

UC San Diego

UC San Diego Electronic Theses and Dissertations

Title

Specific Solutions to General Problems in Data Science and Ecology

Permalink

<https://escholarship.org/uc/item/5mf8b18k>

Author

Saberski, Erik

Publication Date

2024

Supplemental Material

<https://escholarship.org/uc/item/5mf8b18k#supplemental>

Peer reviewed|Thesis/dissertation

UNIVERSITY OF CALIFORNIA SAN DIEGO

Specific Solutions to General Problems in Data Science and Ecology

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

in

Oceanography

by

Erik Saberski

Committee in charge:

Professor George Sugihara, Chair
Professor Jeff Bowman
Professor Ian Eisenman
Professor Jack Gilbert
Professor Art Miller

2024

Copyright

Erik Saberski, 2024

All rights reserved.

The Dissertation of Erik Saberski is approved, and it is acceptable in quality and form for publication on microfilm and electronically.

University of California San Diego

2024

DEDICATION

To my number one supporter, my biggest fan, my Mom.

TABLE OF CONTENTS

DISSERTATION APPROVAL PAGE.....	iii
DEDICATION	iv
TABLE OF CONTENTS	v
LIST OF FIGURES.....	vii
LIST OF SUPPLEMENTAL FILES.....	ix
ACKNOWLEDGEMENTS	x
VITA	xi
ABSTRACT OF THE DISSERTATION	xiii
Chapter 0 INTRODUCTION	1
Chapter 1 Improved Prediction of Managed Water Flow into Everglades National Park Using Empirical Dynamic Modeling.....	4
Abstract	4
Introduction	5
Managed flows	5
Target Flows and the Tamiami Trail Flow Formula	10
Non-linear, Non-parametric Approaches	13
EDM	14
Models	15
Results and Discussion.....	16
Model Coefficients	16
Model Performance and Causal Inference	18
State-space (Non-linear) relationships	23
Conclusion.....	26
Acknowledgements	27

Chapter 2 The Impact of Data Resolution on Dynamic Causal Inference in Multiscale Ecological Networks	28
Abstract	28
Introduction	29
Methods	33
Results	38
Discussion	45
Conclusion.....	48
Chapter 3 Networks of Causal Linkage Between Eigenmodes Characterize Behavioral Dynamics of <i>Caenorhabditis elegans</i>	50
Abstract	50
Introduction	50
Materials and Methods	55
Results	58
Discussion	64
Acknowledgements	68
Chapter 4 Conclusion	69
REFERENCES.....	72

LIST OF FIGURES

Figure 1.1 Schematic of flow paths in South Florida: (a) predrainage; and (b) modern. In the predevelopment era, the Kissimmee Valley floodplain drained into Lake Okeechobee, which then overflowed its southern rim in a river of grass to the southern peninsula.....	7
Figure 1.2 Schematic of Everglades water control structures and projects. The Tamiami Trail, S-12 and S-333 structures separate upstream water conservation areas (WCAs) from Everglades National Park. (Base map courtesy of South Florida Natural Resources Center.).....	8
Figure 1.3 Inputs and outputs to TTFF. Each week, the current week’s environmental conditions are used to generate a target flow that dictates management for the following week. This data is also used to forecast next week’s target flow using the TTFF, a forecast that implicitly accounts for next week’s environmental conditions..	9
Figure 1.4 (a) Time series of Q^{sum} and the presumed causal variables; and (b) scatter plots of the variables versus Q^{sum}	12
Figure 1.5 (a) Recalculating the TTFF coefficients every 5 years yields different coefficients over time compared with the linear coefficients depending on the time of year reveals seasonal dynamics among the TTFF defined by the TTFF (red lines); and (b) recalculating the coefficients variables.	17
Figure 1.6 Comparison of different predictive models accuracy (mean absolute error between observed and predicted flow change). The TTFF (green) is a linear regression across all historical data. The 5 year predictor (purple) recalculates the TTFF coefficients in a 5 year moving window.....	19
Figure 1.7 Removing variables in S-Map forecasts to measure the impact on forecasts. Note that the only two variables that have significant negative impact on forecasts (increased MAE) when removed are water levels in the WC3A and NESRS regions.	21
Figure 1.8 Comparison of TTFF and S-Map forecasts. S-Maps here do not utilize Rain, Za, or PET. (a) and (c) show that overall, S-Maps outperform the TTFF on the original dataset (1965-2005, (a) as well as on contemporary data spanning 2007-2020 (c). (b) and (d) show average error for both forecasting algorithms as a function of flow	22
Figure 2.1 A model simulation in which resources, primary consumers, and secondary consumers exist on a grid and move randomly. Each timestep, primary consumers eat resources and secondary consumers consume primary consumers. Both primary and secondary consumers have rules determining whether they survive, starve, or reproduce.....	39
Figure 2.2 The number of interactions in each system resolved at a monthly timescale ($\tau = 1$) and annual timescale ($\tau = 12$). Note that in all systems more interactions are resolved at the monthly timescale, but there are still interactions in each system that are exclusively resolved at the annual timescale.	40
Figure 2.3 Fine scale connectance translates to aggregate interaction strength: Logistic models with two aggregates (10 predators and 10 prey) run at varying levels of high-resolution connectance (proportion of non-zero elements joining the two aggregates in the interaction matrix)..	41

Figure 2.4 Real world examples showing the relationship between aggregated and fine-scale linkages. Boxplots show that causally linked aggregates (coarse-scale) have more causal links at the species level (fine-scale). 42

Figure 2.5 A comparison of food webs and aggregate causal webs for the four systems studied. (A) Aggregated causal webs (blue arrows) overlaid with food webs (red arrows). (B) High-resolution causal webs mapping interactions between individual species. 43

Figure 3.1 Identifying relationships between eigenmodes of worm position. (a) During forward locomotion, the static view of the first two eigenmodes form a quadrature pair; any given value of a_1 is likely to be 90 degrees out of phase with a_2 52

Figure 3.2 Interaction profiles represented by time-lag versus normalized CCM correlation coefficient (ρ) for each pair of the first four eigenmodes for 12 foraging worms (grey lines - individuals, red line - average). Note that pairs of eigenmodes interact at different timescales, however, these relationships are relatively consistent across individuals. 59

Figure 3.3 Comparing dynamics of the eigenmodes in worms foraging and exhibiting an escape response. (a) The average CCM values between the first four eigenmodes of worm position plotted against τ_p for foraging (red) and escape response (blue) worms. 60

Figure 3.4 (a) Differences in dynamics between pairs of phenotypically similar mutants. The red line indicates the median distance of all mutants from each other. The boxplots are ordered from top to bottom in increasing median difference between strains, where smaller differences indicate more similar dynamics (less variance) between profiles. 63

LIST OF SUPPLEMENTAL FILES

Saberski_description.zip

ACKNOWLEDGEMENTS

I would like to acknowledge Professor George Sugihara for being my advisor. Together we did some really cool stuff.

Chapter 1, in full, is a reprint of the material as it appears in *The Journal of Water Resources Planning and Management*. Saberski, Erik, Joseph Park, Troy Hill, Erik Stabenau, and George Sugihara. "Improved Prediction of Managed Water Flow into Everglades National Park Using Empirical Dynamic Modeling." *Journal of Water Resources Planning and Management* 148, no. 12 (2022): 05022009. The dissertation author was the primary investigator and author of this paper.

Chapter 3, in full, is a reprint of the material as it appears in *PLoS Computational Biology*. Saberski, Erik, Antonia K. Bock, Rachel Goodridge, Vitul Agarwal, Tom Lorimer, Scott A. Rifkin, and George Sugihara. "Networks of causal linkage between eigenmodes characterize behavioral dynamics of *Caenorhabditis elegans*." *PLoS computational biology* 17, no. 9 (2021): e1009329. The dissertation author was the primary investigator and author of this paper.

VITA

- 2017 Bachelor of Science in Math and Physics, Bates College
- 2020 Master of Science in Marine Biology, Scripps Institution of Oceanography
- 2024 Doctor of Philosophy in Oceanography, Scripps Institution of Oceanography

PUBLICATIONS

Saberski, Erik, Antonia K. Bock, Rachel Goodridge, Vitul Agarwal, Tom Lorimer, Scott A. Rifkin, and George Sugihara. "Networks of causal linkage between eigenmodes characterize behavioral dynamics of *Caenorhabditis elegans*." *PLoS computational biology* 17, no. 9 (2021): e1009329.

Saberski, Erik, Joseph Park, Troy Hill, Erik Stabenau, and George Sugihara. "Improved Prediction of Managed Water Flow into Everglades National Park Using Empirical Dynamic Modeling." *Journal of Water Resources Planning and Management* 148, no. 12 (2022): 05022009.

Saberski, Erik T., Julia Daisy Diamond, Nathaniel Fath Henneman, and Daniel A. Levitis. "Post-reproductive parthenogenetic pea aphids (*Acyrtosiphon pisum*) are visually identifiable and disproportionately positioned distally to clonal colonies." *PeerJ* 4 (2016): e2631.

Merz, Ewa, **Erik Saberski**, Luis J. Gilarranz, Peter DF Isles, George Sugihara, Christine Berger, and Francesco Pomati. "Disruption of ecological networks in lakes by climate change and nutrient fluctuations." *Nature Climate Change* 13, no. 4 (2023): 389-396.

Park, Joseph, **Erik Saberski**, Erik Stabenau, and George Sugihara. "Dynamics of Florida milk production and total phosphate in Lake Okeechobee." *Plos one* 16, no. 8 (2021): e0248910.

Orenstein, Eric C., **Erik Saberski**, and Christian Briseño-Avena. "Discovery and dynamics of a cryptic marine copepod-parasite interaction." *Marine Ecology Progress Series* 691 (2022): 29-40.

Choi, Emma S., **Erik Saberski**, Tom Lorimer, Cameron Smith, Unduwap Kandage-Don, Ronald S. Burton, and George Sugihara. "The importance of making testable predictions: A cautionary tale." *PloS one* 15, no. 12 (2020): e0236541.

Lorimer, Tom, Rachel Goodridge, Antonia K. Bock, Vitul Agarwal, **Erik Saberski**, George Sugihara, and Scott A. Rifkin. "Tracking changes in behavioural dynamics using prediction error." *Plos one* 16, no. 5 (2021): e0251053.

Giron-Nava, Alfredo, Stephan B. Munch, Andrew F. Johnson, Ethan Deyle, Chase C. James, **Erik Saberski**, Gerald M. Pao, Octavio Aburto-Oropeza, and George Sugihara. "Circularity in fisheries data weakens real world prediction." *Scientific reports* 10, no. 1 (2020): 6977.

FIELD OF STUDY

Major Field: Biological Oceanography
Studies in Data Science
Professor George Sugihara

ABSTRACT OF THE DISSERTATION

Specific Solutions to General Problems in Data Science and Ecology

by

Erik Saberski

Doctor of Philosophy in Oceanography

University of California San Diego, 2024

Professor George Sugihara, Chair

Nature is hard to predict. Rules and relationships you discover about a system today may be totally different tomorrow. These relationships do not change randomly over time; rather, they change as the *state* of the system evolves. In a deterministic view of the world, similar states lead to similar outcomes. In this thesis, I leverage this principle to better understand and ultimately predict, complex systems.

In chapter 1, I work closely with the National Parks Service to understand variables that influence flow target values through the Everglades National Park. The Tamiami Trail Flow Formula, a linear (not state-dependent) was previously developed to predict such values. In this

chapter I show that with only minor adjustments to their linear approach, a non-linear (state-dependent) predictor can be made with significant prediction improvement.

Chapter 2 focuses on the role of scale in understand ecosystem relationships. Using both models and real world examples I show that not one scale can capture all of the dynamics of a real world system: for example, some relationships are better resolved at an annual timescale while others are best resolved monthly.

In chapter 3 I develop a new method for classifying systems based on the delay in their dynamic relationships. This method is applied to study the behavioral states of the nematode *Caenorhabditis elegans*. By analyzing the causal relationships between eigenvectors that represent the worm's posture (“eigenworms”), I am able to classify the behavioral state of the worm (foraging or reacting to a harmful stimulus). Additionally, I demonstrate that this technique can identify genetic mutations in these worms solely through analysis of their bodily movements.

This work demonstrates the that powerful models and non-linear relationships can be extrapolated directly from data without the need for assumptions or fixed equations.

Chapter 0 INTRODUCTION

The natural world is rich with complex systems where relationships between variables change over time (Deyle *et al.* 2022). In this thesis, I explore approaches to gain a deeper mechanistic understanding of these systems with the overall aim to (1) make forecasts on their future values and (2) understand how variables influence each other.

I focus on Empirical dynamic modelling (EDM) – an evolving suite of time-series analysis techniques for understanding relationships between variables and making predictions on non-linear systems (Sugihara & May 1990; Sugihara 1994; Sugihara *et al.* 2012; Deyle *et al.* 2016b; Cenci, Sugihara & Saavedra 2019). Unlike other methods that rely on theoretical equations or assumptions about system dynamics, EDM looks directly at data to build models that encapsulate changing relationships between variables.

A common thread through this thesis is the concept of a “*dynamic causal relationship*” – where a change in the state of one variable causes a change in the state of another variable (Sugihara *et al.* 2012). Such relationships in non-linear systems may not be identified with classical approaches such as path analysis and structural equation modeling that depend on linear correlation (Spirtes, Glymour & Scheines 2000). Understanding causal relationships in natural systems can help predict future states and understand the consequences of human (or other) interventions (Saberski *et al.* 2022).

In Chapter 1, I focus on predicting target flows throughout the Tamiami Trail in the Everglades National Park. This serves as a demonstration of how relationships between variables are not static – rather, they can change over the course of seasons or decades, and can be state-dependent, changing with evolving variable values (e.g., precipitation, upstream water levels).

Managers of this system created the Tamiami Trail Flow Formula, a linear equation to predict weekly target flows based on a few hand-chosen variables including precipitation and upstream and downstream water levels. This linear solution inherently assumes relationships are static over time and thus only achieves moderate success. Managers of this system also state that they understand that nonlinear solutions may perform better but are far too complex to implement. I show that it is possible to achieve significant prediction improvement with only minimal changes to their current linear setup to create an improved, state-dependent (non-linear) formula. Further, I show that using tests to detect causal relationships many of their hand-picked variables had little-to-no influence on target flows.

While chapter 1 shows examples of how dynamics can change over time, chapter 2 explores how relationships in ecosystems change as a function of scale. This is a subtle but crucial property to consider when managing ecosystems in the context of causal relationships. I show through both model and real-world examples that as systems are analyzed at varying scales, new sets of causal relationships emerge. From this work, I urge managers to consider multiple scales (spatial, temporal, species aggregation, etc.) when constructing and implementing causal ecosystem networks.

In Chapter 3 I introduce a novel implementation of causal relationships by using the delay in dynamic relationships as a fingerprint to classify systems. This chapter focuses on behavioral states of the nematode *Caenorhabditis elegans*. Using the causal relationships between *eigenvectors* – abstract variables that define the worm’s position – I classify the behavior of the worm. Specifically, by measuring how different regions of the worm’s body are influencing each other, I discern whether the worm is either foraging for food or escaping a

noxious stimulus. I further show that this same technique can be used to identify genetic mutations in worms from just their body movements alone.

Collectively, this thesis underscores the significance of understanding causal relationships in natural systems. It highlights the need for adaptable, nuanced models and approaches that reflect the inherent complexity and dynamic nature of ecological systems. This work serves as a foundation for future research and management strategies, guiding us towards more effective stewardship of the natural world. Through a combination of theoretical innovation and practical application, it contributes profoundly to our understanding of the intricate tapestry of life and its myriad interactions.

Chapter 1 Improved Prediction of Managed Water Flow into Everglades National Park Using Empirical Dynamic Modeling

Abstract

Alteration of natural surface flow paths across South Florida has been detrimental to the environmental health and sustainability of the Everglades and surrounding ecosystems. As part of the Comprehensive Everglades Restoration Plan (CERP), predicting flows into Everglades National Park (ENP) is a central concern of effective management strategies. Management efforts have established weekly target flows into Everglades National Park through optimization of numerically intensive hydrological models. These target flows are focused specifically on flows across US Highway 41, also known as the Tamiami Trail. To aide in timely management assessments in response to current or predicted hydrologic conditions, the Tamiami Trail Flow Formula (TTFF) was developed previously to predict weekly target flows based on linear regression of six theorized flow drivers. It is known that these drivers exhibit nonlinear dynamics, suggesting that there is room for improvement in relation to the strictly linear TTFF. We used empirical dynamic modeling (EDM), a nonparametric modeling paradigm for forecasting and analyzing nonlinear time series, to show that prediction accuracy is improved when nonlinearity is accounted for. This method relies on weighted linear regressions that depend on specific environmental conditions or system states, i.e., in which the regression gives greater weight to input variables that have values that are more similar to the current state. Surprisingly, we found that only two of the six standard TTFF variables are required in the nonlinear weekly forecast model (upstream and downstream water levels), and thus the EDM model not only improves accuracy but also greatly simplifies hydrologic forecasting.

Introduction

Many analytic approaches use equation-based models as approximations of real-world systems to test hypothesized mechanisms or to predict future outcomes. However, real world systems are often nonlinear and multidimensional, which can render explicit parametric approaches intractable. Empirical approaches, which extract information from the data instead of relying on hypothesized equations, represent a natural and flexible approach to modeling complex, nonlinear systems such as managed water resources.

Empirical dynamic modeling (EDM) is a non-parametric framework for modeling nonlinear systems based on the mathematical theory of reconstructing attractors (vector fields that can show how variables interact though time) from time series data (Takens 1981). EDM was initially aimed to address problems in ecology (Sugihara & May 1990; Sugihara 1994; Dixon, Milicich & Sugihara 1999; Sugihara *et al.* 2012; Ye & Sugihara 2016), however its applications have extended to many areas such as climate change (Van Nes *et al.*), atmospheric sciences (Sugihara *et al.*), neuroscience (Segundo *et al.*), studying the dynamics of infant heart rhythms (Sugihara *et al.*), identifying the drivers of influenza outbreaks (Deyle *et al.*), and classifying complex behaviors in the nematode *C. Elegans* (Lorimer *et al.* 2021; Saberski *et al.* 2021). To our knowledge, EDM has not yet been used specifically to map hydrologic dynamics. Here, we introduce the use of EDM as a tool for forecasting managed water flows in Everglades National Park as a component of the Comprehensive Everglades Restoration Plan (CERP). A lucid and accessible introduction to EDM is provided by (Chang, Ushio & Hsieh 2017).

Managed flows

The Florida Everglades originally consisted of 3 million acres of marsh draining the Kissimmee River Basin and Lake Okeechobee southward into Florida Bay. Starting in the late

19th Century, ambitious plans to “drain” the Everglades to produce arable and habitable lands were initiated, eventually coalescing in 1948 under the Congressionally authorized Central and Southern Florida Project under auspices of the United States Army Corps of Engineers (USACE). Design goals were to provide flood control and agricultural sustainability, with major features including the Herbert Hoover dike impounding Lake Okeechobee, creation of a large agricultural area along the southern lake border, a levee along the eastern boundary of the Everglades, and impoundment of three water conservation areas (WCAs) linking Lake Okeechobee to Everglades National Park and the southern coast (Council).

The result of these water control efforts was a fundamental alteration of the natural flow paths and hydroperiods (figure 1.1), which was eventually recognized as detrimental to environmental health and sustainability of the Everglades and its ecosystem services. Recognition of these changes led to the Congressionally mandated Comprehensive Everglades Restoration Plan (CERP) in 2000, a framework for restoring, preserving, and protecting the South Florida ecosystem. CERP was originally designed with 68 project components expected to take 30 years at an estimated cost of \$8 billion. Over the last two decades, it has been recognized that CERP and state restoration efforts must encompass an adaptive management approach. As such, the restoration today is a complex, adaptive collaboration continuing to evolve (National Academies of Sciences 2019).

