comprising a large purchasable odor set and a natural odor library, the latter providing many chemicals that are safe and already approved for human use by both United States and European governmental safety organizations. We select four odors from the 150 natural odor library predictions and demonstrate that all four are able to cause strong avoidance in both *D. melanogaster* and *A. aegypti*. Due to the large number of candidates we are able to select new repellents with ideal properties; safe for human consumption, do not dissolve plastics and mild and pleasant aroma. These results suggest that our integrative approach of computational predictions and behavioral analysis can revolutionize the discovery of safe and effective repellent odors that could be very useful in our struggle against the increasing spread of insect disease vector species. Although it has been several years since odor receptors were identified in vertebrates and invertebrates, very few odor receptor targets have been identified for known behavior modifying odorants. This study therefore has broader implications since the approach presented can be used for identification of improved behavior modifying odorants for any organism, even if the target odor receptor is unknown. Moreover upon identification of target receptors, the same methodology can be easily adapted for a receptor-activity based approach.

RESULTS AND DISCUSSION

Contributions of Olfaction and Gustation in DEET Avoidance in *Drosophila*

In order to select a reliable behavior assay, we first tested for repellency in the well-established T-maze. We found that *Drosophila* do not avoid the arm containing DEET (10%) in this assay (Figure 4.1A), perhaps due to its low volatility. By introducing a well known repellent CO₂ (0.66%) in the other arm, we asked whether we could observe an effect by forcing flies to enter the DEET arm closer to the DEET source. We find a marginal, but not significant, reduction in avoidance to CO₂, suggesting that DEET is not an effective olfactory repellent for *Drosophila* in this assay (Figure 4.1A).

Drosophila do however show strong avoidance when they are allowed to come in close contact to DEET in a plate-based 2-choice trap assay as described previously (Syed et al., 2011), where they can sense DEET-treated filter paper positioned at the entrance of a trap with both the olfactory and gustatory system (Figure 4.1B). *Orco* mutant *Drosophila* that lack functional odor receptors belonging to the *Or* family show ~30% reduction in avoidance to DEET (Figure 4.1B). The simplest interpretation of these results is that although olfactory receptor pathways participate in the avoidance behavior, the response mediated by the gustatory system is sufficient to generate strong avoidance of *Drosophila* to DEET.

Aedes aegypti Detect and Avoid DEET Primarily using Olfaction

To test whether the gustatory system of mosquitoes plays a substantial role in DEET repellency as well, we developed a modified hand-in-glove assay (described in Supplementary Methods) (Figures 4.1C, 4.1D, and 4.2). The assay allows us to

quantitatively analyze the repellent effects of DEET on mosquitoes attracted to a human arm, without being able to bite. Female *Aedes aegypti* mosquitoes show an equally strong avoidance behavior to DEET in both the contact and non-contact versions of the assay (Figure 4.1D). For rare landings, the time spent on the net before escape is marginally lower, but not significantly different, when direct contact with DEET was permitted (Figure 4.1E). In order to test whether the reduction in attraction is caused by a direct detection or a masking/scrambling of skin-odor detection, we measure attraction of female *Aedes* to humidity and warmth and demonstrate a significant dose-dependent reduction in presence of DEET (Figure 4.1F). These results indicate that the female *Aedes* primarily use the olfactory system to sense DEET directly and avoid it. This is consistent with previous observations that *Culex quinquefasciatus* (Syed and Leal, 2008) and *A. aegypti* (Turner et al., 2011) avoid DEET presented as a volatile stimulus directly. In fact a DEET-sensitive neuron type has been identified in *C. quinquefasciatus* (Syed and Leal, 2008) and *A. aegypti* (Stanczyk et al., 2010), however it is not known whether these neurons are responsible for repellency, or which odor receptors they express. A broadly tuned larval odor receptor AgOr40 has been shown to respond to DEET (Liu et al., 2010; Xia et al., 2008), however its role in avoidance behavior in adults has not been determined.

A Computer can be trained to Predict Repellent Behavior from chemical structure

Since behavior assays in *Drosophila* and mosquitoes appear to differ in requirements for gustatory and olfactory exposure for repellency (Figure 4.1G), we would require use of both systems to reliably identify improved DEET substitutes. However volatile chemical space that can be exploited to find such improved DEET

substitutes is vast (potentially >400,000) thus using behavior assays alone is unfeasible from the perspective of time and cost of chemical purchase. Moreover, since a protein target for DEET action is unknown, high-throughput screening, or sophisticated computational protein-ligand docking based approaches to identify new ligands are also not possible.

To circumvent these problems so as to enable selection of compounds for behavior assays, we developed a novel chemical informatics approach. We hypothesized that the unknown target protein is recognizing specific structural features of DEET and other known structurally related repellents. By identifying structural features that are shared amongst DEET and the known repellent compounds, we can utilize them to rapidly screen an extremely large number of compounds *in-silico* to identify novel repellents, thus greatly reducing both the cost and time required. We assembled a training set of known repellents that included: the two approved ones DEET (a carboxamide) and picaridin; and also 34 N-acyl piperidines (Katritzky et al., 2008) that are structurally related to picaridin; eucalyptol, linalool, alpha-thujone, beta-thujone (Kline et al., 2003; Klocke et al., 1987; Syed and Leal, 2008) and a structurally diverse panels of odors that are not expected to elicit repellency via similar target receptors (Carey et al., 2010; Hallem and Carlson, 2006). The study where the 34 n-acyl piperidines were identified also showed that a chemical-structure-based approach could be successfully applied to predicting repellency (Katritzky et al., 2008). For our analysis, compounds from different sources were approximated into a single metric of “protection duration” as a rough indicator of repellency. The non-repellent diversifying training set of odors were assigned protection times of zero, while the approved repellents DEET and Picaridin were assigned the highest value since these would have the properties

important for regulatory approval. Compounds were clustered using Euclidean distance and hierarchical clustering based on differences in repellency values, and a set of 5 compounds with the highest activity that clustered together was classified as “training repellents”.

We expect that only specific structural features of the repellent odors will interact with target proteins to elicit repellent responses, and not the entire molecule. We assumed that identification of structural features that are shared across repellent odors would enable a search for these features within a large chemical space, potentially identifying novel repellents (Maldonado et al., 2006). We decided to focus on a descriptor-based computational approach that is effective in structure analysis and is efficient in subsequently searching a large chemical space rapidly. We calculated mathematical values for 3,242 molecular descriptors, that describe the 3-D structure of a chemical, for our 201 compound training set and using a Sequential-Forward-Selection method (Haddad et al., 2008) we incrementally identified a unique subset of 18 descriptors that were highly correlated with repellency (correlation of 0.912)(Figures 4.3A, 4.3B, and 4.3C). As expected the repellent odors cluster together in the training set if the optimized descriptor subset is used to calculate Euclidean distance amongst them (Figure 4.3B). These 18 molecular descriptors represent a collection of predominately 2D and 3D descriptor types. Inspection suggests that 6 member rings, carbon-nitrogen distances, tertiary amides, and oxygen placement are prominent in the optimized subset. Interestingly, although repellents are topically applied chemicals, the Ghose-Viswanadhan-Wendoloski drug-like index is selected, which is an aggregate descriptor that usually suggests similarity of chemical features important for drugs (Ghose et al., 1999).

In order to improve the predictive ability of the chemical informatics approach we used the optimized descriptor set to train a Support Vector Machine (SVM), which is a well-known supervised learning approach (Cortes and Vapnik, 1995)(Figure 4.3A). To validate the predictive ability of our approach, we performed 5-fold cross-validation using SVMs on the training set. Each cross-validation run excluded ~20% of the repellents as a test-set, while the remaining repellents were used to train a SVM. We predicted repellency for the withheld test-set odors using the trained SVM. This operation was repeated 5 times, each trial performed excluding a different subset of the training set, and the whole process was repeated 20 times for consistency. A mean Receiver-Operating-Characteristic (ROC) analysis curve representing the prediction accuracy was generated, and the Area-Under-Curve (AUC) value was determined to be 0.994, indicating that the *in-silico* approach was extremely effective in predicting repellents from compounds that are excluded from the training set (Figure 4.3D).

High-Throughput *In Silico* Screen identifies Safe Natural Odors as repellents

We use the 18 optimized-descriptors and SVM method to screen *in-silico* a large chemical library consisting of >440,000 chemicals from a database called eMolecules of putative volatiles. Upon inspection, we find the top 1,000 predicted compounds (0.23% of hits) represent a structurally diverse group of chemicals that retain some structural features of the known repellents (Figures 4.4A and 4.4B). We computed logP values of the 1,000 compounds to identify ones that are predicted to be lipophilic (logP >4.5) thus allowing for selection of other compounds that are less likely to pass through the skin barrier in topical applications (Walker et al., 2003) (Figure 4.4B). In addition, we

computed the predicted vapor pressures of these chemicals since volatility may predict ability of long-term protection vs. increased spatial domain of action (Figure 4.4B). Taken together the results of the screen present a very large collection of novel predicted repellents with desirable properties, identified via a computationally guided search of odor space.

Although the *in-silico* screen bypasses the challenge of not knowing the protein target, the most significant challenge lies in identifying effective repellent substitutes for DEET that are affordable and safe and that can be rapidly approved for human use. In order to identify compounds that fit these criteria, we applied our *in-silico* screen to an assembled natural odor library consisting of >3,000 chemicals identified as either originating from plants, insects, or vertebrate species or compounds already approved for human use as fragrances, cosmetics or flavors (See Methods). Using the trained SVM and optimized descriptor set on the natural library, we identified the top 150 ranked predicted repellent compounds. Predicted repellents share similarity in some parts of the structure while providing a diverse set of compounds (Figure 4.4C). For example, several anthranilates and pyrazines were identified that represent novel groups of safe and natural compounds, although such compounds were absent from the training set. These 150 compounds were arranged by predicted logP and vapour pressure values to provide a high-priority list of candidates for behavioral testing (Figure 4.4C).

Candidate Natural Odors Strongly Repel *Drosophila* and mosquitoes

In order to test for repellency we first used *Drosophila melanogaster* in the 2-choice trap assay that previously showed DEET repellency (Reeder et al., 2001; Syed et al., 2011). We selected 4 affordable compounds from the list, (Methyl N,N-Dimethyl

anthranilate, Ethyl anthranilate, Butyl anthranilate, 2,3-dimethyl-5-isobutyl pyrizine) the first 3 of which have been thoroughly tested for human use, have excellent safety profiles, a very mild and pleasant aroma like grapes, and have already been approved for human consumption/oral inhalation by the FDA, World Health Organization and European Food Safety Authority (EFSA, 2008; JECF, 2007) (Figure 4.4D). The pyrizine is a natural compound produced by ants as a trail pheromone (Tentschert et al., 2000). All 4 of the predicted compounds showed a very strong repellent effect on *Drosophila* across multiple doses tested using the trap-based 2-choice assay (Reeder et al., 2001; Syed et al., 2011) (Figure 4.4E). The effect was consistent at two different time points, 48 hours and 72 hours after the start of the assay.

One of the major disadvantages of DEET is its property of solubilizing plastics and synthetic materials (Krajick, 2006), which affects its usefulness by the military, and at home. We tested the ability of the 4 repellents to dissolve a 3 x 3mm square of vinyl. The plastic square disappeared completely in DEET within 6 hrs; however in the 4 predicted repellents there was no significant difference in the weight of the vinyl squares after 6 hrs, and even at the 30hr time point (Figure 4.5A).

In order to test whether these promising candidates were in fact repellent to mosquitoes, we performed behavior trials using the modified hand-in-a-cage assay. Notably, we find that all 4 compounds applied at 10% concentration demonstrate substantial repellency. Both the fraction of mosquitoes present on the net window over the arm throughout the duration of the assay, and the cumulative number of mosquitoes present on the net window were substantially lower in the test compounds as compared to solvent controls (Figures 4.5B, 4.5C, 4.6). For the mosquitoes that did land on the

repellent treatment, the escape index, as measured by the frequency of take off, was substantially higher as well (Figure 4.7).

A New Generation of Safe and Effective Repellent Odors

Taken together we have identified a number of affordable and safe compounds with repellent properties, which establishes a significant advance toward the identification of DEET substitutes that are excellent candidates for regulatory approval for human and animal use (Figure 4.5D). This was made possible by a computational screening strategy that identified shared structural features of existing repellents to use for rapid identification of structurally related novel repellent compounds. Apart from the 4 compounds we have behaviorally identified as repellents, we have also used the screen to identify ~1000 novel compounds and ~146 additional natural compounds, many approved for use in human food and cosmetics, that may lead to several additional effective repellents for deadly insects. The repellency strategy can also have great potential if used in combination with other behavior control strategies such as masking of CO₂-mediated attraction behavior or population control by trapping as a part of an integrated pull-mask-push approach (Turner et al., 2011; Turner and Ray, 2009). Since several of these compounds are affordable, repellent for fruit flies and approved for human consumption, they may have implications for control of agricultural pest insects as well that cause billions of dollars in crop loss as well as in protecting animals and pets that have tendencies to lick their skin. Moreover such substitutes may also have implications for control of DEET-resistant insects. We expect such repellents that are safe, affordable and do not solubilize plastics to pave the way for formulations that can

be used in control of insect-human contact in disease endemic areas of the world and to provide an important line of defense against deadly diseases.

Figure 4.1: Contribution of olfaction and gustation in DEET avoidance

(A) Preference index of *Drosophila* in a T-Maze to DEET (10%), CO₂ (0.66%), and opposing gradient of both compounds. N = 8-39 trials, ~40 flies/trial, error bars = s.e.m., ** = p-value < 10⁻⁴. **(B)** Preference index of *Drosophila* wild-type and *Orco*^{-/-} flies to DEET (10%) in a two choice trap assay measured at 48 hrs and 72 hrs after start. N = 5-10, 10 flies/trial, error bars = s.e.m. ** = p-value < 10⁻³. **(C)** Photograph of the hand-in-glove assay to measure repellency in mosquitoes. **(D)** Relative attraction of female *Aedes aegypti* in hand-in-glove assay measured as a percentage of mosquitoes spending >5-sec on the net covered window of glove, measured at 1 min intervals. (Left) DEET (10%) treated netting placed at the top level allowing contact, or (Right) at the lower level not allowing contact. N=5 trials for each, 10 mosquitoes/trial. **(E)** Average time spent per mosquito on net for each landing event. **(F)** Schematic of assay and mean number of female *Aedes* mosquitoes at net probing a heat and humidity source >5-secs. Mean calculated across 10 time-points, every 30-sec interval. DEET (10%) treated or Acetone (Solvent) treated netting is placed at a distance that mosquitos can not contact. N=3 trials for each, 10 mosquitoes/trial, error bars = s.e.m. p-value ** = < 10⁻³, *=<0.05. **(G)** DEET mediated repellency in *Drosophila* is mediated by both olfaction and gustation in *Drosophila* while primarily by olfaction in *Aedes aegypti*.

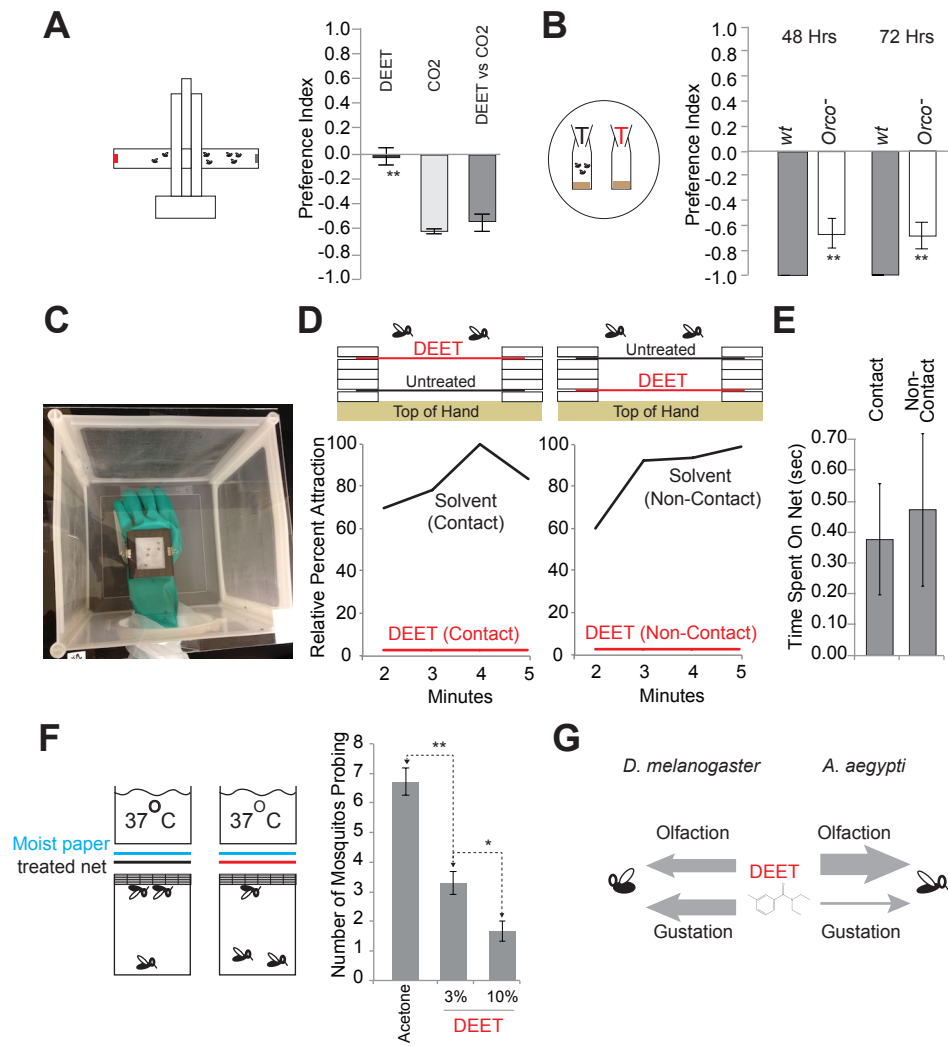


Figure 4.1

Figure 4.2: Mosquito behavioral assay glove setup

The assay glove is assembled in the following order: rubber glove with window cut into the hand, magnet glued around the cut window, control or test odor treatment mesh, three spacer magnets that prevent mosquitoes from biting through to the hand, untreated mesh to prevent mosquitoes from touching the treated mesh, and finally a top magnet. One metal clip is then used on each side of the stack to further reinforce the arrangement of magnets and mesh.

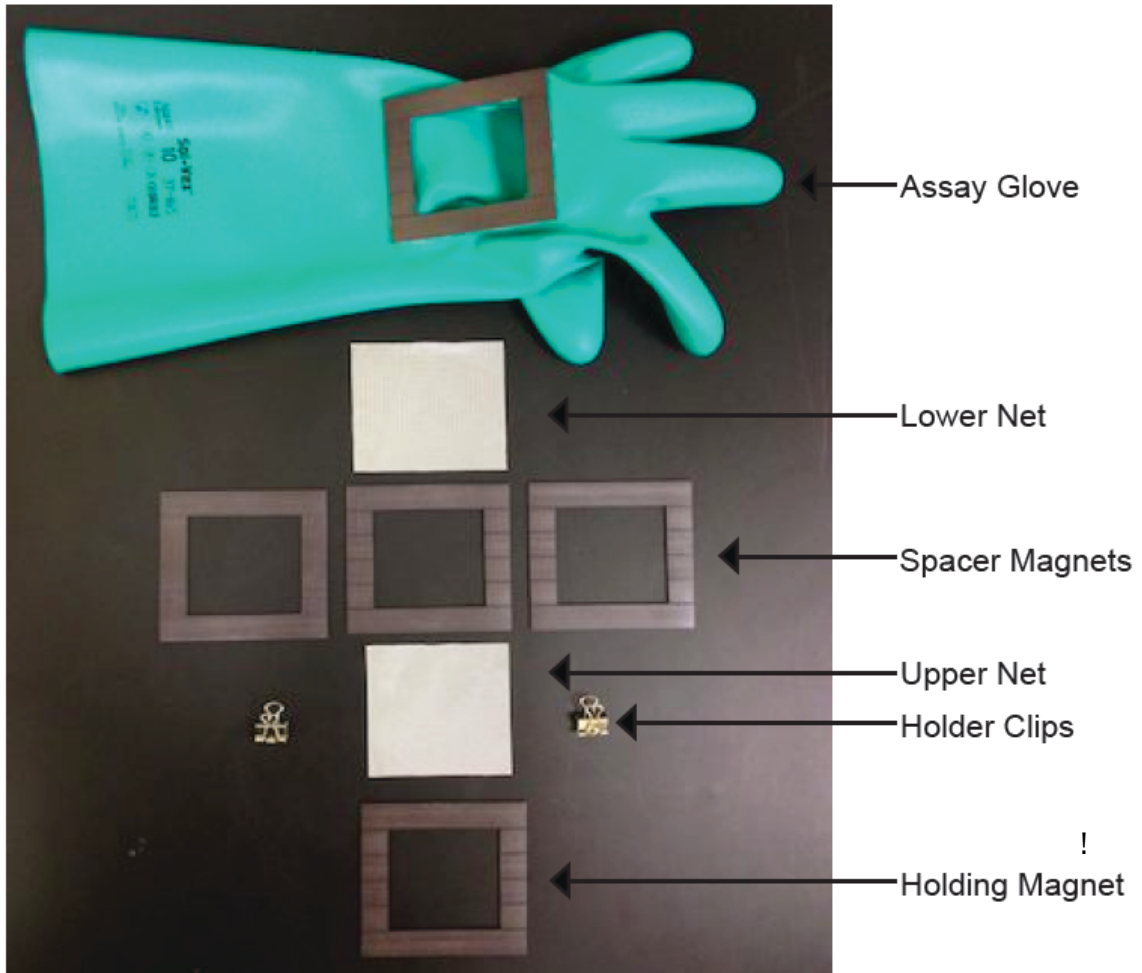


Figure 4.2

Figure 4.3: A chemical informatics method to predict repellency

(A) Overview of the cheminformatics pipeline used to identify novel DEET-like ligands from a larger chemical space. **(B)** Hierarchical cluster analysis of the 201 odorants of the training set using the optimized descriptor set to calculate distances in chemical space. **(C)** Repellency-optimized descriptor symbols and brief descriptions arranged according to order in which they were selected for the optimized set. **(D)** Receiver-operating-characteristic curve (ROC) representing computational validation of repellent predictive ability. The mean true-positive value from 20 independently run 5-fold cross validations is plotted, where ~20% of the dataset was left out of training-set as a test-set for each run. The mean area under the curve (AUC) is provided.