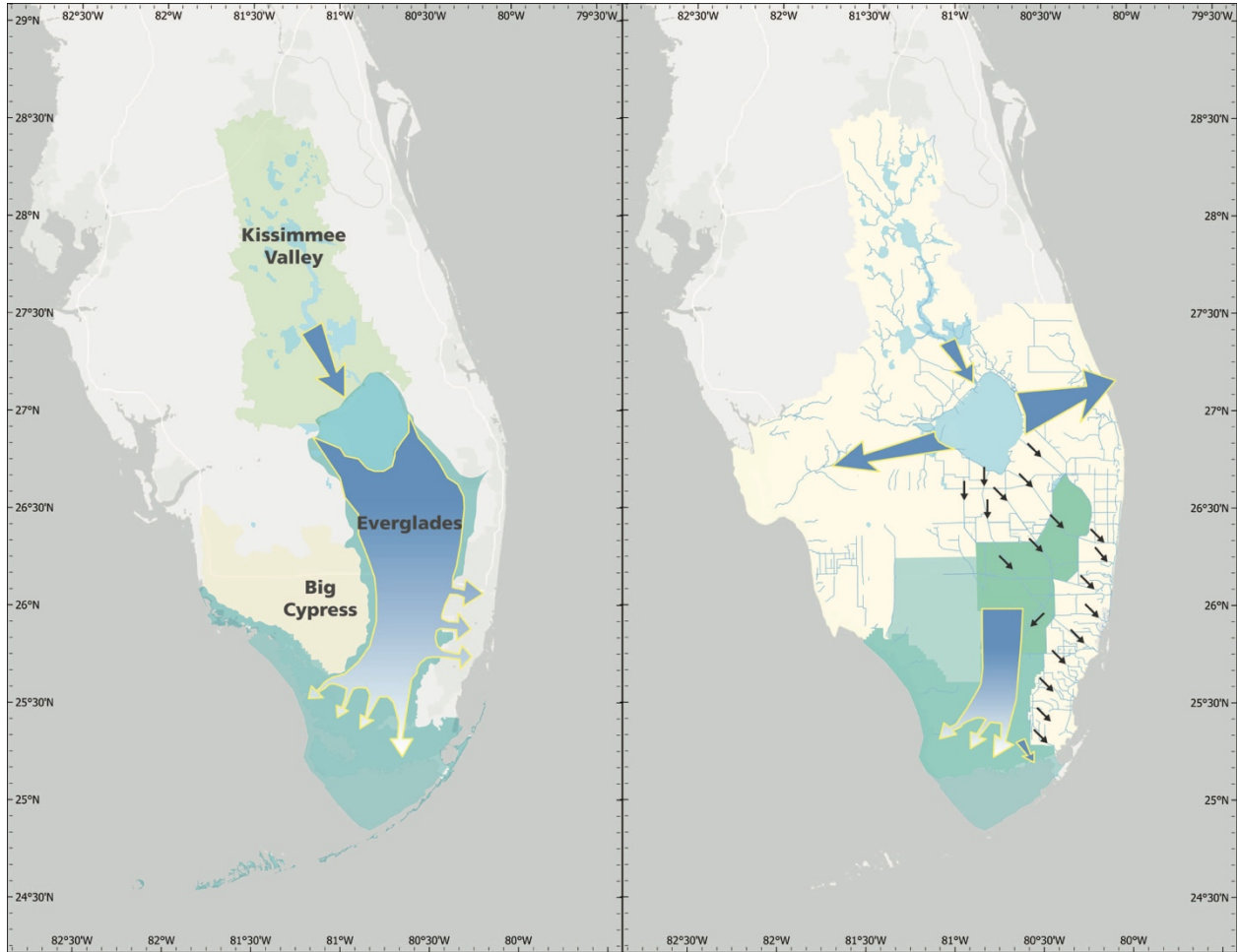


Figure 1.1 Schematic of flow paths in South Florida: (a) predevelopment; and (b) modern. In the predevelopment era, the Kissimmee Valley floodplain drained into Lake Okeechobee, which then overflowed its southern rim in a river of grass to the southern peninsula. Postdevelopment, flow paths were channelized, represented by arrows, and the remaining segments of the Everglades were impounded with levees and canals. (Base map courtesy of South Florida Natural Resources Center.)

A central tenant of CERP is to increase water flows and hydroperiods within Everglades National Park. A fundamental barrier to this was construction of the Tamiami Trail (U.S. Highway 41) in the early 20th Century. The trail acts as a levee preventing natural flow from the upstream WCAs and natural areas. Flows from the upstream WCAs are managed primarily through the S-12 gated weirs (figure 1.2) with both upstream and downstream regulation limits.

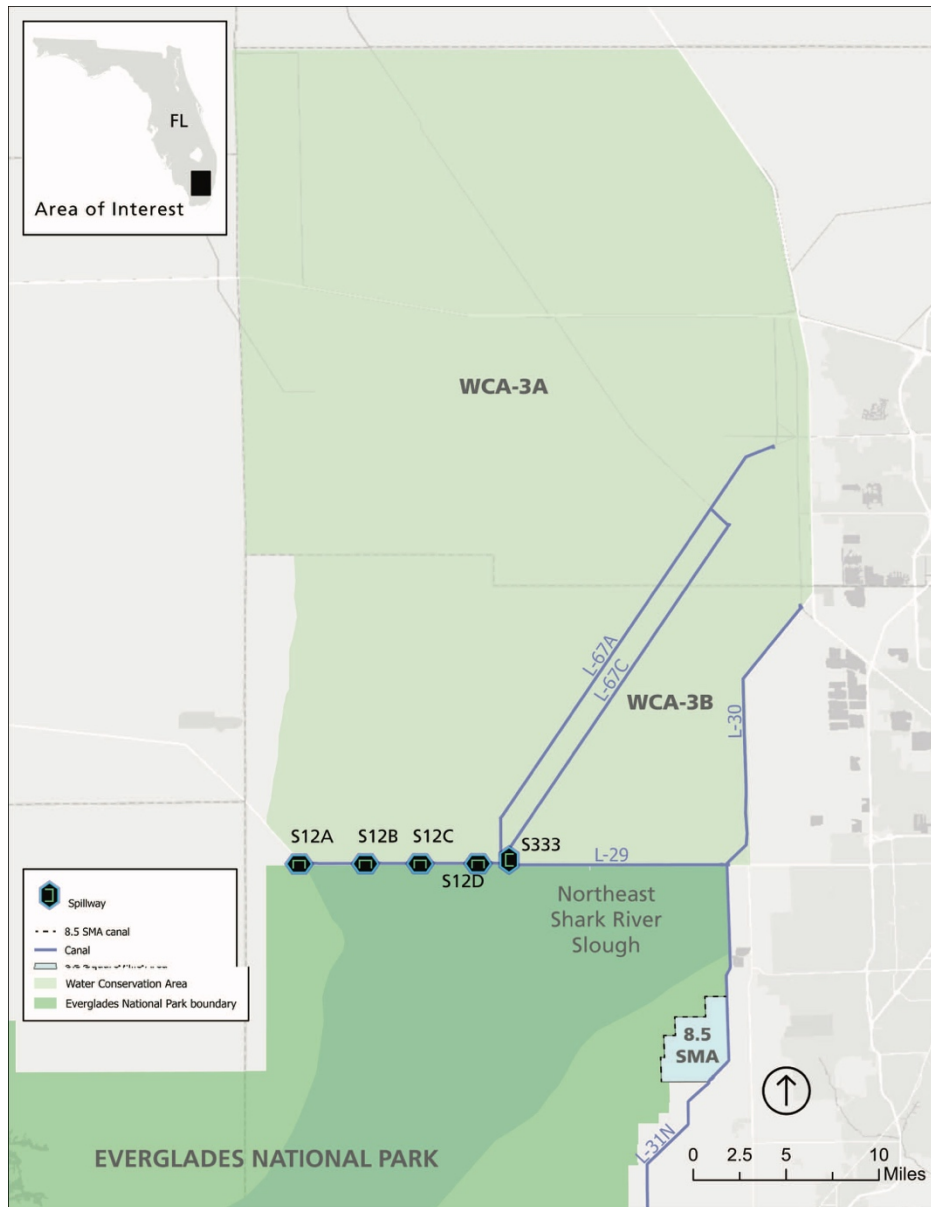


Figure 1.2 Schematic of Everglades water control structures and projects. The Tamiami Trail, S-12 and S-333 structures separate upstream water conservation areas (WCAs) from Everglades National Park. (Base map courtesy of South Florida Natural Resources Center.)

A recent adaptation of restoration water management is the redevelopment of flow targets for releases into Everglades National Park as part of the Combined Operational Plan (COP) (U.S. Army Corps of Engineers 2017). These targets serve as goals for maintaining healthy water levels throughout the greater Everglades system based on multiple environmental variables such

as weekly rainfall and estimated evaporation and are recalculated on a weekly basis. Because these flow targets are derived from up-to-date environmental conditions, future flow targets cannot be exactly determined without knowing the future environmental state. In order to best prepare for these weekly targets, the Tamiami Trail Flow Formula (TTFF) (SFWMD 2020) was recently developed to forecast future flow targets. A diagram of the inputs and outputs to these models is represented in figure 1.3. This figure emphasizes that the goal of the TTFF is not to predict the following week's flow into the ENP; rather, it is to forecast what the following week's *target flow* will be. This subtle difference has huge implications: although the following week's target will guide what the managed flow into the system will be, predicting the target is fundamentally different than predicting the raw flow.

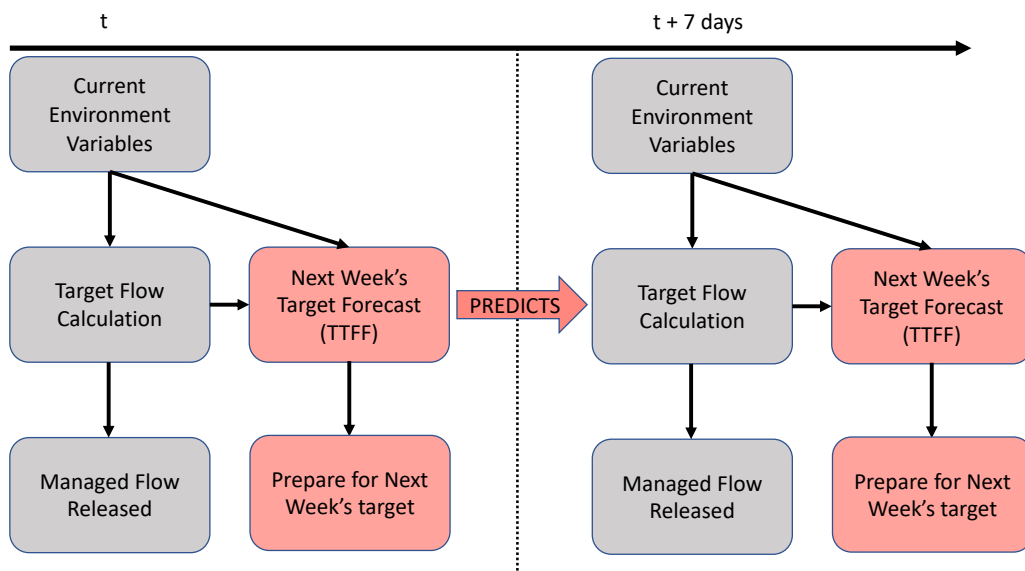


Figure 1.3 Inputs and outputs to TTFF. Each week, the current week's environmental conditions are used to generate a target flow that dictates management for the following week. This data is also used to forecast next week's target flow using the TTFF, a forecast that implicitly accounts for next week's environmental conditions. Overall, this produces both a target to guide management for the upcoming week and a forecast to predict next week's new target, giving time for managers to prepare.

The quality of the forecasts made by the TTFF are particularly critical as the system must be managed as a whole, taking into account both the desired downstream conditions and the upstream storage capacity. Upstream basins are large and respond slowly to changes in operational efforts. Thus, considerable lead time is needed to adjust basin water levels. Forecasting next week's target flow is extremely important as it can provide valuable lead time for managers to appropriately prepare the upstream basin level to effectively meet the future desired target.

Setting the incoming flow volumes appropriately is critical to make projected deliveries through Tamiami Trail and not create adverse impacts due to flooding. Improvements in forecast skill will reduce the likelihood of ecologically adverse conditions within the WCAs or, simply put, having too much or too little water to manage the system properly. A primary aim of the present work is to examine the TTFF to ascertain the completeness of its information content, and to compare it to, and determine the potential benefits of, forecasts made using EDM.

Target Flows and the Tamiami Trail Flow Formula

Target flows were determined over the 1965-2005 period using the Regional Simulation Model (RSM) (SFWMD 2020b) and an inverse modeling tool, identifying optimal flows in response to hydrologic constraints (SFWMD 2020). The resultant time series is referred to as $Q^{\text{sum}}(t)$ representing cumulative, weekly target flows across Tamiami Trail into Everglades National Park. To model these target flows in response to current or future conditions, the Tamiami Trail Flow Formula (TTFF), a linear model, was developed (SFWMD 2020). TTFF developers recognized the nonlinear nature of the problem, but decided that a linear formulation performed adequately and was simpler and easier to understand than a nonlinear or machine

learning model. The TTFP presumes that precedent values of rain, evapotranspiration, upstream and downstream water levels, and flow, are required to best predict target flows.

The TTFP is formulated as:

$$\widehat{Q}^{sum}(t) = \beta_1 S^{WCA}(t) + \beta_2 S^{ENP}(t) + \beta_3 Q^{sum}(t-1) + \beta_4 R(t) + \beta_5 PET(t) + \beta_6 ZA(t)$$

where $\widehat{Q}^{sum}(t)$ is the predicted target flow release for the coming week (sum of S-12A, S-12B, S-12C, S-12D and S-333, see figure 1.2). $S^{WCA}(t)$ is the spatial average of observed water levels in WCA-3A at the start of the current week (start of a week is Sunday and the end of a week is Saturday). $S^{ENP}(t)$ is observed water level in Everglades National Park, Northeast Shark River slough (NESRS), for the current week. $Q^{sum}(t-1)$ is the daily average of target flow releases for the previous week. $R(t)$ is areal average of total weekly rainfall for WCA-3A, PET is the total weekly potential evapotranspiration at the 3AS3WX station, and ZA is the zone A regulation water level of the current week in WCA-3A. When water levels in WCA-3A are above ZA , flood control water releases are authorized across Tamiami Trail. β are linear regression fit coefficients.

Plotting the raw variables against target flows (figure 1.4) reveals that the highest correlation among the variables is the previous week's target flow (autocorrelation). This makes sense because water flows relatively slowly through the Everglades, giving the system large inertia. We therefore expect flows to have relatively low change from the prior week and to exhibit temporal autocorrelation. Upstream and downstream water levels are also noticeably correlated with future flows. The other two variables that show positive correlations with flow are upstream (WCA-3A) and downstream (NESRS) water-levels. Despite the positive correlations, the data suggests a nonlinear fit may be more suitable for these variables (e.g.,

exponential relationship between flow and NESRS level). The remaining variables, rain, PET, and ZA, show no clear indication of linear relationship or covariation (correlation) of any kind; a circumstance that often occurs with nonlinear dynamics and overlapping effects from explanatory variables (Sugihara *et al.* 2012). These variables are also likely coupled with each other (e.g., upstream water level influencing downstream water level, rainfall influencing water levels), creating a complex web of dynamics that may be difficult to define with parametric models. Taken together this suggests that predictions might be improved when the system is viewed through a nonlinear, non-parametric lens.

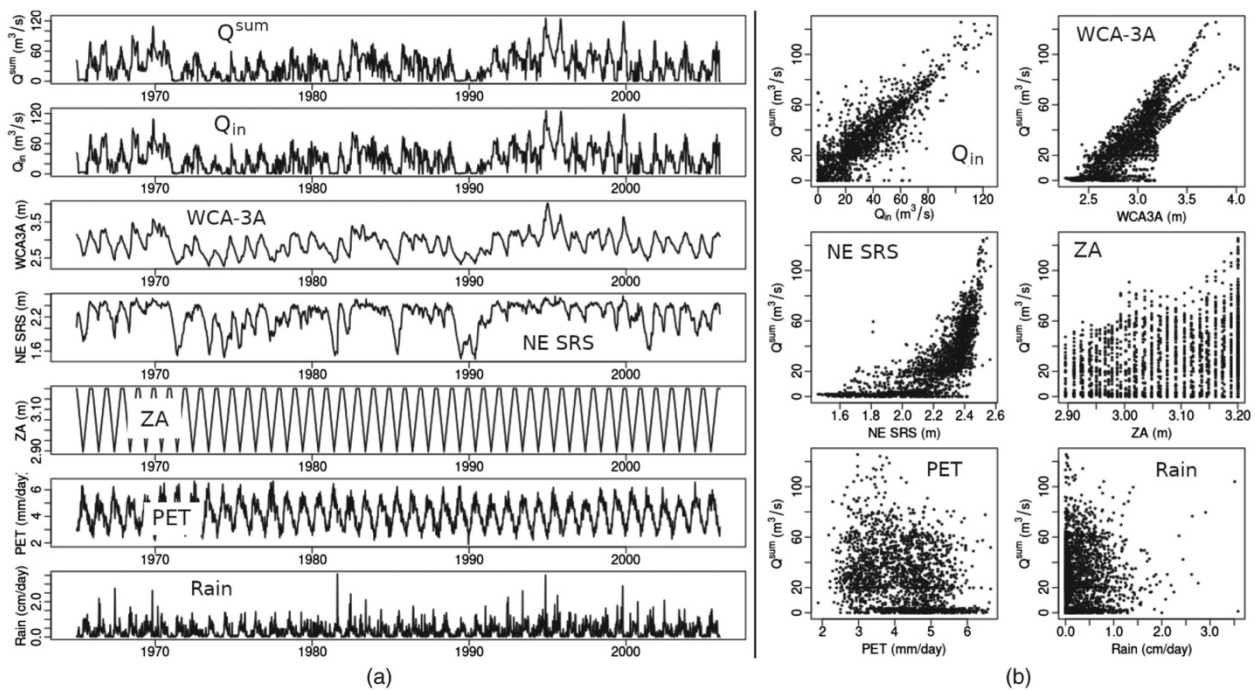


Figure 1.4 (a) Time series of Q^{sum} and the presumed causal variables; and (b) scatter plots of the variables versus Q^{sum} .

Despite the strictly linear nature of the TTF, it has seemingly impressive short-term predictive accuracy, achieving a correlation between observed and predicted weekly values of $\rho = 0.90$. However, this is largely due to the significant amount of autocorrelation in the data on a weekly time scale: a constant predictor (predicting the value next week will be the

same as the current) achieves a predictive accuracy of $\rho = 0.88$. Other metrics for predictive accuracy show that the formula has much room for improvement: it only correctly predicts the directional change in flow 60% of the time and prediction accuracy on changes in flow ($\Delta Q = Q_{t+1} - Q_t$) achieves a predictive accuracy of $\rho = 0.45$.

Non-linear, Non-parametric Approaches

As noted above, the TTFP was generated through a generalized linear model of six variables hypothesized to be influential to flow using data collected from 1965-2005. Because the TTFP was generated from a single best-fit solution on the entire data record, the model is implicitly stationary: resolved coefficients of the TTFP are fixed constants reflecting the global nature of the statistical regression. This is fundamentally distinct from dynamic nonlinear models where relationships among variables can change. In fact, nonlinear models can be constructed piece-wise from segmented linear models to address how relationships among variables change as the system state evolves.

For example, if one assumes the dynamics are slowly changing over time, the linear solution can be recalculated every few years to find a new set of coefficients specific to recent data. Similarly, if dynamics are theorized to change seasonally one may calculate coefficients for each month of the year. Such partitions are known as "similar states", where a state refers to a set of conditions associated with a specific set of dynamics. Similar states exhibit similar dynamics.

Typically, the "state" of a natural system depends on multiple factors (so-called *state variables*). In the seasonal model which calculated coefficients for each month (described above), for example, the time of year would be considered one state variable. Just as the seasonal model recalculates coefficients depending on the month, similar nonlinear models can be built to

account for other state-variables. For example, in the case of the TTFF it may be sensible to re-derive coefficients by partitioning the data into subsets with similar flow rates: high flow rates may have different dynamics than low flow rates.

EDM

Empirical dynamic modelling (EDM) focuses on reconstructing a system's *state-space*: a multi-dimensional representation of system variables as a function of time (Sugihara & May 1990; Sugihara 1994; Sugihara *et al.* 2012). Forecasting leverages the fact that points localized in state-space (nearest neighbors) exhibit similar dynamics (Sugihara & May 1990; Sugihara 1994). While the examples above describe ways to define the state space based on one variable (e.g. month, flow regime), EDM considers multiple variables together to identify similar states without presuming specific relationships (Deyle & Sugihara 2011; Deyle *et al.* 2016b); rather, dynamics are derived directly from the data.

EDM can be used to screen the available time series data and identify which variables are relevant to usefully include in a nonlinear forecasting model. Additionally, EDM involves the use of a causality test, convergent cross-mapping (CCM, (Sugihara *et al.*)) which identifies nonlinear coupling between variables directly from time series data (figure S1). This contrasts with the normal modelling procedure of TTFF where the specific variables used are asserted or hypothesized to be relevant.

Here, we utilize a state-space forecasting technique within the EDM framework called "S-Maps" (sequential locally weighted global linear maps) (Sugihara): At each point in time, coefficients are recalculated based on a linear-fit that maps state-variables onto a target variable, similar to the TTFF formulation. However, each fit at time t is weighted towards similar states to

that of time t . This is analogous to the non-linear methods described above, however instead of rigid cut-offs defined by the partition (e.g., partitioning data strictly by months; either weights of 1 or 0), weights are applied smoothly based on proximity with an exponential kernel applied to all points in the state-space. A nonlinearity parameter, θ , can be adjusted to change how state-specific forecasts are: θ of zero gives all states equal weight regardless of state-similarity, equating to a global autoregressive model. Higher values of θ , however, localize the forecast to more state-specific conditions, accounting for how system dynamics change over time.

We note that the S-Map method only adds one more step into the process of how the TTFF was originally formulated. Both the TTFF and S-Map forecasts utilize regressive maps from driving variables onto flow targets, S-Maps just also include weights with these regressions such that nearest neighbors are weighted more heavily in each prediction. It is worth noting that typically when performing S-Map (or other state-space based methods such as those described in the following section), the value being predicted from is left out of local linear regressions in order to obtain an unbiased predictions (i.e., leave-one-out-cross-validation).

Models

- System state can be defined in many ways. For example, the state can be simply defined by the value of one variable (e.g., a high flow state vs. a low flow state), or combinations of multiple (e.g., high/low flow in the winter vs. high/low flow in the summer). Here, we analyze the performance of four different ways to define system state:
- Interannual Predictor: Re-calculates linear regressions on the six variables from equation 1 using the previous 5 years of data (time points $t-260:t$) to predict flow at time $t+1$.

- Seasonal Predictor: Re-calculates coefficients using historical data within 6 weeks of the current year-day. For example, if a forecast is predicting flow in the first week of March, all historical data between mid-January and mid-April is used to generate linear coefficients.
- Variable-specific Predictor: Re-calculates coefficients at each time (t) using the 10% of values in the dataset that have the closest values of the given variable (v); i.e. the timepoints (t*) with the smallest values $|v_t - v_{t^*}|$. This was performed on the 5 theorized drivers of flow.
- EDM S-Maps: Contrasting the variable-specific model in which data is partitioned based on the value of a single variable, S-Maps use all input variables to define the "state space". At each point in time, the nearest neighbors are defined as the other points in time that also have similar set of all of the variables (as measured by Euclidean distance in state-space), not just one. Coefficients are re-calculated at each point in time just as in the models above, however, regressions are weighted towards state-space coordinates that have similar states (similarly valued state-variables at a particular time). The variables included in these S-Map forecasts are the six TTF variables as well as a sine and cosine term each with a 1-year period to represent the time of year.