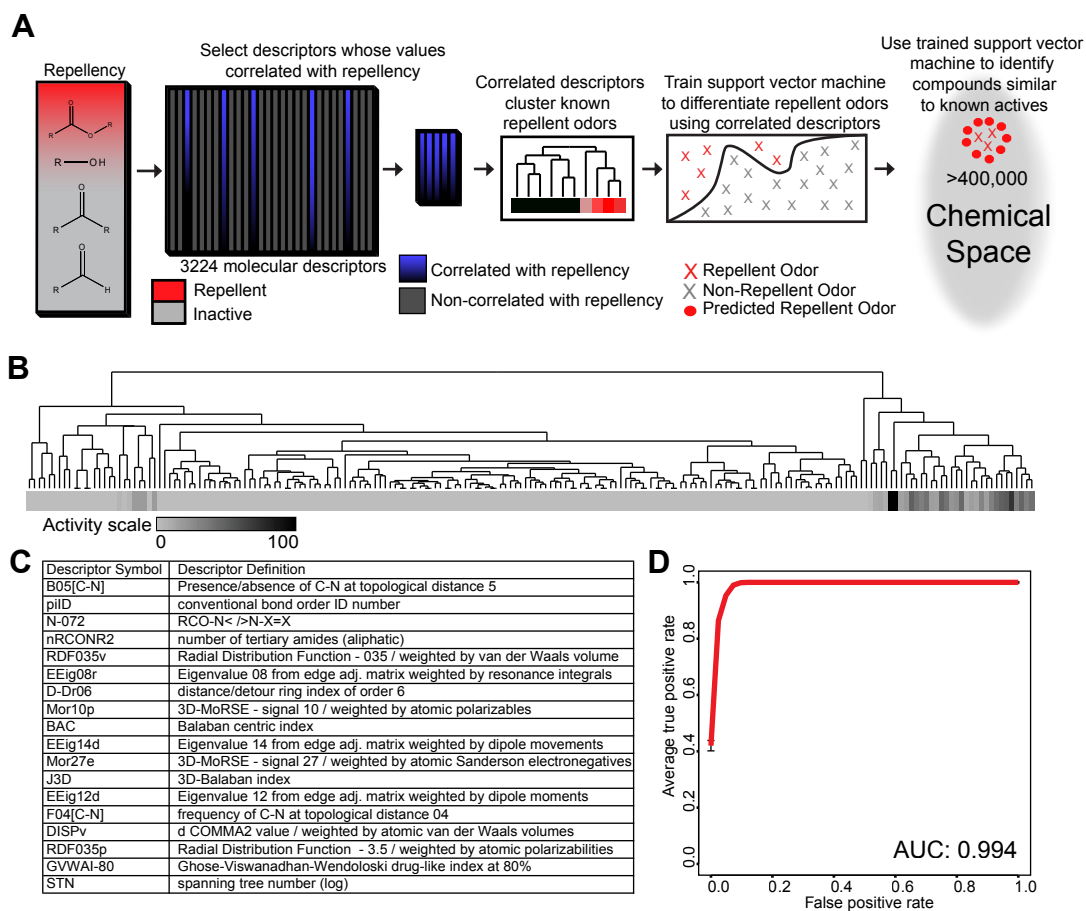


Figure 4.3

Figure 4.4: Identification of repellents using *in-silico* screening of a large chemical space

(A) Examples of two approved repellents DEET and Picaridin, and two unapproved repellents (Katritzky et al., 2008). **(B)** (Left) Representative predicted repellent odors from the odor library of >400,000. Computationally determined (Middle) LogP values and (Right) vapor pressure values for the top ranked 1000 predicted repellent compounds. **(C)** (Left) Representative structures from the top 150 predicted repellent compounds from the natural odor library of >3,000. Computationally determined (Middle) LogP values and (Right) vapour pressure values for the top ranked 150 predicted repellent compounds. Color arrowheads indicate values for DEET and odors selected for behavior experiments from the natural library indicated in **D**. **(D)** Preference index of *Drosophila* adults to predicted repellents at three different concentrations in a two choice trap assay measured after 24 hrs and after 48 hrs. N = 7-10 trials each treatment at 48 hrs and 3-10 at 24 hrs (trials with <30% participation were excluded), 10 flies/trial, error bars = s.e.m

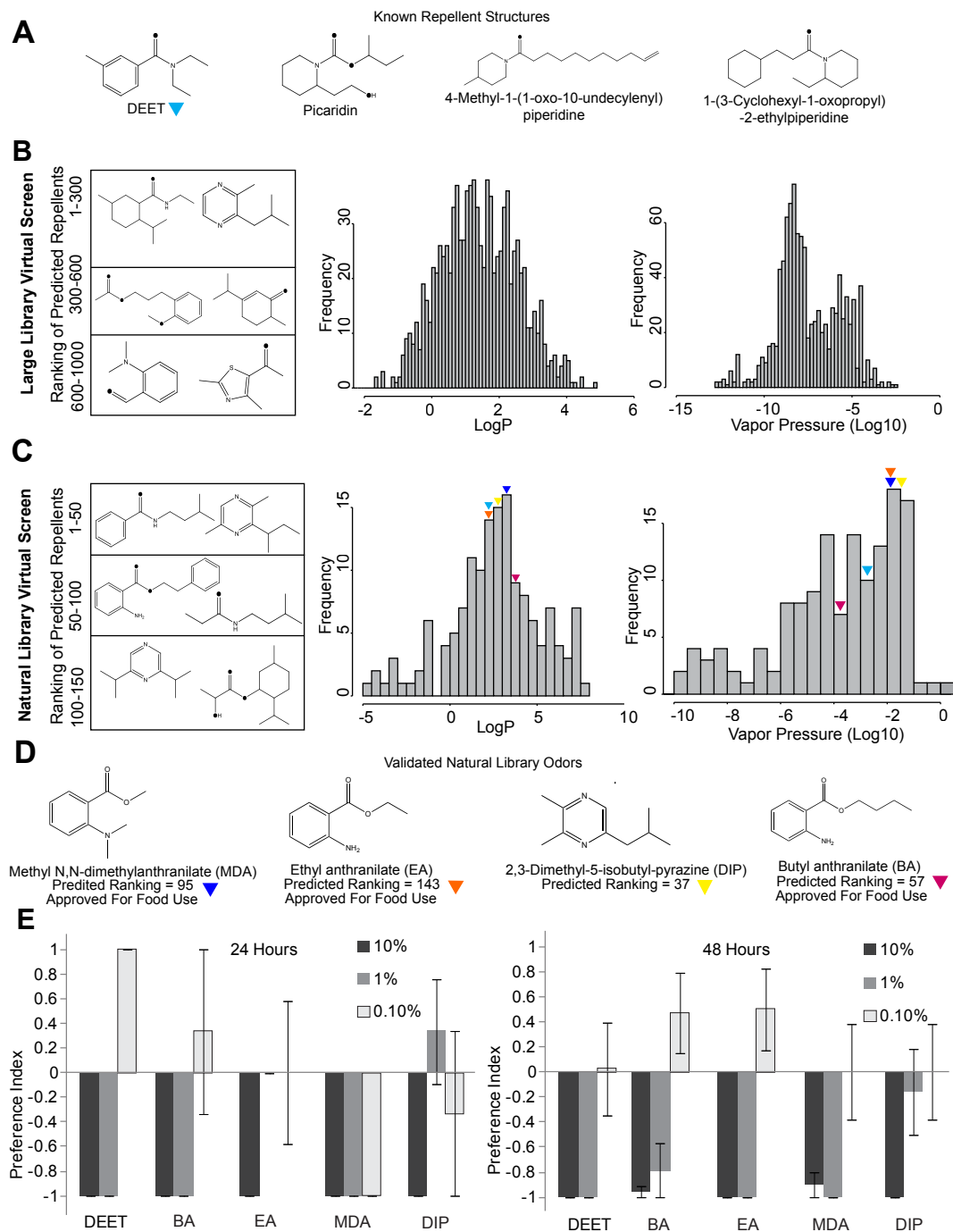


Figure 4.4

Figure 4.5: A new class of mosquito repellents with desirable safety profiles

(A) Mean weight of vinyl pieces following submersion in indicated compounds for indicated amount of time. N = 3, error bars = s.e.m., *** = p-value < 10⁻⁵. **(B)** Cumulative percentage of repellency across minutes 2,3,4 and 5 of indicated treatment (10%), in comparison to the appropriate solvent control. N=5 trials/treatment, 40 mosquitoes/trial. **(C)** Mean percentage landing as measured by mosquitoes spending at least-5 secs on the protective window of glove, measured at different time-points in the 5-min assay. **(D)** Summary of desirable properties of new insect repellents reported in this study.

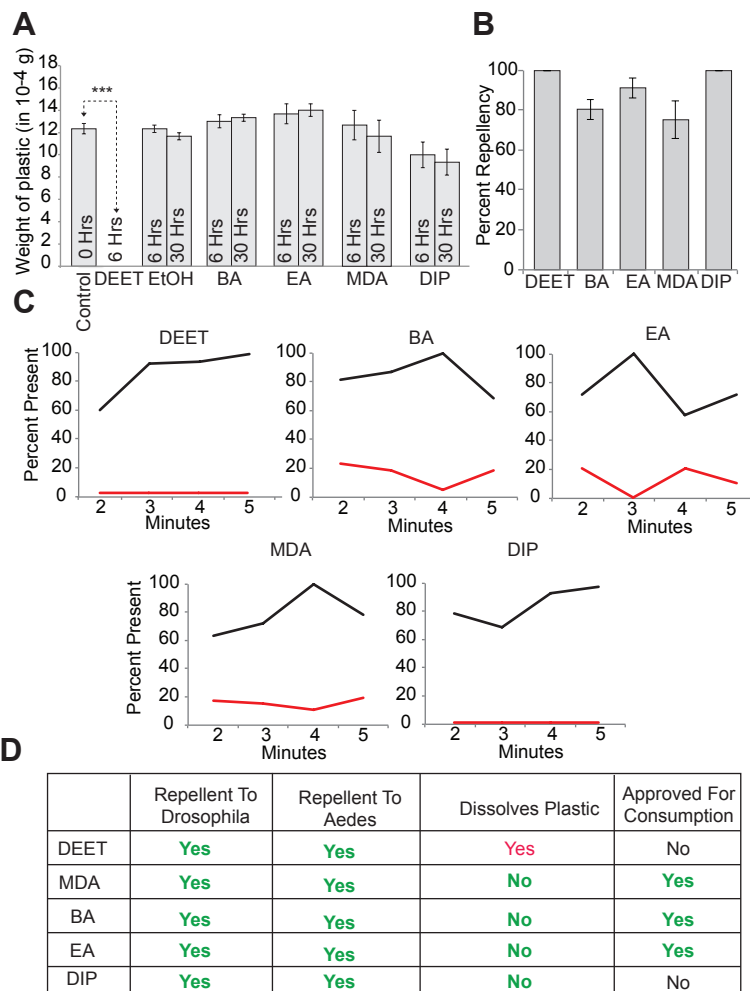


Figure 4.5

Figure 4.6: Mosquito escape index

An index representing the number of mosquitoes that touched and left the behavioral assay glove mesh vs the number that touched the mesh.

Escape Index = (Average Number of mosquitoes that touched and left the mesh during a five second window over the following time points: 1 minute, 2 minutes, 3 minutes, 4 minutes, 5 minutes) / ((Average Number of mosquitoes that touched and left the mesh during a five second window over the following time points: 1 minute, 2 minutes, 3 minutes, 4 minutes, 5 minutes) + (Average Number of mosquitoes that stayed in the mesh during a five second window over the following time points: 1 minute, 2 minutes, 3 minutes, 4 minutes, 5 minutes)).

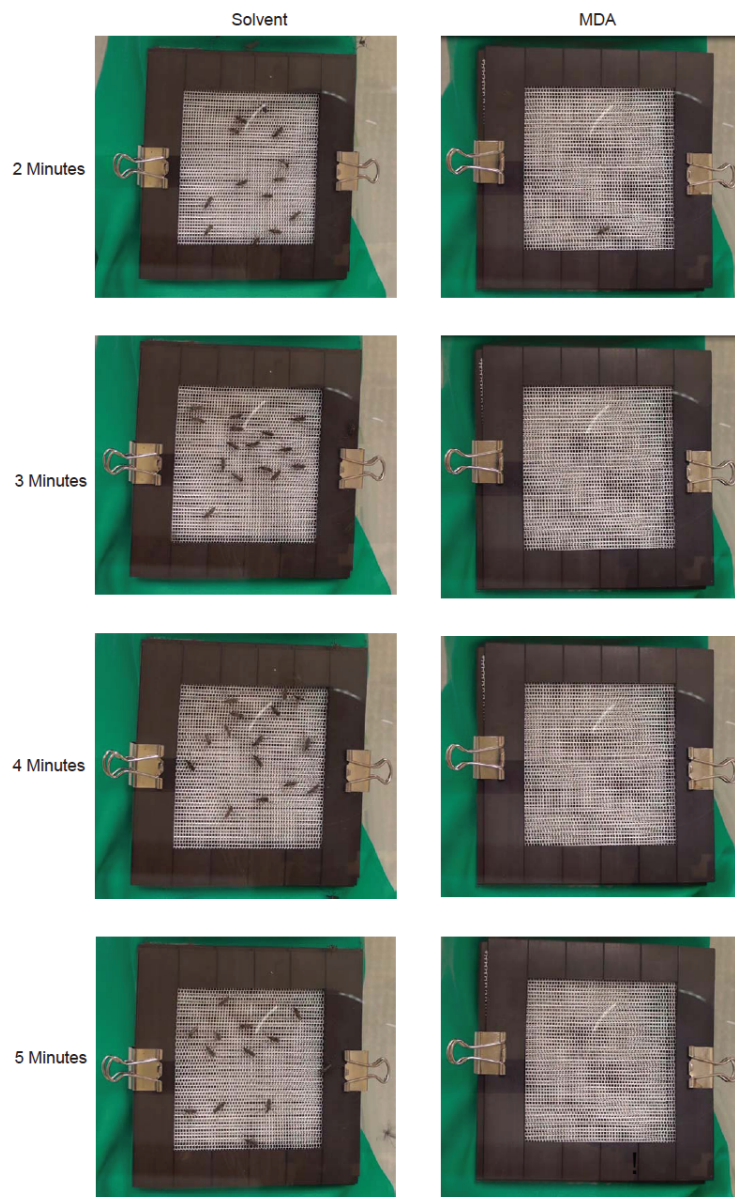


Figure 4.6

Figure 4.7: Natural compounds are effective at repelling *aedes aegypti* in the hand in glove assay

Representative still photographs from specific time-points of video assaying landing of female *Aedes aegypti* on solvent treated and MDA treated netting in the hand-in-glove assay.

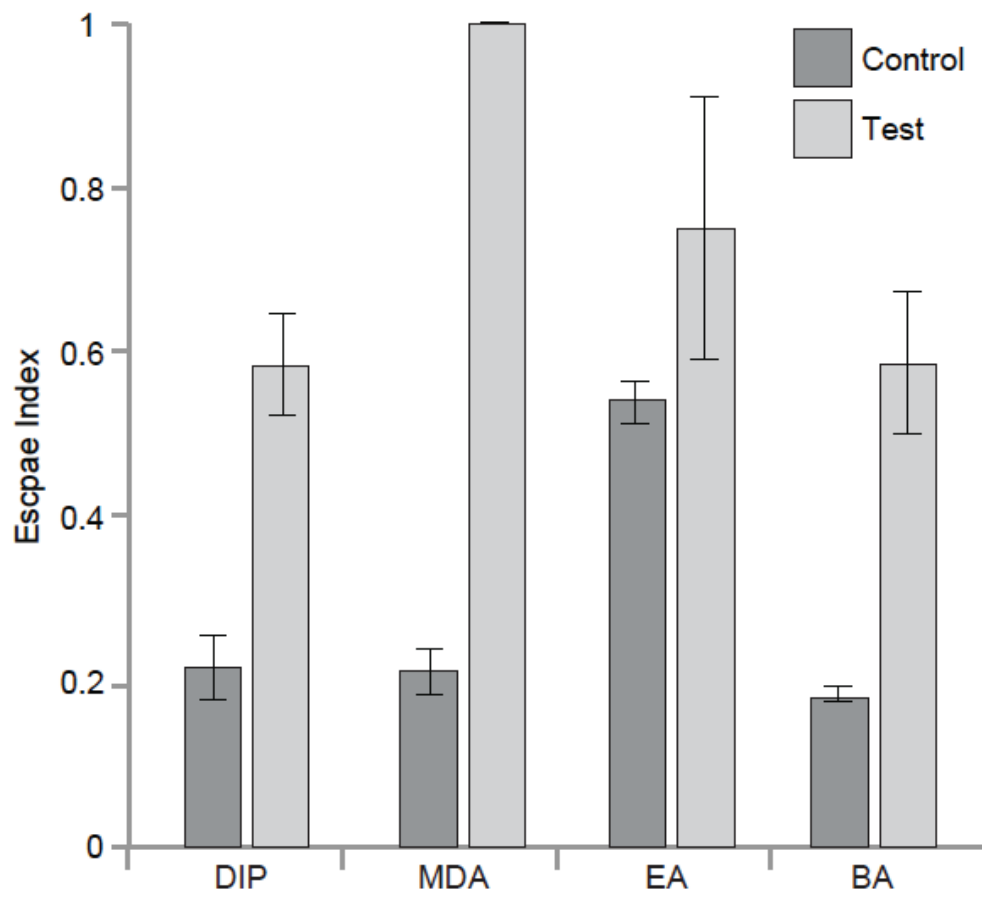


Figure 4.7

CHAPTER V:

Analyzing Termination Dynamics of Prolonged Activating Odors and their Effects on Receptor-Mediated Olfactory Behavior

INTRODUCTION

Successful interpretation of a volatile chemical environment allows an organism to make more informed decisions based upon what could be critically important information. The olfactory system is often utilized by insects and mammals for tasks such as finding food, determining oviposition sites, avoiding predators, and identifying mates (van Naters and Carlson, 2006). These high priority tasks are achieved in part through recognition of odors by large families of Odor receptor (Or) genes that often encode 7-transmembrane proteins that are expressed in the Olfactory Receptor Neurons (ORNs) (Buck and Axel, 1991; Clyne et al., 1999; Dahanukar et al., 2005; de Bruyne and Baker, 2008; Vosshall et al., 1999). Generally, a single functional Or gene, along with an obligate co-receptor Orco, is expressed per odor receptor neuron class, causing neuronal responses to be a direct result of the activity of a single odor receptor protein (Couto et al., 2005; Fishilevich and Vosshall, 2005; Larsson et al., 2004; Vosshall and Hansson, 2011). In this fashion, neurons are able to send environmental information from the periphery to higher brain centers for interpretation and behavioral decisions. Using this system an organism can decode and analyze either single compounds, such as pheromones, or complex blends, such as rotting fruit in the environment.

Activation of ORNs can lead to distinct behavioral responses, attraction and repellency, which can be robustly and reproducibly measured in *Drosophila* larvae (Aceves-Pina and Quinn, 1979; Ayyub et al., 1990; Fishilevich and Vosshall, 2005;

Kreher et al., 2008; Monte et al., 1989). While the degree of behavioral response varies across odors, it is dependent on both the odor stimuli and odor concentration (Kreher et al., 2008). Recently, large numbers of Ors have been decoded using the empty neuron system, providing a wonderful knowledgebase to draw upon for analysis (Hallem and Carlson, 2006; Kreher et al., 2008). By testing a panel of 27 odors against all 21 Ors expressed in *Drosophila* larvae one group was able to analyze the relationship between Or activity and behavioral responses on a systems level for the first time (Kreher et al., 2008). This analysis found that attraction and avoidance can be caused by activation of either multiple Ors or, perhaps more interestingly, by single Ors. It was also shown that odors that activated similar sets of Ors are able to behaviorally mask the presence of each other.

One Or in *Drosophila* in particular, Or42b, has been identified as an odor receptor whose response is explicitly linked to an innate behavioral response (Semmelhack and Wang, 2009). Activation of this Or in the adult, and thus its associated ORN ab1A and glomerulus DM1, by the well-known *Drosophila* attractant apple cider vinegar causes a robust attractive response (Semmelhack and Wang, 2009). It has since been demonstrated that the Or42b associated ORN is a starvation modulated neuron (Root et al., 2011). Starvation increases transcription of short neuropeptide F Receptor (sNPFR1) in ab1A, increasing the neurons response to odor stimuli (Root et al., 2011). This cellular response to starvation was shown to be necessary and more importantly sufficient, for increased food search behavior to apple cider vinegar, suggesting that Or42b and associated ORN ab1A may be hardwired for appetitive choices and critically important for *Drosophila* food searching (Figure 5.1A).

This system provides a powerful tool to investigate whether periphery stimulation by odors alone is sufficient to modify the strengthened food search behavior of starved flies.

Traditionally, odors have been thought to interact with receptors causing one of two unique responses, either activation or inhibition of the neuron in which it is expressed (Dobritsa et al., 2003; Hallem et al., 2004). When considered at a systems level, these two categories of responses allow some odors to activate and others to shut down odor receptor neurons in a combinatorial fashion. Very recently, a third type of response to a receptor-odor interaction has been discovered wherein a receptor is activated for a greatly extended period of time (Montague et al., 2011; Turner et al., 2011). In one analysis a special class of odors called pyrazines were found to cause prolonged activation with a maximum response of ~160 spikes/second that lasted for minutes before returning to baseline in DmOr33b and DmOr59b (Montague et al., 2011). In the other, our lab identified CO₂ detecting ORNs that express members of the gustatory receptor gene family in mosquitoes that demonstrate ultra-prolonged termination kinetics, which can disrupt detection of CO₂ and behavior for several minutes (Turner et al., 2011). Additionally, in a few instances, ORNs have been shown to have a tonic response that lasts up to 10 sec beyond the end of the 0.5 sec stimulus period (Hallem and Carlson, 2006; Hallem et al., 2004). The currently understood effects of ORN activation on behavioral responses have emerged from studying the relationship between a very limited set of receptor-odor interactions (~0.05% of putative and ~5% of known volatiles tested) and associated behavioral responses. As a result, there is potential for new behavioral responses to be observed, especially in regard to the newly identified UP activators, as very little is known about what role differences in temporal aspects of receptor activity will have on behavior.

Here we identify the first UP activators of several *Drosophila* odor receptors, including for the behaviorally important ab1A ORN, and demonstrate that these odors have profound behavioral effects on *Drosophila* attraction. We first applied a novel chemical informatics approach to identify structural features that are important for activation of the ab1A neuron. We then utilized these features to perform an *in silico* screen of a large chemical space, predicting a set of odors which are active on ab1A. We validated the predictive approach computationally and tested a number of compounds using single unit electrophysiology, both showing high success rates. This large number of new ligands enabled us to identify a small set of odors that represent the UP mode of signaling. Odors that cause this UP signaling leave a short-term memory trace, where for minutes after exposure of ab1A to an UP activator the neuronal response to other ab1A ligands is dramatically changed. We demonstrate that even a brief exposure of ab1A to these odorants trigger activity in the neuron in a manner that it loses the ability to subsequently detect other ligands in the environment, which ultimately effects the ability of *Drosophila* to track towards highly attractive ab1A activating odors. In addition to UP activators, we also discover that an inhibitory odor of ab1A can cause an avoidance behavior that is able to overcome the response to highly attractive odors.

RESULTS

Chemical informatics can be applied to describe activity of the ab1A neuron

Only a small number of odors have been tested on the ab1A associated receptor Or42b, however the testing of these odors represents a significant body of work. In the

largest analysis to date 47 odors were tested on the ab1A neuron in the *Drosophila* antenna using single unit electrophysiology (de Bruyne et al., 2001). In two additional analyses, 27 odors were tested on Or42b using single unit electrophysiology in the “empty neuron” expression system (Dobritsa et al., 2003; Kreher et al., 2008).

Volatile chemical space is vast (potentially >400,000) and can be exploited to find ligands of ab1A, however the cost of purchase and time required to test a large number of compounds to identify UP activators using traditional wet lab techniques could prove to be unfeasible. To circumvent these challenges, we have implemented a novel chemical informatics approach to identify UP activators of ab1A (Figure 5.1B). Our approach allowed us to sidestep these challenges by only purchasing and testing odors that share important structural features with known activating odors of ab1A, thus greatly reducing both the cost and time required for identification. While we did not directly predict UP activators, we greatly increased our efficiency at screening active compounds, some of which we hoped would show UP termination kinetics.

Since we wanted to train our pipeline from the largest and most diverse set of compounds that had been tested for activity on ab1A, we used the deBruyne analysis of 47 odors as our training set. To identify important structural features our chemical informatics pipeline selected molecular descriptors, which are mathematical values that describe the structural features of a chemical, that were highly correlated with the activity of a set of training odors (Figure 5.1B). Molecular descriptors can describe a great deal of features about a molecule, ranging from simple information such as molecular weight to much more complex characteristics such as the relationships between atoms in a 3D space. Using the commercially available Dragon suite we calculated 3,224 molecular descriptors for each of the 47 training odors, producing a broad range of information for

each odor. Calculated descriptors consist of 0, 1, 2, and 3 dimensional information, examples of which are vapor pressure, atom type counts, the number and type of atoms that connect each of the atoms in a molecule, and the 3D distances between atoms in energy minimized molecules, respectively. We then implemented a Sequential Forward Selection (SFS) approach that incrementally grows an optimal subset of molecular descriptors best describing the important structural features of activating odors (Whitney, 1971), which for descriptor optimization we considered to be only the most strongly activating odors (>113 spikes/sec) Figure 5.1B, See Supplemental Methods). The SFS approach began by selecting the single descriptor that is most highly correlated with the activity of the 47 training odors. The approach then iteratively built upon this single descriptor, identifying the second best descriptor when combined with the first. The process continued until a single iteration failed to identify a descriptor that further improved the correlation between activity and the growing descriptor set. The resulting list contained 13 optimized descriptors (Table 5.1).

Interestingly, the SFS approach selected four 1D descriptors, one 2D descriptor and eight 3D descriptors. The optimized descriptor set containing the 1D descriptors focusing on the number of aliphatic esters, number of acceptor atoms for H-bonds, number of hydroxyl groups, and the number of carbon atoms attached to heteroatom, as well as containing a large number of 3D descriptors, suggested that the most important structural features dependent for activity in ab1A are the absence of hydroxyl groups, the presence aliphatic esters, and their placement in the odor molecules 3D space. As expected from descriptor optimization, the ab1A activators cluster tightly together when the training is clustered using Euclidean distances calculated between odors using the optimized descriptor values (Figure 5.1C).

Next we used the optimized set of molecular descriptors to train a Support Vector Machine (SVM) to perform predictions, which is a well known and highly effective machine learning technique (Cortes and Vapnik, 1995) (Figure 5.1B, See Supplemental Methods). We then performed 100 independent 4-fold cross validations to computationally validate the predictive success, which is an established computational approach, followed by a Receiver-Operating-Characteristic (ROC) analysis (See Supplemental Methods). A mean ROC curve representing the prediction accuracy across all 100 iterations was generated and the Area UnderCurve (AUC) value was determined to be 0.999, indicating that our *in silico* approach to describe the activity of odors on the ab1A neuron using structural features of the odors themselves was almost perfect (Figure 5.1D).

Screening a large untested chemical space *in silico*

Since the Or-optimized descriptor sets are effective at grouping activating odors based on their structural features and we had successfully trained a SVM to perform predictions based on those descriptors, we next applied them to screen a large chemical space *in silico*. We assembled and calculated molecular descriptors for both an eMolecules library, which contains ~440,000 compounds with similar molecular weight and atom type constraints to known volatiles, and a natural odor library, which contains 3,197 naturally occurring odors (See Methods). We then applied the trained SVM to predict activity for each of the compounds in both libraries using the optimized molecular descriptor set, resulting in a ranked list of candidate agonists for Or42b (ab1A).