Results and Discussion

Model Coefficients

When the TTF was formulated, the static coefficients were associated with physical processes. For example, rain had a positive coefficient that was interpreted as more rainfall should increase overall flow. However, depending on the data used in the linear regression to formulate the

TTFF, one can obtain either a negative coefficient, positive coefficient or a coefficient near 0 (figure 1.5). For example, in model 1 (Interannual Predictor), we find that the linear regression calculated using only data from 1985-1990 yields a negative coefficient for rain (figure 1.5). This highlights that it can be dangerous to make physical interpretations based on linear coefficients if the results change depending on the data used. In this case, we suggest that the influence of rain on flow targets can in fact change from positive to negative depending on the state of the system. For example, certain states may cause rainfall to increase downstream water level more so than upstream, which may in turn reduce the overall flow.

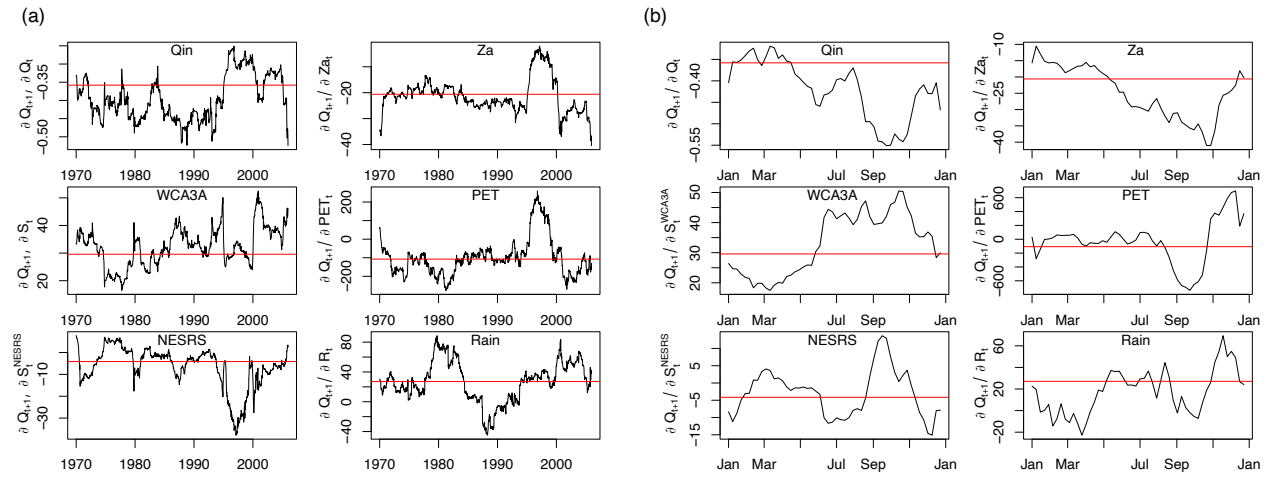


Figure 1.5 (a) Recalculating the TTFF coefficients every 5 years yields different coefficients over time compared with the linear coefficients depending on the time of year reveals seasonal dynamics among the TTFF defined by the TTFF (red lines); and (b) recalculating the coefficients variables.

Each model resulted in coefficients that change over time. For example, figure 1.5 shows how coefficients change for the interannual predictor (a) and seasonal predictor (b). The interannual predictor reveals coefficients with significant temporal variation, exhibiting dynamics reflective of external interactions. We also note a large excursion in NESRS, ZA and PET coefficients from the mid 1990's to 2000. This was a period of high water levels in the

upstream WCAs, with accordingly negative influence of NESRS (downstream water levels) and positive forcing associated by upstream water availability (ZA).

The seasonal predictor clearly recovers dynamics reflective of the South Florida summer monsoon, with a dry season from November through April, and wet season May to October. Here, WCA-3A and ZA (upstream water supply) closely reflect these monsoon patterns with distinct shift from positive to negative coefficients in April and November. Further, the downstream NESRS exhibits a delayed response consistent with water management releases.

Model Performance and Causal Inference

Figure 1.6 compares the model accuracy of all models tested. We use the mean absolute error (MAE) as our main metric for model accuracy because it is a meaningful value for managers implementing these predicted target values. The TTFF achieves a MAE of 7.2 m³/s. The inter-annual predictor performs the worst with a MAE of flow forecasts of 7.3 m³/s. Three of the variable-specific predictors (using Rain, ZA, and PET to define system state) also perform poorly (worse than the TTFF). The other three variable-specific predictors (using NESRS, WCA-3A and flow) outperform the TTFF, achieving a MAE under 7.2 m³/s. The second best model was the seasonal model, achieving a MAE of 6.9 m³/s. S-Map forecasts, however, provided the highest fidelity achieving a MAE of 6.5 m³/s.

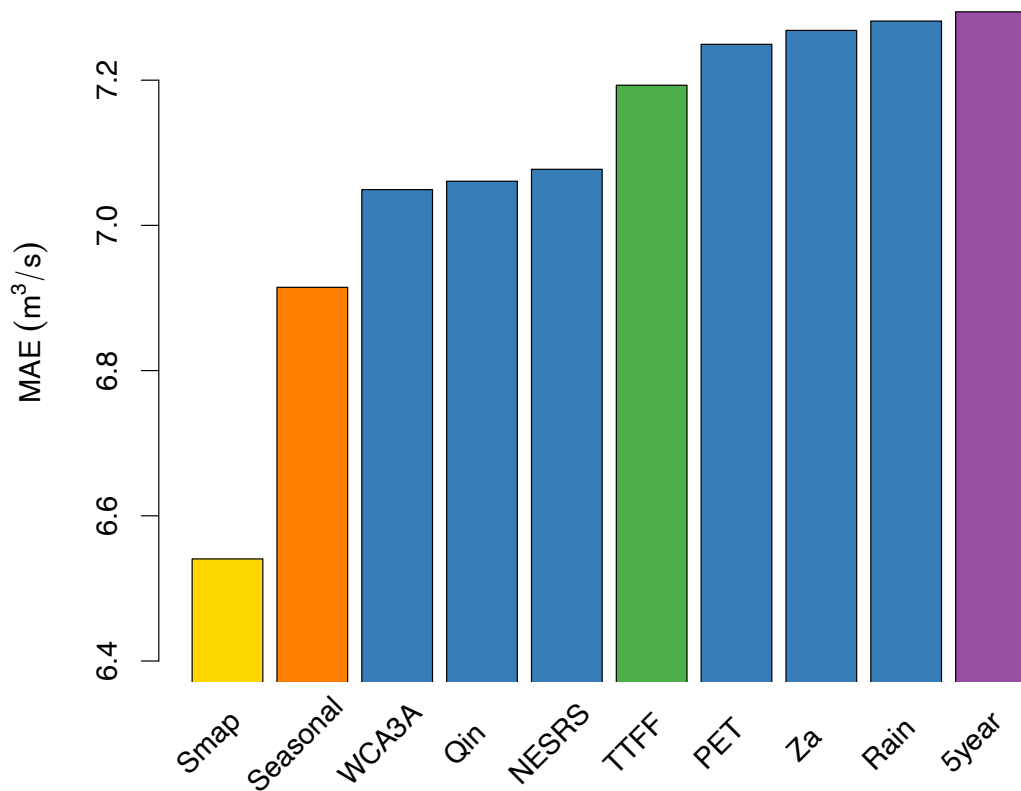


Figure 1.6 Comparison of different predictive models accuracy (mean absolute error between observed and predicted flow change). The TTFF (green) is a linear regression across all historical data. The 5 year predictor (purple) recalculates the TTFF coefficients in a 5 year moving window. The seasonal predictor (orange) recalculates the coefficients using historical data within 6 weeks of the given year date (see Methods). Each variable (v) predictor (blue) makes forecasts from time t using only historical data with similar values of v_t (see methods). S-Maps (gold) account for all of the nonlinearities among the variables to make a smooth "state-space" that is not specific to the state of any one variable.

The six variables selected as independent variables of the TTFF model make complete sense from the perspective of a hydrologic system. However, all variables may not provide significant information for improving forecasts.

To verify that the five hypothesized variables are indeed causal drivers of flow targets, we performed the EDM nonlinear causality test convergent cross-mapping (CCM), (Sugihara *et al.*). Despite the limited correlation between flow targets and these variables (figure 1.4), CCM reveals evidence for nonlinear coupling between all five variables and flow targets (figure S1). For an introduction to CCM, see the video at <http://tinyurl.com/EDM-intro>. Variables that show weak coupling (low CCM values, e.g. Rain and PET) do not necessarily provide useful information for improving predictions beyond the information gained from the strong drivers (e.g. upstream and downstream water levels). To further evaluate whether the five theorized driving variables of the TTFF are important for making predictions, we measured the performance of the S-Map predictor with variables removed one at a time (figure 1.7). Three of these variables (ZA, PET, and Rain) showed little-to-no negative impact on overall predictions when removed. This suggests that these variables, although shown to be weak causal drivers with CCM, may not be important for defining the state space of the system: no matter their values, the dynamics of Q^{sum} do not appreciably change. As a further check we performed an exhaustive assessment of state-space variable combinations using the EDM "multiview" algorithm (figure S2, (Ye & Sugihara)). The multiview approach tests the predictive accuracy of using different combinations of variables (with varying time delays, see supplement S2) to reconstruct the state-space. This gives a more complete measure of how important variables are for making predictions (see supplement S2). Combined with the CCM results, these analyses confirm that the variables ZA, PET, and Rain in the historical data, although potentially important, are not in themselves historically important for the overall goal of predicting integrated water flows on a weekly timescale across the Tamiami Trail.

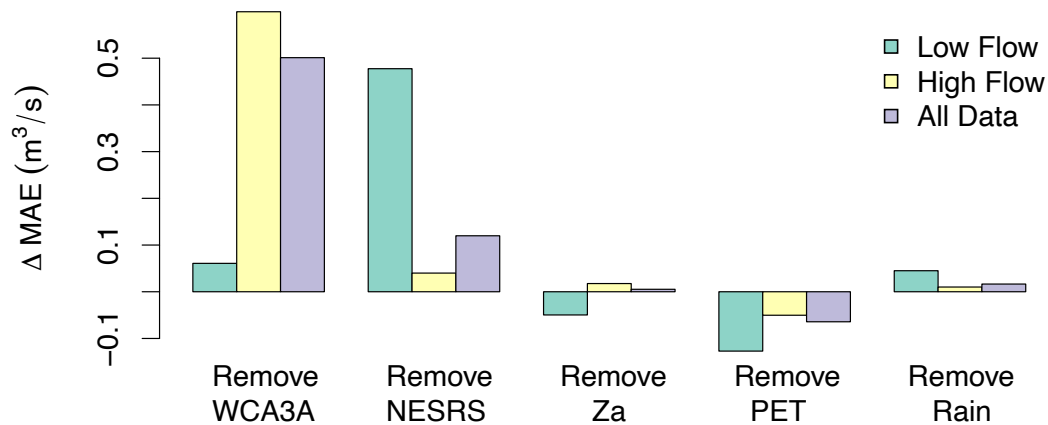


Figure 1.7 Removing variables in S-Map forecasts to measure the impact on forecasts. Note that the only two variables that have significant negative impact on forecasts (increased MAE) when removed are water levels in the WC3A and NESRS regions.

In accordance with figure 1.6, we find that predictions are significantly hindered when WCA-3A and NESRS are removed. Physically, this aligns with the fact that upstream (WCA-3A) and downstream (NESRS) water levels are the primary variables determining weir flow, whereas rain and PET accumulated over one week are integrated drivers of these upstream and downstream water levels. Further, focusing on periods of high flow and low flow reveals that the WCA-3A water stage is more important for making forecasts during high flow while the stage in NESRS is more important when predicting low flow regimes.

The accuracy of the S-Map forecasts is further improved to a MAE of 6.3 m³/s when ZA, PET, and Rain were excluded from the embedding. Figure 1.8a shows this improvement compared to the performance of the TTFF on the test data set of weekly sampled data spanning 1965 - 2005. Further, Figure 1.8b shows that the prediction improvement varies depending on the flow: S-Maps outperform the TTFF during all flow regimes, however periods of lower flow show the greatest improvement. Figure 1.8c&d show that this finding is also true when looking at a contemporary data set (weekly data from 2007 - 2020).

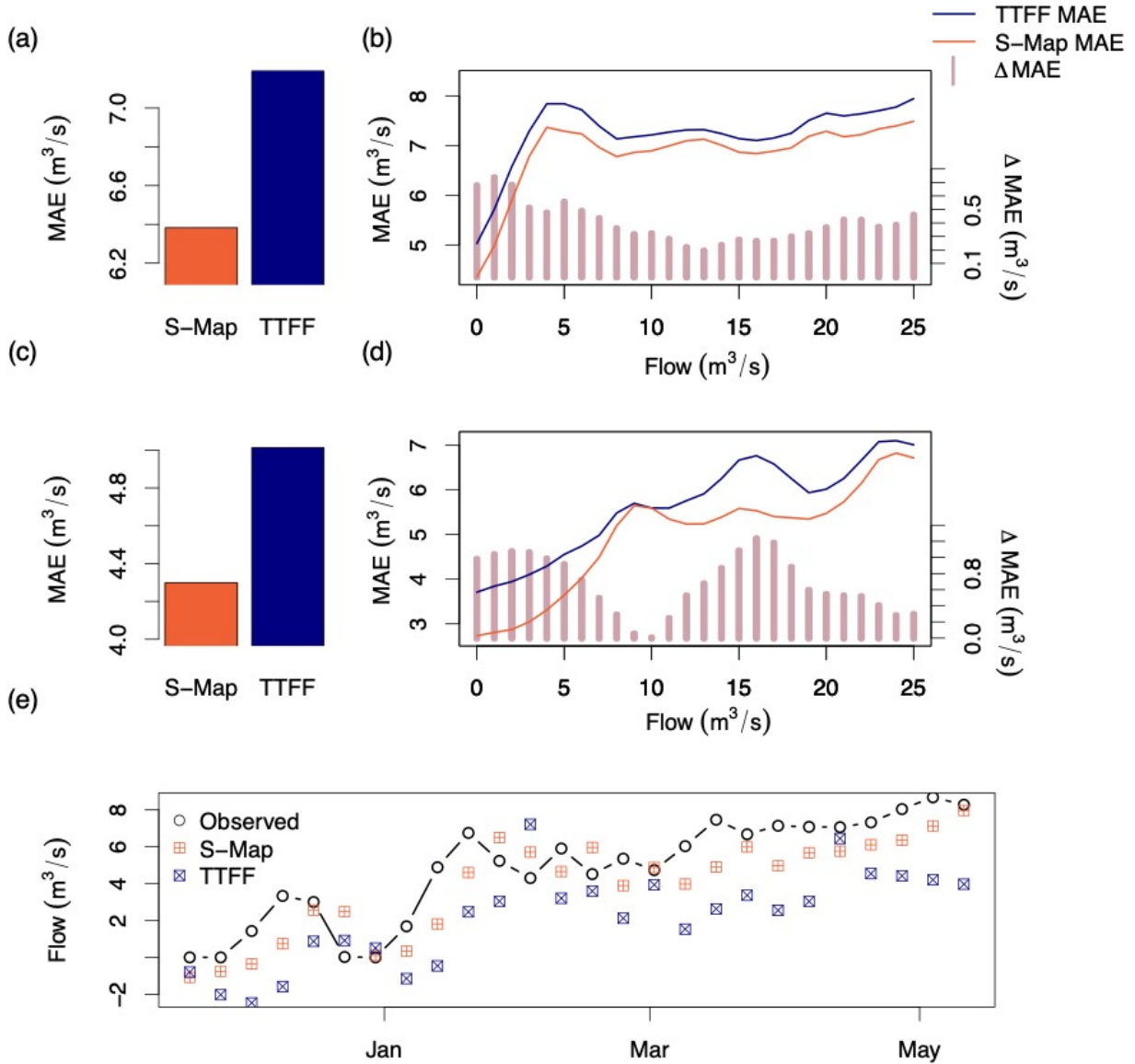


Figure 1.8 Comparison of TTF and S-Map forecasts. S-Maps here do not utilize Rain, Za, or PET. (a) and (c) show that overall, S-Maps outperform the TTF on the original dataset (1965-2005, (a)) as well as on contemporary data spanning 2007-2020 (c). (b) and (d) show average error for both forecasting algorithms as a function of flow, as well as the difference between the errors (maroon). Note that the S-Map forecasts significantly outperform the TTF during low-flow regimes. (e) shows an example of a 6-month period with relatively low flow with corresponding predictions made by the TTF (blue) and S-Maps (gold).

The TTF achieves a seemingly significant predictive accuracy with correlation between observed and predicted target flows of 0.90. However, upon inspection, it becomes apparent that such accuracy is not hard to achieve: simply predicting next week's flow will be the same as this

week achieves a comparable correlation of 0.88. By removing variables one at a time from the TTFF, model performance stays essentially constant. This suggests that relationships presumed by TTFF may not be as fully informative for forecasting dynamics of the system as one might presume.

Because correlation between observed and predicted values is obscured by the high level of autocorrelation in the system, correlation is not the best metric to determine the significance of predictions. Here, we focus on mean absolute error as a standard for measuring predictive accuracy. Using S-Maps, we find an average improvement of 0.9 m³/s per weekly prediction (from a MAE of 7.2 to 6.3 m³/s). This translates to a predicted flow of over 500,000 m³/s of water over the course of a week. Still, without a point of reference, the relative magnitude of this improvement is difficult to assess. As additional metrics of predictive accuracy, we find that predicting the correct directional change (higher or lower next week than the current week) increased from 60% with the TTFF to 70% with S-Maps. We determined a null standard for this metric to be 55% by predicting next week's change will be the same as the previous week's (i.e. if the flow target increased last week, it will increase again next week). Thus, an improvement from 60% to 70% corresponds to an improvement from 5% above the null to 15%. Further, correlation between predicted and observed changes in target flows from the prior week ($\Delta Q = Q_{t+1} - Q_t$) improved from $\rho = 0.45$ with the TTFF to $\rho = 0.58$ with S-Maps.

State-space (Non-linear) relationships

If a real-world system exhibits purely linear dynamics, reducing the amount of data used in the best-fit solution should hinder predictive accuracy because it reduces the signal-to-noise ratio (assuming equal amounts of noise throughout the timeseries). However, if partitioning the

data into state-dependent subsets leads to improved predictions, the system dynamics are in fact different within each partition (i.e. the system is nonlinear). We find that certain partitions lead to increased predictability over the general linear solution (TTF), suggesting that this system is indeed nonlinear. Specifically, we find that the seasonal partitions perform relatively best (aside from S-Maps), suggesting that the dynamics of this system are highly dependent on seasonal forcing.

A 5-year moving window performed the worst out of the models tested, obtaining a MAE of 7.3 m³/s (figure 1.6). This suggests that the system is not significantly changing on a year-to-year basis. However, that is not to say that dynamics do not change inter-annually at all; rather, the potential nonlinearity accounted for does not improve predictions more than the negative impact of using fewer data-points, reducing the signal-to-noise ratio. We did find that a 5-year window outperformed all other window sizes tested (ranging from 2 - 20 years, figure S3). This suggests window sizes that are too small have a large signal-to-noise ratio while window sizes too large obscure the nonlinear effects.

It is interesting to note that ZA seems to potentially have a relatively strong causal influence on flow targets: it has the third highest CCM value (figure S1) and is the third most important variable in multiview embeddings (figure S2). This is likely due to the strong seasonal forcing in this system: ZA represents the Zone A regulation, a waveform with constant annual periodicity. While this value may influence managed flows in the region, it more likely contributes in predictive models as a variable that helps define the time of year (season). This is affirmed by predictive accuracy being hardly diminished when it is removed from S-Map embeddings (figure 1.7), which already include sine and cosine terms to provide seasonal information.

When partitioning these results specifically into periods of high and low target flows, we find that water levels in WCA-3A are more important when making predictions during high flow periods while water levels in NESRS are relatively more important during low flow periods (1.7). This may be explained by water management operations in this region: when upstream water levels (WCA-3A) are high, water is available for release into ENP. Conversely, when downstream (NESRS) water levels are low, there is a need to release water to mitigate drought, likely at lower flow values.

The TTFF was specifically formulated on weekly data collected from 1965-2005. It is instructive to measure whether the predictive improvement obtained with S-Maps constructed using data from the same 1965-2005 period is consistent in contemporary data. Figure 1.8c\&d show that S-Maps still outperform TTFF on data spanning 2007-2020.

Although EDM outperforms the TTFF over the course of the entire timeseries on average, it may still be possible that the TTFF outperforms S-Maps during specific flow regimes. An important management concern of flows from the WCAs into ENP are low-flow regimes during dry season and drought conditions. Figure 1.8b\&d shows that S-Maps outperform the TTFF during low flow (0-25 m³/s), especially during flows close to 0 m³/s. Figure 1.8e shows an example of the significant improvement gained through using S-Maps.

The S-Map forecasts are fundamentally similar to that of the TTFF. The main differences being 1) S-Maps utilize fewer variables (they do not include Rain, ZA, or PET) and, most importantly, 2) S-Maps solve for linear fits only on *similar states* rather than on all historical data. Yet, just these two changes significantly improve forecasts. This demonstrates that nonlinear forecasting does not need to be complex; rather, it can be implemented with nearly the same ease as linear formulations while providing insight about nonlinear relationships.

Conclusion

A guiding principle of the Comprehensive Everglades Restoration Plan is to "get the water right". This refers to restoring the quantity, quality, timing and distribution of water throughout the greater Everglades system. This work focuses on the "quantity" aspect of this plan. A core component of this objective is management of water delivered across the Taimiami Trail from the upstream water conservation areas into Everglades National Park. This management is highly constrained by competing interests of hydroperiod and water depths for ecologic benefit, flood control for agricultural and urban interests, and water quality. These issues become particularly acute during the seasonal dry periods and droughts. While efforts continue to remove barriers to natural "sheetflow" across the Trail, the active management of this complex, nonlinear objective is an fundamental lever in the water managers toolbox towards Everglades restoration.

This work highlights the importance of model selection when dealing with real-world systems. In cases where the system is multidimensional and dynamic, it is ambitious to assume a single linear equation can describe the dynamics of a system. Despite this, such linear models are often favored due to their simplicity. However, significantly improved, nonlinear approaches do not necessitate significantly complicated models. Here, we used the same linear regressive approach as used to formulate the TTFF; however, we added a nonlinear perspective by partitioning the data into similar states. This effectively changes the question the models address from "what is the single set of rules that define this system?" to "what are the rules of this system when it looks like it does *right now*"? This nonlinear perspective significantly improves predictions of weekly integrated flows from the WCAs into ENP, while also revealing dynamical

truths about the system. Given that nonlinear dynamics are ubiquitous in nature, such nonlinear approaches should also be ubiquitous in management efforts.

Acknowledgements

Chapter 1, in full, is a reprint of the material as it appears in *The Journal of Water Resources Planning and Management*. Saberski, Erik, Joseph Park, Troy Hill, Erik Stabenau, and George Sugihara. "Improved Prediction of Managed Water Flow into Everglades National Park Using Empirical Dynamic Modeling." *Journal of Water Resources Planning and Management* 148, no. 12 (2022): 05022009. The dissertation author was the primary investigator and author of this paper.

Chapter 2 The Impact of Data Resolution on Dynamic Causal Inference in Multiscale Ecological Networks

Abstract

While it is commonly accepted that ecosystem dynamics are nonlinear, what is often not acknowledged is that nonlinearity implies scale-dependence. With the increasing availability of high-resolution ecological time series, there is a growing need to understand how scale and resolution in the data affect the construction and interpretation of causal networks – specifically, networks mapping how changes in one variable drive changes in others as part of a shared dynamic system (“dynamic causation”). We use Convergent Cross Mapping (CCM), a method specifically designed to measure dynamic causation, to study the effects of varying temporal and taxonomic/functional resolution in data when constructing ecological causal networks. Our analyses involve time-series data from mean-field logistic models, multi-timescale individual-based automata models, and observational data from several aquatic ecosystems. As the system is viewed at different scales relationships will appear and disappear. The relationship between data resolution and interaction presence is not random: the temporal scale at which a relationship is uncovered identifies a biologically relevant scale that drives changes in population abundance. Further, causal relationships between taxonomic aggregates (low-resolution) are shown to be influenced by the number of interactions between their component species (high-resolution). Both phenomena are observed in models and real-world data. Because no single level of resolution captures all the causal links in a system, a more complete understanding requires multiple levels when constructing causal networks. This approach provides a more nuanced view of how ecosystems operate that can improve our ability to predict and manage ecosystems.

Introduction

One of the fundamental goals of ecology is to understand causal interactions as they occur within naturally evolving ecosystems. Here causation can be direct or transitive, span multiple mechanisms (e.g., trophic, competition, mutualism, etc.), and change with ecosystem state. All of this ultimately determines how effects (natural or managed) propagate, and travel in ways that can sometimes lead to unintended consequences. Although, controlled experiments can be important for establishing direct causal links *in principle*, in practice, because interactions in nature tend to change with the evolving ecosystem state (Deyle *et al.* 2016b; Ushio *et al.* 2018; Deyle *et al.* 2022; Liu & Gaines 2022) static single-factor assessments fail to translate into predictive understanding. This is a challenge that can be met with a data-driven approach for inferring causal effects between ecosystem components using observational time series (Dixon, Milicich & Sugihara 1999; Brookshire & Weaver 2015; Deyle *et al.* 2016a; Matsuzaki *et al.* 2018; Yang, Peng & Huang 2018; Liu & Gaines 2022; Orenstein, Saberski & Briseño-Avena 2022).

How data resolution impacts perception is fundamental. Indeed, the basic notion of what constitutes the variables of study in real ecological applications is necessarily tied up in the scales of observation. For example, in some lakes we might measure chlorophyll-a every hour but only at the surface, while in others we might measure and define each known species of chlorophyte and diatom at various depths, but only once per summer. Although both observe something of dynamics underlying primary production, such differences in scale and aggregation will determine what we see as the causal factors shaping the dynamics of those observations. Accounting for these differences is the focus of this work.

While approaches that rely on a statistical framework do not assume an underlying deterministic dynamical system (Glymour *et al.* 1997; Spirtes, Glymour & Scheines 2000; Pearl 2009), here we take the position that dynamics are an essential part of the machinery. Through a dynamical systems lens, causality can be regarded as explicitly deterministic, mechanistic, and dynamic. This contrasts with statistical definitions of causality where relationships are independent of changing system states. Thus, we are interested in whether a change in one variable produces a change in another due to their mechanistic coupling in a shared dynamic system (i.e., “dynamic causation”).

Here, we revisit the role of scale and aggregation in causal pattern and process using a common data-driven approach specifically aimed at measuring dynamic causation in ecosystems: convergent cross-mapping (CCM) (Sugihara *et al.* 2012; Brookshire & Weaver 2015; Matsuzaki *et al.* 2018; Yang, Peng & Huang 2018; Chang *et al.* 2020; Liu & Gaines 2022). CCM infers causal relationships from time series data by exploiting Takens’ Theorem (Takens 1981), which states as a generic property that, quite remarkably, any one variable in a coupled dynamic system will contain information about the other variables in the network. This means that links inferred using CCM are not simply binary and independent but include transitive effects across multiple components of the full natural system. Thus, causal interaction webs produced by CCM provide a comprehensive picture of causal interdependence that can be used, for example, to effectively study direct and indirect consequences of interventions, and, in principle, it should be able to do so using readily available monitoring data.

Unlike classical structural modeling approaches for detecting causal association (Glymour *et al.* 1997; Spirtes, Glymour & Scheines 2000; Pearl 2009), CCM is specifically designed to detect nonlinear relationships that are invisible to correlation-based methods. Adding

to its practical significance in ecological network analysis is the fact that CCM does not require all relevant causal variables to be observed—a consequence of Takens' Theorem. Other commonly used methods to construct causal networks from time series data include Granger causality and structural causal models informed by machine learning (Runge *et al.* 2019). None of these methods are both explicitly nonlinear and dynamic, and many are based on conditional probabilities that require all relevant causal variables to be observed. This is a constraint that makes them less practical for ecological applications. Further, these methods focus on direct linkages which allows for easier conceptualization, yet precisely for that reason does not capture the true level of interdependence in ecosystems. CCM can address the kinds of problems that are directly relevant to conservation and management efforts (Saberski *et al.* 2022) —for example, how small changes in one variable can propagate and push a system toward or pull it away from collapse (Cenci, Sugihara & Saavedra 2019; Medeiros *et al.* 2023).

Indeed, whichever causal inference method is chosen, it is an unavoidable fact that data will be aggregated through primary observations and subsequent processing over some spatial, taxonomic, and/or temporal scale in constructing any kind of ecological network. Food webs are often constructed in terms of functional groups by pooling species into taxonomic aggregates (Dunne, Williams & Martinez 2002) as well as trophic equivalence classes (Sugihara 1983) and pollination networks have been analyzed at spatial and temporal scales that span many orders of magnitude (Bascompte & Jordano 2007; Rasmussen *et al.* 2013). Some have argued that aggregated data can reveal robust patterns and valuable insight (Sugihara, Schoenly & Trombla 1989; Martinez 1993; Sugihara, Bersier & Schoenly 1997; Dunne, Williams & Martinez 2004), while others suggest that aggregation can mask important ecosystem dynamics that arise more coherently at finer scales (Abarca-Arenas & Ulanowicz 2002; Allesina, Bondavalli & Scharler

2005; Pinnegar *et al.* 2005). Taking in both points of view suggests that one can be intentional about the often-unacknowledged choice of scale and even take full advantage of it to improve understanding of ecological dynamics on scales relevant to management.

Already as early as 1992, in his MacArthur Award keynote address Simon Levin stated: “Applied challenges, such as the prediction of the ecological causes and consequences of global climate change, require interfacing phenomena that occur on very different scales of space, time, and ecological organization” (Levin 1992). Indeed, the recognition by ecologists (Iwasa, Andreasen & Levin 1987; Allen & Starr 2017) that dynamic processes occur simultaneously at multiple spatial and temporal scales has arisen in many fields including economics (e.g., Lange-Hicks Condition (Lange 1944)) and neuroscience (Van de Ville, Britz & Michel 2010), where it is often termed the “aggregation problem”. Here the focus has been on investigating conditions under which a coarse-grained “macrosystem” view (where dynamics occur between aggregated macroscopic variables like functional groups) and a fine-grained “microsystem” view (where dynamics occur between disaggregated or less aggregated variables like population abundances of individual species) give different results (Sugihara *et al.* 1984; Sugihara, Schoenly & Trombla 1989). The simple answer given in the Lange-Hicks condition is that unless the dynamics are linear or can be separated (i.e., where fast components can be treated as if they are in equilibrium and slow components as if they are constant), scale matters. Indeed, the overwhelming evidence that ecological dynamics are nonlinear and state-dependent means that analyses at different scales will present different portraits of the functional relationships – a potential liability if ignored, but when accounted for can become a substantial asset to support the understanding and management of these systems.

Because nonlinearity implies scale-dependence and ecological dynamics are nonlinear (Sugihara 1994; Hsieh *et al.* 2005; Clark & Luis 2020; Munch *et al.* 2022), it is not surprising that identifying dynamic causal linkages will necessarily depend on the scale and resolution of the data used. However, how this plays out in practice is not known. Here we aim to provide a better understanding of the implications of scale and resolution when constructing and interpreting dynamic causal networks for ecosystems.

Methods

1) Model and Field Data

Multiple Time Scales in a 3-Component Individual-based Automata Model

To understand how dynamics spanning multiple time scales can be accommodated we construct a simple game-of-life analogue that incorporates trophic activity on multiple time scales. It is an individual-based ecological automata model with three components intended to simulate species dynamics in three trophic levels, each operating on a different time scale (Fig. 2.1): “resources”, “primary consumers”, and “secondary consumers”. The components (individuals) exist in a 2-dimensional (1000 x 1000) grid. Each individual moves randomly and follows the simple trophic rules described below to determine whether it survives or reproduces in subsequent timesteps. The system is initialized arbitrarily with 15,000 resources, 1,500 primary consumers, and 20 secondary consumers placed randomly on the grid. Each is allowed to move randomly from their position in both x and y directions with a speed (distance traveled in one timestep) s (s_R and $s_{PC} = 10$, $s_{SC} = 25$). Primary consumers eat resources within a radius $r_{PC} = 10$ and secondary consumers consume primary consumers within a radius $r_{SC} = 100$. At each timestep, 500 new resources are introduced randomly on the grid. This pedagogical example is intended to be a simple template to demonstrate scale effects, and the qualitative results are robust to the specific parameters chosen.

If an individual primary consumer consumes at least 1 resource at time t , it survives to time $t+1$; if it consumes at least 2 resources at time t it will have an offspring at time $t+1$. However, if a primary consumer consumes 0 resources at time t , it does not survive to time $t+1$. Secondary consumers follow similar rules but operate on much larger scales: it survives if it has consumed at least a minimum number of 3,000 of primary consumers in the prior 500 timesteps and has one offspring if it consumes at least 5,000 primary consumers in the last 500 timesteps (limited to 1 offspring per 200 timesteps per secondary consumer).

The simulation ran for 15,000 timesteps to generate three distinct abundance time series, one for each trophic category (Fig. 2.1). We performed CCM between each time series using $E = 4$, $\tau = 1$, $\tau_p = -1$ (high-resolution web, Fig. 2.1 middle) and $E=4$, $\tau = 500$, $\tau_p = -500$ (low-resolution web, Fig. 2.1 right). Only CCM linkages having rho values greater than the cross correlation were taken to show nonlinear causal connection. Subtracting out the linear cross-correlation is a simple way to measure whether there are causal dynamics beyond the linear correlation (Deyle *et al.* 2013).

Aggregation in a Logistic Model

We used a simple logistic model to examine the relationship between species-resolved connectance and resolved causal influence between two aggregated functional groups. This was accomplished by simulating 20 timeseries loosely representing 10 predators and 10 prey. A randomized interaction matrix defined the one-timestep relationships of each timeseries on each other. Predators had negative influences on prey and prey had positive influences on predators. For simplicity, the model simulations did not include any intra-aggregate interactions. The abundances were constrained to $[0,1]$ by taking the reciprocal of any abundance if it exceeded 1.

A connectance parameter (C) determined the number of non-zero elements in the interaction matrix. The main diagonal of the interaction matrix was set to -0.15 for all timeseries.

We performed 500 model simulations with C ranging between 0.3 and 0.9. After each simulation completed, we took the sum of the 10 predators and 10 prey to generate two aggregate timeseries. CCM was then performed to measure the influence of the predators on prey using $E = 5$, $\tau = 1$, and $\tau_p = -1$. The resolved interaction strength was measured as the correlation coefficient between observations and predictions from this CCM analysis. This was used to explore the effect of connectance on aggregated CCM values.

Observational Data from Four Aquatic Sites

We focus on four exceptional long-term ecological monitoring studies containing highly resolved time series for the individual taxa located in 1) The North Sea (from the Survey of the Marine Biological Association, formerly the Sir Alister Hardy Foundation for Ocean Science (SAHFOS), 2) Port Erin Bay (MetaBase), 3) Lake Zurich (Pomati *et al.* 2020) and 4) A kelp forest system of San Nicholas Island (Kenner *et al.* 2013). Sites 1-3 were sampled monthly while the kelp forest system was analyzed as annual averages. The four studies were chosen based on data quality: their time-series data have a high degree of continuity and overall length, and there is sufficient knowledge about the ecosystems to construct a credible food web from systematic literature searches and expert knowledge. To ensure quality and uniformity in the analyses for each dataset, taxa whose time series contained less than 35 non-zero data points or were known to be inconsistently monitored were removed from the analysis.

Food web construction

Systematic literature surveys were used to construct food webs for each system. Data extracted for each taxon included: genus, species, prey, predators, trophic role (autotroph, heterotroph, mixotroph, primary consumer, secondary consumer), and additional ecological or biological notes of interest (including, but not limited to competition, known defenses and size). In certain cases where no information could be found for a taxon, researchers from the long-term monitoring sites were solicited to provide expert opinion to fill in remaining gaps.

Final food-web constructions represent the collation of taxon-specific relationships into functional group aggregates (Fig. 2.5). Functional groups are defined by both trophic level and taxonomic criteria (e.g., Omnivorous copepods). Food-web interactions were drawn between functional groups if any member within one group had a direct trophic interaction with any member within another group. Food webs for Port Erin Bay and the North Sea datasets were reviewed by plankton experts at SAHFOS.

2) Analysis

Convergent Cross Mapping

Convergent cross-mapping (CCM) measures dynamic causation by using cross-map prediction to assess how well one time series can be used to predict another. If time series Y has been influenced by a driver X, it contains information that can be extracted (using Takens' Theorem (Takens 1981)) to predict ("map onto") values of time series X (Sugihara *et al.* 2012). Thus, in CCM the recipient time series has information that allows one to recover states of the driver, where predictions are based on time-indexed nearest-neighbors (Sugihara & May 1990) in time-lagged embeddings. For the monthly sampled systems, we use an embedding dimension (E) of 12 and for the Kelp Forest system we use an E of 4 with the prediction horizon (tp) set to

0. A “CCM value” is defined as the Pearson’s correlation between observed and predicted values.

Accommodating Seasonal Synchrony

The abundances of species within each monthly-sampled system show high levels of seasonal synchrony. As described in Sugihara et al (2012), when time series are synchronized, if applied uncritically CCM can return a false positive result (Cobey & Baskerville 2016; Sugihara, Deyle & Ye 2017). To address this, “seasonal surrogates” can be constructed that maintain the seasonal relationship between two time series but shuffle the other properties of the time series. This is accomplished here by randomly shuffling the time series values within each month (e.g., shuffling all January values, then all February values, etc.). CCM is then performed on the resulting null surrogate time series to measure how accurate cross-mapping is on samples having only this seasonal property but whose dynamics are otherwise randomized. We repeat this 100 times for each possible interaction to get a distribution of null surrogate CCM values. An observed CCM value that is more accurate than at least 95 out of the 100 surrogate values is considered significant and the causal link is included.

Causal Web Construction

To compare food webs and causal webs at the same aggregate resolution, we create aggregate time series by normalizing each species-abundance time series between 0 and 1 then add their abundances at each point in time. This normalization procedure gives species equal contribution to the aggregate time series which prevents a single, highly abundant taxa from dominating the aggregate. We then perform CCM as described above between each aggregate of

each system. High-resolution webs are constructed by performing CCM between each individual species (when the system is viewed at highest resolution).

Results

First, we investigate how varying timescales can play a role in resolving ecosystem dynamics using the individual-based automata (IBA) model. By adjusting the time-lag (τ) in reconstructed embeddings to capture causal relationships at different timescales, we find that high-frequency causal webs resolve bidirectional influences between resources and primary consumers, but only show unidirectional effects of secondary consumers on primary consumers and resources (Fig. 2.1). The high-frequency causal web does not show any influence of primary consumers or resources on secondary consumer abundance. However, these dynamics are well-resolved at a 500-timestep scale (Fig. 2.1).

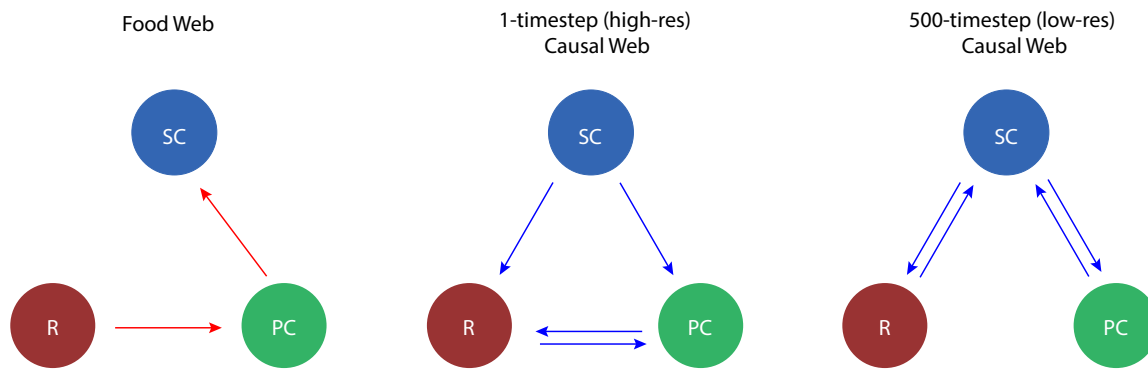
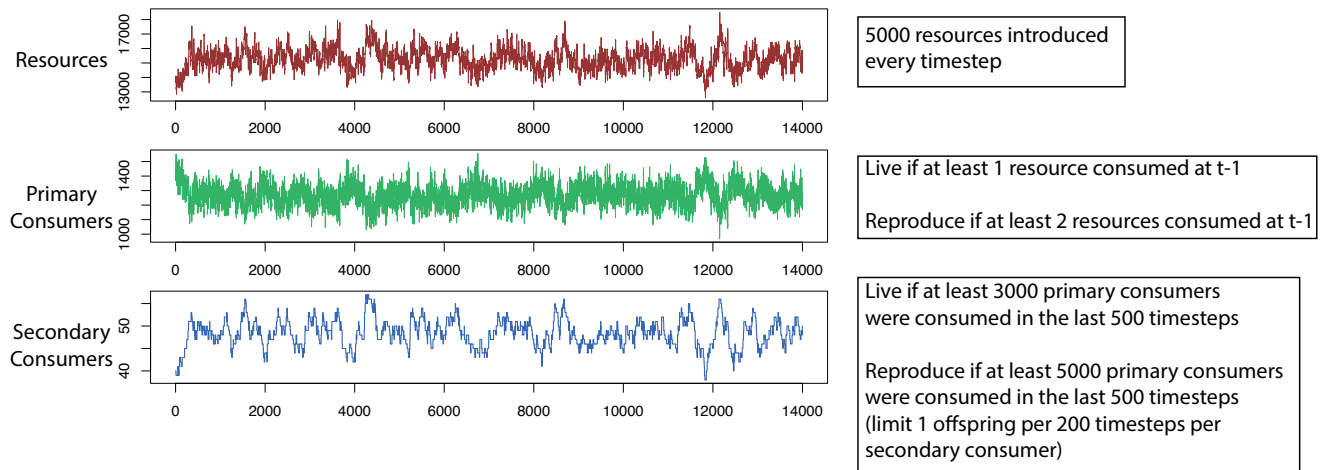


Figure 2.1 A model simulation in which resources, primary consumers, and secondary consumers exist on a grid and move randomly. Each timestep, primary consumers eat resources and secondary consumers consume primary consumers. Both primary and secondary consumers have rules determining whether they survive, starve, or reproduce. Primary consumers and resources interact on a 1-timestep timescale, and secondary consumers consume primary consumers at a 1-timestep scale as well; however, primary consumers influence secondary consumers at a 500-timestep scale. Thus, causal webs constructed using a 1-timestep frequency do not resolve the influence of primary consumers on secondary consumers (middle), but the causal webs constructed using a 500-timestep frequency do (right).