Identification of Or42b agonists through an electrophysiological validation of *in silico* screen

In order to identify agonists of Or42b, some of which we hoped would demonstrate prolonged termination kinetics, as well as to further validate our computational approach, we obtained 15 high ranking odors from our natural odor library predictions. We tested each odor for activity against ab1A and its associated receptor Or42b using single unit electrophysiology (Figure 5.2A). We observed that all tested odors were agonists of Or42b. 13 Odors were highly active (>100 spikes/sec above spontaneous activity), 1 odor was active (>50 spikes/sec above spontaneous activity), and 1 odor validated as an inhibitor (>50% reduction in spontaneous activity). In order to further verify that the response was unequivocally from the Or42b receptor expressed in the ab1A neuron and not from neighboring neurons housed in the same ab1 sensilla, we performed single unit electrophysiology on Or42b *-/-* flies (Figure 5.2B). We observed a substantial reduction in response, indicating that Or42b is indeed the target receptor for these odors.

Identification of odors for ab1A with unusual ultra-prolonged termination kinetics

It has been proposed that not just frequency of activity, but the temporal properties of a response can potentially be a rich source of information encoded by sensory neurons. For a majority of odors, ORNs show rapid activation and termination kinetics, turning off abruptly shortly after the end of the odor stimulus (de Bruyne et al., 2001; Hallem and Carlson, 2006; Hallem et al., 2004). Now that we had identified several additional activating ligands, we could systematically inspect each odor for large differences in kinetics of response termination. Interestingly, 2 of the *in silico* predicted

odors that we tested using long-term electrophysiology recordings showed an ultra-prolonged response that lasted for several minutes beyond the 0.5 sec odor stimulus (Figure 5.3A). In both instances exposure of ab1A to a 0.5-second pulse of an ultra-prolonged (UP) activator elicits a strong initial response followed by decrease in the ORN firing rate to approximately half maximal frequency in about 8-10 seconds, after which the ORN has a relatively steady firing rate of ~30-70 spikes/sec above spontaneous activity for several minutes (Figure 5.3A -green line). This unusual pattern of activation is quite different from other equally strong activators that return close to baseline between 2-6 seconds after odor exposure (Figure 5.3A -grey line).

This form of a prolonged response elicited by a short 0.5 second stimulus appears even stronger than the activity pattern evoked by a continuous odor pulse of 30 secs, in which the adaptation of the initial strong response of the neuron occurs rapidly, followed by a low frequency of action potentials throughout the duration of the stimulus (de Bruyne et al., 2001). These observations suggest that the class of UP odors we have identified by cheminformatics have an atypical property to activate ORNs for periods >500 times the duration of the initial 0.5 sec stimulus.

A mechanism for short-term memory of odor exposure

Surprisingly, a brief 0.5-second pre-exposure to an UP activator renders the neuron unresponsive to changes in concentration of other activating odorants that the ORN normally responds to. The UP-activated ORNs are unable to respond to odorants normally for more than 100 seconds, in some cases up to the duration of the entire recording (300 seconds) (Figure 5.3B, green bars). This short-term memory effect is specific to the UP-activators and pre-exposure to regular odorants that activate to

comparable levels does not affect the ability of the ORN to detect subsequent repeated 0.5 second stimuli of an activating odorant (Figure 5.3B, grey bars). This form of change in coding capacity of a peripheral sensory level after ultra-prolonged activation is a novel phenomenon that has not been reported for neurons expressing Odor receptor proteins.

Although termination kinetics has been shown to be a property associated with the receptor (Hallem et al., 2004), it remains to be tested whether this novel type of short-term memory effect could be influenced by odor binding proteins (Wang et al., 2007) found in the sensillum lymph or other factors involved in signal transduction. In order to test directly whether the short-term memory effect depends upon the odor receptor protein, we expressed *Or42b* in the well-established “empty neuron” decoder system (Dobritsa et al., 2003). We find that the UP-activation effect is partially transferred to the Δ halo ab3a neuron upon ectopic expression of *Or42b* (Figure 5.4). Surprisingly, the short-term memory effect is transferred to the acceptor ORN simply by the ectopic misexpression of *Or42b* (Figure 5.4B, blue bars), suggesting that the integration of the two inputs and retention of the memory of prior exposure is maintained at the level of the receptor.

The ability of an animal to respond behaviorally to changes in concentration of a ligand along an odor gradient or navigate along an odor plume in the environment is presumed to depend on two important coding properties of the ORN: sensitivity and rapidity of detecting incremental changes in concentration. Since both these properties are severely compromised upon exposure to UP-odorants, we expect a long-term effect on behavior even after the animal is removed from the vicinity of the UP-odorant. The

Drosophila larvae are ideal for behavioral testing of the short-term memory effects to prior exposure to UP-activators since they have a simpler olfactory system and behavior assays are robust and quantitative (Fishilevich and Vosshall, 2005; Kreher et al., 2005; Kreher et al., 2008; Louis et al., 2008). An antennal receptor Or42b, for which we have identified UP-activators, is expressed in the larval system as well (Fishilevich and Vosshall, 2005; Kreher et al., 2005; Kreher et al., 2008). It has also been demonstrated that Or42b is the exclusive receptor for detection of ethyl acetate at 10^{-4} dilution, which causes attractive behavior (Kreher et al., 2008). We performed a two-choice assay where larvae were given 90 seconds to choose between an ethyl acetate (10^{-4}) side and a solvent side on a plate after they were pre-exposed to either the UP-activator (methyl propionate), a regular activator of comparable initial strength with normal response termination (methyl isobutyrate), or solvent (paraffin oil) (Figure 5.5A). In solvent pre-exposed and activator pre-exposed larvae a robust attractive behavior towards ethyl acetate was observed, as expected (Figure 5.5B). However a 10-second pre-exposure to an UP-activator almost completely disrupted attraction. Even a brief 1-second pre-exposure to an UP-activator was sufficient to cause a significant reduction of attraction towards ethyl acetate (Figure 5.5B). Taken together these results demonstrate that pre-exposure to a super activator not only disrupts the ability of ab1A to detect odors, but also disrupts behavior elicited by the odorant.

In order to demonstrate that the change in behavior is caused directly by UP-activation of Or42b, we performed a similar pre-exposure experiment on larvae that were genetically manipulated to have only the *Or42b* expressing ORN pair functional. This was achieved by specifically rescuing expression of the obligate co-receptor *Orco* in only

the *Or42b* expressing neurons in an *Orco* mutant background. In both solvent pre-exposed and activator pre-exposed larvae robust attractive behavior towards ethyl acetate was observed, as expected (Figure 5.5C). However, pre-exposure to the UP-activator completely disrupted attraction towards ethyl acetate (Figure 5.5C). This shows that the olfactory system has the ability to integrate two odor stimuli, separated in time, directly at the level of the peripheral ORN, and compute two dramatically different outputs (change or no-change in firing frequency) to transmit to the CNS depending upon the identity of the first odorant.

Disruption of attraction towards apple cider vinegar in starved larvae

In the natural environment attractive cues are often more complex consisting of blends of odorants and insects. Moreover, the olfactory ability to find food odors may be enhanced by the NPF pathway-mediated starvation response (Root et al., 2011). For example, *Or42b* neurons were shown to be necessary for the starvation-induced food-search behavior in the adults (Root et al., 2011). In order to test whether exposure of *Or42b*, which was been demonstrated to be hardwired for behaviorally important appetitive choices, to an UP-activator was sufficient to interfere with the attraction of starved larvae towards complex food odor, we tested them in behavior assays. Surprisingly, a brief pre-exposure of 2-hour starved larvae (Koon et al., 2011) to an *Or42b* UP-activator was sufficient to dramatically reduce the ability to find an attractive source of apple cider vinegar (Figure 5.5D). This result suggests that The UP-activators may provide a powerful tool for behavior modification to protect against harmful pest species, many of which cause significant loss of agricultural produce and crops globally.

Inhibition of appetitive choice ORN causes strong avoidance response

While we have demonstrated that UP activators disrupt attraction and activators have been previously shown to cause a strong attractive response, the behavioral effect of inhibitors tuned to this receptor alone has not been previously reported. We performed a two-choice assay where larvae were given 90 seconds to choose between the inhibitor Isopentyl formate (10^{-2}), which was found in this study, and a solvent side (paraffin oil) on a plate. Interestingly, larvae behaviorally avoided the inhibitor, with a preference index of -0.2 (Figure 5.6). These results demonstrate that inhibitors are able to repel insects on their own.

Identification of odors that display Ultra Prolonged activation in additional ORNs

In order to test whether we could identify UP activators of additional Ors, we performed similar chemical informatics analysis for Or22a and Or85b. Here we trained our pipeline using the largest panel available for these receptors, which consists of 109 previously tested odors (Hallem and Carlson, 2006). In chapter 2 we predicted activating odors and validated several additional odors that displayed delayed termination kinetics reminiscent of UP activators (Figure 5.7A). We identified 5 UP activators for Or22a and 1 for 85b (Figures 5.7B, 5.8A), further suggesting that such ultra-prolonged termination kinetics are not properties specific only a select few receptors. The prolonged activation of Or85b by UP activator 2-heptanol elicits a strong initial response similar to the one caused by previously identified Or42b up activators, however the decrease in the ORN firing rate to approximately half maximal frequency requires a greatly extended time of 30 seconds, indicating that activation decay rates can vary significantly (Figure 5.7B).

Visual inspection of the 5 UP-activators of Or22a does not reveal any features that easily distinguish this class from that described for other activators of the same receptor (Figure 5.8B). To investigate whether we could identify an UP-optimized descriptor set that separates UP-activators from a large set of active compounds, we identified a subset of descriptors that cluster the UP-activators close together in a chemical space containing activators of Or22a (Figure 5.7C). We then ranked odors for UP-activation for Or22a, tested one using electrophysiology (Butyl propionate) and found that it had the hallmarks of an UP-activator (Figure 5.7D). These results suggest that the UP-optimized descriptors may be useful to distinguish odorants that cause prolonged activation even from structurally similar odorants.

We found instances where the UP-activator did not show prolonged activation of other ORNs it activates, suggesting that the prolonged response is not caused simply due to a property of the odor (Figure 5.8C). However, we also found that 2-heptanol can cause UP-activation of another ORN that expresses Or85a. These results suggest that although UP-activation is probably caused by specific interaction between odor and receptor, some odors may have properties that can cause prolonged activation in at least two different Or containing neurons.

Once again we find that a brief 0.5-second pre-exposure of these receptors to an UP activator renders the neuron unresponsive to changes in concentration of other activating odorants that the ORN normally responds to (Figure 5.9A). The UP-activated Or22a and Or85b expressing ORNs ab3A and ab3B are unable to respond to odorants normally for more 100 seconds for the case of ab3A and for the duration of the entire recording for ab3B (300 seconds) (Figure 5.9B, green bars). This short-term memory effect again contrasts from pre-exposure to regular odorants where activation to

comparable levels does not affect the ability of the ORN to detect subsequent repeated 0.5 second stimuli of an activating odorant (Figure 5.9B, grey bars).

DISCUSSION

We have designed a computational approach that allows us to identify UP activating odors and investigate their effects on odor receptors, including subsequent odor detection. Our computational pipeline builds an optimal set of molecular descriptors that explain odor activity and predicts Or42b activity to over 400,000 odors by applying them to a trained SVM. All experimentally tested odors either activate or inhibit Or42b, additionally, a number of validated odors exhibited prolonged termination kinetics. Through examination of their effects, we have identified that exposure to these UP odors disrupts the ability of *Drosophila* larvae to track towards highly attractive odors. We have also demonstrated the repellent properties of inhibitory odors of Or42b.

No analyses have been previously performed to computationally identify UP activating odors. Computational approaches have been made at predicting activating ligands, however if compared at the same activity threshold applied in our study (≥ 50 spikes/second), their success rates were far lower than ours at predicting ligands (25%)(Schmucker et al., 2007). We find that our computational approach was 100% efficient in predicting ligands of Or42b, as all odors tested with electrophysiology were either activators or inhibitors. Additionally, through rational prioritization of odors, the likelihood of identifying odors with unique properties becomes far higher. 13% (2/15) of validated odors displayed UP activation characteristics for Or42b. Without our initial computational prioritization to identify activating odors, this rate would reduce

significantly. For example, if we were screening for UP activators using the same ligand identification success rate of the hand selected training set at 19% (9/47) (de Bruyne et al., 2001), we could roughly expect our accuracy at identifying UP activators to be 13% of the 19%, or 2%. Clearly, a computational approach is advantageous.

The new ligands presented in this study include a small number of novel UP-activators whose discovery expands the olfactory code from activation and inhibition to now include ultra-prolonged activation. Most of these UP-odors, exposure to which is retained as a short-term memory trace by the odor-receptor, are present in natural substances suggesting that they may play a significant role in the natural olfactory behavior of animals. Functional insect odor receptors are thought to act as ion-channels (Sato et al., 2008) as well as GPCRs (Wicher et al., 2008) and have a unique inside-out 7-transmembrane structure (Benton et al., 2006) acting along with a heteromeric partner Orco (Benton et al., 2006; Larsson et al., 2004). Although outside the scope of this study, understanding mechanistically how these odors bind and cause prolonged activation will be of great interest.

We have demonstrated two very important characteristics of these UP activating odors: a brief exposure is sufficient to functionally silence the responses of Or42b, Or85b, and Or22a to normally excitatory odors for greater than 2 minutes and that a brief exposure is sufficient to disrupt robust attractive behavior. Olfactory signals are very important for many biological processes. Identification of subtle changes in odor plume identity and concentration allows for successful navigation by insects to odor sources. For example, casting, which involves identifying the edge of an odor plume and course correcting to towards it's center, is utilized by many species to travel up odor plumes. Our prolonged activators cause the insect to be functionally unable to determine these

concentration changes, leading to a loss of casting ability. By the time the effects of the UP activator wear off the insect could easily have flown far off course. In a case where they do reacquire the plume, they would also reacquire the UP activator, once again causing them to drift off course.

We have also demonstrated a very important characteristic of our inhibitory odor: inhibition of Or42b causes repellent behavior. Additionally, if the inhibitory odor is introduced next to a highly attractive odor for the same Or, it negates attraction. These results suggest that Or42b is not solely hardwired for innate attraction and is instead hardwired for behavior, with activating odors causing attraction and inhibitory odors causing avoidance. Further exploration on the effects of other inhibitory odors on this receptor as well as those for additional receptors should be an important focus.

An emerging area of research is the identification of odors that can modify host-seeking behavior in insect disease vectors, either by virtue of their ability to inhibit ORNs that detect host-seeking cues (Montague et al., 2011; Turner et al., 2011; Turner and Ray, 2009), by activating ORNs that cause avoidance behavior (Syed and Leal, 2008). UP-activators and inhibitors provide a novel paradigm for behavior modification strategies targeting important odor receptors (Figure 5.10). *In silico* screens can thus provide a rational foundation for identification of novel insect repellents and lures that are environmentally safe and can aid in the fight against insect-borne diseases.

Figure 5.1: Activity of a behaviorally important neuron, Ab1A, can be described by an odor receptor neuron-optimized chemical informatics approach

(A) Activation of Odor Receptor Neuron (ORN) ab1A, and thus glomerulus DM1, leads to a strong attractive response in *Drosophila melanogaster*. **(B)** Schematic of the cheminformatics pipeline used to identify novel ligands from a larger chemical space. **(C)** Receiver-operating-characteristic curve (ROC) representing computational validation of ligand predictive ability of the ORN-optimization approach. The mean true-positive value from 100 independent 4-fold cross validation runs of the Support vector machine (SVM) approach is plotted. **(D)** Hierarchical cluster analysis of the 46 training set odorants using ORN optimized descriptor sets to calculate distances in chemical space for ab1A. Training set odorant activity is indicated using a color gradient scale. **(E)** Pharmacophore of odors that are highly active against ab1A.

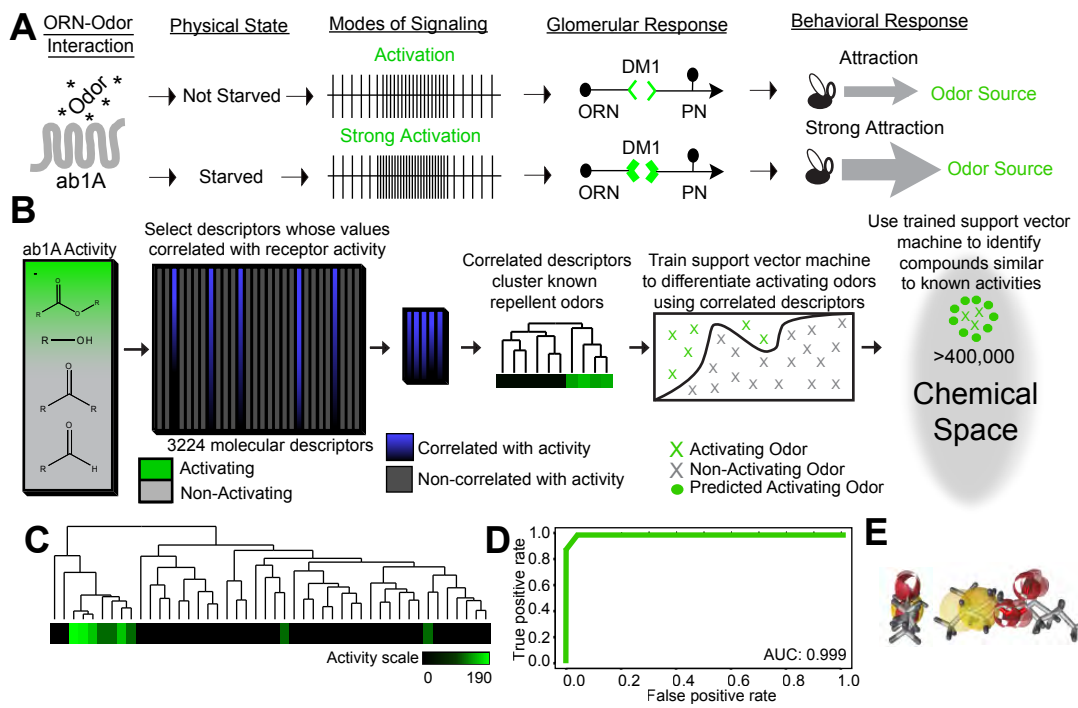


Figure 5.1

Figure 5.2: Electrophysiology validates that ORN-optimized molecular descriptors can successfully identify new ligands for ab1A

Mean increase in response of neuron to 0.5-sec stimulus of indicated predicted odors in **(A)** ab1A (Or42b) ORN; **(B)** and ab1A (Or42b^{-/-}) (10⁻² dilution). N=3, error bars=s.e.m.

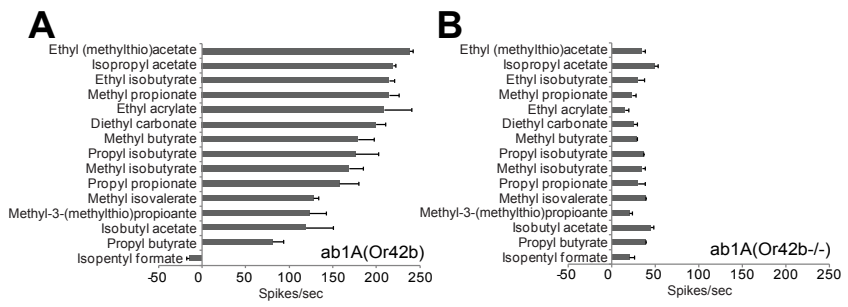


Figure 5.2

Figure 5.3: Functional identification of Ultra-Prolonged activators and analysis of their long-term effects on ab1A using electrophysiology

(A) Mean long-term response of ab1A to a 0.5- sec stimulus of indicated odor at t=0.

Each response curve is depicted in 3 separate graphs with different time windows, 10

sec, 90 sec and 300 sec. N=3, error bars=s.e.m **(B)** Mean increase in frequency of

response of the ab1A to the indicated odor applied at indicated time points after pre-

exposure to 0.5-sec odor stimulus indicated (grey=activator, green=UP-activator,

blue=UP-activator response in UAS-Or42b expressing ab3A neurons in $\Delta Halo; Or22a-$

$G4, UAS-Or42b$ flies). N=5, error bars= s.e.m.

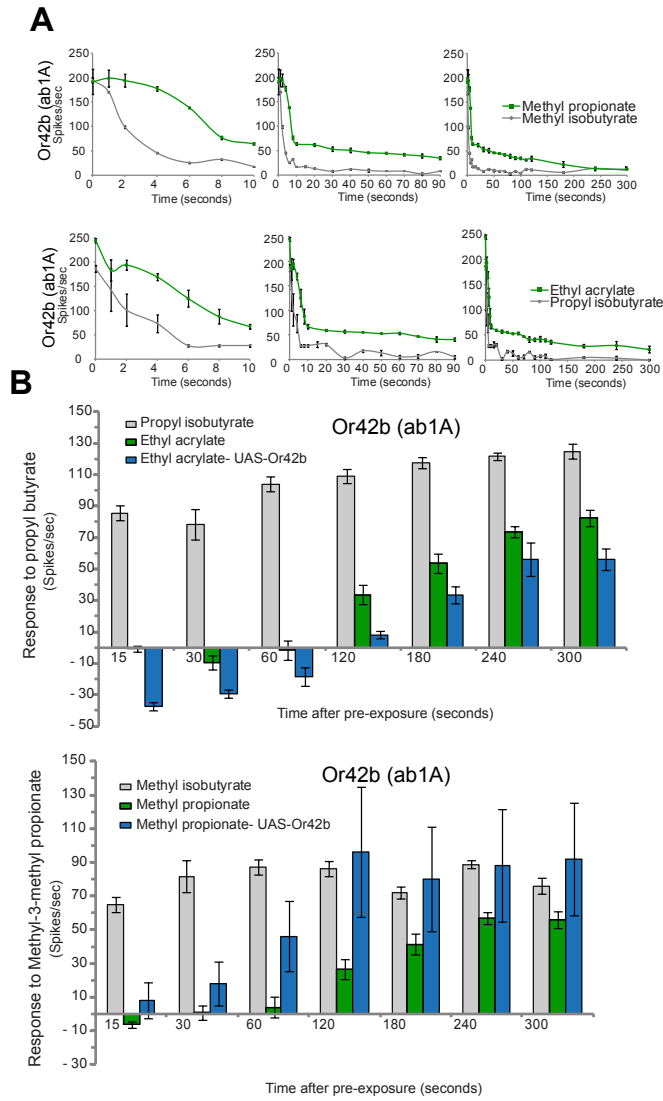


Figure 5.3

Figure 5.4: Ectopic expression confers Or42b UP activation Mean long-term response of Or42b expressing neuron (ab1A) to a 0.5- sec stimulus of indicated odor. Each response curve is depicted in 2 separate graphs with different time windows, 10 sec and 90 sec. UAS-42b recordings were performed in Δ Halo; Or22a-G4,UAS-Or42b flies where ab1A neurons express Or42b. N=3, error bars=s.e.m

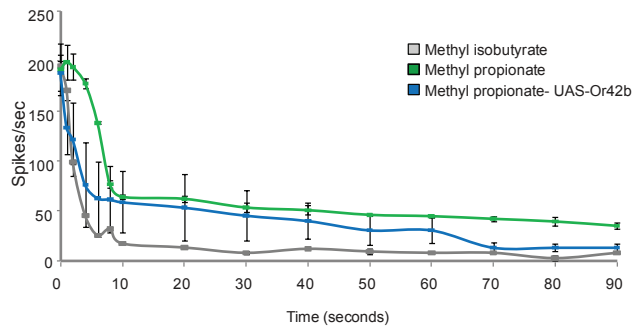
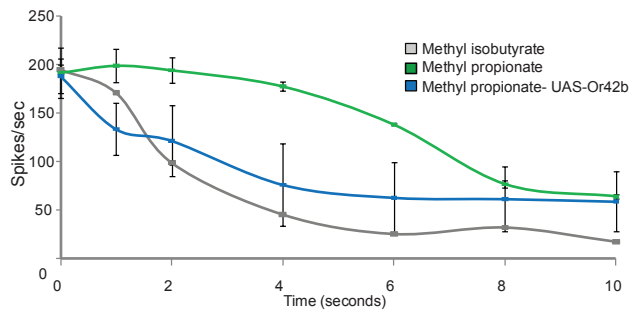


Figure 5.4

Figure 5.5: Behavioral effects of Ultra-Prolonged activators on odor detection in ab1A

(A) Overview of pre-exposure integration into larval two-choice preference assay.

(B) Preference index of *Drosophila* larvae to Ethyl acetate 10^{-4} . Larvae were pre-exposed to the indicated odors for either 10-sec (left) or 1-sec (right) immediately prior to the preference assays. N=10 (~40 larvae/ trial), error bars=s.e.m. **(C)** Preference index of *Drosophila* larvae (*w;Or42b-GAL4/UAS-Orco; ΔOrco/ΔOrco*) to Ethyl acetate 10^{-4} .

Larvae were pre-exposed to the indicated odors for either 10-sec immediately prior to the preference assays. N=10 (~40 larvae/ trial), error bars=s.e.m. **(D)** Similar experiment as (b) performed on 2-hr starved larvae given a choice between Apple Cider Vinegar (5%) and water. N=20 trials (~40 larvae/trial), error bars=s.e.m. P-values= *<0.05, **<0.001.

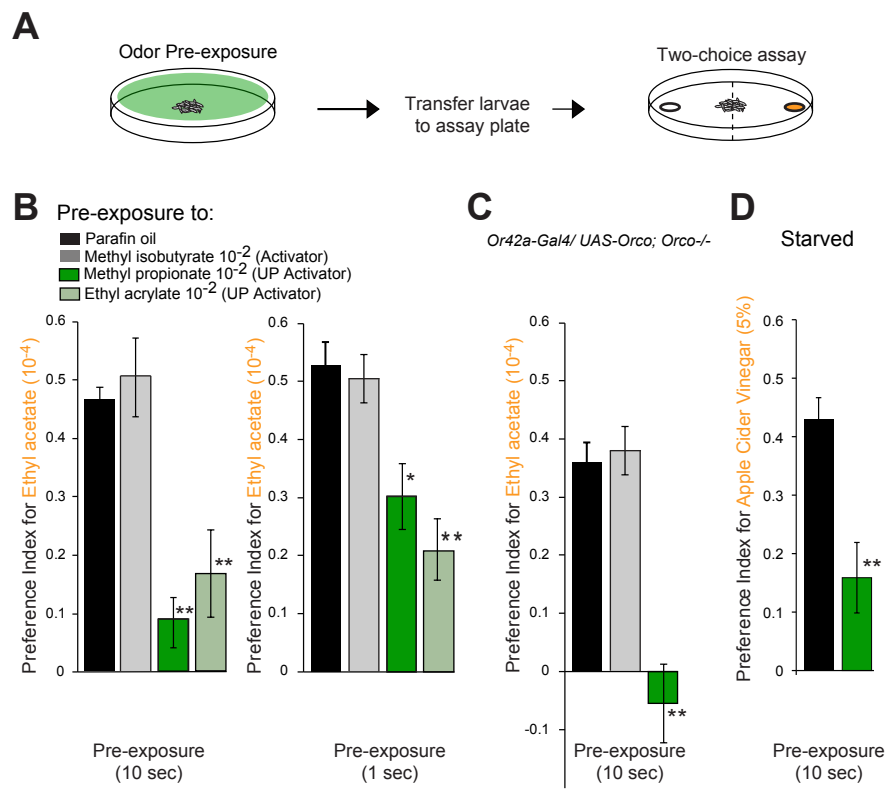


Figure 5.5

Figure 5.6: Behavioral effects of inhibitory odor on odor detection in ab1A

(A) Overview of larval two-choice preference assay. Red circle denotes inhibitor (iso-amyl formate) placement and empty circle represents solvent (paraffin oil).

(B) Preference index of *Drosophila* larvae to iso-amyl formate 10^{-2} . Larvae were given a choice between solvent or iso-amyl formate 10^{-2} for the preference assay. N=8 (~40 larvae/ trial), error bars=s.e.m.

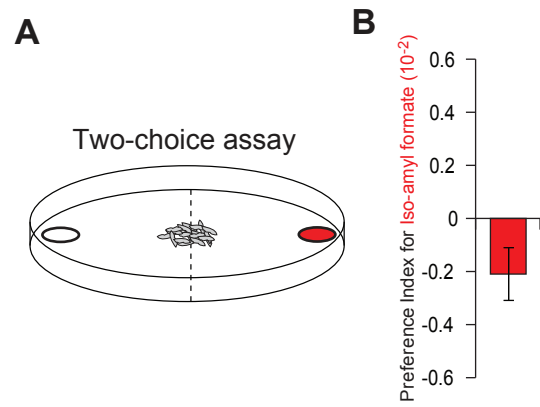


Figure 5.6

Figure 5.7: Functional identification of Ultra-Prolonged activators of additional odor receptors using electrophysiology.

(A) Long-term response of a single Or22a expressing ab3A neuron to a brief 0.5-sec stimulus of ethyl valerate (indicated as a red bar). Action potentials from 0.25-sec windows shown from indicated regions of the response. **(B)** Mean long-term response of indicated receptor expressing neuron to a 0.5- sec stimulus of indicated odor at t=0. Each response curve is depicted in 3 separate graphs with different time windows, 10 sec, 90 sec and 300 sec. For ease of spike counting ab3B (Or85b) recordings were performed in a Δ Halo (Δ H) mutant background (Dobritsa et al., 2003) where ab3A neuron is unresponsive. N=3, error bars=s.e.m **(C)** Schematic representing the identification of prolonged activator-specific molecular descriptors that can cluster prolonged activators together in a tree containing all other Or22a activators. **(D)** Mean long-term response from an ab3A neuron expressing Or22a to a 0.5-sec stimulus of a predicted UP-activator depicted in 2 separate graphs with different time windows, 10 sec and 90 secs. N=3, error bars=s.e.m

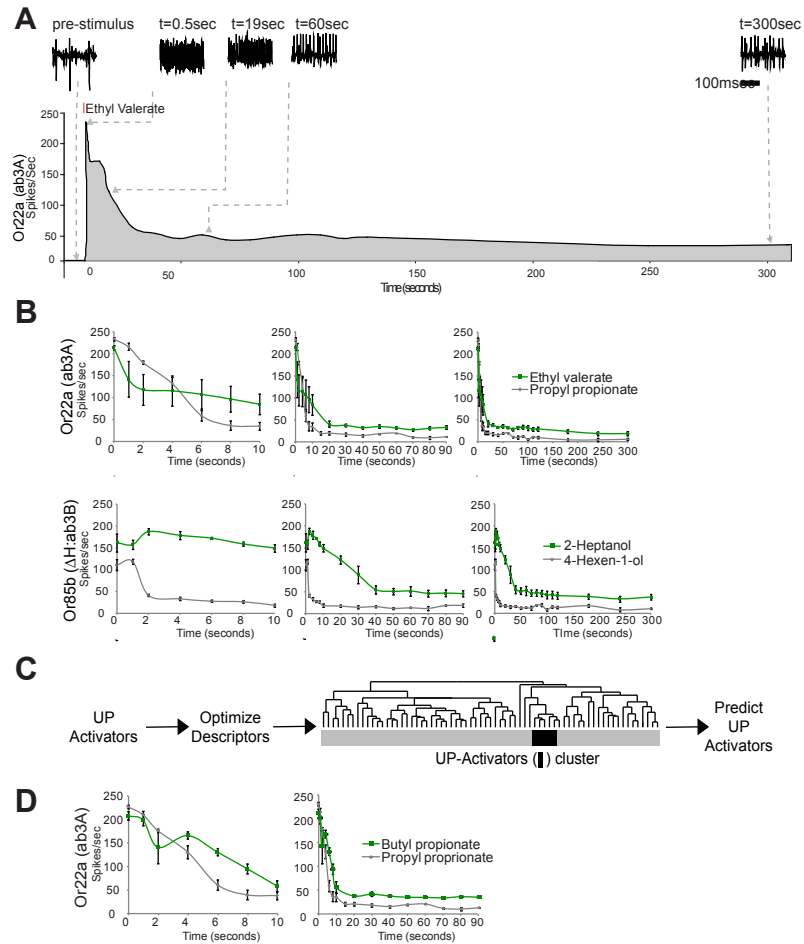


Figure 5.7

Figure 5.8: Identification of ultra-prolonged activating odors for additional Ors

(A) Mean responses from ab3A neurons expressing Or22a across 10 seconds to a 0.5 second odor (10^{-2}) stimulus, red bar. **(B)** Pharmacophores of Or22a prolonged activators and Or22a activator pharmacophore. **(C)** Mean responses of indicated receptor expressing neurons to a 0.5 second stimulus of indicated odors (10^{-2}), red bar. For (A) and (C) N=3, error bars=s.e.m.

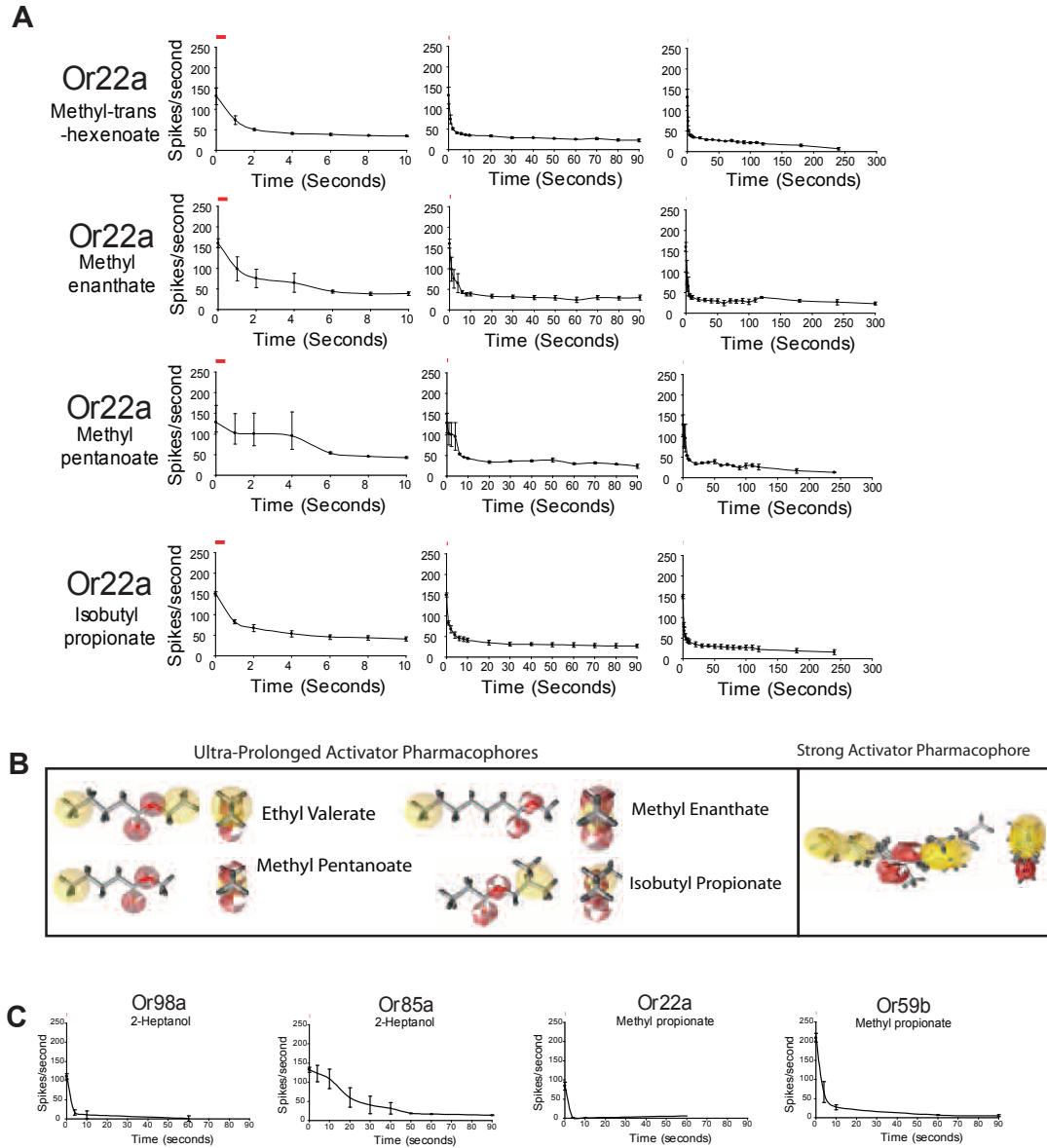


Figure 5.8

Figure 5.9: Long-term effects of Ultra-Prolonged activators on odor detection

(A) Representative electrophysiology traces of Or85b expressing ab3B neuron starting ~13 sec after a 0.5-sec pre-exposure to an activator (4-hexen-1-ol) or an UP-activator (2- heptanol) with repeated exposures to another activator (5-hexen-1-ol). Boxed area denotes 2-sec regions at 15 sec and 30 sec after pre-exposure where 0.5-sec of 5-hexen-1-ol (black bars) was applied, which are magnified above and below. Ab3B (Or85b) recordings performed in a Δ Halo (Δ H) mutant background (Dobritsa et al., 2003). **(B)** Mean increase in frequency of response of the indicated neuron to the indicated odor applied at indicated time points after pre-exposure to 0.5-sec odor stimulus indicated (grey=activator, green=UP-activator). N=5, error bars= s.e.m.

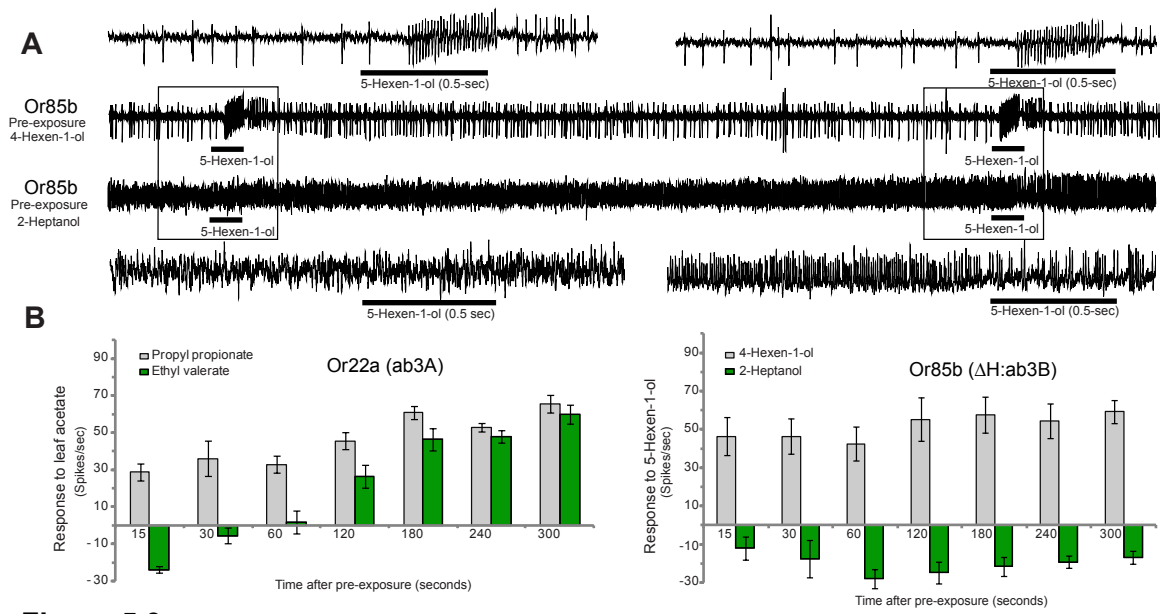


Figure 5.9

Figure 5.10: Modes of signaling and their behavioral responses

Expansion of the known behavioral responses from modes of signaling: inhibition of ab1A leads to repulsion and UP activation of ab1A leads to confusion. Flies are repulsed by inhibitory odors that specifically affect ab1A. Exposure to an UP activator causes flies to be unable to respond to normally attractive odors that activate ab1A.

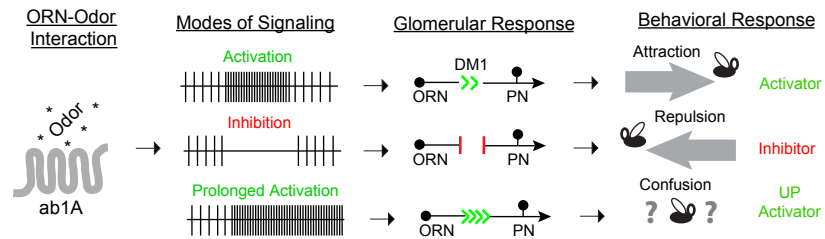


Figure 5.10

Table 5.1: Optimized descriptor sets for ab1A (Or42b) Optimized descriptor symbol, brief description, class, and dimensionality are listed. Descriptors are listed in ascending order of when they were selected into the optimized set.

Optimized Descriptor List			
Symbol	Description	Class	Dimensionality
nRCOOR	number of esters (aliphatic)	Functional group counts	1
Mor10u	signal 10 / unweighted	3D-MoRSE descriptors	3
Mor04m	signal 04 / weighted by mass	3D-MoRSE descriptors	3
R1e+	R maximal autocorrelation of lag 1 / weighted by Sanderson electronegativity	GETAWAY descriptors	3
Mor27m	signal 27 / weighted by mass	3D-MoRSE descriptors	3
nHAcc	number of acceptor atoms for H-bonds (N,O,F)	Functional group counts	1
E1m	1st component accessibility directional WHIM index / weighted by mass	WHIM descriptors	3
GATS5m	Geary autocorrelation of lag 5 weighted by mass	2D autocorrelations	2
nROH	number of hydroxyl groups	Functional group counts	1
R5v	R autocorrelation of lag 5 / weighted by van der Waals volume	GETAWAY descriptors	3
Mor10p	signal 10 / weighted by polarizability	3D-MoRSE descriptors	3
C-006	CH2RX	Atom-centred fragments	1
Mor11e	signal 11 / weighted by Sanderson electronegativity	3D-MoRSE descriptors	3

Table 5.1

CHAPTER VI:

Applying Structure-Based Virtual Screening to Identify Potent Inhibitors of the Tuberculosis Target EthR

INTRODUCTION

While virtual screening is just now becoming an effective tool for modeling and screening odors in mammalian ORs, it is not yet an option for insect Ors, as the structure of receptors for insects have not yet been solved. However, as technologies advance the structures of insect and mammalian odor receptors will undoubtedly be determined in the near future. Once the structures are available, application of structure-based virtual screening approaches to predict ligand binding will be effective tools for predicting receptor-odor interactions. Additionally, application of these approaches will also provide valuable insight into how odors are functionally interacting with the receptors. Structure-based virtual screening can be applied as a standalone method or, as was done with the goldfish Or, it can be applied in combination with ligand-based virtual screening (Triballeau et al., 2008). Comparing the structures of activating or inhibiting compounds using ligand-based approaches, such as molecular descriptors, can quickly and effectively screen millions of compounds. A second screen can then be applied where potential ligands that share important structural features with known actives can be docked into the binding site of a target protein, which is a more time consuming process.

In order to learn the technique of structure-based virtual screening, I performed an external research project under the guidance of Professor Sir. Tom Blundell at the University of Cambridge, UK. Under the guidance of specialists in his laboratory, including Dr. Willaim Pitt, I successfully applied structure-based virtual screening to

identify potent inhibitors of the enzyme EthR, a promising second line drug target currently under investigation for Tuberculosis (TB) treatment.

TB is a major health challenge around the globe. Roughly 1/3 of the world population is infected with *Mycobacterium tuberculosis*, which is responsible to leading to the infectious disease Tuberculosis that is responsible for approximately 2 million deaths each year. The currently approved course of treatment, titled Directly Observed Therapy, Short-Course (DOTS), comprises 2 months of treatment using isoniazide, rifampicin, ethambutol, and pyrazinamide followed by a 4-6 month treatment with only 2 of the aforementioned drugs (Raviglione, 2003). Unfortunately as has been seen repeatedly in bacteria, a multidrug resistant strain has been emerging, with which it is estimated that roughly 50 million people are currently infected (Chan and Iseman, 2002; Mahmoudi and Iseman, 1993). As a response the efficacy of several “second line” drugs have been under investigation and considered increasingly important.

One of these second line drugs ethionamide (ETH) has shown good efficacy, but only at relatively high doses, leading to high rates of side effects (Frenois et al., 2004). ETH functions as a pro-drug, since it needs to be activated by a mycobacterial enzyme prior to being able to perform its antimicrobial effect (Baulard et al., 2000; DeBarber et al., 2000). The activator of ETH has been determined to be EthA protein, which is a monooxygenase that oxidizes ETH, thus activating the compound (Figure 6.1) (Baulard et al., 2000; Vannelli et al., 2002). A repressor gene EthR, which has been shown to transcriptionally regulate expression of EthA through a 55bp operator located in the ethA-R intergenic region, regulates EthA (Aramaki et al., 1995; Engohang-Ndong et al., 2004). Inhibition of EthR leads to increased ETH efficacy, as the prodrug is more readily activated.

Several small molecule inhibitors of EthR have been identified over the last decade. Initial crystal structures of the repressor identified a potential drugable cavity between helices 4 and 9, which contained two dioxane molecules (Dover et al., 2004) (Figure 6.2A). Shortly thereafter another group discovered a naturally bound ligand hexadecyl octanoate which was co-purified and co-crystallized with the repressor (Frenois et al., 2004)(Figure 6.2B). The same group also identified and tested several additional ketone based compounds, identifying benzacetone as a likely inhibitor of the EthA repressor EthR. Over the years since, a number of groups have identified and co-crystallized additional EthR inhibitors, all of which function within the same long cylindrical cavity (Flipo et al., 2011; Flipo et al., 2012; Willand et al., 2010; Willand et al., 2009) (Figure 6.2C). The efficacy of these compounds ranges from micro to nanomolar.

In an attempt to identify more effective inhibitors of EthR we have applied a combination of both ligand-based and structure-based virtual screening to predict additional putative inhibitors of EthR. As we are able to use our approach to screen a vastly larger chemical space than has previously been inspected, we hoped to identify more potent EthR inhibitors.

RESULTS AND DISCUSSION

Analyzing the ligand binding pocket of EthR

As previously mentioned, a number of inhibitors have formerly been identified and their bound orientations within the EthR cavity had been determined. In addition to these published compounds, several additional molecules have been identified and co-crystallized by Dr. Sachin Surade, a member of the Blundel lab. Analysis of the EthR

structure displays a long narrow cavity that runs between helices 4, 5, 7, 8, and 9 (Figure 6.3). The cavity entrance is open and resides between helices 4, 5, and 7.

By Integrating all solved ligand binding locations and analyzing their orientations in EthR, it appears there are certainly two and possibly up to three distinct binding regions within the EthR cavity (Figures 6.4, 6.5). Very few ligands span the entire length of the cavity. The majority of compounds bind close to the cavity entrance at site 2 and a select few reside deep within the pocket in Site1. A few longer compounds extend into what could be considered site 3, which is an extension of site 2, and leave the binding pocket. From previous analysis it appears that site 2 is the most potent of these cavity positions, however it is important to note that the number of ligands identified for site 1 are limited, making a more thorough analysis required before a final decision can be made. This led us to focus our efforts on identifying and virtually docking ligands into sites 1 and 2 of the protein target, which we will refer to as small site 1 and small site 2. As a select few known ligands reside across both sites, we also performed a separate analysis where we attempted to dock ligands across the entire cavity, which we will call the combined site.

Preparing a Virtual Screening Library

We assembled a virtual library for inhibitory compound screening that contained only commercially available chemicals that were purchasable. Our library contained compounds from the Enamine (advanced library: 204,772 compounds), Maybridge (screening library: 56,000 compounds), and Asinex (merged library: 439,946 compounds) libraries containing 700,930 unique structures. Upon removal of duplicate

molecules shared across libraries (9,661 removed) and molecules containing salts or unappealing atom types, we remained with a sizable chemical library of 588,723 compounds (See Methods).

Compounds identified to bind within small site 1 were duplicates of those bound within small site 2 (Sachin Unpublished). Interestingly, binding site 1 was only occupied when site 2 was filled, raising the possibility that this site may be more challenging to fill. As a result, a single library was docked within both site 1 and site 2 and a separate and larger compound set was created for docking into the combined site.

As it would take a great deal of time to computationally dock >550,000 chemicals into each of our binding sites, we analyzed previously identified ligands to reduce the size of our screening library. We determined structural features that may be important for fitting within the cavity and for binding to EthR. By analyzing known ligands, it appears that the inhibitory cavity of EthR is long and fairly thin, thus all known compounds bound to either site are also relatively thin and unbranched. We also noticed that compounds that bound only into sites 1 or 2 were small, contained at least one aromatic ring, and often contained a ketone within the structure that formed a hydrogen bond within each of the cavities (Asparagine 179 – Site 2, Glutamine 125 - Site 1). Additionally, we identified 2 potential hydrogen bonding sites that are rarely satisfied by any of the known ligands (Asparagine 176 – Site2, Threonine 121 – Site 1). Considering these criteria, we assembled our two separate libraries, one to individually dock into binding sites 1 and 2 and another for the cavity spanning combined site.

We selected compounds for our small ligand library from the full set of >550,000 compounds using the following criteria: ≤ 3 H-bond acceptors/donors, at least 1 aromatic ring, a molecular weight of less than 250, and had no rings of size greater than

7. As the compounds in this library were relatively small, we were not concerned with removing compounds containing side chain branches that would be too wide for the cavity. These criteria resulted in selection of 38,893 compounds for a small site library that could be docked into binding sites 1 and 2.