We use the same technique on the three monthly-sampled systems to quantify causal relationships at different time scales. Fig. 2.2 shows the number of network interactions resolved by CCM at each scale. All three systems had more interactions resolved at the monthly scale than at the annual scale ($\tau = 12$).

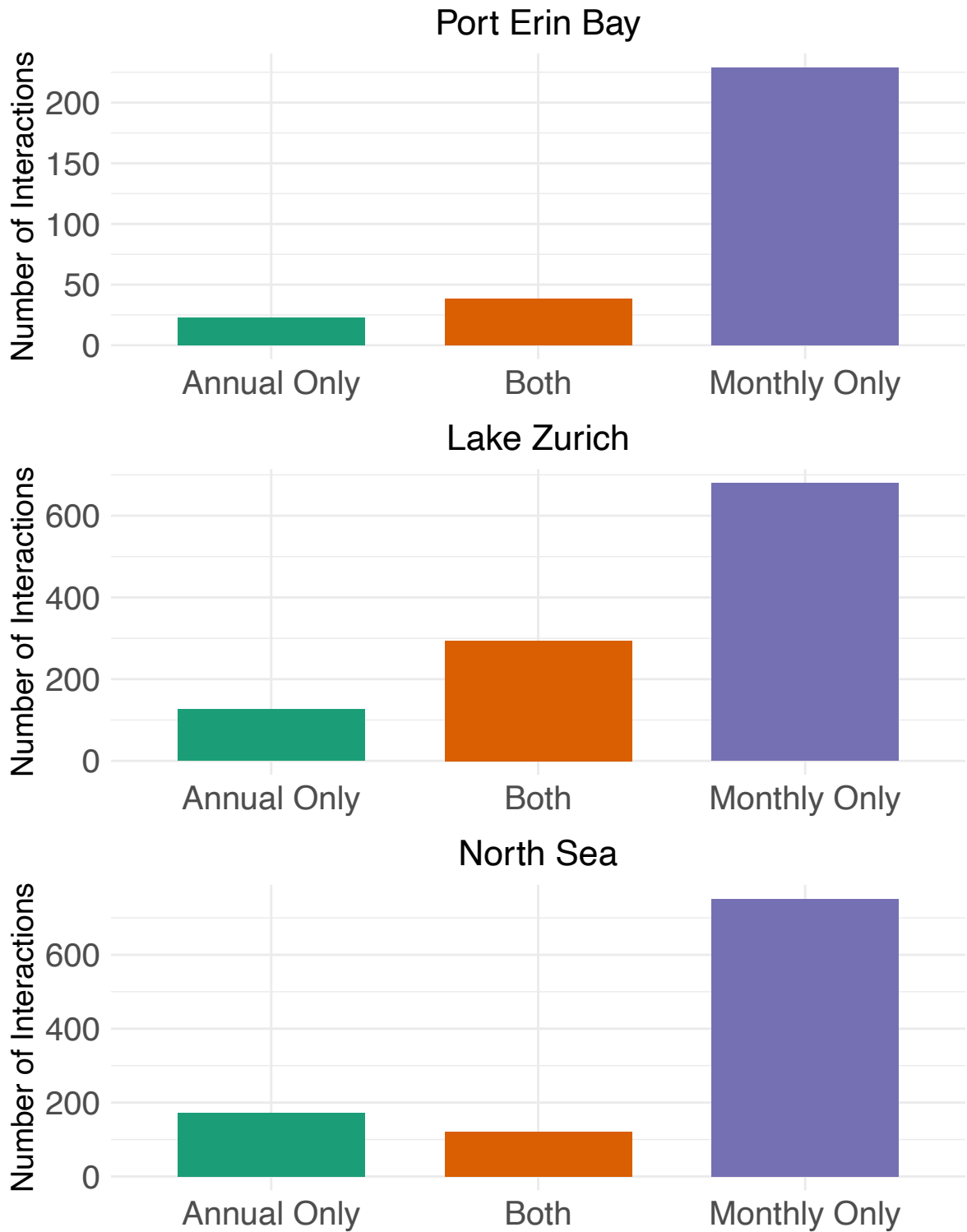


Figure 2.2 The number of interactions in each system resolved at a monthly timescale ($\tau = 1$) and annual timescale ($\tau = 12$). Note that in all systems more interactions are resolved at the monthly timescale, but there are still interactions in each system that are exclusively resolved at the annual timescale.

By varying the connectance of the interaction matrix of the logistic models, we find a positive association between model connectance and resolved interaction strength between the predator and prey aggregates (Fig. 2.3). That is, the more links there are connecting individual species across aggregates, the more likely there will be a resolved link between the aggregates in a lower-resolution network. A similar pattern is shown in real-world data (Fig. 2.4).

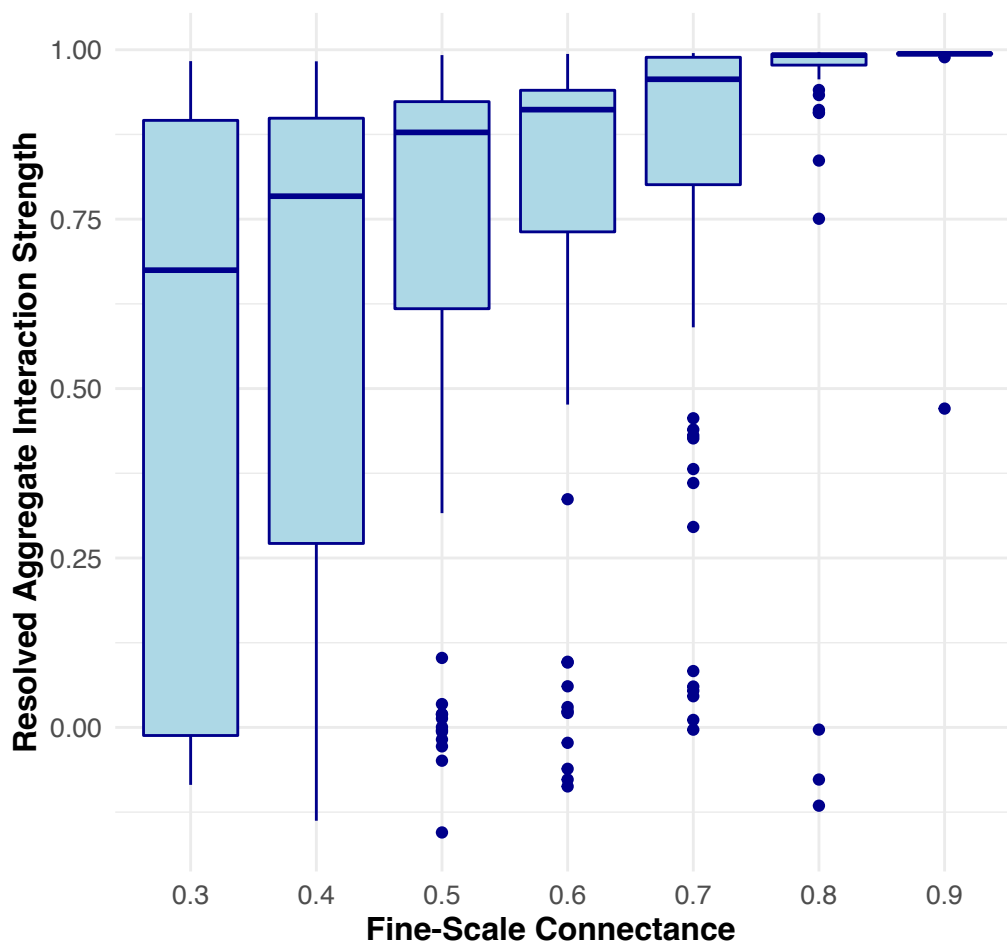


Figure 2.3 Fine scale connectance translates to aggregate interaction strength: Logistic models with two aggregates (10 predators and 10 prey) run at varying levels of high-resolution connectance (proportion of non-zero elements joining the two aggregates in the interaction matrix). Higher connectance at the fine scale creates stronger aggregated interaction strength resolved between the functional groups.

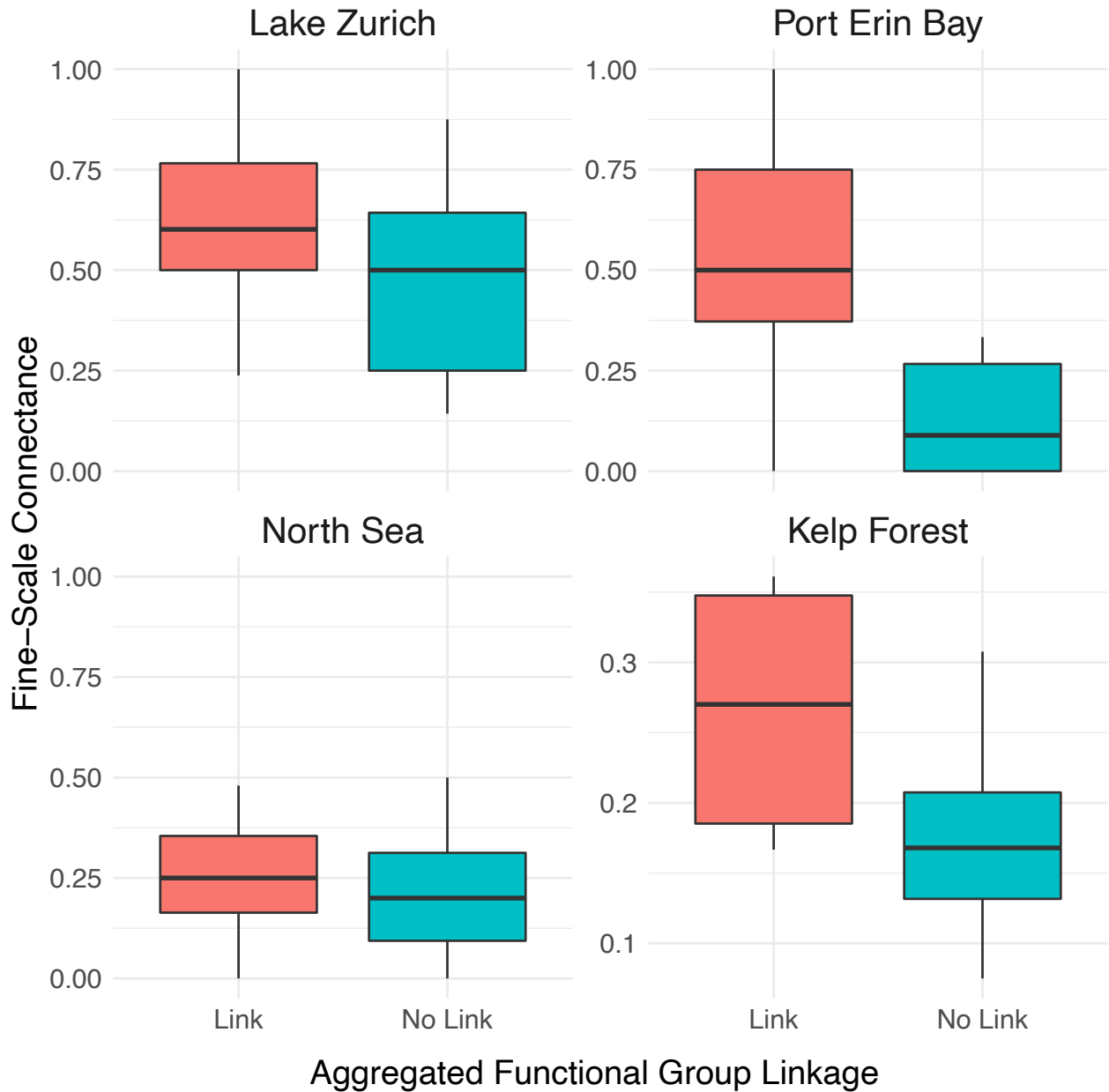


Figure 2.4 Real world examples showing the relationship between aggregated and fine-scale linkages. Boxplots show that causally linked aggregates (coarse-scale) have more causal links at the species level (fine-scale).

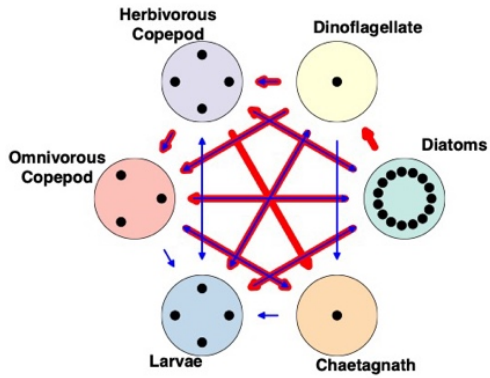
Finally, we compared the low-resolution aggregate webs with food webs constructed at the same taxonomic resolution (Fig. 2.5). We find that most food web links are resolved to be causal, but there are also causal interactions that are non-trophic, and trophic interactions that are not resolved as causal. The lack of a resolved causal link does not mean that interaction is absent; rather, there is no dynamic dependence between abundances detected in the data at that scale.

Figure 2.5 A comparison of food webs and aggregate causal webs for the four systems studied. (A) Aggregated causal webs (blue arrows) overlaid with food webs (red arrows). (B) High-resolution causal webs mapping interactions between individual species. Each large circle represents an aggregated function group of species, and each dot represents an individual species. Note that although there are more causal links (blue arrows) than food web links (red arrows) in A, not every food web link is detectably causal as might be expected from scale considerations. Indeed, when the systems are views through a high-resolution (species-resolved) lens, there is always at least one link between all trophically linked aggregates (B).

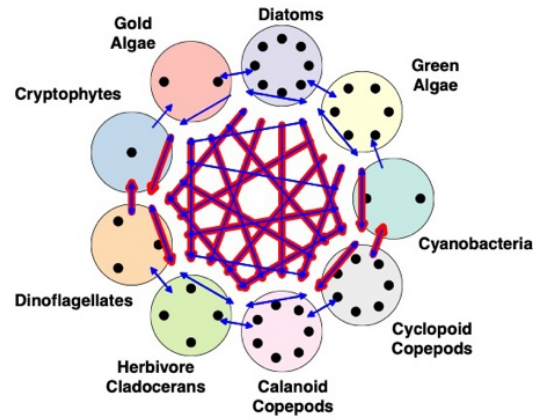
Causal Link: A → B B is influenced by A
 Trophic Link: A → B B is consumed by A
 Causal + trophic: →

A

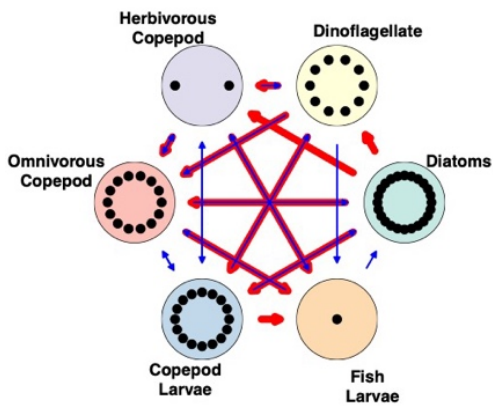
Port Erin Bay



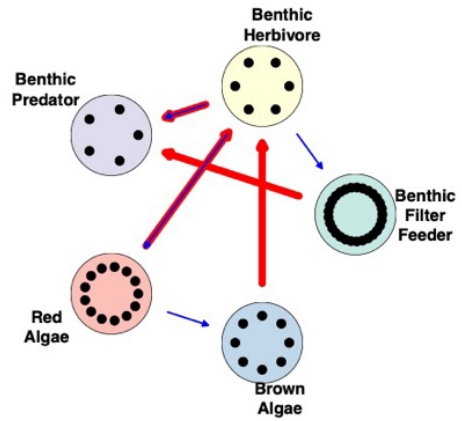
Lake Zurich



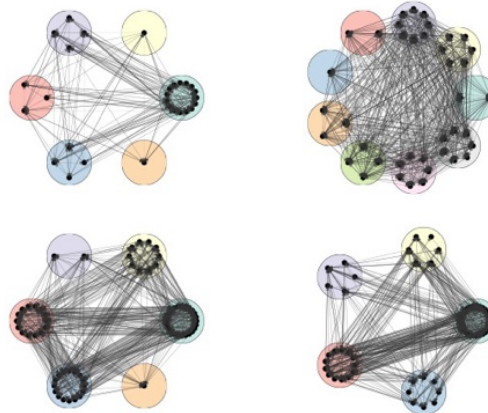
North Sea



Kelp Forest



B



Discussion

The individual-based automata (IBA) model clearly shows that different causal relationships appear when data of different temporal resolution is used. Although primary consumers (and resources, indirectly) influence secondary consumer abundance, it makes sense that the causal relationship will not be resolved using high-frequency data, as the influence of primary consumers on secondary consumers occurs over 500 timesteps. When dynamics are measured at a 500-timestep scale, the influences of resources and primary consumers on secondary consumers are resolved (Fig. 2.1).

If this were a real system and causal linkages were only measured at a high-resolution timescale, one may inaccurately conclude that the primary consumers have no influence on the abundance of secondary consumers. In a management situation, one may then incorrectly deem it is safe to alter the abundance of the primary consumers (e.g., increasing fishing, removing habitat, etc.) without any predicted consequences on the abundance of the secondary consumer. Of course, this would be a mistake since the secondary consumers are entirely dependent on the primary consumers, but at a much slower temporal scale.

The causal influences measured depend on the scales defined by the data and by parametric decisions (e.g., embedding dimension chosen for the analysis). If the data is aggregated, the analysis will quantify relationships between aggregates. Similarly, if time series are sampled at a specific frequency (e.g., monthly), then relationships that occur at that time scale will be measured (assuming time-lagged embeddings are made with 1-timestep). Interactions occurring at higher (or lower) frequencies than those captured by the sampling frequency will not be properly resolved.

We see similar patterns in real-world data. Fig. 2.2 shows interactions from three distinct ecosystems, each evaluated annually and monthly. This comparison illuminates the dynamic

nature of ecological interactions that manifest differently across these timescales: certain interactions become apparent only in the annual data, whereas others are evident exclusively in the monthly data. Thus, the temporal scale chosen will influence the resolved interaction network. Specifically, the annual data reveal interactions that may be obscured monthly due to the finer temporal resolution capturing more transient dynamics. Conversely, the monthly data can unveil more short-term interactions that might be averaged out or diluted when viewed through the broader lens of yearly assessments.

As a crude heuristic, it has been argued that scaling state-space-reconstruction to the generation times of species may help resolve dynamics (Munch *et al.* 2023). For example, In Lake Zurich the animal species with the smallest ratio of number of annual/monthly resolved influences were Cyclopoida C1-C3 (34 monthly, 13 annual), nauplii (32 monthly, 13 annual) and eggs (24 monthly, 10 annual) which either die or grow into their next life-stage on the order of days to weeks. The animal species with the largest ratio were adult Cyclops (9 monthly, 14 annual) and Eudiaptomus Gracilis (13 monthly, 16 annual) which has generation times on the order of months (WÆRVÅGEN & Nilssen 2010) (highlights in Table S1). However, because species can exhibit dynamics that span many orders of magnitude (e.g., as illustrated in Fig. 2.1) it is not surprising to find exceptions that do not follow this rule-of-thumb.

Similar to temporal scale, taxonomic aggregation can also be a double-edged sword that can obscure or clarify interaction patterns. Some have argued that aggregating abundances across multiple individuals, populations, or habitats can reveal the emergent structure of ecosystem networks by smoothing over small scale stochasticity (Sugihara, Schoenly & Trombla 1989; Sugihara, Bersier & Schoenly 1997; Dunne, Williams & Martinez 2004). Aggregation can have practical advantages in terms of increased efficiency of data collection, ease of visualization etc.

However, others argue that aggregation can also lead to oversimplification that can potentially mask important network properties (Abarca-Arenas & Ulanowicz 2002; Allesina, Bondavalli & Scharler 2005; Pinnegar *et al.* 2005).

The model simulations within our study provide a context for the debate on aggregation in ecological research (Fig. 2.3). In the simulations, timeseries were aggregated to analyze the relationship between model connectance and the aggregate relationship. The model shows that higher connectance at a fine scale translates to higher connectance at a larger aggregate scale: the more links there are connecting individual species across aggregates, the more likely there will be a resolved link between the aggregates in a lower-resolution network. This underscores the sensitivity of CCM to underlying structural parameters of ecological networks and suggests that higher fine-scale connectance, indicative of more densely interconnected ecosystems, can amplify detectable interactions in aggregate data.

Fig. 2.4 shows a similar association in four real-world systems: higher connectance between individual species across aggregates is associated with an increased likelihood of resolving a significant association between the aggregates. Thus, high-resolution connectivity within ecosystems can significantly influence the detectability of interactions among aggregated groups. However, it is worth noting that the correlation between fine-scale connectance and aggregate interaction strength may not be generalizable to other systems, and is likely sensitive to how timeseries are aggregated and other parametric choices.

In the context of ecosystem management, the absence of detected influence at an aggregate level might overlook vital high-resolution (species-resolved) interactions that are crucial for ecosystem functioning. Such an oversight could lead to management decisions that

inadvertently destabilize ecological relationships. This emphasizes the need to consider high-resolution interactions within broader ecosystem management strategies.

Fig. 2.5 shows that while most of the food web links are measurably causal, there are also causal interactions that are non-trophic (blue arrows), and trophic interactions that are non-causal (red arrows). While finding non-trophic causal interactions (blue arrows) is not surprising (competition, mutualism etc.), the converse (known trophic links that measure up as non-causal) is more surprising (red arrows) (Kawatsu *et al.* 2021). At first glance, this looks like a mistake. After all, if a predator consumes a prey, one should expect the predator to have a causal influence on the abundance of the prey. However, we emphasize that lack of a resolved causal link does not indicate a lack of interaction; rather, it reveals there is no resolved influence at that scale at which it was measured. It is likely that the non-causal trophic links may be resolved as causal links when the system is viewed at a higher taxonomic resolution or different timescale.

Conclusion

Emerging tools that allow for the collection of higher-resolution ecological data (Merz *et al.* 2021) should enable deeper insights into how ecosystems operate. These results show that beyond enabling a fine-scale view, a major advantage of high-resolution data is that it allows viewing the system at multiple time scales. The same principle applies to data collected at high taxonomic resolution. Species-level networks together with networks based on coarser functional aggregates can provide a more robust picture of ecosystem functioning than a high-resolution network can achieve alone, and critically better reflect emerging quantitative paradigms for ecosystem-based management.

Statements about causality (e.g., “species A influences species B”) are usually made in absolute terms without considering the data resolution context. This work can be a reminder that

such statements be kept pragmatic, acknowledging that relationships may appear, disappear, or change as systems are viewed at different scales.

Chapter 3 Networks of Causal Linkage Between Eigenmodes Characterize Behavioral Dynamics of *Caenorhabditis elegans*

Abstract

Behavioral phenotyping of model organisms has played an important role in unravelling the complexities of animal behavior. Techniques for classifying behavior often rely on easily identified changes in posture and motion. However, such approaches are likely to miss complex behaviors that cannot be readily distinguished by eye (e.g. high-dimensional dynamics). To explore this issue, we focus on the model organism *Caenorhabditis elegans*, where behaviors have been extensively recorded and classified. Using a dynamical systems lens, we identify high-dimensional, non-linear causal relationships between four basic shapes that describe worm motion (eigenmodes, also called “eigenworms”). We find relationships between all pairs of eigenmodes, but the timescales of the dynamics vary between pairs and across individuals. Using these varying timescales, we create “interaction profiles” to represent an individual’s behavioral dynamics. These profiles show consistent patterns among individuals in similar well-known behavioral states: i.e., the profiles for foraging individuals are distinct from those of individuals exhibiting an escape response. More importantly, we find that interaction profiles can distinguish high dimensional behaviors among mutant strains previously classified as phenotypically similar and can detect differences not previously identified in strains related to dysfunction of hermaphrodite-specific neurons.

Introduction

Caenorhabditis elegans has long been an important model species for studying the drivers of behavioral dynamics, in part due to its ease of lab culture, completely mapped nervous system, and sequenced genome (White *et al.* 1986; Waterston 1998). This body of knowledge is the basis for discoveries regarding the locomotory phenotypic consequences of changes in

genotype and neural structure (e.g. (Krajacic *et al.* 2012; Brown *et al.* 2013; Koren *et al.* 2015; Lin & Chuang 2017; Yan 2017; Javer, Ripoll-Sánchez & Brown 2018; McDiarmid, Yu & Rankin 2018)). To this end, immense work has been dedicated to generating and analyzing high-resolution recordings of *C. elegans* in motion under experimental conditions of interest (e.g. (Feng *et al.* 2004; Stephens *et al.* 2008)).

Based on such recordings, Stephens *et al.* (Stephens *et al.* 2008) found that linear combinations of just four static vectors, or “eigenworms” (also referred to as “eigenmodes”), can account for 92% of the variance of body poses formed by *C. elegans*. Thus, the posture of any worm can be represented with only a few variables at high precision. This seminal work further found relationships between the eigenworms. For example, when a worm is moving in a straight path foraging for food, the first two coefficients (a_1 and a_2 , defining sinusoidal oscillations of the body shape) form a quadrature pair: for any value of one of these coefficients during forward locomotion, the other value is predictably out-of-phase by 90 degrees ((Stephens *et al.* 2008), Fig 3.1a).

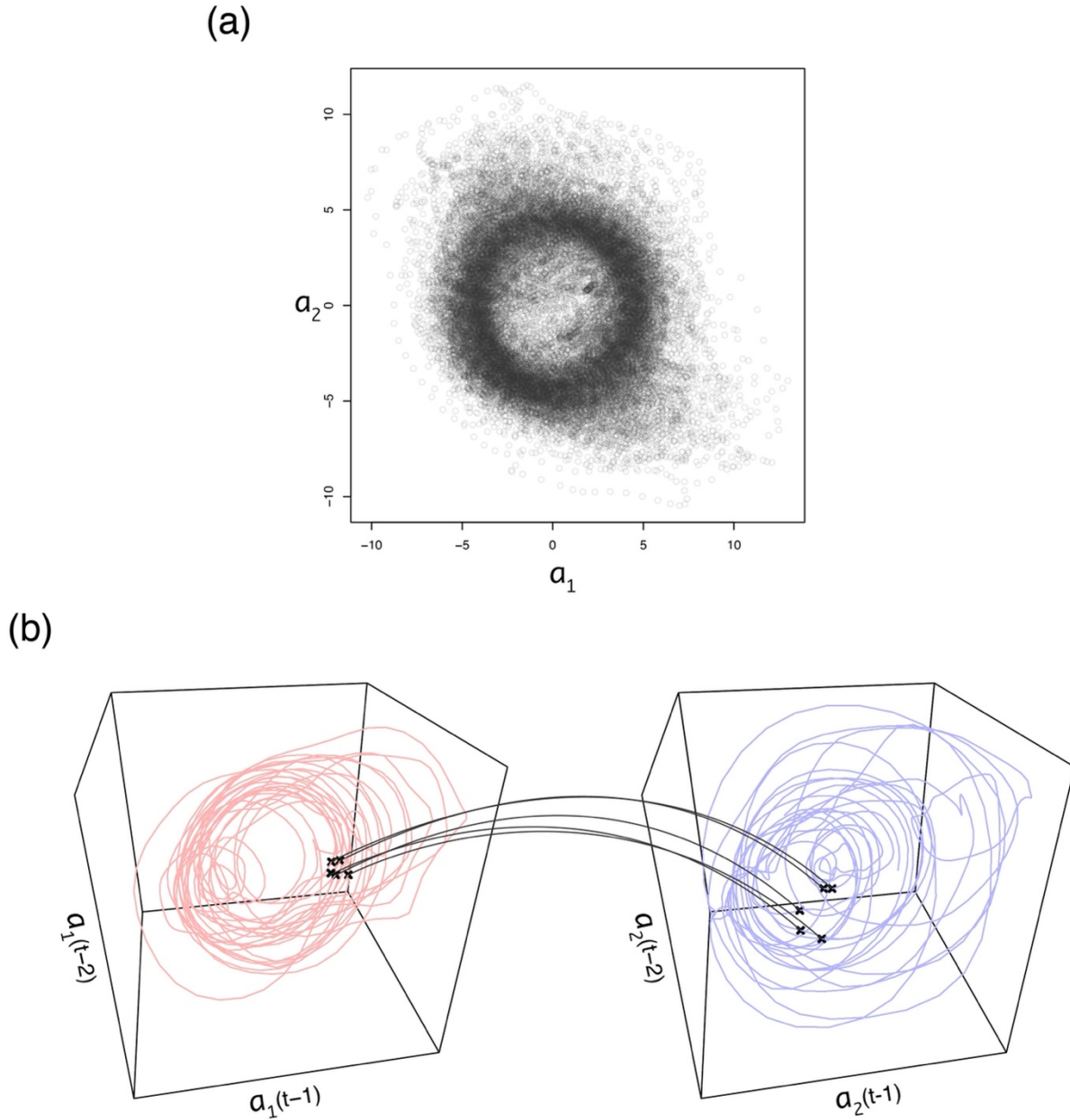


Figure 3.1 Identifying relationships between eigenmodes of worm position. (a) During forward locomotion, the static view of the first two eigenmodes form a quadrature pair; any given value of a_1 is likely to be 90 degrees out of phase with a_2 . (b) CCM identifies complex relationships between dynamically coupled variables by using time-lagged embeddings to measure the extent to which nearest neighbors on one attractor (e.g. black crosses, left) map to nearest neighbors on the other (e.g. black crosses, right).

Because eigenmode dynamics describe patterns of an individual's motion, measured changes in the eigenmodes can be used to identify changes in factors driving behavior. For

example, changes in the relationship between a_1 and a_2 have been used to quantify the effect of simulated damage to the neural network (Kunert, Maia & Kutz 2017) and changes in a_1 and a_4 (the fourth eigenmode coefficient) have been used to describe effects of locomotory neural ablation (Yan 2017). Furthermore, recent studies have used patterns within the eigenmode projections to classify distinct behaviors (Ahamed, Costa & Stephens 2019; Costa, Ahamed & Stephens 2019). Although there is a wealth of information encoded in the eigenmodes, identifying systematic relationships can be challenging. There may be many high-dimensional, complex relationships between eigenmodes that have not yet been identified. Here, we seek to uncover such hidden relationships between eigenmodes and explore how these can be used to quantitatively characterize and distinguish complex suites of behavior.

To develop an intuition for how dimensionality and dynamics interrelate, consider again the quadrature pair relationship between eigenmode coefficients a_1 and a_2 (Fig 3.1a). This static relationship describes a manifold-like structure on which a_1 and a_2 vary through time, but it does not discern dynamics along that structure. Specifically, knowing the value of a_1 and a_2 at any moment does not tell us whether the next value of a_1 and a_2 should follow the manifold in a clockwise, counterclockwise, or some other direction. In fact, forward and reverse worm motion correspond to opposite directions of rotation along this structure (Stephens *et al.* 2008). Resolving this uncertainty requires additional information. In this case, phase velocity would be sufficient, but more generally, identifying the required information (dimensions) will be challenging. Fortunately, a powerful theorem from dynamical systems theory can help resolve this problem. Takens' Embedding Theorem (Takens 1981) tells us that by taking time lags of just one of the component variables of a dynamical system, we can create an *embedding* that allows important dynamic properties of the entire system to be preserved and observed. For an

accessible introduction to this concept, and how it can be practically exploited, see the video at <http://tinyurl.com/EDM-intro>. Indeed, looking at a time-lagged coordinate representation of a_1 (Fig 3.1b) shows us that the quadrature-pair structure of the a_1, a_2 manifold is preserved, and that the third time-lag dimension reveals additional subtleties (more time-lag dimensions may be beneficial, but cannot be visualised).

To reveal hidden relationships between eigenmodes, we use convergent cross-mapping (CCM, (Sugihara *et al.* 2012)), a method based on Takens-style time-series embeddings. CCM detects nonlinear relationships between time series by tracking the temporal correspondence of nearest neighbors (Fig. 3.1b). If the nearest neighbors in the time-lag embedding formed from one variable (Fig. 3.1b, left), map temporally to nearby neighbors in the time-lag embedding formed from another variable (Fig 3.1b, right), then the state of the second variable can be predicted from the first, indicating a dynamic causal relationship from the second variable to the first (see (Sugihara *et al.* 2012) for further details).

Because several timescales may be relevant, we measure the strength of interaction across a range of delays between cause and effect, to produce an *interaction profile*, for a given ordered pair of eigenmode coefficients (see Methods). The collection of these pairwise interaction profiles between all ordered pairs of the first four eigenmode-coefficient timeseries then provides a full quantitative characterization of the worm motion dynamics. The complexity seen in *C. Elegans* behavior can (at least partially) be attributed to eigenmode relationships occurring at varying timescales; at any given time, each pair of eigenmodes may be in a different ‘phase’ of their relationship allowing for many possible combinations of active causal influences (see S1 Video). Here we assess the stability, robustness and utility of this novel characterization

by examining how it differs within and between broad categories of phenotypic behavior, in several strains of worm.

Previous studies have used a variety of methods to extract phenotypic features or behavioral motifs to cluster mutant strains, such as biomechanical profiling (e.g. (Krajacic *et al.* 2012)), construction of dictionaries of features (e.g. (Brown *et al.* 2013; Koren *et al.* 2015; Javer, Ripoll-Sánchez & Brown 2018) empirical mode decomposition of body curvature of the worm (e.g. (Lin & Chuang 2017)), and machine learning (Javer *et al.* 2018). Because nonparametric models have successfully classified *C. elegans* behaviors (Ahamed, Costa & Stephens 2019; Costa, Ahamed & Stephens 2019), we expect methods specifically designed to identify relationships in nonlinear systems (e.g. (Sugihara *et al.* 2012; Deyle *et al.* 2016b)) to excel in exploring differences between behaviors as well.

Materials and Methods

To analyze the dynamics of 12 foraging worms (Broekmans *et al.* 2016), we performed convergent cross-mapping (CCM (Sugihara *et al.* 2012)) between all pairs of the first four eigenmodes. To do this, 200 time indices (t) were randomly selected from the time series (of 33,600 values) such that none were blank (NA or NaN) and used as our library in CCM [16]. We embedded these 200 values of t in ten-dimensions, taking time lags to make points in the form $[x_t, x_{t-1}, \dots, x_{t-9}]$. This embedding was then used to make predictions on the target time series (at time $t - tp$) as described in (Sugihara *et al.* 2012). We repeated this 50 times, selecting a different sample of 200 time indices each time, and the average correlation coefficient between observed and predicted values was used as the resolved CCM value for the time delay (tp) tested. This was repeated for tp values from -32 to 0 (equivalent to -2 to 0 seconds) to quantify how the strength of the interactions change with different delays (see (Ye *et al.* 2015)). These CCM values were then normalized between 0 and 1 by subtracting the minimum value and dividing by the

difference between the maximum and minimum values. This process was repeated for all pairs of the first 4 eigenmodes ($4 \times 3 = 12$ total) across all foraging individuals (12), creating 12 plots of tp vs. CCM for each pair for each of the 12 individuals (Fig 3.2b). The average interaction profile (red lines Figs 3.2, 3.3a) is the normalized average across all individuals for each pair of eigenmodes.

To calculate the optimal embedding dimensions (E) for each individual, a similar process was done, however tp was set to the value that gave the highest CCM value in the average interaction profile for the pair of eigenmodes tested with ten dimensions. E (number of lags in the embedding) was tested from 2 to 40 with the optimal E being the number of lags that produced the highest CCM value.

For the 91 escaping worms tested (Broekmans *et al.* 2016), time series were 600 values long and sampled at 20 Hz with a laser pulse to the head occurring at 10 seconds (200th time index). Because we are only interested in escape response behavior, the first 200 values are removed from each time series, making the remaining time series 400 values long—much shorter than the 33,600 values of foraging behavior. The analysis follows similarly as that of the foraging analysis; however, to account for the shorter time series, only 100 random time indices were selected. Also, an “exclusion radius” of 10 values was also implemented, making nearest neighbors at least 10 indices (half a second) apart when making each prediction (see (Sugihara & May 1990)). tp values were then tested from -40 to 0 (still equivalent to -2 to 0 seconds due to sampling frequency) to make the interaction profile for each escape response individual.

The difference (D) between two interaction profiles was calculated by taking the average of the absolute difference of the all CCM values for each pair of eigenmodes (Fig 3.3c). This can be summarized by equation 1 below:

$$(1) \quad D = \frac{1}{12} \sum_{i=1}^4 \sum_{j \neq i}^4 |w_{1ij}^{\rightarrow} - w_{2ij}^{\rightarrow}|$$

Because the two datasets were sampled at different frequencies, values of the foraging dataset were interpolated in order to compare the difference between foraging and escape response individuals (S7 Fig). This allows for values to line up in time between foraging and escaping interaction profiles. To measure the significance of these differences, we compared these values to that of random surrogates. Surrogates were generated by randomly shuffling normalized CCM values within each interaction profile before differences were calculated. This process simulates interaction profiles that have the same distribution of values of real interaction profiles but without any temporal correspondence between them. Thus, this is what one might expect in an extreme case of two completely different interaction profiles of two individuals.

Mutant strain data was downloaded from the Worm Behavior Database from the Schafer Lab: <https://www2.mrc-lmb.cam.ac.uk/groups/wschafer/WormBehaviorDatabase.tmp.html>. Time series were filtered to only include those which had less than 25% NA's in the data and over 200 indices. This left 6,376 individuals encompassing 287 distinct strains. The analysis to generate interaction profiles for individuals of specific strains follows analogously from that of the foraging worms. However, because this data set was sampled at 30 Hz, tp was set to -60 to 0, only testing every sixth value to reduce computation time. To compare the differences between each strain, the interaction profiles for genetically identical individuals were averaged and differences between strains were taken between the averaged interaction profiles. This gave us a 287 x 287 symmetric distance matrix of differences between all pairs of strains. These values

were then used to generate the boxplots shown in Fig 3.4a,b by considering the differences between strains within specific categories defined by (Brown *et al.* 2013).

This distance matrix was also used to test whether the groups were distinct (Fig 3.4c). To measure whether two groups were distinct from each other, we considered every pairwise combination between two groups (group A and group B, Fig 3.4c). For each strain in group A, we considered which other strain had the most similar interaction profile (smallest D) out of all other strains in group A and B. After doing this for each strain in group A, we calculated the percentage of strains that had their most similar strain also in group A. This percentage was then recalculated 1000 times, however, the associations of which group each strain belonged to (A or B) was randomized each time. The level of distinction was measured as the fraction of randomized percentages that were less than that of the non-randomized percentage. In Fig 3.4c, the lowest level of significant distinction (lightest shade of green) corresponds to 80% of the randomized data showing less distinction than the real data and the darkest shade corresponds to 100%.

Results

Looking first at N2 wildtype worms during foraging (data from (Broekmans *et al.* 2016)), we find a high degree of consistency in the interaction profiles across individual worms (Fig 3.2). These profiles indicate that there exist characteristic timescales of interaction specific to particular eigenmode pairs that remain consistent throughout the foraging behavior. Indeed, contrasting these foraging interaction profiles with those obtained from the same type of worm during an escape response elicited by an aversive stimulus (data also from (Broekmans *et al.* 2016)), shows that the interaction profiles are not only consistent, but also specific: a clear qualitative difference in shape is evident in each profile between the foraging and escaping worms (Fig 3.3a).

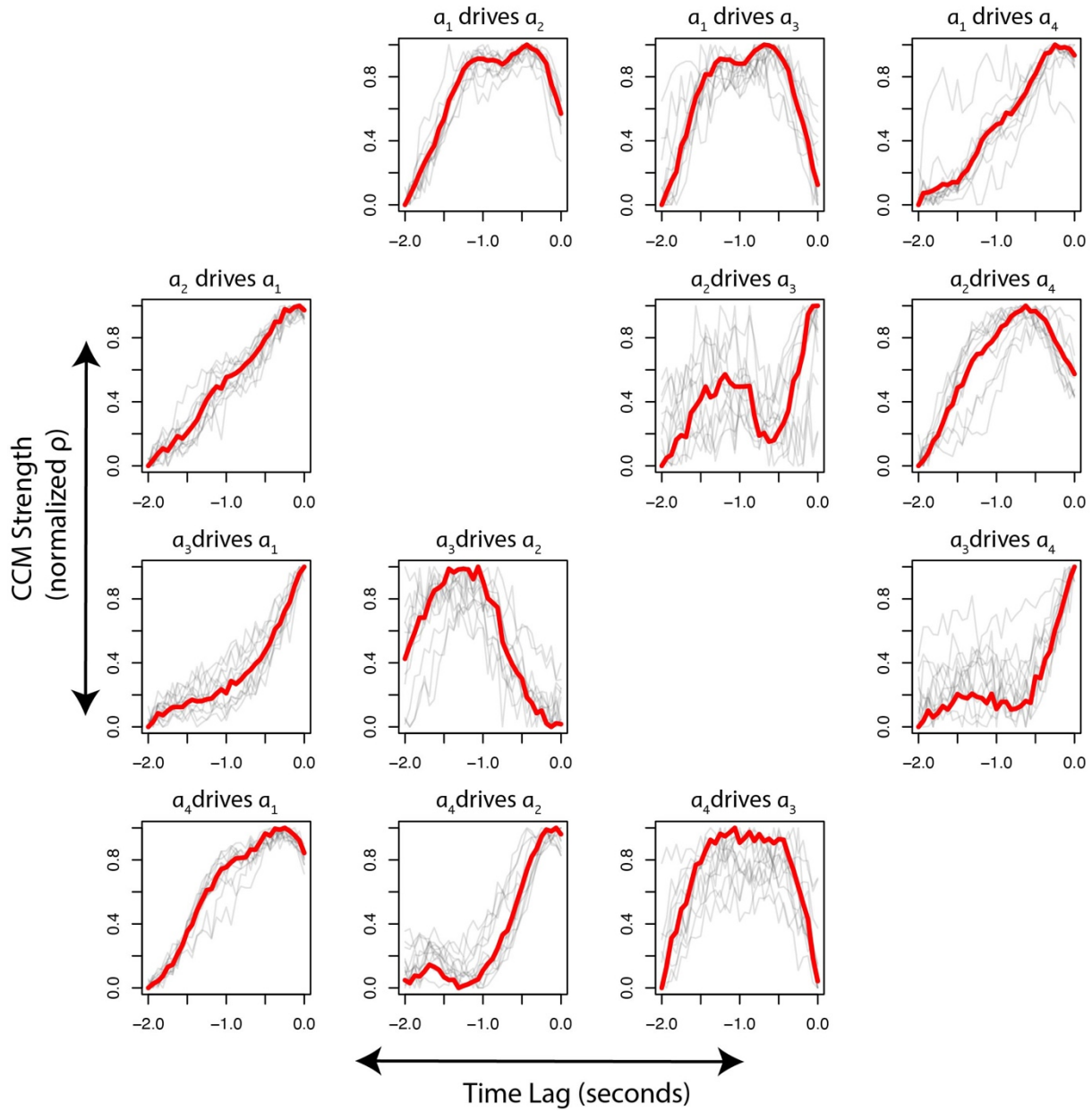


Figure 3.2 Interaction profiles represented by time-lag versus normalized CCM correlation coefficient (ρ) for each pair of the first four eigenmodes for 12 foraging worms (grey lines - individuals, red line - average). Note that pairs of eigenmodes interact at different timescales, however, these relationships are relatively consistent across individuals.