We selected our large ligand library by application of both molecular descriptors and machine learning. While the combined site library was also selected from the full set of >550,000 compounds, it needed to contain much larger structures. As there were many compounds with sizes that would span the entire cavity in the full library, a more meticulous criterion was used in its selection. We first attempted to identify a single molecular descriptor from either the commercially available Dragon descriptor suite or pipeline pilot that defined the length and narrowness of a compound. Unfortunately, we were unable to identify such a descriptor. As a result, we compared three different sets of general shape descriptors (Balaban, Chi, and Kappa) for their ability to separate out long thin compounds. At the same time we also compared several machine learning approaches including random forests, support vector machines, and linear discriminant analysis for their ability to separate out the ideal compound shapes using each of the three descriptors. Each of these methods requires both a positive and a negative training set, which in this case would be a set of compounds that inhibit EthR and another set that do not. We assembled a training compound set containing ~20 previously known ligands (positive training set) and, as no negative ligands were known, 100 randomly selected ligands from the Asinex library (negative training set) in order to compare each method and descriptor for the ability to separate known binding from non-binding chemicals. Application of the random forest approach trained on balaban descriptor values was the most effective at differentiating the training set with the largest Area

Under the Curve (AUC) of the Receiver Operating Characteristic (ROC) of 0.95. We next applied the random forest technique to discriminate between similar and dissimilar compounds based on their balaban descriptor values, reducing the full original library from 588,723 to 112,692 compounds. As it requires a great deal longer to dock compounds into the long combined binding site, the library was still far too large. We further reduced the size of the large ligand library by randomly selecting 6,000 chemicals meeting the following requirements that were determined based upon the structures of previously identified compounds that bound into the combined cavity: MW > 400, ≥ 1 aromatic ring, and no ring sizes greater than 6 atoms.

Selecting Docking Structures of EthR

We next selected the crystal structures and binding pockets that would be used to perform the virtual screen of our final compound libraries. In total, 32 previously determined EthR crystal structures were available, consisting of 24 unpublished (Surade Unpublished) and 8 previously published (3G1M, 1T56, 1U9N, 1U9O, 3G1L, 3G1O, 3O8H, and 3O8G) structures. While initial attempts were made to use ensemble docking, which uses multiple crystal structures of the same protein to allow for protein backbone and R group movements, the computational docking time greatly exceeded our time constraints. As a result, we selected a single structure that best defined each binding site for docking.

We identified the single structure that underwent the largest conformational shift in the DNA binding region, which in theory would also cause the largest reduction in repression of EthA, and contained the widest ligand cavity. All protein structures were aligned to their largest alpha helix using Pymol. Each structure was compared for R

group shifts leading to increase or decrease in binding cavity shape. After careful inspection, it appeared that MB1047 and NT06 had the most open configurations in site 2 and site 1, respectively, and were selected for docking (Surade Unpublished). NT06 was also selected for combined site docking as this structure had the largest cavity spanning both binding sites.

Selecting Hydrogen Bond Locations Important for Protein-Ligand Interface

The hydrogen bonding between a ligand and its protein-binding site is critically important for describing the mechanism of attachment. As such, it is essential that we identify which binding locations appear to be important for EthR. Visual analysis of all previously identified ligands, their binding sites, and the protein cavity itself uncovered Asparagine 179 for site 1 and Glutamine 125 for site 2, as being highly important to hydrogen bonding at the protein-ligand interface. Asparagine 176 and Threonine 149 in site 1 and Threonine 121 in site 2 may also play a less substantial role, as evidence for their interactions with known ligands is present, but less pronounced.

Self-Docking and Cross-Docking Validations

We next compared how well the Gold software suite self-docked, which involves docking a ligand back into its own crystal structure, and cross-docked, which involves docking other solved ligands for the same receptor into a crystal structure, for the structures NT06 and MB1047. As we know exactly where each of these compounds should bind, they provide good indicators as to how well the docking software and protein system will work.

For self-Docking, ligands for both NT06 and MB1047 were re-docked into their own solved structures using the GOLD suite (Verdonk et al., 2003). H-Bond weighted constraints for Asparagine 179 in site 1 and Glutamine 125 in site 2 were compared to docking runs where no constraints were implemented. Self-docking runs were often able to successfully reinsert their own ligands back into the correct position and orientation in the receptor (Figures 6.6, 6.7). For our cross-docking analysis, we docked a variety of ligands known to bind to EthR. Cross-docking is a more stringent test of binding as a structure that is used for docking was not solved with a ligand that is attempted to be docked creating a test to see how “accepting” the binding site is compounds with different shapes that are known to bind as well. The analyses yielded promising results, as many of the ligands were positioned very similarly to the natural binding positions of their own solved structures.

We next compared the effectiveness of both the presence and absence of H-bond constraints. Interestingly, docking with and without H-bond constraints provided similar results, indicating that even without constraints this interaction is likely energetically preferred. However, as a large number of known ligands contain H-bond interactions, we feel they will be highly important in discriminating between the tens of thousands of compounds in our virtual screen, many of which may not be able to form this H-bond. As such, we have decided to include constraints for all docking.

Performing structure-based virtual screening on EthR

We performed 3 independent docking experiments. Each experiment attempted to dock a single library into a binding site using the Gold software suite. While the many of the same settings were applied across experiments, several were unique to a binding

site, such as H-bond constraints. Each compound was docked multiple times into a binding site to increase the likelihood of identifying the correct location and orientation in the binding pocket. The structure of the EthR protein was kept rigid for all 3 docking scenarios.

We first docked our small library into the EthR binding site 1 of the solved structure NT06. Each of the 38,893 compounds from the small site library were independently docked into binding site 1, 10 times. H-bond constraints were applied for both Asparagine 179 and Asparagine 176.

We next docked our small library into the EthR binding site 2 of the solved structure MB1047. Each of the 38,893 compounds from the small site library were independently docked into binding site 2 10 times. An H-bond constraint was applied for Glutamine 125.

Lastly, we docked our combined library into the full EthR cavity of the solved structure NT06, which spans nearly the entire cavity including from deep within the pocket at binding site 1 to binding site 2 and beyond to the opening of the tunnel. Each of the 6,000 structures was independently docked into the cavity 3 times. We had reduced the number of docking iterations to 3 due to computational constraints. As the size of both the ligands and docking cavity was much larger for the combined site, a far longer time was required for each individual docking. All H-bond constraints applied in both sites 1 and 2 were implemented for combined site docking (Asparagine 179, Asparagine 176, Glutamine 125).

Selecting docked ligands for experimental validation

We analyzed Gold docking results using Goldmine, which is an integrated portion of the Gold software suite. The Goldmine package allows for visual inspection and thresholding of docking results based on H-bond constraints and fitness values. Fitness values and H-bond constraints were individually applied to each docking set, resulting in a reduction to 708 compounds for site 2, 619 compounds for site 1, and 69 compounds for the combined site. We then sorted the remaining compounds by H-bond score and visually inspected each docking, selecting the most promising leads. We ultimately decided on 14 compounds for purchase and experimental validation. 6 of the selected compounds originated from binding site 1, 3 from binding site 2, and 5 from the combined site.

Experimental validation of predicted ligands

All experimental validation was performed by Dr Sachin Surade. The 14 compounds were assayed for their inhibition of EhtR binding affinity to the promoter sequence of EthA. Surface Plasmon Resonance (SPR) assays were performed to determine the inhibition caused by each compound at a 50 micro molar concentration. SPR assay activity determined for the compounds ranged from 95% inhibition to 159% activation (Table 6.1). We believe activation values may be a product of interference with the assay itself.

Analysis of predicted ligand binding and experimentally validated activity

We next inspected each of the docking conformations and compared their validated inhibitory values to orientations within their associated binding pocket. It is important to note that predicted binding orientations may not be exactly as they are

when truly bound. Our screening approach only considered the rigid structure of a single enzyme. Since protein structures in cellular conditions are able to shift and flex and many amino acid side chains are able to rotate, the true binding conditions will likely allow for more ligand movement and a better orientation in the binding pocket.

Compound BAS118136, which was docked into site 1, had the highest validated inhibition of 95% (Figure 6.8 (Top)). From analysis of the predicted binding location it appears that two hydrogen bonds may be strengthening the interaction. It appears that the ligand may actually be shifted forward and rotated forming interactions with Glutamine 125 and Glutamine 180, however this is only a hypothesis would need to be determined with a crystal structure. The wide shape of the structure placed into binding site 1 should cause a large shift in DNA binding region of EthR, thus reducing the repression of EthA and increasing ETH effectiveness.

Compound T6977356, which is docked into site 2 had the second highest validated inhibition of 92% (Figure 6.8 (Middle)). From analysis of the predicted binding location it appears that a hydrogen bond with Asparagine 179 is formed. The ligand shape is long and flat, which fits very well into the cavity and spans from the edge of binding site 1 to the cavity opening. This long addition in the cavity should cause a large shift to the backbone of EthR.

Compound T6779777, which is also docked into site 2 had the third highest validated inhibition of 86% (Figure 6.8 (Bottom)). This compound has a very similar orientation in the binding cavity to T6977356. There appears to be a hydrogen bond forming with Asparagine 179, however this molecule contains two additional H-bond accepting oxygen atoms that do not appear to be forming specific interactions. It is possible that the ligand has a slightly different orientation in its naturally bound state.

The SPR analysis determined that 4 additional compounds have inhibitory values that are between 25% and 75% (Figure 6.9). These moderate inhibitors range in shape and size. One was predicted to interact with binding site 1, one is predicted to interact at binding site 2, and two are predicted to interact across the whole combined cavity. All 4 molecules appear to form H-bonds with at least one amino acid. The two compounds from the combined docking set may form multiple H-bond interactions. 3 of the molecules have bi or tri-cyclic rings that are likely effective in expanding the cavity and thus effectively reducing the ability of EthR to bind to the DNA promoter region upstream of EthA.

Finally, 6 of the compounds had validated inhibitory values lower than 25% and were poor inhibitors of EthR (Figure 6.10). Poor inhibitory compounds originated from all three binding sets. Several compounds contained very large side groups including an adamantane group and a bicyclo[2.2.2]octane. Interestingly, several of the compounds validated as activators of EthR. While it is possible these compounds increased the binding affinity of EthR to DNA, it is probably more likely they have some off target effect on the assay itself.

Future directions in validation

Preliminary results presented here show great promise and we will continue validation of these compounds in several ways. Firstly, It will be important to determine IC50 values so that they can be compared to the current best in class inhibitors of EthR. Each of these compounds has been tested for activity at a 50 micro molar concentration. Since two of our compounds caused very high inhibition (>90%), it is possible that these compounds could be as effective as the best in class compounds. Secondly, it will be

important to test for inhibition across multiple assays. The responses of EthR to a compound may change in a different assay. Finally, it will be very important to determine the true binding orientations of the best compounds within the binding pocket. In order to achieve this we will need to co-crystallize compounds into EthR. As EthR has been co-crystallized with a number of other compounds, this should be an achievable goal. It will be very interesting to compare the predicted with the actual ligand orientations in EthR.

The high success rate in identifying highly inhibitory ligands is a testament to the utility that structure-based virtual screening can have on ligand identification. In a period of 6 weeks we assembled a large untested library of >550,000 compounds, applied ligand-based virtual screening to analyze structural features that are important for ligand binding to limit the library sizes for individual binding sites, and performed structure-based virtual screens for 3 independent binding sites on EthR. From our very small experimental validation we find an accuracy of 36% at identifying inhibitors (>50% inhibition), with 14% being very strong inhibitors (>90% inhibition). The ligand specificity of the EthR receptors is almost certainly vastly more selective than for insect odor receptors, as many insect odor receptors are activated or inhibited by an array of structural features and sizes. While the accuracy of this analysis was lower than what we have observed by applying ligand-based virtual screening in insect species (~72% in *Drosophila* and ~65% in *Anopheles*), we would expect a far higher success rate in odor receptors by application of a two stage (ligand-based and then structure-based) screen. As the first structures of insect odor receptors and additional mammalian odor receptors become available, application of structure-based virtual screening should be of high priority.

Figure 6.1: Schematics of Ethionamide Activation Pathway

The Ethionamide activation pathway. EthR represses the expression of EthA. EthA is responsible for activating Ethionamide causing the antibacterial effect. By inhibiting our target EthR, a larger amount of Ethionamide is converted into active form, thus increasing the effectiveness of the drug.

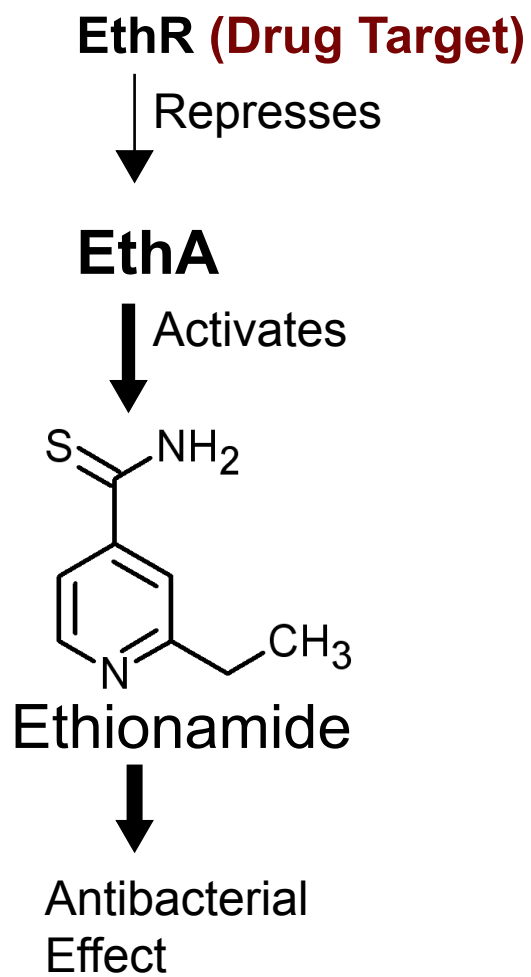


Figure 6.1

Figure 6.2: Solved EthR structures

(A) The first solved structure of EthR. Fortuitously, the receptor was co-crystallized with two dioxane molecules, allowing for identification of a cavity within the protein. This structure was previously determined, obtained from PDB, and visualized in PyMol (Dover et al., 2004). **(B)** The first identified structure of EthR bound to a naturally occurring ligand. EthR is bound to hexadecyl octanoate, which spans a very large cavity in EthR. This structure was previously determined, obtained from PDB, and visualized in PyMol (Frenois et al., 2004). **(C)** An overlay of several bound EthR ligands. Individually co-crystallized structures have been aligned by their largest alpha helix, producing a clear representation of the binding cavity. All structures were previously determined, obtained from PDB, and visualized by PyMol (Dover et al., 2004; Flipo et al., 2011; Frenois et al., 2004; Willand et al., 2010; Willand et al., 2009).

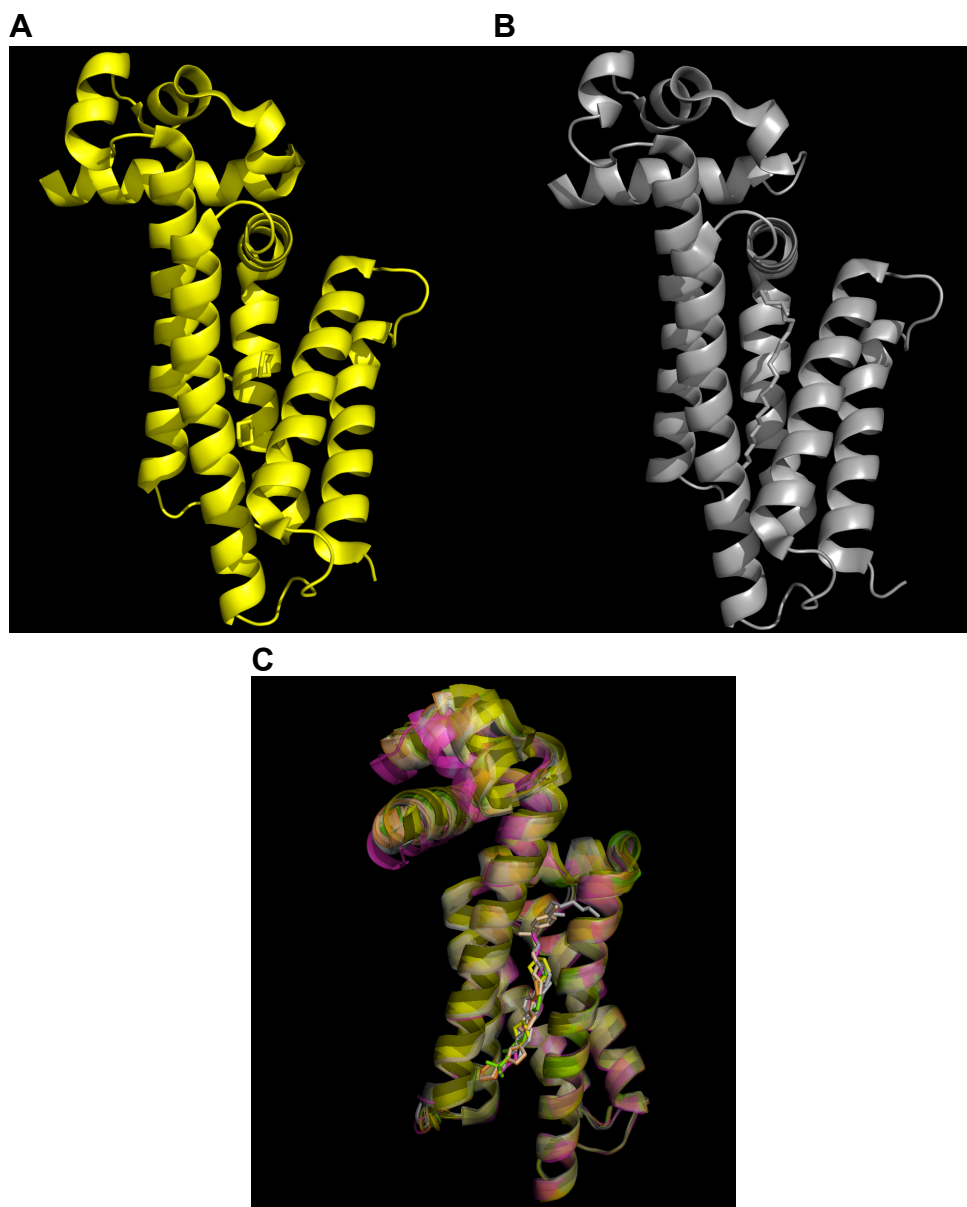


Figure 6.2

Figure 6.3: The shape of the EthR cavity

A structure of EthR (NT06), which was recently solved by members of the Blundell lab, was used along with PyMol to visualize a long and unbranched cylindrical cavity that traverses nearly the entire length of EthR .

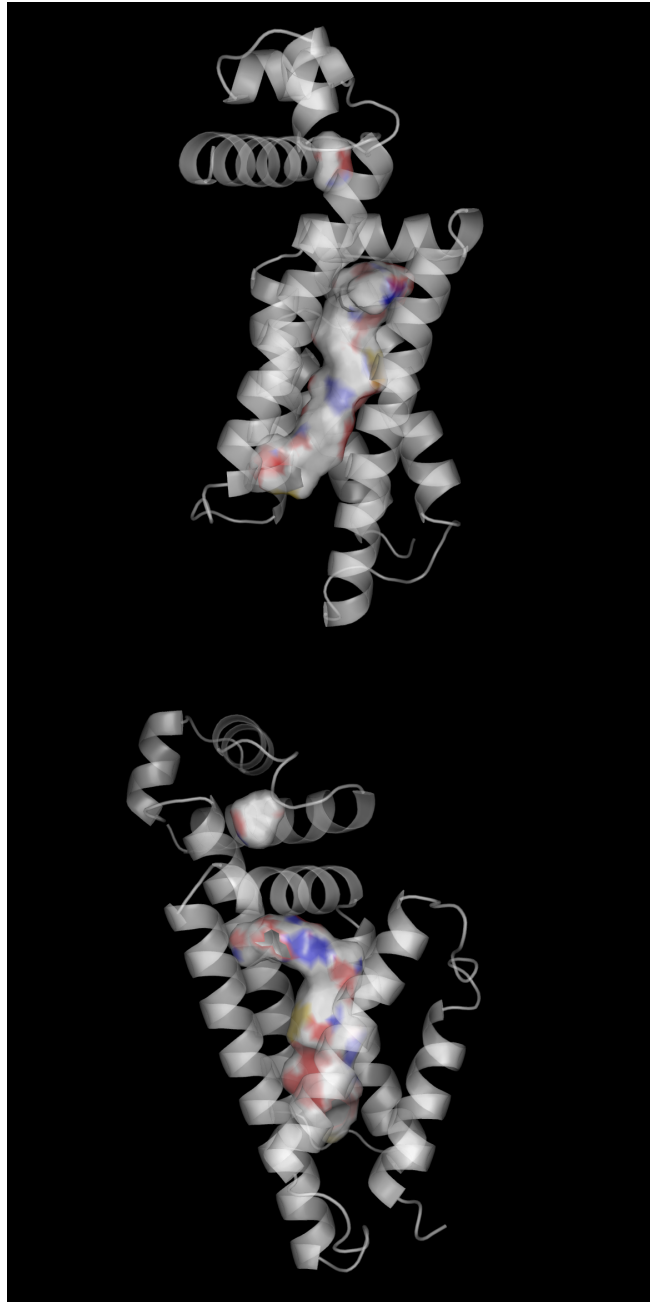


Figure 6.3

Figure 6.4: Previously identified EthR ligands bind within the proposed cavity

Previously determined ligands were oriented into the EthR binding cavity. Three distinct binding regions can be observed. Similar compounds bind into sites 1 and 2.

Compounds that extend from the site 2 occupy a hypothetical site 3.

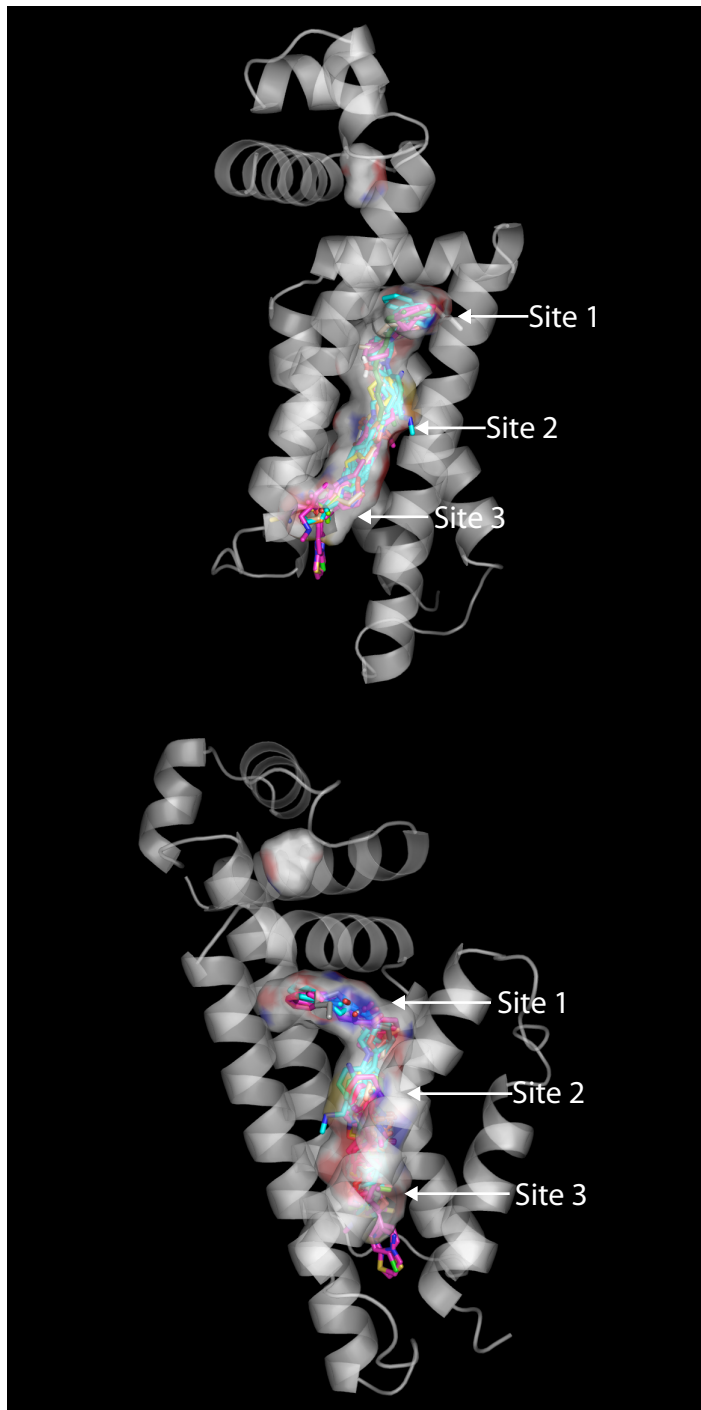


Figure 6.4

Figure 6.5: Visualizing previously identified EthR ligands

Three different views of previously identified EthR ligands residing within the cavity.