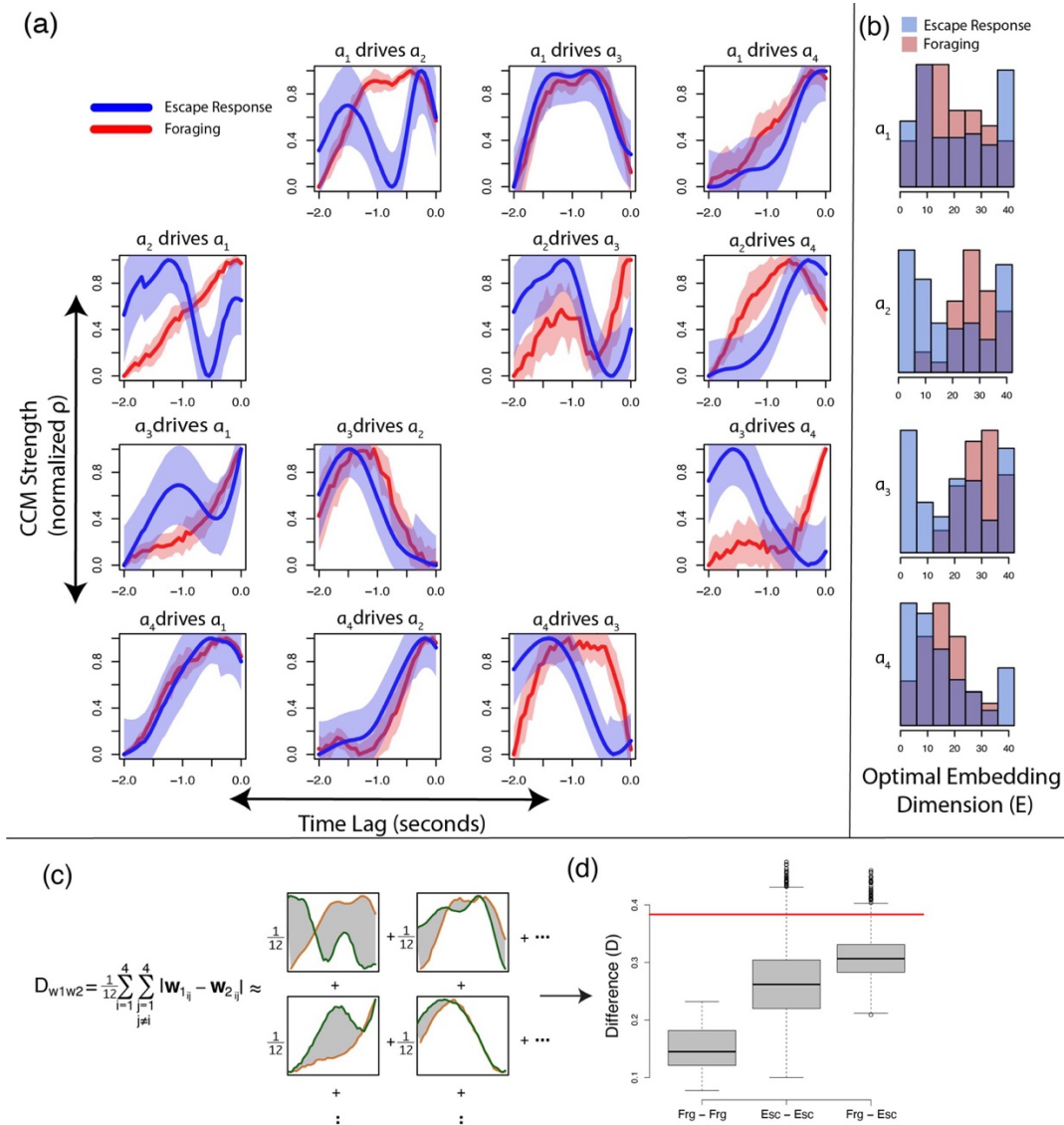


Figure 3.3 Comparing dynamics of the eigenmodes in worms foraging and exhibiting an escape response. (a) The average CCM values between the first four eigenmodes of worm position plotted against τ_p for foraging (red) and escape response (blue) worms. Shaded regions represent one standard deviation on either side of the mean. (b) The optimal embedding dimension to resolve driving dynamics in foraging (red) and escape response (blue) worms. (c) The difference between two individuals' dynamics can be quantified by adding the absolute difference between their respective CCM values versus τ_p (interaction profiles). (d) Boxplots showing the differences in dynamics between all pairs of worms. Each pair falls into one of three categories: two foraging worms (Frg - Frg), two escape response worms (Esc - Esc), or one foraging and one escape response (Frg - Esc). The red line represents the average distance between randomly shuffled surrogate profiles (see Methods).

To quantify this difference, we introduce a simple statistic: we denote by D the average absolute difference between each pair of curves across the 12 eigenmode comparisons (Fig 3.3c). Calculating D between all pairs of individuals (12 foraging and 91 escaping) shows that the average distance between pairs of foraging worms is less than that between pairs of escaping worms (t-test $p < 10^{-6}$), implying that there is more variability in escape response behaviors (Fig 3.3d). Furthermore, i) pairs of worms exhibiting the same category of behavior have a smaller difference than those exhibiting different categories of behavior (Fig 3.3d); and ii) the differences between categories of behavior are smaller than those between random surrogate data (Fig 3.3d; see Methods). Taken together, this suggests that the interaction profiles provide a quantitative characterization of these complex suites (categories) of behavior that is consistent and specific at the individual category level, and distinctive and meaningful in between-category comparisons.

The interaction profiles so far discussed were generated using 10-dimensional embeddings (see Methods), but we find that the structural consistency of these profiles is also robust to embedding dimension, across a very broad range (S2 Fig). Despite this robustness to embedding dimension, examining the effect of embedding dimension on the strength of interaction independently, may still yield additional insight into dynamical differences between behaviors. By keeping the time delay for each ordered pair of eigenmodes fixed, we examine which embedding dimensions reveal the strongest causal interaction (see Methods). Interestingly, we find that the average optimal embedding dimension for eigenmode dynamics driven by a_2 and a_3 (other eigenmodes mapped onto a_2 and a_3 by CCM) is significantly lower (t-test $p < 10^{-6}$) for worms exhibiting an escape response than for those that are foraging (Fig 3.3b). Furthermore, for foraging individuals, the driving dynamics of a_1 and a_4 can typically be

resolved in fewer dimensions than those of a_2 and a_3 , which show a strong left skew (t-test $p < 10^{-6}$ for all pairs except $[a_2, a_3]$ and $[a_1, a_4]$).

A number of genetic mutations are known to affect locomotion, and previous studies have made great strides in quantifying and categorizing the locomotory phenotypes caused by thousands of mutations (Brown *et al.* 2013; Koren *et al.* 2015; Javer, Ripoll-Sánchez & Brown 2018). From recordings of over 6,000 individuals encompassing 287 distinct mutations, we generated average interaction profiles (see Methods) for each mutation and calculated pairwise distances (D) between all pairs of mutations. Each mutation falls into one of nine categories of “phenotypically or functionally similar” as defined in (Brown *et al.* 2013). Calculating D between distinct mutant strains showed that phenotypically similar mutations had significantly smaller differences in dynamics (lower variance) than the average across all strains for 6 out of the 9 categories (t-test $p < 10^{-6}$ for 5 out of 6, $p = 0.01$ for TRP Channel, Fig 3.4a). The “egg laying defective” and “uncoordinated” groups showed average within-group differences significantly greater than the average difference between all mutants (t-test $p < 10^{-6}$). This implies that strains within these two categories have variable and dissimilar patterns of behavior. When we split the “egg laying defective” category into mutants with hypothesized or confirmed effects on hermaphrodite specific motor neurons (HSNs), neurons essential for normal reproduction (Desai & Horvitz 1989), and strains without known effects on HSNs (see S3 Table), we find that strains affecting HSNs exhibit more similar behavior than those that do not (t-test $p = 10^{-4}$, Fig 3.4b). Both subgroups, however, still show relatively high variance when compared to the differences seen within other groups of mutations.

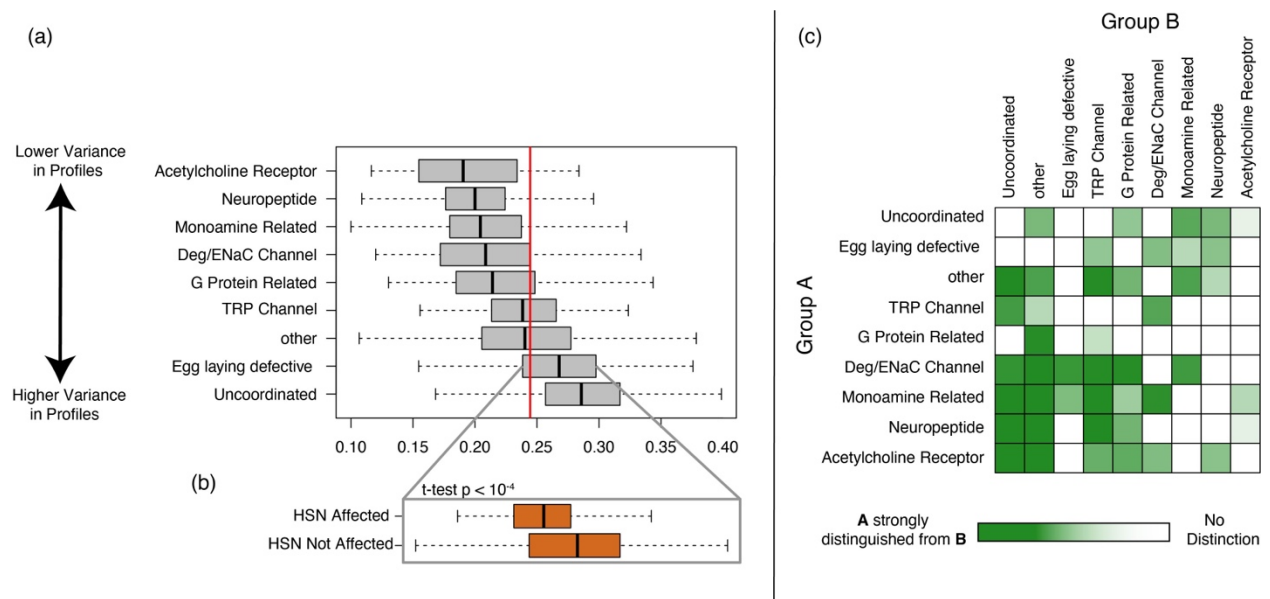


Figure 3.4 (a) Differences in dynamics between pairs of phenotypically similar mutants. The red line indicates the median distance of all mutants from each other. The boxplots are ordered from top to bottom in increasing median difference between strains, where smaller differences indicate more similar dynamics (less variance) between profiles. (b) Differences between strains of two subgroups of “egg laying defective” strains: those that affect hermaphrodite specific motor neurons (HSNs) and those that do not. Note that this difference is not identified by (Brown *et al.* 2013) (S4 Fig). (c) Groups that show significantly different interactions profiles from each other (see Methods). Note that although some groups show lower variance, they are not necessarily discernable from each other.

Although some groups show lower variance than others, they are not necessarily distinct from each other. To measure how distinct two groups are from each other, we identify the most similar strain (strain with the smallest difference D) to each strain in the two groups. If strains of one group tend to have a most similar strain that is also in that group, then that group is distinct from the other (see Methods). This process allows for directionality in these distinctions: group A may be distinct from group B, but group B may not be distinct from A. Mutations affecting acetylcholine receptors show the lowest variance, however their interaction profiles are not significantly distinct from *egl* or monoamine related mutations (Fig 3.4c). Interestingly, monoamine related mutations are, however, moderately distinct from acetylcholine receptor-affecting mutations (Fig 3.4c). In fact, monoamine related mutations are the most distinct group of mutations (distinct from all except neuropeptide-affecting mutations). Analogously, G-protein

related mutations are the least distinct group (only distinct from “other” and TRP channel-affecting mutations). Performing this analysis on the *egl* mutation subgroups (affecting/not affecting HSNs) shows that the two groups are indeed distinct from each other (see Methods).

Discussion

Worms exhibiting different categories of phenotypic behavior show significantly greater differences in their interaction profiles than those of the same behavior. Within a single strain, the interaction profiles of escaping individuals show higher variation than those of foraging worms. This may be due to slight differences in the stimuli triggering the escape response (e.g. exactly where on the worm the stimulus hit, or the position the worm was in when the stimulus occurred) and future work could explore the exact drivers of escape response behavior. Still, there are consistent patterns in the relationships between eigenmodes of escape response individuals (S5 Fig). Furthermore, when comparing differences of the interaction profiles between foraging and escaping individuals, differences are greater than those for individuals exhibiting the same behavior, however the differences are still less than that of randomly sampled surrogates, which represent a theoretical average difference between two entirely different worm behaviors (see Methods). This implies that although they are distinct behaviors, both foraging and escaping dynamics show behaviors consistent with an underlying structure that is not maintained in the randomly shuffled surrogate data.

Different classifications of mutations show different levels of variance between their component strains’ interaction profiles (Fig 3.4a). For example, mutations affecting acetylcholine receptors show relatively low variance in their interaction profiles while those of *egl* and *unc* mutations show high variance. The high variance in the *egl* and *unc* mutations can be attributed to the known variety of phenotypes caused by genes in multiple pathways within these two general groupings, as well as overlap between these categories (e.g. (Trent, Tsung & Horvitz

1983; Bany, Dong & Koelle 2003)). Although these two groups do not seem to describe one specific set of dynamics, there may be patterns among specific strains within these groups. For example, *egl* mutants with known effects on HSNs (Desai & Horvitz 1989; Garriga 1995; Jacob & Kaplan 2003; Kwok *et al.* 2006; Ringstad & Horvitz 2008; Díaz-Balzac *et al.* 2015) have less variation in their profiles than the rest of the *egl* group on average (Fig 3.4b). Although there is less variance in interaction profiles among *egl* mutations with known effects on HSNs than those without, there is still relatively high variance in both groups. In addition to controlling vulval and uterine muscles, HSNs are known to help regulate feeding (Lee *et al.* 2017) and some HSN mutations cause uncoordinated locomotion (Desai *et al.* 1988). Furthermore, the mutations classified as having a potential effect on HSNs disrupt a variety of functions, ranging from HSN cell migration in development to G-protein coupled receptors expressed in HSNs (e.g. (Desai & Horvitz 1989; Garriga 1995; Forrester, Kim & Garriga 2004)), and can play roles in the function of additional neuron types and pathways (e.g. (Jacob & Kaplan 2003)). The multiple roles of HSNs and of mutations classified as affecting HSNs could explain some of the variability among the worms in this group.

Mutations with effects on acetylcholine receptors showed the lowest variance in their interaction profiles; however, most other strains cannot be distinguished from them (Fig 3.4c). This can be imagined graphically as two clusters of points with one cluster tightly grouped within the other. If a point lies outside the inner cluster, it can be confidently determined that it is not a member of the inner cluster. However, if a point lies within the inner cluster, it cannot be confidently distinguished from the outer cluster. This may imply that the dynamics shown in acetylcholine receptor affecting groups are somewhat representative of the overall average dynamics seen in all strains. This result was not seen in all groups with low variance. For

example, the interaction profiles of strains affecting Deg/ENaC channels are largely distinguished both from other groups and other groups from them. Analogously, this can be graphically imagined as two non-overlapping clusters of points. It is possible that mutations affecting Deg/ENaC channels cause distinct changes in worm movement which manifests here as a distinct difference in their interaction profiles.

Our results resonate with that of Brown et al. (Brown *et al.* 2013), who used a nonparametric approach to cluster strains into phenotypic groups based on similarities in frequencies of repeated positions. The two methods broadly agree (S6 Fig), however, interaction profiles can show greater sensitivity in making some distinctions. We find that both methods identify high variance in the *unc* group and low variance among individuals with mutations affecting acetylcholine receptors. Also, both methods find that although individuals with mutations affecting acetylcholine receptors show low phenotypic variance, other groups cannot be confidently distinguished from them (i.e. they cluster within other groups). There are some groups that were found to be distinct only with one method and not the other. For example, mutations affecting neuropeptides and those affecting G-protein coupled receptors were not distinct from each other in (Brown *et al.* 2013) but were found to be distinct based on their interaction profiles (Fig 3.4c). Further, the differences identified here in HSN-affecting strains is not identified in Brown et al. (S4 Fig).

Still, there are a few differences identified in (Brown *et al.* 2013) that are missed with the interaction profiles described here. For example, Brown et al. found that monoamine related mutations are phenotypically distinct from those affecting neuropeptides while their interaction profiles were not significantly different (Fig 3.4c). However, it is possible that some dynamics identified by Brown et al. are not resolved with the methods described here. For example, certain

behaviors (motifs) identified in their work may be best resolved in lower-dimensional spaces. Interaction profiles here strictly measure how well dynamics are resolved in 10 dimensions and consequently may not resolve lower-dimensional dynamics. It is encouraging that even with a fixed embedding dimension, interaction profiles can identify novel phenotypic differences between strains. Future work should explore the potential of exploring interaction profiles without fixed parameters.

Interestingly, when examining the effect of dimensionality on foraging and escaping N2 worm behavior representation, the dynamics of a_2 and a_3 show lower optimal embedding dimensions for escape than for foraging. We note that reduction in dimensionality of dynamics is also observed in other systems under atypical or stressful conditions, such as in brain activity preceding an epileptic seizure (Scheffer *et al.* 2009).

Beyond exploring interaction profiles in other embedding dimensions, it would be interesting to consider other distance metrics between these profiles. Perhaps certain relationships (different panels in Figs 3.2, 3.3a) are more indicative of certain changes in worm behavior. The difference metric used here is minimalistic in that it makes very few assumptions and is consequently far from optimized; differences between distinct strains may become further resolved with improved distance metrics.

Nonetheless, we find interaction profiles reveal novel distinctions between groups of mutations without exhaustively testing parameters (distance metric, embedding dimension, eigenmodes tested, etc.). These findings shed light on the potential of using these complex relationships between eigenmodes as a classifier of worm behavior. Future work should explore the extent to which further distinctions between strains can be made with different parametric considerations.

Acknowledgements

Chapter 3, in full, is a reprint of the material as it appears in PLoS Computational Biology. Saberski, Erik, Antonia K. Bock, Rachel Goodridge, Vitul Agarwal, Tom Lorimer, Scott A. Rifkin, and George Sugihara. "Networks of causal linkage between eigenmodes characterize behavioral dynamics of *Caenorhabditis elegans*." *PLoS computational biology* 17, no. 9 (2021): e1009329. The dissertation author was the primary investigator and author of this paper.

Chapter 4 Conclusion

We do not live in a linear world with static relationships. Behaviors that we observe today may not be seen tomorrow. Similarly, patterns we see in one part of the world may not exist in another part of the world. We live in a state-dependent world, where the rules change depending on the conditions at play. In order to best understand and ultimately predict natural, complex systems, we must consider how such rules and relationships change with respect to time, space, resolution, and more generally, **state**.

As shown in chapter 1, it is all too common in management situations to choose the simple solution over the complex. However, if you are to take anything away from this thesis, I hope it is that non-linearity does not necessitate complexity. With only a small amount of work a linear solution can become a much-improved non-linear solution by simply considering, for example, how relationships change with the season (Saberski *et al.* 2022).

If you are to dissect the steps taken in chapter 1 starting from the original formulation of the TTFF to the much-improved non-linear predictor it would only include one extra step.

TTFF Linear Formulation:

- 1) Collect historical data.
- 2) Run a multi-variate linear regression on the data.

TTFF Non-linear Formulation:

- 1) Collect historical data.
- 2) Group the data into similar-looking chunks.
- 3) Run a multi-variate linear regression on each chunk of data.

This one extra step immediately takes what was a linear solution that is *never* a true representation of the system but rather a global average of its dynamics and makes it something that can be a true representation of the system at a given point in time.

Beyond considering how relationships may potentially change over time, chapter 1 also highlights the importance of what it means for a relationship to be *causal*. A causal relationship is one in which a change in one variable induces a change in another. Thus, if X has a causal influence on Y then knowing the current state of X can help predict a future state of Y (note: the directionality stated here is reversed from how the causality test CCM is performed). The variables chosen for the original TTFB formulation were picked due to their hypothesized influence on target flows. However, after testing for causality we found that two of the variables (precipitation and the Zone A regulation) had little-to-no causal relationship on target flows. In fact, including them made predictions worse.

Data can tell us a lot about how systems are structured and how they function. We do not need to rely on hypothesized relationships to build models. Rather, as shown in chapter 1, we can both identify which relationships are causal, and construct non-linear state-dependent models based entirely on what the data shows. However, building models entirely from data is a double-edged sword: when you use data to construct models your models will be held to the same scaling constraints as your data. Chapter 2 explores how this plays out both in model and real-world systems.

Data is always confined by some scale. For example, data can be confined to some spatial region, a specific temporal resolution, or an aggregation across some population of individuals. Relationships and dynamics obtained from the data will be inherently specific to the scale of the data. This makes mapping system dynamics tricky especially since one system can exhibit dynamics across many scales.

In figure 2.1, we construct a simple model that has three species exhibiting dynamics at varying temporal scales. Specifically, blue dots influence green dots at a 1-timestep scale, but

green dots influence blue dots at a 500-timestep scale. When the data is resolved to a 1-timestep resolution, no significant causal relationship is identified for green dots influencing blue dots. In a management setting, this could lead to potentially catastrophic outcomes: if one ran their causality test and found no influence of one species on another, they may deem it safe to alter or remove that species with no predicted influence on the other. However, there may in fact be an important relationship at a different scale. In this model example, when causality was measured at a 500-timestep scale there was a clear association between the two.