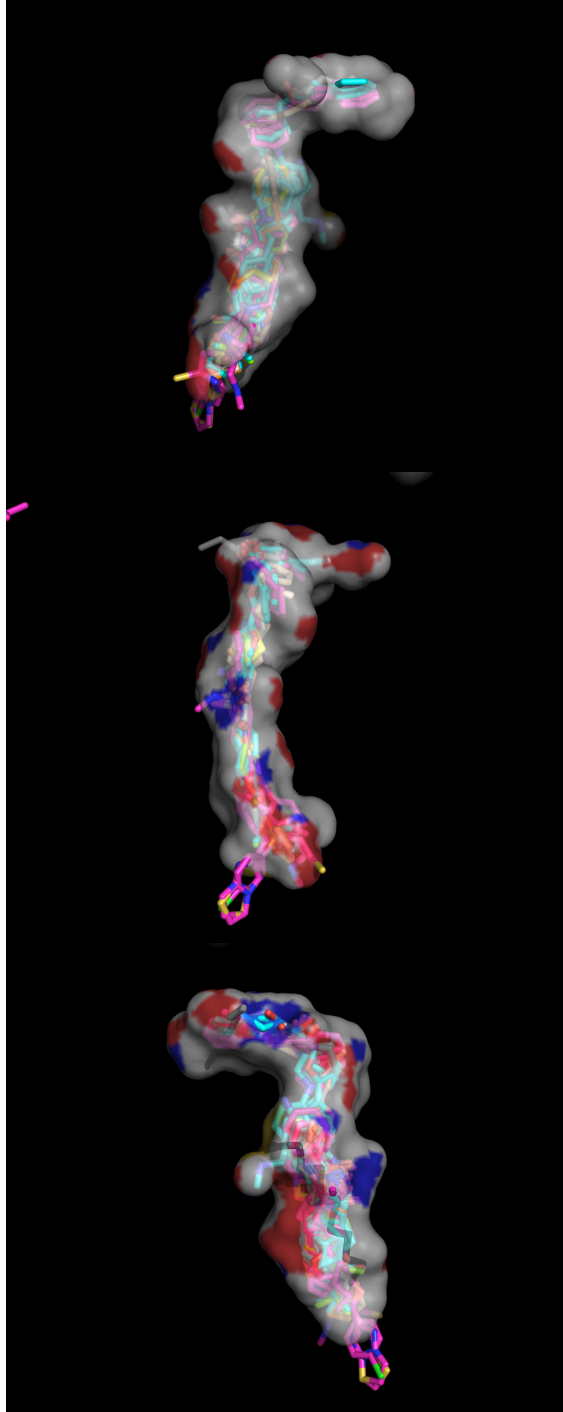


Figure 6.5

Figure 6.6: Single example of successful EthR self-docking

The EthR ligand (brown ligand) has been successfully self-docked back into the correct orientation (yellow ligand) in the EthR binding site 1.



Figure 6.6

Figure 6.7: The results of many EthR site 1 self-docking runs

The Gold software suite is highly successful at self-docking ligands back into EthR site 1. Each overlaid ligand represents an individual self-docking into binding site 1. The majority of ligands reside in the correct orientation.



Figure 6.7

Figure 6.8: The predicted orientation of our strongest identified inhibitors in EthR

The resulted docking of our top three strongest validated inhibitors bound into EthR. The percent inhibition represents the inhibition observed by SPR experimental analysis. An individual image is provided for each inhibitor.

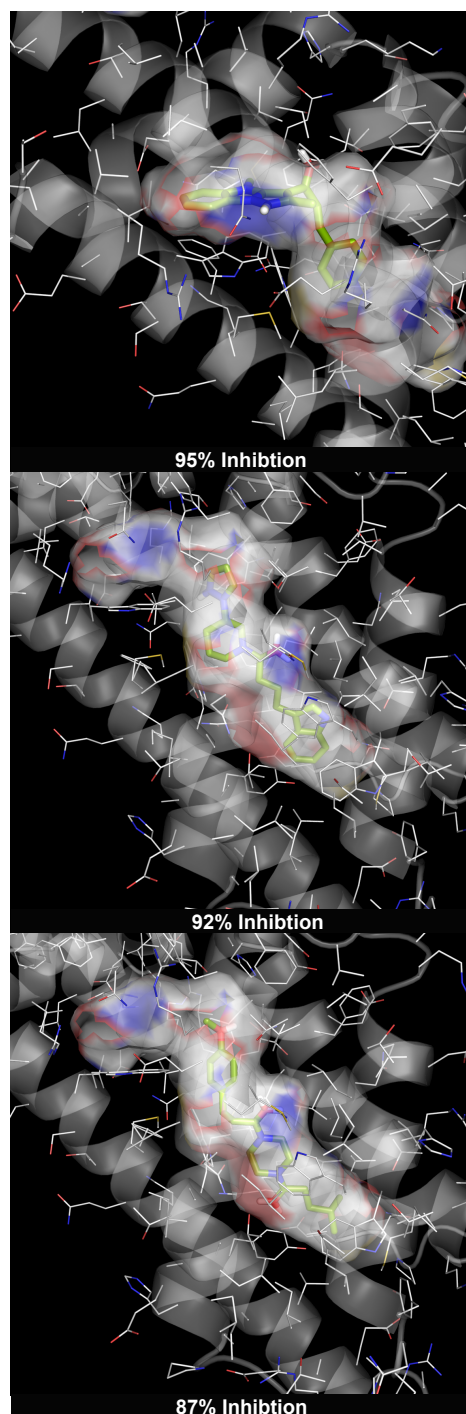


Figure 6.8

Figure 6.9: The predicted orientations of 4 modest inhibitors of EthR

The resulted dockings of our modest validated inhibitors bound into EthR. The percent inhibition represents the inhibition observed by SPR experimental analysis. An individual image is provided for each inhibitor.

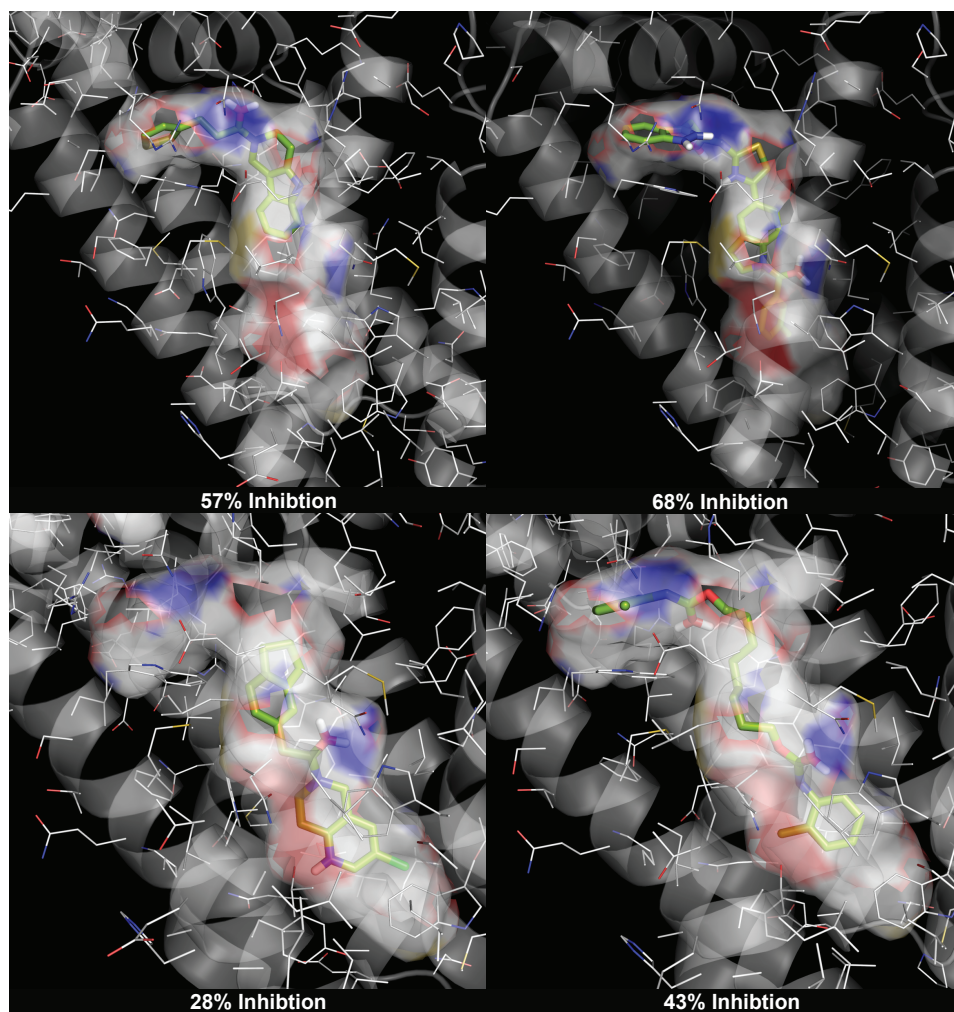


Figure 6.9

Figure 6.10: The predicted orientations of the poor inhibitors of EthR

The resulted dockings of our poor inhibitors bound into EthR. The percent inhibition or activation represents either the inhibition or activation observed by SPR experimental analysis. An individual image is provided for each compound.

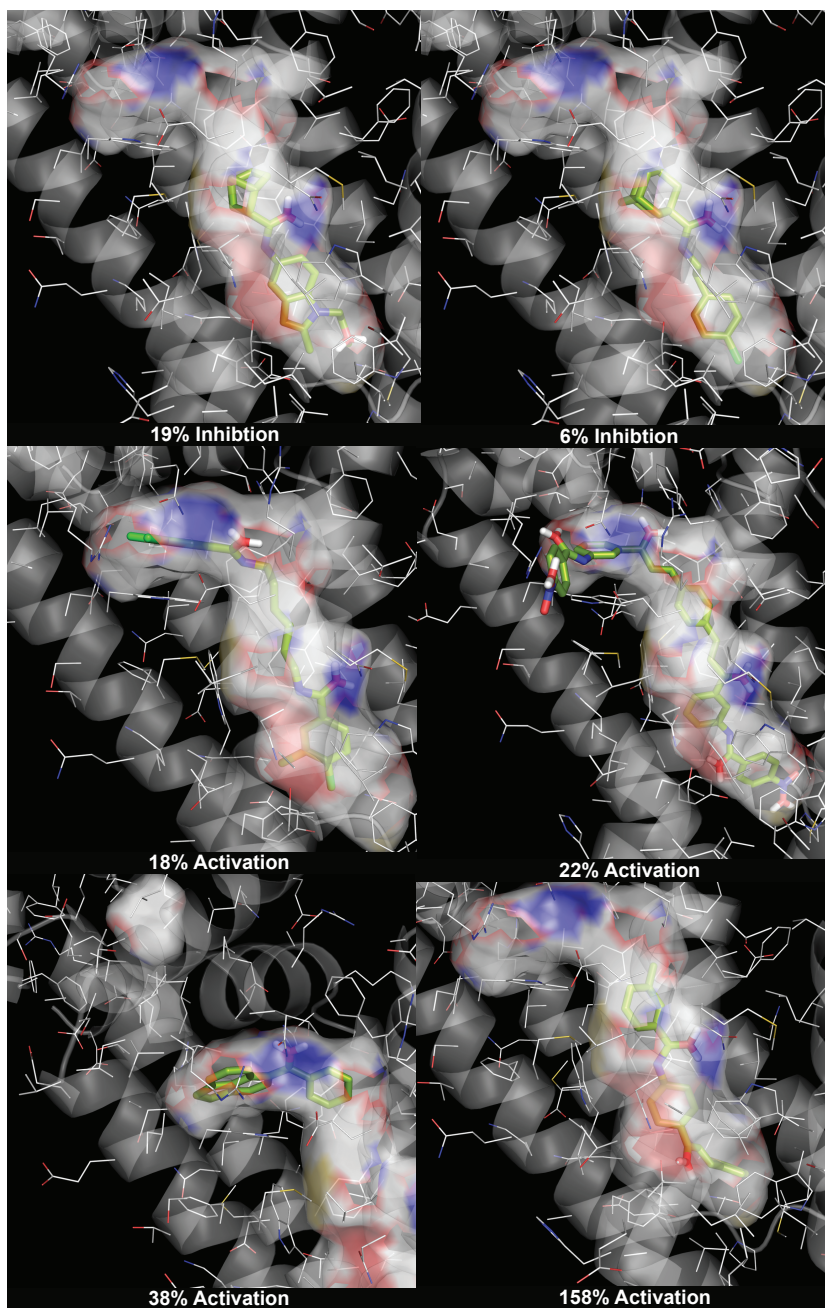


Figure 6.10

Table 6.1: A breakdown of our structure-based virtual screening accuracy for EthR

The accuracy of our virtual screening approach was validated by SPR. The number of compounds that were determined to inhibit EthR are listed by percent inhibition.

Inhibition Criteria	Number of Compounds
> 90% Inhibition	2/14
> 75% Inhibition	3/14
> 50% Inhibition	5/14
> 25% Inhibition	8/14
< 25% Inhibition	6/14

Table 6.1

CHAPTER VII:

Methods

Chemical informatics

Energy minimized 3-D structures were generated using Omega2 software (OpenEye)(Hawkins et al., 2010). Optimized descriptor subsets were iteratively identified from the large Dragon (Talete) and Cerius2 (Accelrys) molecular descriptor sets using a Sequential Forward Selection (SFS) method (Whitney, 1971). The SFS approach functioned by selecting descriptors that increased the correlation between the training set odor distances calculated by activity (spikes/sec) and the growing optimized descriptor set. A 5-fold cross-validation and a receiver-operating-characteristics (ROC) analysis was applied to analyse the performance of the receptor-optimized descriptor sets. Each receptor-optimized descriptor set was then used to rank the combined Natural Odor and Pubchem Libraries (>240,000 compounds) based on each odors distance to a known activator for the corresponding receptor.

Electrophysiology

Extracellular single-sensillum electrophysiology was performed as before (de Bruyne et al., 2001b; Dobritsa et al., 2003; Hallem and Carlson, 2006) with a few modifications. Diagnostic odorants were used to distinguish individual classes of ORNs in sensilla (ab1-ab7) and therefore unequivocally identify the target ORN for testing (de Bruyne et al., 2001b; Hallem et al., 2004). 50ml odor at 10^{-2} dilution in paraffin oil was applied to cotton wool in odor cartridge. Odor stimulus flow = 12ml/second. Due to variability in temporal kinetics of response across various odors, the counting window was shortened

to 250 milliseconds from the start of odor stimulus. A diagnostic panel of odorants were used to distinguish individual classes of sensilla (ab1-ab7) and therefore unequivocally identified the target ORN(de Bruyne et al., 2001a; Hallem et al., 2004).

Natural odor compound library

We assembled a subset of 3197 volatile compounds from annotated origins including plants(Knudsen et al., 2006), insects(El-Sayed, 2009), humans, and a fragrance collection(Sigma-Aldrich, 2007) that may have additional fruit and floral volatiles(Cork and Park, 1996; Curran et al., 2005; Gallagher et al., 2008; Knudsen et al., 2006; Logan et al., 2008; Meijerink et al., 2000; Zeng et al., 1991; Zeng et al., 1996).

Pubchem compound library

We assembled a subset of 241,150 odors from Pubchem, which have similar characteristics to known odor molecules. Compounds met a criteria of MW < 200 and only being composed of the following atoms (C, O, N, H, I, Cl, S, F).

eMolecules compound library

We assembled a subset of >440,000 odors from eMolecules(eMolecules), which have similar characteristics to known odor molecules. Compounds met a criteria of MW < 350 and only being composed of the following atoms (C, O, N, H, I, Cl, S, F).

Calculation of 3D conformations

The 3-Dimensional structures were predicted for compounds through use of the Omega2 software package(Bostrom et al., 2003; Hawkins et al., 2010). The Omega2 software

package functions in three major stages: assembly of an initial 3D structure from a library of fragments; exhaustive enumeration of all rotatable torsions using values drawn from a knowledge-based list of angles, creating a large number of conformations; and sampling of this set by geometric and energy criteria (Hawkins et al., 2010). The lowest energy 3D conformer for each compound in our Pubchem and Natural compound libraries were stored for use in molecular descriptor calculation.

Calculation of molecular descriptors

Commercially available software packages Cerius2 (200 individual descriptors) and Dragon (3224 individual descriptors) from Accelrys and Talete were used to calculate molecular descriptors from 3D molecular structures. Descriptor values were normalized across compounds to standard scores by subtracting the mean value for each descriptor type and dividing by the standard deviation. Molecular descriptors that did not show variation across compounds were removed. Maximum Common Substructures were determined using an existing algorithm (Cao et al., 2008b). Atom Pairs were computed from the version implemented in ChemmineR (Cao et al., 2008a).

Classification of active compounds

Since we were interested in identifying descriptors which best described active compounds, we needed to first determine which compounds to classify as “active” based on their electrophysiology activity for the receptor being studied. All of the training odors were clustered using hierarchical clustering by activity individually for each Or. The resulting tree can then be pruned such that the branch containing the majority of activating odors (>50 spikes/second) are selected. The activity threshold therefore was

set as the lowest spike/second activity of any odor present in the selected branch.

Calculation of Accumulative Percentage of Actives (APoA)

The accumulative percentage of actives is calculated for each descriptor set individually as previously described (Chen and Reynolds, 2002). Compounds are ranked according to their distance from each known active using the optimized descriptor values as distance, resulting in one set of ranked compound distances from each activating odor. Moving down the list for each of these rankings, ratios are calculated for the number of active compounds observed divided by the total number of compounds inspected, or the APoA. APoA values are averaged across all active compound rankings, creating a single set of mean values representing the APoA for a single Or and descriptor set. Using this approach, ApoA mean values are calculated for each of the 24 Ors separately for each descriptor set used, including optimized sets, all Dragon, all Cerius2, Atom Pair, and Maximum Common Substructure. The Area Under the Curve (AUC) scores were calculated by approximation of the integral under each plotted APoA line.

Determination of optimized descriptor subsets

A compound-by-compound activity distance matrix was calculated using training odor activity data for each of the Ors (Hallem and Carlson, 2006). A separate compound-by-compound descriptor distance matrix was calculated using the 3,424 descriptor values for training odors calculated by Dragon and Cerius2. Active compounds for each Or were identified individually through activity thresholds, as described above. The

correlation between the compound-by-compound activity and compound-by-compound descriptor distance matrices were compared for each actively classified compound, considering their distances to all other compounds. The goal was to identify descriptors that correlate most closely with activity. Using a Sequential Forward Selection (SFS) approach, which involves incrementally adding a single best choice item to growing list in an attempt to produce an optimal final set of items (Whitney, 1971), all descriptors are individually compared for their ability to increase this correlation. The descriptor that correlates best is retained and the process is iteratively used to search for additional descriptors. Each iteration aims to further increase correlation values. In this manner, the size of the optimized descriptor set increases by one in each iteration, as the best descriptor set from the previous step is combined with all possible descriptors to find the next best descriptor. This process is halted when all possible descriptor additions in an iteration fails to improve the correlation value from the previous step. This whole process is run independently for each Or resulting in unique descriptor sets that are optimized for each Or.

Clustering Ors by most common descriptors

The first 20 descriptors selected by our optimized descriptor selection algorithm for each Or were used to create an identity matrix. Each row representing an Or and column a specific descriptor. Ors that share common descriptors contain 1s in the same column. This matrix was then converted into an Or-by-Or Euclidean distance matrix and clustered using hierarchical clustering and complete linkage.

Clustering compounds by activity of Or

The responses of each of the Ors that had previously been tested against a panel of compounds were converted into an Or-by-Or Euclidean distance matrix (Hallem and Carlson, 2006). Ors were clustered using hierarchical clustering and complete linkage. Specifically, this was achieved by creating a compound-by-compound distance matrix using the differences in activity between compounds tested on a single Or. Hierarchical clustering using each Or distance matrix and then identifying the sub cluster which contained the most compounds.

Clustering Ors by predicted ligand space

Percentages of overlapping predictions within the top 500 predicted compounds were calculated pair-wise for all Ors. Euclidean distances were calculated from the similarity between Ors.

Calculation of Or prediction distribution frequencies

Initially, all extreme outliers were removed from the dataset for each Or. On average 5.82 compounds were removed for each Or, resulting in a mean dataset reduction of 0.0024%. Next, all compounds whose distance was greater than 3 standard deviations from the strongest activating compound were removed to reduce outliers. Distribution frequencies were produced for each Or. All compound distances were converted into a percentage of the most distant compound for each Or. Frequencies of compounds in the top 15% were plotted.

Or-ligand interaction map

The Or-ligand interaction map was developed using Cytoscape(Shannon et al., 2003).

Each predicted Or-ligand interaction from the top 500 predicted ligands for all of the Ors listed were used to calculate the map (Figure 5A). All predicted interactions are labelled in purple. In addition all interactions identified in this study and the previous study (Hallem and Carlson, 2006) were included and labeled in gray. All compounds are represented as small black circles and Ors are represented as large colored circles. Or names are provided on the upper right corner of each Or.

Computational validation of ligand-based virtual screening approach (non- SVMs)

We performed 5 independent 5-fold cross-validations. For each independent validation the dataset was divided into 5 equal sized partitions containing roughly 22 compounds each. During each run, one of the partitions is selected for testing, and the remaining 4 sets are used for training. The training process is repeated 5 times with each unique odorant set being used as the test set exactly once. For every training iteration, a unique set of descriptors was calculated from the training compound set. These descriptors were then used to calculate distances of the test set compounds to the closest active compound, exactly as we use to predict ligands in our ligand discovery pipeline. Once test set compounds have been ranked by distance from closest to furthest to a known active in the training set, a receiver operating characteristics (ROC) analysis is used to analyze the performance of our computational ligand prediction approach. Using ROC we were able to determine our predictive ability for the 12 receptors. We decided to perform this validation only on receptors for which sufficient training odors had previously been identified. We consider this to consist of at least one very strongly activating known odor (>150 spikes/sec) and at least five strongly activating odors (>100 spikes/sec), thus allowing for at least one activating odor for each of the 5 test sets in the

cross-validation (DmOr7a, DmOr9a, DmOr10a, DmOr22a, DmOr35a, DmOr43b, DmOr12, DmOr59b, DmOr67a, DmOr67c, DmOr85b, DmOr98a). Test set validations for all 12 Ors were combined and a single ROC curve representing an average across all Ors was plotted (Fig. 2A).

Calculation of LogP and vapor pressure values

SMILES structures of the predicted odors were used with EPI Suite (<http://www.epa.gov/oppt/exposure/pubs/episuite.htm>) to calculate predicted LogP and Vapor Pressure values.

Repellency behavior testing

Repellency was tested using *Drosophila melanogaster* 2-choice trap assay as described previously (Reeder et al., 2001; Syed et al., 2011).

Preference Index = number of flies in treated trap / (number of flies in treated + control traps).

The *Drosophila* T-maze assay was conducted as described previously (Turner et al., 2011; Turner and Ray, 2009).

Preference index = (number of flies in test arm - number of flies in control arm) / (number of flies in test arm + number of flies in control arm).

Repellency was tested in mated and starved *Ae. aegypti* females using a hand-in-glove assay. Briefly, a gloved hand with an opening exposing skin odorants protected by 2 layers of netting was presented to mosquitoes for 5 min inside a cage and video taped for landing and avoidance responses. Mosquitoes were unable to bite due to the

outer protective layer of netting and the inner layer of netting was treated with either test compound (10%) or solvent, such that mosquitoes were able to respond to volatiles but unable to make physical contact. The number of mosquitoes present for more than 5 seconds, and the numbers departing during the same period were counted from the videos at minutes 2,3,4, and 5 mins and repellency percentage and escape index calculated by comparing with similar numbers in solvent treated controls.

Percentage repellency = $100 \times [1 - (\text{mean cumulative number of mosquitoes on the window of treatment for 5 seconds at time points 2,3,4,5 min} / \text{mean cumulative number of mosquitoes that remained on window of solvent treatment for 5 seconds at time points 2,3,4,5 min})]$.

Percentage present = average number of mosquitoes on window for 5 seconds at a given time-point across trials. All values were normalized to percentage of the highest value for the comparison, which was assigned a 100 percent present.

Mean Escape Index = $(\text{Average Number of mosquitoes in treatment that landed yet left the mesh during a five second window over the following time points: 2 minutes, 3 minutes, 4 minutes, 5 minutes}) / (\text{Average Number of mosquitoes that landed yet left the mesh during a five second window over the same time points in (treatment + control)})$.

Each time point has N=5 trials, 40 mosquitoes per trial, Except for EA, where N=4.

The humidity and warmth attraction assay *A. aegypti* were placed in a cylindrical clear acrylic cage covered with insect screen mesh on one end. At 0.5 cm above the cage a moist filter paper and heated source was presented to cause attraction. In between the moist filter paper and the cage either DEET or acetone treated nets were placed far enough above the top of the cage so there could be no contact. For 5 second windows at 30 second intervals throughout the 5 minutes assay the number of

mosquitoes present at the top of the mesh cage were measured and an average was calculated for each condition.

Number of mosquitoes probing = Average number of mosquitoes present at top mesh for 5-sec or more, across the 10 time-point measurements/trial, for 3 trials.

Classification of repellent compounds

Training odors were clustered by their repellency values using Euclidean distance and hierarchical clustering. The resulting cluster containing the majority of strongest repellents was then selected and the repellency threshold for identification of descriptors determined as the lowest protection time within the cluster.

Support Vector Machine (SVM) predictions

The R Package e1071 interface with libsvm, a well established program, was used in the analysis(Chang and Lin, 2001; Karatzoglou et al., 2006). The Tune.SVM function was used to determine optimal gamma (0.01) and cost (100) values. The SVM was trained with the determined optimized descriptors for the training compound set using regression and a radial basis function kernel. The trained SVM then ranked both the eMolecules and natural odor libraries by repellency using their optimized descriptors.

Computational validation of ligand-based virtual screening (SVM)

We performed 20 independent 5-fold cross-validations by dividing the dataset into 5 equal sized partitions containing roughly 40 compounds each. During each run, one of the partitions is selected for testing, and the remaining 4 sets are used for training. The optimized descriptor values for the training set were used in training the SVM to

recognize repellent chemicals. The trained SVM is then used to predict repellent activity for the withheld test set. This process is repeated five times, each trial excluding a different subset of compounds as the training set and assigning the remainder as the test set. The whole process is repeated 20 times to improve consistency. A receiver operating characteristics (ROC) analysis is then used to analyze the performance of our computational repellency prediction approach.

Vinyl solubility assay

One 3 x 3 mm square of 4 gauge vinyl was submerged in 1mL of each test compound in a glass container and stirred at a constant rate on a shaker and checked every 30 minutes until the vinyl square in DEET was completely dissolved (6 Hrs). The vinyl pieces in each of the other compounds was removed, rinsed in ethanol and weighed. The process was repeated at 30 Hrs (24 Hrs after the vinyl square completely disappeared in DEET).

Olfactory avoidance assay trap assay for *Drosophila*

Trap Assay experiments were performed as described previously (Reeder et al., 2001) (Syed et al., 2011) with minor modifications. Briefly, traps were made with two 1.5 ml microcentrifuge tubes (USA Scientific) and 200 ul pipette tips (USA Scientific), each cap contained standard cornmeal medium. T-shape piece of filter paper (Whatman #1) was impregnated with 5 ul of acetone (control) or 5 ul of 10%, 1%, 0.10% test odor, diluted in acetone. Traps were placed within a petri dish (100 x 15mm, Fisher) containing 10ml of 1% agarose to provide moisture. Ten flies wCs 4-7days old were used per trial which lasted 48 hours by which time point nearly all flies in the assays had made a choice. For

the 24 hour time point data was considered only if 30% of flies had made a choice, at 48 hours the majority of flies had made choices.

Preference Index = number of flies in treated trap / (number of flies in treated + control traps).

T-maze Assay Methods

Carbon Dioxide and DEET trials were conducted as described previously (Turner and Ray, 2009). Sawyer Jungle Juice 100 Insect Repellent (DEET) was dissolved at 10% in dimethyl sulfoxide (DMSO) and impregnated in a filter paper disc and placed at the bottom of test tube. An equal volume of DMSO impregnated on filter paper disc was used in control arm.

Preference index = (number of flies in test arm - number of flies in control arm) / (number of flies in test arm + number of flies in control arm). For trials with DEET and carbon dioxide, tubes were set up with 10% DEET and DMSO in respective arms. Carbon dioxide was injected into DMSO control tube.

Modified hand-in-glove olfactory repellency assay for mosquitoes

Ae. aegypti mosquitoes (eggs obtained from Benzon Research Inc.) were maintained at ~27 °C and 70% RH on 14h: 10h L: D cycle. Behavioral tests were done with 40 mated, non-blood fed, ~24 hour starved, 4-10 day old females in 30 x 30 x 30 cm cages with a glass top to allow for video recording (Figure 1C, Figure S1). Each test compound solution (500µl) of 10% concentration in acetone solvent was applied evenly to a white rectangular 7 x 6 cm polyester netting (mesh size 26 x 22 holes per square inch) in a glass petri-dish and suspended in the air for 30 minutes to allow solvent evaporation.

The more volatile 2,3-dimethyl-5-isobutyl pyrazine was dissolved in paraffin oil. Acetone or paraffin oil (500µl) served as control. A nitrile glove (Sol-vex) was modified as described in Supplementary Figure 1 such that a 5.8 x 5cm window was present for skin odor exposure. A set of magnetic window frames were designed to secure the treated net ~1.5 mm above skin, and a second untreated netting ~4.5 mm above the treated net in a manner so that mosquitoes were attracted to skin emanations in the open window but unable to contact treated nets with tarsi, or contact and pierce skin. Additionally the test compound had minimal contact with skin. A clean set of glove and magnets were used for every trial. Care was taken that experimenter did not use cosmetics, soap etc on arms. For each trial the arm was first inserted for 5 min and the number landing or escaping test window recorded on video for 5 min period. Solvent controls were always tested prior to treatment, No cage was tested more than once within 1 hour of a testing session and not more than twice on any single day. Videos were analyzed offline on computers and numbers of mosquitoes present for a 5-sec continuous duration were counted every minute. Mosquitoes reliably started accumulating in controls at the 2 min point, and data from this time point was considered for analysis.

Percentage repellency = $100 \times [1 - (\text{mean cumulative number of mosquitoes on the window of treatment for 5 seconds at time points 2,3,4,5 min} / \text{mean cumulative number of mosquitoes that remained on window of solvent treatment for 5 seconds at time points 2,3,4,5 min})]$.

Percentage present = average number of mosquitoes on window for 5 seconds at a given time-point across trials. All values were normalized to percentage of the highest value for the comparison, which was assigned a 100 percent present.

Mean Escape Index = (Average Number of mosquitoes in treatment that landed yet left the mesh during a five second window over the following time points: 2 minutes, 3 minutes, 4 minutes, 5 minutes) / (Average Number of mosquitoes that landed yet left the mesh during a five second window over the same time points in (treatment + control))
Each time point has N=5 trials, 40 mosquitoes per trial, Except for EA, where N=4.

Humidity and warmth attraction assay

Experiments with female *Aedes aegypti* for DEET repellency were conducted at $27\pm 1^{\circ}\text{C}$ and 50-55% under fluorescent lighting from 1400-1700hrs. Each assay lasted 5 minutes. In the humidity-and-heat-attraction assay, 10 females (24hr starved mosquitoes) were held in cylindrical clear acrylic cages (7cm diameter x 5 cm high) covered with insect screen mesh on one end and sealed with clean manila paper on the open end. Mosquitoes were exposed to DEET-treated or acetone(solvent)-treated mesh placed 0.5cm above the screened end of the cage. A 55mm Whatman #1 filter paper disc moistened with 400 μL of water was placed over the net to provide humidity. Heat was provided by placing a 148mL dram plastic vial (Thornton Plastics, Utah USA) containing 75mL of water at 37°C over the filter paper. Mosquito behavior around the treated net was recorded on HD video and the number of probing mosquitoes sitting on the top screen mesh for >5-sec were counted every 30-sec interval of the 5-min duration of the assay.

Number of mosquitoes probing = Average number of mosquitoes present at top mesh for 5-sec or more, across the 10 time-point measurements/trial, for 3 trials.

Larval behavior assays

Behavior assays were conducted as in (Kreher et al., 2008) with some modifications. Odorant solution and paraffin oil solvent were presented in small plastic containers fashioned out of inverted eppendorf tube caps. For pre-exposure experiments ~50 larvae per trial were placed on a 1% agarose base in a 90 mm disposable petri dish. 750 ml of odor solution was evenly spread on the lid and it was placed for indicated times (1 sec or 10 sec) on dish. Larvae were gently removed using a paint brush and placed on a separate petridish and assayed for odor preference in the dark for 1 min 30 sec as described above. The mutant *Orco*^{-/-} (previously called *Or83b*¹) (Larsson et al., 2004) was obtained from Bloomington stock center, and the *Or42b-Gal4* flies were a kind gift from Dr. John Carlson.

Assembling of EthR screening library

Three large libraries were assembled consisting of the advanced library from Enamine (204,772 compound), the screening library from Maybridge (56,000 compounds) and the merged library from Asinex (439,946 compounds). All chemical from each library were combined and duplicates or any compound that contained an atom other than C, O, H, N, S, Br, and F were removed. The resulting library contained 588,723 unique structures.

Calculating molecular descriptors for EthR analysis

Molecular descriptors were calculated independently using both the Dragon software (Taleté) suite and Pipeline pilot (Accelrys). The 3D conformations for each compound were calculated using the Omega application (OpenEye). Only a single best conformation was saved for each compound resulting in 588,723 unique 3D structures.

These structures were then fed into Dragon, which calculated 3,224 unique molecular descriptor values for each of the compounds. These same structures were also fed into pipeline pilot for descriptor calculation (Balaban, Chi, and Kappa values) and machine learning.

Applying machine learning with Pipeline Pilot

The Pipeline pilot suite was used to compare three different machine learning approaches (Random Forests, Support Vector Machines, and linear discriminate analysis) for their ability to identify what shape characteristics separate compounds that bind into EthR from those that do not. As we only knew compounds that have been effectively docked into the binding site and machine learning techniques require both a positive and negative training set, we had to estimate a set of compounds that would not. We randomly selected 100 compounds from the Asinex library, which as EthR is selective in its binding, are unlikely to bind into EthR by chance. Each of the three machine learning methods was then individually trained to separate out inhibitors from our hypothetical non-inhibitors based on molecular descriptor values. Each method had a built in validation approach, such as cross fold validation or out of bag, that was applied along with a Receiver Operating Characteristic (ROC) analysis to compare between methods.

Performing structure-based virtual screening with Gold

The Gold suite was used for all protein-ligand docking and Goldmine was applied for analysis of docking results. An individual docking run was performed for each of the three binding sites. Other than for the following exceptions, default settings were applied

for the gold algorithm. Protein H-bond constraints were individually applied for each binding site with the following values: Asparagine 179 and Asparagine 176 for binding site 1, Glutamine 125 for binding site 2, and Asparagine 179, Asparagine 176, Glutamine 125 for the combined site. We performed 10 independent docking iterations for both binding site 1 and 2 and 3 independent iterations for the combined binding site. Early termination was allowed if 3 independent iterations failed to improve the rms_tolerance value by better than a value of 1.5. The binding cavity for each of the sites was set to a 10 angstrom radius from the natural ligands bound into site 1 or 2 of structures NT06 and MB1047, respectively. The binding cavity was set to be a 10 angstrom radius from either of the natural bound ligands in both sites 1 and 2 in structure NT06 for the combined site, which created a single site spanning nearly the entire cavity. None of the amino acid side chains were allowed to rotate in any of the docking sets. Both goldscore and chemscore scoring functions were calculated for each binding site.

REFERNCES:

- Abuin, L., Bargeton, B., Ulbrich, M.H., Isacoff, E.Y., Kellenberger, S., and Benton, R. (2011). Functional architecture of olfactory ionotropic glutamate receptors. *Neuron* **69**, 44-60.
- Aceves-Pina, E.O., and Quinn, W.G. (1979). Learning in normal and mutant *Drosophila* larvae. *Science* **206**, 93-96.
- Adams, M.D., Celniker, S.E., Holt, R.A., Evans, C.A., Gocayne, J.D., Amanatides, P.G., Scherer, S.E., Li, P.W., Hoskins, R.A., Galle, R.F., *et al.* (2000). The genome sequence of *Drosophila melanogaster*. *Science* **287**, 2185-2195.
- Altschul, S.F., Gish, W., Miller, W., Myers, E.W., and Lipman, D.J. (1990). Basic local alignment search tool. *J Mol Biol* **215**, 403-410.
- Aramaki, H., Yagi, N., and Suzuki, M. (1995). Residues important for the function of a multihelical DNA binding domain in the new transcription factor family of Cam and Tet repressors. *Protein Eng* **8**, 1259-1266.
- Araneda, R.C., Kini, A.D., and Firestein, S. (2000). The molecular receptive range of an odorant receptor. *Nat Neurosci* **3**, 1248-1255.
- Ayyub, C., Paranjape, J., Rodrigues, V., and Siddiqi, O. (1990). Genetics of Olfactory Behavior in *Drosophila-Melanogaster*. *J Neurogenet* **6**, 243-262.
- Baldwin, E., Plotto, A., Manthey, J., McCollum, G., Bai, J.H., Irey, M., Cameron, R., and Luzio, G. (2010). Effect of *Liberibacter* Infection (Huanglongbing Disease) of Citrus on Orange Fruit Physiology and Fruit/Fruit Juice Quality: Chemical and Physical Analyses. *J Agr Food Chem* **58**, 1247-1262.
- Baulard, A.R., Betts, J.C., Engohang-Ndong, J., Quan, S., McAdam, R.A., Brennan, P.J., Locht, C., and Besra, G.S. (2000). Activation of the pro-drug ethionamide is regulated in mycobacteria. *Journal of Biological Chemistry* **275**, 28326-28331.
- Belluscio, L., Gold, G.H., Nemes, A., and Axel, R. (1998). Mice deficient in G(olf) are anosmic. *Neuron* **20**, 69-81.
- Benton, R., Sachse, S., Michnick, S.W., and Vosshall, L.B. (2006). Atypical membrane topology and heteromeric function of *Drosophila* odorant receptors in vivo. *PLoS Biol* **4**, e20.
- Benton, R., Vannice, K.S., Gomez-Diaz, C., and Vosshall, L.B. (2009). Variant Ionotropic Glutamate Receptors as Chemosensory Receptors in *Drosophila*. *Cell* **136**, 149-162.
- Bohacek, R.S., McMartin, C., and Guida, W.C. (1996). The art and practice of structure-based drug design: A molecular modeling perspective. *Med Res Rev* **16**, 3-50.

- Bohbot, J., Pitts, R.J., Kwon, H.W., Rutzler, M., Robertson, H.M., and Zwiebel, L.J. (2007). Molecular characterization of the *Aedes aegypti* odorant receptor gene family. *Insect Mol Biol* *16*, 525-537.
- Bolton, E.E., Wang, Y., Thiessen, P.A., and S.H., B. (2008). PubChem: Integrated Platform of Small Molecules and Biological Activities In Annual Reports in Computational Chemistry (Washington DC, American Chemical Society).
- Bostrom, J., Greenwood, J.R., and Gottfries, J. (2003). Assessing the performance of OMEGA with respect to retrieving bioactive conformations. *J Mol Graph Model* *21*, 449-462.
- Bove, J.M. (2006). Huanglongbing: A destructive, newly-emerging, century-old disease of citrus. *J Plant Pathol* *88*, 7-37.
- Brockhoff, A., Behrens, M., Niv, M.Y., and Meyerhof, W. (2010). Structural requirements of bitter taste receptor activation. *P Natl Acad Sci USA* *107*, 11110-11115.
- Buck, L., and Axel, R. (1991). A Novel Multigene Family May Encode Odorant Receptors - a Molecular-Basis for Odor Recognition. *Cell* *65*, 175-187.
- Butler, D. (2011). Mosquitoes score in chemical war. *Nature* *475*, 19.
- Cao, Y., Charisi, A., Cheng, L.C., Jiang, T., and Girke, T. (2008a). ChemmineR: a compound mining framework for R. *Bioinformatics* *24*, 1733-1734.
- Cao, Y., Jiang, T., and Girke, T. (2008b). A maximum common substructure-based algorithm for searching and predicting drug-like compounds. *Bioinformatics* *24*, i366-374.
- Carde, R.T., and Gibson, G. (2010). Long Distance orientation of mosquitoes to host odours and their clues. in *Ecology of Vector-Borne Diseases. Vol. 2. Olfaction in Vector-Host Interactions* (Wageningen, The Netherlands, Wageningen Academic Publishers).
- Carey, A.F., Wang, G.R., Su, C.Y., Zwiebel, L.J., and Carlson, J.R. (2010). Odorant reception in the malaria mosquito *Anopheles gambiae*. *Nature* *464*, 66-U77.
- Carhart, R.E., Smith, D.H., and Venkataraghavan, R. (1985). Atom Pairs as Molecular-Features in Structure Activity Studies - Definition and Applications. *J Chem Inf Comp Sci* *25*, 64-73.
- Chan, E.D., and Iseman, M.D. (2002). Current medical treatment for tuberculosis. *Brit Med J* *325*, 1282-1286.
- Chang, C., and Lin, C. (2001). Libsvm: A Library for Support Vector Machines.
- Chen, X., and Reynolds, C.H. (2002). Performance of similarity measures in 2D fragment-based similarity searching: comparison of structural descriptors and similarity coefficients. *J Chem Inf Comput Sci* *42*, 1407-1414.

- Clyne, P.J., Warr, C.G., and Carlson, J.R. (2000). Candidate taste receptors in *Drosophila*. *Science* *287*, 1830-1834.
- Clyne, P.J., Warr, C.G., Freeman, M.R., Lessing, D., Kim, J., and Carlson, J.R. (1999). A novel family of divergent seven-transmembrane proteins: candidate odorant receptors in *Drosophila*. *Neuron* *22*, 327-338.
- Corbel, V., Stankiewicz, M., Pennetier, C., Fournier, D., Stojan, J., Girard, E., Dimitrov, M., Molgo, J., Hougard, J.M., and Lapied, B. (2009). Evidence for inhibition of cholinesterases in insect and mammalian nervous systems by the insect repellent deet. *Bmc Biol* *7*, -.
- Cork, A., and Park, K.C. (1996). Identification of electrophysiologically-active compounds for the malaria mosquito, *Anopheles gambiae*, in human sweat extracts. *Medical and Veterinary Entomology* *10*, 269-276.
- Cortes, C., and Vapnik, V. (1995). Support-Vector Networks. *Mach Learn* *20*, 273-297.
- Couto, A., Alenius, M., and Dickson, B.J. (2005). Molecular, anatomical, and functional organization of the *Drosophila* olfactory system. *Curr Biol* *15*, 1535-1547.
- Curran, A.M., Rabin, S.I., Prada, P.A., and Furton, K.G. (2005). Comparison of the volatile organic compounds present in human odor using SPME-GC/MS. *Journal of Chemical Ecology* *31*, 1607-1619.
- Dagraca, J.V. (1991). Citrus Greening Disease. *Annu Rev Phytopathol* *29*, 109-136.
- Dahanukar, A., Foster, K., van Naters, W.M.V.D.G., and Carlson, J.R. (2001). A Gr receptor is required for response to the sugar trehalose in taste neurons of *Drosophila*. *Nat Neurosci* *4*, 1182-1186.
- Dahanukar, A., Hallem, E.A., and Carlson, J.R. (2005). Insect chemoreception. *Curr Opin Neurobiol* *15*, 423-430.
- de Bruyne, M., and Baker, T.C. (2008). Odor detection in insects: volatile codes. *J Chem Ecol* *34*, 882-897.
- de Bruyne, M., Clyne, P.J., and Carlson, J.R. (1999). Odor coding in a model olfactory organ: the *Drosophila* maxillary palp. *J Neurosci* *19*, 4520-4532.
- de Bruyne, M., Foster, K., and Carlson, J. (2001). Odor coding in the *Drosophila* antenna. *Neuron* *30*, 537-552.
- DeBarber, A.E., Mdluli, K., Bosman, M., Bekker, L.G., and Barry, C.E. (2000). Ethionamide activation and sensitivity in multidrug-resistant *Mycobacterium tuberculosis*. *P Natl Acad Sci USA* *97*, 9677-9682.
- Ditzen, M., Pellegrino, M., and Vosshall, L.B. (2008). Insect odorant receptors are molecular targets of the insect repellent DEET. *Science* *319*, 1838-1842.

- Dobritsa, A.A., van der Goes van Naters, W., Warr, C.G., Steinbrecht, R.A., and Carlson, J.R. (2003). Integrating the molecular and cellular basis of odor coding in the *Drosophila* antenna. *Neuron* 37, 827-841.
- Dobson, C.M. (2004). Chemical space and biology. *Nature* 432, 824-828.
- Dover, L.G., Corsino, P.E., Daniels, I.R., Cocklin, S.L., Tatituri, V., Besra, G.S., and Futterer, K. (2004). Crystal structure of the TetR/CamR family repressor *Mycobacterium tuberculosis* EthR implicated in ethionamide resistance. *Journal of Molecular Biology* 340, 1095-1105.
- EFSA (2008). Consideration of Anthranilate derivatives evaluated by JECFA (65th meeting) Opinion of the scientific Panel on Food Additives, Flavourings, Processing Aids and Materials in Contact with Food. *The EFSA Journal*, 1-24.
- El-Sayed, A. (2009). *The Pherobase: Database of Insect Pheromones and Semiochemicals*.
- Engohang-Ndong, J., Baillat, D., Aumercier, M., Bellefontaine, F., Besra, G.S., Locht, C., and Baulard, A.R. (2004). EthR, a repressor of the TetR/CamR family implicated in ethionamide resistance in mycobacteria, octamerizes cooperatively on its operator. *Mol Microbiol* 51, 175-188.
- Fishilevich, E., and Vosshall, L.B. (2005). Genetic and functional subdivision of the *Drosophila* antennal lobe. *Curr Biol* 15, 1548-1553.
- Flipo, M., Desroses, M., Lecat-Guillet, N., Dirie, B., Carette, X., Leroux, F., Piveteau, C., Demirkaya, F., Lens, Z., Rucktooa, P., *et al.* (2011). Ethionamide Boosters: Synthesis, Biological Activity, and Structure-Activity Relationships of a Series of 1,2,4-Oxadiazole EthR Inhibitors. *Journal of Medicinal Chemistry* 54, 2994-3010.
- Flipo, M., Desroses, M., Lecat-Guillet, N., Villemagne, B., Blondiaux, N., Leroux, F., Piveteau, C., Mathys, V., Flament, M.P., Siepmann, J., *et al.* (2012). Ethionamide Boosters. 2. Combining Bioisosteric Replacement and Structure-Based Drug Design To Solve Pharmacokinetic Issues in a Series of Potent 1,2,4-Oxadiazole EthR Inhibitors. *Journal of Medicinal Chemistry* 55, 68-83.
- Floriano, W.B., Vaidehi, N., Goddard, W.A., Singer, M.S., and Shepherd, G.M. (2000). Molecular mechanisms underlying differential odor responses of a mouse olfactory receptor. *P Natl Acad Sci USA* 97, 10712-10716.
- Fox, A.N., Pitts, R.J., Robertson, H.M., Carlson, J.R., and Zwiebel, L.J. (2001). Candidate odorant receptors from the malaria vector mosquito *Anopheles gambiae* and evidence of down-regulation in response to blood feeding. *P Natl Acad Sci USA* 98, 14693-14697.

- Frenois, F., Engohang-Ndong, J., Locht, C., Baulard, A.R., and Villeret, V. (2004). Structure of EthR in a ligand bound conformation reveals therapeutic perspectives against tuberculosis. *Mol Cell* *16*, 301-307.
- Galizia, C.G., Munch, D., Strauch, M., Nissler, A., and Ma, S.W. (2010). Integrating Heterogeneous Odor Response Data into a Common Response Model: A DoOR to the Complete Olfactome. *Chemical Senses* *35*, 551-563.
- Gallagher, M., Wysocki, J., Leyden, J.J., Spielman, A.I., Sun, X., and Preti, G. (2008). Analyses of volatile organic compounds from human skin. *Brit J Dermatol* *159*, 780-791.
- Ghaninia, M., Ignell, R., and Hansson, B.S. (2007). Functional classification and central nervous projections of olfactory receptor neurons housed in antennal trichoid sensilla of female yellow fever mosquitoes, *Aedes aegypti*. *Eur J Neurosci* *26*, 1611-1623.
- Ghose, A.K., Viswanadhan, V.N., and Wendoloski, J.J. (1999). A knowledge-based approach in designing combinatorial or medicinal chemistry libraries for drug discovery. 1. A qualitative and quantitative characterization of known drug databases. *J Comb Chem* *1*, 55-68.
- Gillies, M.T. (1980). The Role of Carbon-Dioxide in Host-Finding by Mosquitos (Diptera, Culicidae) - a Review. *B Entomol Res* *70*, 525-532.
- Goldman, A.L., van Naters, W.V., Lessing, D., Warr, C.G., and Carlson, J.R. (2005). Coexpression of two functional odor receptors in one neuron. *Neuron* *45*, 661-666.
- Goto, S., Okuno, Y., Hattori, M., Nishioka, T., and Kanehisa, M. (2002). LIGAND: database of chemical compounds and reactions in biological pathways. *Nucleic Acids Research* *30*, 402-404.
- Grant, A.J., and Oconnell, R.J. (1996). Electrophysiological responses from receptor neurons in mosquito maxillary palp sensilla. *Ciba F Symp* *200*, 233-253.
- Grosse-Wilde, E., Kuebler, L.S., Bucks, S., Vogel, H., Wicher, D., and Hansson, B.S. (2011). Antennal transcriptome of *Manduca sexta*. *P Natl Acad Sci USA* *108*, 7449-7454.
- Guo, S., and Kim, J. (2007). Molecular evolution of *Drosophila* odorant receptor genes. *Mol Biol Evol* *24*, 1198-1207.
- Guo, S., and Kim, J. (2010). Dissecting the molecular mechanism of *drosophila* odorant receptors through activity modeling and comparative analysis. *Proteins* *78*, 381-399.
- Gupta, R.K.a.B., A.K. (2007). Discovery and Design of New Arthropod/Insect Repellents by Computer-Aided Molecular Modeling. In *Insect Repellents: principles, methods, and uses*, M. Debboun, Frances, S.P., Strickman, D., ed. (Boca Raton, Taylor & Francis Group), pp. 195-228.