Relationships do not only change over time – they can change across scale as well. Consider figure 1.5: on the y axis is the magnitude of some relationship and on the x axis there is time. However, imagine if the x axis was instead temporal resolution, spatial scale, or number of species aggregated. A similar pattern is likely to emerge: relationships that are dynamic.

It can be daunting to quantify interactions in such a dynamic world. Tools like S-Map offer a broad way to scope data to only similar states – however, if your given embedding is not complete even similar states can evolve differently (a so-called “singularity”). Consider the dynamics of the nematode *C. Elegans* described in chapter 3. At any moment a worm foraging may be in a similar eigen position (similarly valued eigenworms) as a worm exhibiting an escape response. However, if we could know in advance to separate these two states out ahead of time, we could avoid this potential singularity.

This logic is similar to that described in chapter 1 (figure 1.6) where we split river flow data into different states manually based on time of year, upstream and downstream water level, precipitation, etc. In chapter 1 however, we knew ahead of time how to separate these distinct states. In chapter 3 we developed a tool to separate states without knowing ahead of time which were distinct. This is analogous to being able to know which season’s dynamics are being

exhibited by the river flow without having any information about what month it currently is. Based on the dynamic fingerprint (the “*interaction profile*”), we can separate the distinct states directly from the data.

Natural systems are complex because they have many interacting parts, where each relationship can change depending on the state. To best understand and predict systems we need to accurately depict these interactions among their components, considering exactly how they change over time (chapter 1) and change with respect to scale (chapter 2). Further, dynamic relationships resolved directly from data can be leveraged to improve our ability to classify distinct system states (chapter 3).

In conclusion, this thesis not only challenges the traditional reliance on linear models in understanding natural systems but also provides a comprehensive framework for approaching their complexity. By embracing non-linearity, acknowledging the importance of scale, and utilizing data-driven insights, we can better understand, predict, and manage the dynamic and intricate systems that govern our natural world. This work sets the stage for future research and practical applications, paving the way for more effective and sustainable management of complex natural systems.

REFERENCES

Abarca-Arenas, L.G. & Ulanowicz, R.E. (2002) The effects of taxonomic aggregation on network analysis. *Ecological Modelling*, **149**, 285-296.

Ahamed, T., Costa, A.C. & Stephens, G.J. (2019) Capturing the Continuous Complexity of Behavior in *C. elegans*. *arXiv preprint arXiv:1905.10559*.

Allen, T.F. & Starr, T.B. (2017) *Hierarchy*. University of Chicago Press.

- Allesina, S., Bondavalli, C. & Scharler, U.M. (2005) The consequences of the aggregation of detritus pools in ecological networks. *Ecological Modelling*, **189**, 221-232.
- Bany, I.A., Dong, M.-Q. & Koelle, M.R. (2003) Genetic and cellular basis for acetylcholine inhibition of *Caenorhabditis elegans* egg-laying behavior. *Journal of Neuroscience*, **23**, 8060-8069.
- Bascompte, J. & Jordano, P. (2007) Plant-animal mutualistic networks: the architecture of biodiversity. *Annual review of ecology, evolution, and systematics*, 567-593.
- Broekmans, O.D., Rodgers, J.B., Ryu, W.S. & Stephens, G.J. (2016) Resolving coiled shapes reveals new reorientation behaviors in *C. elegans*. *Elife*, **5**, e17227.
- Brookshire, E. & Weaver, T. (2015) Long-term decline in grassland productivity driven by increasing dryness. *Nature communications*, **6**, 1-7.
- Brown, A.E., Yemini, E.I., Grundy, L.J., Jucikas, T. & Schafer, W.R. (2013) A dictionary of behavioral motifs reveals clusters of genes affecting *Caenorhabditis elegans* locomotion. *Proceedings of the National Academy of Sciences*, **110**, 791-796.
- Cenci, S., Sugihara, G. & Saavedra, S. (2019) Regularized S-map for inference and forecasting with noisy ecological time series. *Methods in Ecology and Evolution*, **10**, 650-660.
- Chang, C.-W., Ushio, M. & Hsieh, C.-h. (2017) Empirical dynamic modeling for beginners. *Ecological research*, **32**, 785-796.
- Chang, C.W., Ye, H., Miki, T., Deyle, E.R., Souissi, S., Anneville, O., Adrian, R., Chiang, Y.R., Ichise, S. & Kumagai, M. (2020) Long-term warming destabilizes aquatic ecosystems through weakening biodiversity-mediated causal networks. *Global Change Biology*, **26**, 6413-6423.
- Clark, T. & Luis, A.D. (2020) Nonlinear population dynamics are ubiquitous in animals. *Nature ecology & evolution*, **4**, 75-81.
- Cobey, S. & Baskerville, E.B. (2016) Limits to causal inference with state-space reconstruction for infectious disease. *PloS one*, **11**, e0169050.

- Costa, A.C., Ahamed, T. & Stephens, G.J. (2019) Adaptive, locally linear models of complex dynamics. *Proceedings of the National Academy of Sciences*, **116**, 1501-1510.
- Council, N.R. (2008) *Progress toward restoring the Everglades: the second biennial Review-2008*. National Academies Press.
- Desai, C., Garriga, G., McIntire, S.L. & Horvitz, H.R. (1988) A genetic pathway for the development of the *Caenorhabditis elegans* HSN motor neurons. *Nature*, **336**, 638-646.
- Desai, C. & Horvitz, H.R. (1989) *Caenorhabditis elegans* mutants defective in the functioning of the motor neurons responsible for egg laying. *Genetics*, **121**, 703-721.
- Deyle, E.R., Bouffard, D., Frossard, V., Schwefel, R., Melack, J. & Sugihara, G. (2022) A hybrid empirical and parametric approach for managing ecosystem complexity: Water quality in Lake Geneva under nonstationary futures. *Proceedings of the National Academy of Sciences*, **119**, e2102466119.
- Deyle, E.R., Fogarty, M., Hsieh, C.-h., Kaufman, L., MacCall, A.D., Munch, S.B., Perretti, C.T., Ye, H. & Sugihara, G. (2013) Predicting climate effects on Pacific sardine. *Proceedings of the National Academy of Sciences*, **110**, 6430-6435.
- Deyle, E.R., Maher, M.C., Hernandez, R.D., Basu, S. & Sugihara, G. (2016a) Global environmental drivers of influenza. *Proceedings of the National Academy of Sciences*, **113**, 13081-13086.
- Deyle, E.R., May, R.M., Munch, S.B. & Sugihara, G. (2016b) Tracking and forecasting ecosystem interactions in real time. *Proceedings of the Royal Society B-Biological Sciences*, **283**.
- Deyle, E.R. & Sugihara, G. (2011) Generalized theorems for nonlinear state space reconstruction. *Plos one*, **6**, e18295.
- Díaz-Balzac, C.A., Lázaro-Peña, M.I., Ramos-Ortiz, G.A. & Bülow, H.E. (2015) The adhesion molecule KAL-1/anosmin-1 regulates neurite branching through a SAX-7/L1CAM–EGL-15/FGFR receptor complex. *Cell reports*, **11**, 1377-1384.
- Dixon, P.A., Milicich, M.J. & Sugihara, G. (1999) Episodic fluctuations in larval supply. *Science*, **283**, 1528-1530.

- Dunne, J.A., Williams, R.J. & Martinez, N.D. (2002) Food-web structure and network theory: the role of connectance and size. *Proceedings of the National Academy of Sciences*, **99**, 12917-12922.
- Dunne, J.A., Williams, R.J. & Martinez, N.D. (2004) Network structure and robustness of marine food webs. *Marine Ecology Progress Series*, **273**, 291-302.
- Feng, Z., Cronin, C.J., Wittig, J.H., Sternberg, P.W. & Schafer, W.R. (2004) An imaging system for standardized quantitative analysis of *C. elegans* behavior. *BMC bioinformatics*, **5**, 115.
- Forrester, W.C., Kim, C. & Garriga, G. (2004) The *Caenorhabditis elegans* Ror RTK CAM-1 inhibits EGL-20/Wnt signaling in cell migration. *Genetics*, **168**, 1951-1962.
- Garriga, G. (1995) Genetic Analysis of Neuronal Migration in the Nematode *Caenorhabditis elegans*. *Neural Cell Specification*, pp. 105-110. Springer.
- Glymour, C., Madigan, D., Pregibon, D. & Smyth, P. (1997) Statistical themes and lessons for data mining. *Data mining and knowledge discovery*, **1**, 11-28.
- Hsieh, C.-h., Glaser, S.M., Lucas, A.J. & Sugihara, G. (2005) Distinguishing random environmental fluctuations from ecological catastrophes for the North Pacific Ocean. *Nature*, **435**, 336-340.
- Iwasa, Y., Andreasen, V. & Levin, S. (1987) Aggregation in model ecosystems. I. Perfect aggregation. *Ecological Modelling*, **37**, 287-302.
- Jacob, T.C. & Kaplan, J.M. (2003) The EGL-21 carboxypeptidase E facilitates acetylcholine release at *Caenorhabditis elegans* neuromuscular junctions. *Journal of Neuroscience*, **23**, 2122-2130.
- Javer, A., Brown, A.E., Kokkinos, I. & Rittscher, J. (2018) Identification of *C. elegans* strains using a fully convolutional neural network on behavioural dynamics. *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 0-0.
- Javer, A., Ripoll-Sánchez, L. & Brown, A.E. (2018) Powerful and interpretable behavioural features for quantitative phenotyping of *Caenorhabditis elegans*. *Philosophical Transactions of the Royal Society B: Biological Sciences*, **373**, 20170375.

- Kawatsu, K., Ushio, M., van Veen, F.F. & Kondoh, M. (2021) Are networks of trophic interactions sufficient for understanding the dynamics of multi-trophic communities? Analysis of a tri-trophic insect food-web time-series. *Ecology Letters*, **24**, 543-552.
- Kenner, M.C., Estes, J.A., Tinker, M.T., Bodkin, J.L., Cowen, R.K., Harrold, C., Hatfield, B.B., Novak, M., Rassweiler, A. & Reed, D.C. (2013) A multi-decade time series of kelp forest community structure at San Nicolas Island, California. *Ecological Archives E*, **94**, 244.
- Koren, Y., Sznitman, R., Arratia, P.E., Carls, C., Krajacic, P., Brown, A.E. & Sznitman, J. (2015) Model-independent phenotyping of *C. elegans* locomotion using scale-invariant feature transform. *PLoS One*, **10**, e0122326.
- Krajacic, P., Shen, X., Purohit, P.K., Arratia, P. & Lamitina, T. (2012) Biomechanical profiling of *Caenorhabditis elegans* motility. *Genetics*, **191**, 1015-1021.
- Kunert, J.M., Maia, P.D. & Kutz, J.N. (2017) Functionality and robustness of injured connectomic dynamics in *C. elegans*: linking behavioral deficits to neural circuit damage. *PLoS computational biology*, **13**, e1005261.
- Kwok, T.C., Ricker, N., Fraser, R., Chan, A.W., Burns, A., Stanley, E.F., McCourt, P., Cutler, S.R. & Roy, P.J. (2006) A small-molecule screen in *C. elegans* yields a new calcium channel antagonist. *Nature*, **441**, 91-95.
- Lange, O. (1944) *Price flexibility and employment*. Principia Press Bloomington, IN.
- Lee, K.S., Iwanir, S., Kopito, R.B., Scholz, M., Calarco, J.A., Biron, D. & Levine, E. (2017) Serotonin-dependent kinetics of feeding bursts underlie a graded response to food availability in *C. elegans*. *Nature communications*, **8**, 1-11.
- Levin, S.A. (1992) The problem of pattern and scale in ecology: the Robert H. MacArthur award lecture. *Ecology*, **73**, 1943-1967.
- Lin, L.-C. & Chuang, H.-S. (2017) Analyzing the locomotory gaitprint of *Caenorhabditis elegans* on the basis of empirical mode decomposition. *PloS one*, **12**, e0181469.
- Liu, O.R. & Gaines, S.D. (2022) Environmental context dependency in species interactions. *Proceedings of the National Academy of Sciences*, **119**, e2118539119.

- Lorimer, T., Goodridge, R., Bock, A.K., Agarwal, V., Saberski, E., Sugihara, G. & Rifkin, S.A. (2021) Tracking changes in behavioural dynamics using prediction error. *Plos one*, **16**, e0251053.
- Martinez, N.D. (1993) Effects of resolution on food web structure. *Oikos*, 403-412.
- Matsuzaki, S.i.S., Suzuki, K., Kadoya, T., Nakagawa, M. & Takamura, N. (2018) Bottom-up linkages between primary production, zooplankton, and fish in a shallow, hypereutrophic lake. *Ecology*, **99**, 2025-2036.
- McDiarmid, T.A., Yu, A. & Rankin, C. (2018) Beyond the response—High throughput behavioral analyses to link genome to phenome in *Caenorhabditis elegans*. *Genes, Brain and Behavior*, **17**, e12437.
- Medeiros, L.P., Allesina, S., Dakos, V., Sugihara, G. & Saavedra, S. (2023) Ranking species based on sensitivity to perturbations under non-equilibrium community dynamics. *Ecology Letters*, **26**, 170-183.
- Merz, E., Kozakiewicz, T., Reyes, M., Ebi, C., Isles, P., Baity-Jesi, M., Roberts, P., Jaffe, J.S., Dennis, S.R. & Hardeman, T. (2021) Underwater dual-magnification imaging for automated lake plankton monitoring. *Water Research*, **203**, 117524.
- Munch, S.B., Rogers, T.L., Johnson, B.J., Bhat, U. & Tsai, C.-H. (2022) Rethinking the Prevalence and Relevance of Chaos in Ecology. *Annual Review of Ecology, Evolution, and Systematics*, **53**, 227-249.
- Munch, S.B., Rogers, T.L., Symons, C.C., Anderson, D. & Pennekamp, F. (2023) Constraining nonlinear time series modeling with the metabolic theory of ecology. *Proceedings of the National Academy of Sciences*, **120**, e2211758120.
- National Academies of Sciences, E., and Medicine (2019) *Progress Toward Restoring the Everglades: The Seventh Biennial Review-2018*. National Academies Press.
- Orenstein, E.C., Saberski, E. & Briseño-Avena, C. (2022) Discovery and dynamics of a cryptic marine copepod-parasite interaction. *Marine Ecology Progress Series*, **691**, 29-40.
- Pearl, J. (2009) *Causality*. Cambridge university press.

- Pinnegar, J.K., Blanchard, J.L., Mackinson, S., Scott, R.D. & Duplisea, D.E. (2005) Aggregation and removal of weak-links in food-web models: system stability and recovery from disturbance. *Ecological Modelling*, **184**, 229-248.
- Pomati, F., Shurin, J.B., Andersen, K.H., Tellenbach, C. & Barton, A.D. (2020) Interacting temperature, nutrients and zooplankton grazing control phytoplankton size-abundance relationships in eight Swiss lakes. *Frontiers in microbiology*, **10**, 3155.
- Rasmussen, C., Dupont, Y.L., Mosbacher, J.B., Trøjelsgaard, K. & Olesen, J.M. (2013) Strong impact of temporal resolution on the structure of an ecological network. *PLoS one*, **8**, e81694.
- Ringstad, N. & Horvitz, H.R. (2008) FMRFamide neuropeptides and acetylcholine synergistically inhibit egg-laying by *C. elegans*. *Nature neuroscience*, **11**, 1168.
- Runge, J., Nowack, P., Kretschmer, M., Flaxman, S. & Sejdinovic, D. (2019) Detecting and quantifying causal associations in large nonlinear time series datasets. *Science advances*, **5**, eaau4996.
- Saberski, E., Bock, A.K., Goodridge, R., Agarwal, V., Lorimer, T., Rifkin, S.A. & Sugihara, G. (2021) Networks of Causal Linkage Between Eigenmodes Characterize Behavioral Dynamics of *Caenorhabditis elegans*. *PLoS computational biology*, **17**, e1009329.
- Saberski, E., Park, J., Hill, T., Stabenau, E. & Sugihara, G. (2022) Improved Prediction of Managed Water Flow into Everglades National Park Using Empirical Dynamic Modeling. *Journal of Water Resources Planning and Management*, **148**, 05022009.
- Scheffer, M., Bascompte, J., Brock, W.A., Brovkin, V., Carpenter, S.R., Dakos, V., Held, H., Van Nes, E.H., Rietkerk, M. & Sugihara, G. (2009) Early-warning signals for critical transitions. *Nature*, **461**, 53-59.
- Segundo, J., Sugihara, G., Dixon, P., Stiber, M. & Bersier, L.-F. (1998) The spike trains of inhibited pacemaker neurons seen through the magnifying glass of nonlinear analyses. *Neuroscience*, **87**, 741-766.
- Spirtes, P., Glymour, C.N. & Scheines, R. (2000) *Causation, prediction, and search*. MIT press.
- Stephens, G.J., Johnson-Kerner, B., Bialek, W. & Ryu, W.S. (2008) Dimensionality and dynamics in the behavior of *C. elegans*. *PLoS Comput Biol*, **4**, e1000028.

- Sugihara, G. (1983) *Niche hierarchy: structure, organization and assembly in natural communities*. Princeton University.
- Sugihara, G. (1994) Nonlinear forecasting for the classification of natural time series. *Philosophical Transactions of the Royal Society of London. Series A: Physical and Engineering Sciences*, **348**, 477-495.
- Sugihara, G., Allan, W., Sobel, D. & Allan, K.D. (1996) Nonlinear control of heart rate variability in human infants. *Proceedings of the National Academy of Sciences*, **93**, 2608-2613.
- Sugihara, G., Bersier, L.-F. & Schoenly, K. (1997) Effects of taxonomic and trophic aggregation on food web properties. *Oecologia*, **112**, 272-284.
- Sugihara, G., Casdagli, M., Habjan, E., Hess, D., Dixon, P. & Holland, G. (1999) Residual delay maps unveil global patterns of atmospheric nonlinearity and produce improved local forecasts. *Proceedings of the National Academy of Sciences*, **96**, 14210-14215.
- Sugihara, G., Deyle, E.R. & Ye, H. (2017) Reply to Baskerville and Cobey: Misconceptions about causation with synchrony and seasonal drivers. *Proceedings of the National Academy of Sciences*, **114**, E2272-E2274.
- Sugihara, G., Garcia, S., Platt, T., Gulland, J., Rachor, E., Lawton, J., Rothschild, B., Maske, H., Ursin, E. & Paine, R. (1984) Ecosystems Dynamics: Group Report. *Exploitation of Marine Communities: Report of the Dahlem Workshop on Exploitation of Marine Communities Berlin 1984, April 1-6*, pp. 131-153. Springer.
- Sugihara, G., May, R., Ye, H., Hsieh, C.H., Deyle, E., Fogarty, M. & Munch, S. (2012) Detecting Causality in Complex Ecosystems. *Science*, **338**, 496-500.
- Sugihara, G. & May, R.M. (1990) Nonlinear forecasting as a way of distinguishing chaos from measurement error in time series. *Nature*, **344**, 734-741.
- Sugihara, G., Schoenly, K. & Trombla, A. (1989) Scale invariance in food web properties. *Science*, **245**, 48-52.
- Takens, F. (1981) Detecting strange attractors in turbulence. *Dynamical systems and turbulence, Warwick 1980*, pp. 366-381. Springer.

- Trent, C., Tsung, N. & Horvitz, H.R. (1983) Egg-laying defective mutants of the nematode *Caenorhabditis elegans*. *Genetics*, **104**, 619-647.
- Ushio, M., Hsieh, C.-h., Masuda, R., Deyle, E.R., Ye, H., Chang, C.-W., Sugihara, G. & Kondoh, M. (2018) Fluctuating interaction network and time-varying stability of a natural fish community. *Nature*, **554**, 360-363.
- Van de Ville, D., Britz, J. & Michel, C.M. (2010) EEG microstate sequences in healthy humans at rest reveal scale-free dynamics. *Proceedings of the National Academy of Sciences*, **107**, 18179-18184.
- Van Nes, E.H., Scheffer, M., Brovkin, V., Lenton, T.M., Ye, H., Deyle, E. & Sugihara, G. (2015) Causal feedbacks in climate change. *Nature Climate Change*, **5**, 445-448.
- WÆRVÅGEN, S.B. & Nilssen, J.P. (2010) Life histories and seasonal dynamics of common boreal pelagic copepods (Crustacea, Copepoda) inhabiting an oligotrophic Fennoscandian lake. *Journal of Limnology*, **69**, 311.
- Waterston, R. (1998) Genome sequence of the nematode *C. elegans*: a platform for investigating biology. The *C. elegans* Sequencing Consortium. *Science*, **282**, 2012-2018.
- White, J.G., Southgate, E., Thomson, J.N. & Brenner, S. (1986) The structure of the nervous system of the nematode *Caenorhabditis elegans*. *Philos Trans R Soc Lond B Biol Sci*, **314**, 1-340.
- Yan, G., Vértés PE, Towlson EK, Chew YL, Walker DS, Schafer WR, Barabási AL (2017) Network control principles predict neuron function in the *Caenorhabditis elegans* connectome. *Nature*, **550**, 519-523.
- Yang, A.C., Peng, C.-K. & Huang, N.E. (2018) Causal decomposition in the mutual causation system. *Nature communications*, **9**, 1-10.
- Ye, H., Deyle, E.R., Gilarranz, L.J. & Sugihara, G. (2015) Distinguishing time-delayed causal interactions using convergent cross mapping. *Scientific Reports*, **5**.
- Ye, H. & Sugihara, G. (2016) Information leverage in interconnected ecosystems: Overcoming the curse of dimensionality. *Science*, **353**, 922-925.