- Haddad, R., Khan, R., Takahashi, Y.K., Mori, K., Harel, D., and Sobel, N. (2008). A metric for odorant comparison. *Nat Methods* 5, 425-429.
- Halbert, S.E., and Manjunath, K.L. (2004). Asian citrus psyllids (Sternorrhyncha : Psyllidae) and greening disease of citrus: A literature review and assessment of risk in Florida. *Fla Entomol* 87, 330-353.
- Hallem, E.A., and Carlson, J.R. (2006). Coding of odors by a receptor repertoire. *Cell* 125, 143-160.
- Hallem, E.A., Ho, M.G., and Carlson, J.R. (2004). The molecular basis of odor coding in the *Drosophila* antenna. *Cell* 117, 965-979.
- Hastie, T., Tibshirani, R., and Friedman, J.H. (2001). *The elements of statistical learning : data mining, inference, and prediction : with 200 full-color illustrations* (New York, Springer).
- Hawkins, P.C.D., Skillman, A.G., Warren, G.L., Ellingson, B.A., and Stahl, M.T. (2010). Conformer Generation with OMEGA: Algorithm and Validation Using High Quality Structures from the Protein Databank and Cambridge Structural Database. *Journal of Chemical Information and Modeling* 50, 572-584.
- Hendrickson, J.B. (1991). Concepts and Applications of Molecular Similarity - Johnson, Ma, Maggiora, Gm. *Science* 252, 1189-1189.
- Hill, S.R., Hansson, B.S., and Ignell, R. (2009). Characterization of antennal trichoid sensilla from female southern house mosquito, *Culex quinquefasciatus* Say. *Chem Senses* 34, 231-252.
- Ishimoto, H., Takahashi, K., Ueda, R., and Tanimura, T. (2005). G-protein gamma subunit 1 is required for sugar reception in *Drosophila*. *Embo J* 24, 3259-3265.
- JECF (2007). Safety evaluation of certain food additives and contaminants. Sixty-fifth meeting of the joint FAO/WHO Expert Committee of Food Additives. WHO Food Additives Series: 56 IPCS, WHO Geneva.
- Jefferis, G.S., Potter, C.J., Chan, A.M., Marin, E.C., Rohlfsing, T., Maurer, C.R., Jr., and Luo, L. (2007). Comprehensive maps of *Drosophila* higher olfactory centers: spatially segregated fruit and pheromone representation. *Cell* 128, 1187-1203.
- Jones, W.D., Cayirlioglu, P., Kadow, I.G., and Vosshall, L.B. (2007). Two chemosensory receptors together mediate carbon dioxide detection in *Drosophila*. *Nature* 445, 86-90.
- Kain, P., Chakraborty, T.S., Sundaram, S., Siddiqi, O., Rodrigues, V., and Hasan, G. (2008). Reduced odor responses from antennal neurons of G(q)alpha, phospholipase C beta, and rdgA mutants in *Drosophila* support a role for a phospholipid intermediate in insect olfactory transduction. *Journal of Neuroscience* 28, 4745-4755.

- Kaluza, J.F., and Breer, H. (2000). Responsiveness of olfactory neurons to distinct aliphatic aldehydes. *J Exp Biol* 203, 927-933.
- Karatzoglou, A., Meyer, D., and Hornik, K. (2006). Support Vector Machines in R. *J Stat Softw* 15.
- Katada, S., Hirokawa, T., Oka, Y., Suwa, M., and Touhara, K. (2005). Structural basis for a broad but selective ligand spectrum of a mouse olfactory receptor: Mapping the odorant-binding site. *Journal of Neuroscience* 25, 1806-1815.
- Katritzky, A.R., Wang, Z., Slavov, S., Tsikolia, M., Dobchev, D., Akhmedov, N.G., Hall, C.D., Bernier, U.R., Clark, G.G., and Linthicum, K.J. (2008). Synthesis and bioassay of improved mosquito repellents predicted from chemical structure. *Proc Natl Acad Sci U S A* 105, 7359-7364.
- Keiser, M.J., Setola, V., Irwin, J.J., Laggner, C., Abbas, A.I., Hufeisen, S.J., Jensen, N.H., Kuijer, M.B., Matos, R.C., Tran, T.B., *et al.* (2009). Predicting new molecular targets for known drugs. *Nature* 462, 175-U148.
- Khafizov, K., Anselmi, C., Menini, A., and Carloni, P. (2007). Ligand specificity of odorant receptors. *J Mol Model* 13, 401-409.
- Kim, J., and Carlson, J.R. (2002). Gene discovery by e-genetics: *Drosophila* odor and taste receptors. *J Cell Sci* 115, 1107-1112.
- Kitchen, D.B., Decornez, H., Furr, J.R., and Bajorath, J. (2004). Docking and scoring in virtual screening for drug discovery: Methods and applications. *Nat Rev Drug Discov* 3, 935-949.
- Kline, D.L., Bernier, U.R., Posey, K.H., and Barnard, D.R. (2003). Olfactometric evaluation of spatial repellents for *Aedes aegypti*. *Journal of Medical Entomology* 40, 463-467.
- Klocke, J.A., Darlington, M.V., and Balandrin, M.F. (1987). Biologically-Active Constituents of North-American Plants .3. 1,8-Cineole (Eucalyptol), a Mosquito Feeding and Ovipositional Repellent from Volatile Oil of *Hemizonia-Fitchii* (Asteraceae). *J Chem Ecol* 13, 2131-2141.
- Klun, J.A., Strickman, D., Rowton, E., Williams, J., Kramer, M., Roberts, D., and Debboun, M. (2004). Comparative resistance of *Anopheles albimanus* and *Aedes aegypti* to N,N-diethyl-3-methylbenzamide (Deet) and 2-methylpiperidinyl-3-cyclohexen-1-carboxamide (AI3-37220) in laboratory human-volunteer repellent assays. *J Med Entomol* 41, 418-422.
- Knudsen, J.T., Eriksson, R., Gershenzon, J., and Stahl, B. (2006). Diversity and Distribution of Floral Scent. *The Botanical Review* 72, 1-120.

- Koon, A.C., Ashley, J., Barria, R., DasGupta, S., Brain, R., Waddell, S., Alkema, M.J., and Budnik, V. (2011). Autoregulatory and paracrine control of synaptic and behavioral plasticity by octopaminergic signaling. *Nat Neurosci* 14, 190-U275.
- Krajick, K. (2006). Medical entomology - Keeping the bugs at bay. *Science* 313, 36-38.
- Kreher, S.A., Kwon, J.Y., and Carlson, J.R. (2005). The molecular basis of odor coding in the *Drosophila* larva. *Neuron* 46, 445-456.
- Kreher, S.A., Mathew, D., Kim, J., and Carlson, J.R. (2008). Translation of sensory input into behavioral output via an olfactory system. *Neuron* 59, 110-124.
- Kurland, M.D., Newcomer, M.B., Peterlin, Z., Ryan, K., Firestein, S., and Batista, V.S. (2010). Discrimination of Saturated Aldehydes by the Rat I7 Olfactory Receptor. *Biochemistry* 49, 6302-6304.
- Kwon, J.Y., Dahanukar, A., Weiss, L.A., and Carlson, J.R. (2007). The molecular basis of CO₂ reception in *Drosophila*. *P Natl Acad Sci USA* 104, 3574-3578.
- Lai, P.C., Bahl, G., Gremigni, M., Matarazzo, V., Clot-Faybesse, O., Ronin, C., and Crasto, C.J. (2008). An olfactory receptor pseudogene whose function emerged in humans: a case study in the evolution of structure-function in GPCRs. *J Struct Funct Genomics* 9, 29-40.
- Larsson, M.C., Domingos, A.I., Jones, W.D., Chiappe, M.E., Amrein, H., and Vosshall, L.B. (2004). Or83b encodes a broadly expressed odorant receptor essential for *Drosophila* olfaction. *Neuron* 43, 703-714.
- Lee, Y., Kim, S.H., and Montell, C. (2010). Avoiding DEET through Insect Gustatory Receptors. *Neuron* 67, 555-561.
- Liu, C., Pitts, R.J., Bohbot, J.D., Jones, P.L., Wang, G., and Zwiebel, L.J. (2010). Distinct olfactory signaling mechanisms in the malaria vector mosquito *Anopheles gambiae*. *PLoS Biol* 8.
- Logan, J.G., Birkett, M.A., Clark, S.J., Powers, S., Seal, N.J., Wadhams, L.J., Mordue, A.J., and Pickett, J.A. (2008). Identification of human-derived volatile chemicals that interfere with attraction of *Aedes aegypti* mosquitoes. *Journal of Chemical Ecology* 34, 308-322.
- Louis, M., Huber, T., Benton, R., Sakmar, T.P., and Vosshall, L.B. (2008). Bilateral olfactory sensory input enhances chemotaxis behavior. *Nat Neurosci* 11, 187-199.
- Lu, T., Qiu, Y.T., Wang, G., Kwon, J.Y., Rutzler, M., Kwon, H.W., Pitts, R.J., van Loon, J.J.A., Takken, W., Carlson, J.R., *et al.* (2007). Odor coding in the maxillary palp of the malaria vector mosquito *Anopheles gambiae*. *Curr Biol* 17, 1533-1544.

- Mahmoudi, A., and Iseman, M.D. (1993). Pitfalls in the Care of Patients with Tuberculosis - Common Errors and Their Association with the Acquisition of Drug-Resistance. *Jama-J Am Med Assoc* 270, 65-68.
- Maldonado, A.G., Doucet, J.P., Petitjean, M., and Fan, B.T. (2006). Molecular similarity and diversity in chemoinformatics: From theory to applications. *Mol Divers* 10, 39-79.
- Martin, Y.C., Kofron, J.L., and Traphagen, L.M. (2002). Do structurally similar molecules have similar biological activity? *Journal of Medicinal Chemistry* 45, 4350-4358.
- Meijerink, J., Braks, M.A.H., Brack, A.A., Adam, W., Dekker, T., Posthumus, M.A., Van Beek, T.A., and Van Loon, J.J.A. (2000). Identification of olfactory stimulants for *Anopheles gambiae* from human sweat samples. *Journal of Chemical Ecology* 26, 1367-1382.
- Montague, S.A., Mathew, D., and Carlson, J.R. (2011). Similar Odorants Elicit Different Behavioral and Physiological Responses, Some Supersustained. *Journal of Neuroscience* 31, 7891-7899.
- Monte, P., Woodard, C., Ayer, R., Lilly, M., Sun, H., and Carlson, J. (1989). Characterization of the Larval Olfactory Response in *Drosophila* and Its Genetic-Basis. *Behav Genet* 19, 267-283.
- Nikolova, N., and Jaworska, J. (2004). Approaches to measure chemical similarity - A review. *Qsar Comb Sci* 22, 1006-1026.
- Olsen, S.R., Bhandawat, V., and Wilson, R.I. (2007). Excitatory interactions between olfactory processing channels in the *Drosophila* antennal lobe. *Neuron* 54, 89-103.
- Olsen, S.R., and Wilson, R.I. (2008). Lateral presynaptic inhibition mediates gain control in an olfactory circuit. *Nature* 452, 956-U953.
- Onagbola, E.O., Rouseff, R.L., Smoot, J.M., and Stelinski, L.L. (2011). Guava leaf volatiles and dimethyl disulphide inhibit response of *Diaphorina citri* Kuwayama to host plant volatiles. *J Appl Entomol* 135, 404-414.
- Pellegrino, M., Steinbach, N., Stensmyr, M.C., Hansson, B.S., and Vosshall, L.B. (2011). A natural polymorphism alters odour and DEET sensitivity in an insect odorant receptor. *Nature*.
- Pelz, D., Roeske, T., Syed, Z., de Bruyne, M., and Galizia, C.G. (2006). The molecular receptive range of an olfactory receptor in vivo (*Drosophila melanogaster* Or22a). *J Neurobiol* 66, 1544-1563.
- Pitts, R.J., Rinker, D.C., Jones, P.L., Rokas, A., and Zwiebel, L.J. (2011). Transcriptome profiling of chemosensory appendages in the malaria vector *Anopheles gambiae* reveals tissue- and sex-specific signatures of odor coding. *Bmc Genomics* 12.
- Raviglione, M.C. (2003). The TB epidemic from 1992 to 2002. *Tuberculosis* 83, 4-14.

- Reeder, N.L., Ganz, P.J., Carlson, J.R., and Saunders, C.W. (2001). Isolation of a DEET-insensitive mutant of *Drosophila melanogaster* (Diptera: Drosophilidae). *J Econ Entomol* *94*, 1584-1588.
- Robertson, H.M., and Kent, L.B. (2009). Evolution of the gene lineage encoding the carbon dioxide receptor in insects. *Journal of Insect Science* *9*.
- Robertson, H.M., Warr, C.G., and Carlson, J.R. (2003). Molecular evolution of the insect chemoreceptor gene superfamily in *Drosophila melanogaster*. *Proc Natl Acad Sci U S A* *100 Suppl 2*, 14537-14542.
- Rodriguez-Gil, D.J., Treloar, H.B., Zhang, X.H., Miller, A.M., Two, A., Iwema, C., Firestein, S.J., and Greer, C.A. (2010). Chromosomal Location-Dependent Nonstochastic Onset of Odor Receptor Expression. *Journal of Neuroscience* *30*, 10067-10075.
- Root, C.M., Ko, K.I., Jafari, A., and Wang, J.W. (2011). Presynaptic Facilitation by Neuropeptide Signaling Mediates Odor-Driven Food Search. *Cell* *145*, 133-144.
- Root, C.M., Semmelhack, J.L., Wong, A.M., Flores, J., and Wang, J.W. (2007). Propagation of olfactory information in *Drosophila*. *P Natl Acad Sci USA* *104*, 11826-11831.
- Ruta, V., Datta, S.R., Vasconcelos, M.L., Freeland, J., Looger, L.L., and Axel, R. (2010). A dimorphic pheromone circuit in *Drosophila* from sensory input to descending output. *Nature* *468*, 686-U106.
- Saito, H., Chi, Q., Zhuang, H., Matsunami, H., and Mainland, J.D. (2009). Odor coding by a Mammalian receptor repertoire. *Sci Signal* *2*, ra9.
- Saito, H., Kubota, M., Roberts, R.W., Chi, Q., and Matsunami, H. (2004). RTP family members induce functional expression of mammalian odorant receptors. *Cell* *119*, 679-691.
- Sanger, F., and Coulson, A.R. (1975). A rapid method for determining sequences in DNA by primed synthesis with DNA polymerase. *J Mol Biol* *94*, 441-448.
- Sato, K., Pellegrino, M., Nakagawa, T., Vosshall, L.B., and Touhara, K. (2008). Insect olfactory receptors are heteromeric ligand-gated ion channels. *Nature* *452*, 1002-1006.
- Sato, K., Tanaka, K., and Touhara, K. (2011). Sugar-regulated cation channel formed by an insect gustatory receptor. *P Natl Acad Sci USA* *108*, 11680-11685.
- Schmucker, M., de Bruyne, M., Hahnel, M., and Schneider, G. (2007). Predicting olfactory receptor neuron responses from odorant structure. *Chem Cent J* *1*, -.
- Schneider, G., Schneider, P., and Renner, S. (2006). Scaffold-hopping: How far can you jump? *Qsar Comb Sci* *25*, 1162-1171.

- Schwarz, D., Robertson, H.M., Feder, J.L., Varala, K., Hudson, M.E., Ragland, G.J., Hahn, D.A., and Berlocher, S.H. (2009). Sympatric ecological speciation meets pyrosequencing: sampling the transcriptome of the apple maggot *Rhagoletis pomonella*. *Bmc Genomics* 10.
- Semmelhack, J.L., and Wang, J.W. (2009). Select *Drosophila* glomeruli mediate innate olfactory attraction and aversion. *Nature* 459, 218-U100.
- Shang, Y.H., Claridge-Chang, A., Sjulson, L., Pypaert, M., and Miesenbock, G. (2007). Excitatory local circuits and their implications for olfactory processing in the fly antennal lobe. *Cell* 128, 601-612.
- Shannon, P., Markiel, A., Ozier, O., Baliga, N.S., Wang, J.T., Ramage, D., Amin, N., Schwikowski, B., and Ideker, T. (2003). Cytoscape: A software environment for integrated models of biomolecular interaction networks. *Genome Res* 13, 2498-2504.
- Sigma-Aldrich (2007). Flavors and fragrances 2007-2008 catalog. Sigma-Aldrich Fine Chemicals Company, Milwaukee, WI.
- Singer, M.S. (2000). Analysis of the molecular basis for octanal interactions in the expressed rat I7 olfactory receptor. *Chemical Senses* 25, 155-165.
- Singer, M.S., and Shepherd, G.M. (1994). Molecular modeling of ligand-receptor interactions in the OR5 olfactory receptor. *Neuroreport* 5, 1297-1300.
- Smadja, C., Shi, P., Butlin, R.K., and Robertson, H.M. (2009). Large Gene Family Expansions and Adaptive Evolution for Odorant and Gustatory Receptors in the Pea Aphid, *Acyrtosiphon pisum*. *Mol Biol Evol* 26, 2073-2086.
- Stanczyk, N.M., Brookfield, J.F., Ignell, R., Logan, J.G., and Field, L.M. (2010). Behavioral insensitivity to DEET in *Aedes aegypti* is a genetically determined trait residing in changes in sensillum function. *P Natl Acad Sci USA* 107, 8575-8580.
- Stensmyr, M.C., Giordano, E., Balloi, A., Angioy, A.M., and Hansson, B.S. (2003). Novel natural ligands for *Drosophila* olfactory receptor neurones. *The Journal of experimental biology* 206, 715-724.
- Syed, Z., and Leal, W.S. (2008). Mosquitoes smell and avoid the insect repellent DEET. *P Natl Acad Sci USA* 105, 13598-13603.
- Syed, Z., Pelletier, J., Flounders, E., Chitolina, R.F., and Leal, W.S. (2011). Generic insect repellent detector from the fruit fly *Drosophila melanogaster*. *PLoS ONE* 6, e17705.
- Takken, W. (1996). Synthesis and future challenges: The response of mosquitoes to host odours. *Ciba F Symp* 200, 302-320.
- Tan, P.-N., Steinbach, M., and Kumar, V. (2006). Introduction to data mining, 1st edn (Boston, Pearson Addison Wesley).

- Tentschert, J., Bestmann, H.J., Holldobler, B., and Heinze, J. (2000). 2,3-dimethyl-5-(2-methylpropyl)pyrazine, a trail pheromone component of *Eutetramorium mocquersyi* Emery (1899) (Hymenoptera: Formicidae). *Naturwissenschaften* 87, 377-380.
- Triballeau, N., Van Name, E., Laslier, G., Cai, D., Paillard, G., Sorensen, P.W., Hoffmann, R., Bertrand, H.O., Ngai, J., and Acher, F.C. (2008). High-Potency Olfactory Receptor Agonists Discovered by Virtual High-Throughput Screening: Molecular Probes for Receptor Structure and Olfactory Function. *Neuron* 60, 767-774.
- Turner, S.L., Li, N., Guda, T., Githure, J., Carde, R.T., and Ray, A. (2011). Ultra-prolonged activation of CO₂-sensing neurons disorients mosquitoes. *Nature* 474, 87-91.
- Turner, S.L., and Ray, A. (2009). Modification of CO₂ avoidance behaviour in *Drosophila* by inhibitory odorants. *Nature* 461, 277-281.
- van der Goes van Naters, W., and Carlson, J.R. (2006). Insects as chemosensors of humans and crops. *Nature* 444, 302-307.
- van Naters, W.V.G., and Carlson, J.R. (2007). Receptors and neurons for fly odors in *Drosophila*. *Curr Biol* 17, 606-612.
- Vannelli, T.A., Dykman, A., and de Montellano, P.R.O. (2002). The antituberculosis drug ethionamide is activated by a flavoprotein monooxygenase. *Journal of Biological Chemistry* 277, 12824-12829.
- Verdonk, M.L., Cole, J.C., Hartshorn, M.J., Murray, C.W., and Taylor, R.D. (2003). Improved protein-ligand docking using GOLD. *Proteins-Structure Function and Genetics* 52, 609-623.
- Vosshall, L.B., Amrein, H., Morozov, P.S., Rzhetsky, A., and Axel, R. (1999). A spatial map of olfactory receptor expression in the *Drosophila* antenna. *Cell* 96, 725-736.
- Vosshall, L.B., and Hansson, B.S. (2011). A unified nomenclature system for the insect olfactory coreceptor. *Chemical senses* 36, 497-498.
- Walker, J.D., Rodford, R., and Patlewicz, G. (2003). Quantitative structure-activity relationships for predicting percutaneous absorption rates. *Environ Toxicol Chem* 22, 1870-1884.
- Wang, J.W., Wong, A.M., Flores, J., Vosshall, L.B., and Axel, R. (2003). Two-photon calcium imaging reveals an odor-evoked map of activity in the fly brain. *Cell* 112, 271-282.
- Wang, P., Lyman, R.F., Shabalina, S.A., Mackay, T.F.C., and Anholt, R.R.H. (2007). Association of polymorphisms in odorant-binding protein genes with variation in olfactory response to benzaldehyde in *Drosophila*. *Genetics* 177, 1655-1665.
- Watson, J.D., and Crick, F.H. (1953). Molecular structure of nucleic acids; a structure for deoxyribose nucleic acid. *Nature* 171, 737-738.

Weiss, L.A., Dahanukar, A., Kwon, J.Y., Banerjee, D., and Carlson, J.R. (2011). The molecular and cellular basis of bitter taste in *Drosophila*. *Neuron* 69, 258-272.

Whitney, A.W. (1971). Direct Method of Nonparametric Measurement Selection. *IEEE T Comput C* 20, 1100-&.

Wicher, D., Schafer, R., Bauernfeind, R., Stensmyr, M.C., Heller, R., Heinemann, S.H., and Hansson, B.S. (2008). *Drosophila* odorant receptors are both ligand-gated and cyclic-nucleotide-activated cation channels. *Nature* 452, 1007-1011.

Willand, N., Desroses, M., Toto, P., Dirie, B., Lens, Z., Villeret, V., Rucktooa, P., Locht, C., Baulard, A., and Deprez, B. (2010). Exploring Drug Target Flexibility Using in Situ Click Chemistry: Application to a Mycobacterial Transcriptional Regulator. *ACS Chem Biol* 5, 1007-1013.

Willand, N., Dirie, B., Carette, X., Bifani, P., Singhal, A., Desroses, M., Leroux, F., Willery, E., Mathys, V., Deprez-Poulain, R., *et al.* (2009). Synthetic EthR inhibitors boost antituberculous activity of ethionamide. *Nat Med* 15, 537-544.

Wilson, R.I., and Laurent, G. (2005). Role of GABAergic inhibition in shaping odor-evoked spatiotemporal patterns in the *Drosophila* antennal lobe. *Journal of Neuroscience* 25, 9069-9079.

Wilson, R.I., Turner, G.C., and Laurent, G. (2004). Transformation of olfactory representations in the *Drosophila* antennal lobe. *Science* 303, 366-370.

Wishart, D.S., Knox, C., Guo, A.C., Cheng, D., Shrivastava, S., Tzur, D., Gautam, B., and Hassanali, M. (2008). DrugBank: a knowledgebase for drugs, drug actions and drug targets. *Nucleic Acids Research* 36, D901-D906.

World Health Organization. (2011). World malaria report (Geneva, Switzerland, World Health Organization), pp. v.

Xia, Y., Wang, G., Buscariollo, D., Pitts, R.J., Wenger, H., and Zwiebel, L.J. (2008). The molecular and cellular basis of olfactory-driven behavior in *Anopheles gambiae* larvae. *Proc Natl Acad Sci U S A* 105, 6433-6438.

Yao, C.A., and Carlson, J.R. (2010). Role of G-proteins in odor-sensing and CO₂-sensing neurons in *Drosophila*. *J Neurosci* 30, 4562-4572.

Yao, C.A., Ignell, R., and Carlson, J.R. (2005). Chemosensory coding by neurons in the coeloconic sensilla of the *Drosophila* antenna. *Journal of Neuroscience* 25, 8359-8367.

Zeng, X.N., Leyden, J.J., Lawley, H.J., Sawano, K., Nohara, I., and Preti, G. (1991). Analysis of Characteristic Odors from Human Male Axillae. *Journal of Chemical Ecology* 17, 1469-1492.

Zeng, X.N., Leyden, J.J., Spielman, A.I., and Preti, G. (1996). Analysis of characteristic human female axillary odors: Qualitative comparison to males. *Journal of Chemical Ecology* 22, 237-257.

Zhan, S., Merlin, C., Boore, J.L., and Reppert, S.M. (2011). The Monarch Butterfly Genome Yields Insights into Long-Distance Migration. *Cell* 147, 1171-1185.

Zhang, X., and Firestein, S. (2002). The olfactory receptor gene superfamily of the mouse. *Nat Neurosci* 5, 124-133.

Zwiebel, L.J., and Takken, W. (2004). Olfactory regulation of mosquito-host interactions. *Insect Biochem Molec* 34, 645-652.