



Published in final edited form as:

Electron Struct. 2022 December ; 4(4): . doi:10.1088/2516-1075/acad51.

Convergence in determining enzyme functional descriptors across Kemp eliminase variants

Yaoyukun Jiang¹, Sebastian L Stull¹, Qianzhen Shao¹, Zhongyue J Yang^{1,2,3,4,5,*}

¹Department of Chemistry, Vanderbilt University, Nashville, TN 37235, United States of America

²Center for Structural Biology, Vanderbilt University, Nashville, TN 37235, United States of America

³Vanderbilt Institute of Chemical Biology, Vanderbilt University, Nashville, TN 37235, United States of America

⁴Data Science Institute, Vanderbilt University, Nashville, TN 37235, United States of America

⁵Department of Chemical and Biomolecular Engineering, Vanderbilt University, Nashville, TN 37235, United States of America

Abstract

Molecular simulations have been extensively employed to accelerate biocatalytic discoveries. Enzyme functional descriptors derived from molecular simulations have been leveraged to guide the search for beneficial enzyme mutants. However, the ideal active-site region size for computing the descriptors over multiple enzyme variants remains untested. Here, we conducted convergence tests for dynamics-derived and electrostatic descriptors on 18 Kemp eliminase variants across six active-site regions with various boundary distances to the substrate. The tested descriptors include the root-mean-square deviation of the active-site region, the solvent accessible surface area ratio between the substrate and active site, and the projection of the electric field (EF) on the breaking C–H bond. All descriptors were evaluated using molecular mechanics methods. To understand the effects of electronic structure, the EF was also evaluated using quantum mechanics/molecular mechanics methods. The descriptor values were computed for 18 Kemp eliminase variants. Spearman correlation matrices were used to determine the region size condition under which further expansion of the region boundary does not substantially change the ranking of descriptor values. We observed that protein dynamics-derived descriptors, including $\text{RMSD}_{\text{active_site}}$ and $\text{SASA}_{\text{ratio}}$, converge at a distance cutoff of 5 Å from the substrate. The electrostatic descriptor, $\text{EF}_{\text{C-H}}$, converges at 6 Å using molecular mechanics methods with truncated enzyme models and 4 Å using quantum mechanics/molecular mechanics methods with whole enzyme model. This study serves as a future reference to determine descriptors for predictive modeling of enzyme engineering.

* Author to whom any correspondence should be addressed. zhongyue.yang@vanderbilt.edu.

Conflict of interest

The authors declare no competing financial interest.

Keywords

convergence test; enzyme kinetics; mutation effect; Kemp eliminase

1. Introduction

Enzymes have been widely used as biocatalysts for chemical synthesis [1–3], biomass conversion [4–7], polymer upcycling [8–11], drug functionalization [12–15], and food allergy treatment [16–18]. Wild-type enzymes usually exhibit low specificity for converting non-native substrate and feeble activity for catalyzing new-to-nature reactions. Experimental strategies of enzyme engineering, such as random mutagenesis [19–21], gene shuffling/recombination [22, 23], CASTing [24, 25], and directed evolution [26–29], have been leveraged to optimize enzymes' capability for accommodating non-native substrates or catalyzing new-to-nature reactions. These strategies require extensive efforts for screening and selecting mutants to achieve desired functions. To accelerate biocatalytic discovery, molecular simulations [30–34] have been augmented with the campaign of biocatalytic discovery. The catalytic actions of enzyme catalysis can be elucidated and quantified using descriptors, including folding stability [35], binding affinity [36–38], activation barriers [39, 40], protein dynamics and correlated motions [41–49], electric field (EF) [50–55], charge transfer [54, 56], and more. These descriptors, derived from quantum mechanical (QM) or molecular mechanical (MM) simulations, have guided the search for beneficial mutants [57, 58]. They also serve as critical features for data-driven enzyme engineering.

For example, protein dynamics-derived descriptors and EFs have been extensively studied, because they were found to correlate with enzyme catalytic efficiency [50, 58–60]. Additionally, their computation is more efficient than that of activation barriers, whose convergence requires intensive conformational sampling and electronic structure calculations. A common descriptor for protein dynamics is the root-mean-square deviation of the active-site region ($\text{RMSD}_{\text{active_site}}$). $\text{RMSD}_{\text{active_site}}$ quantifies the structural fluctuation of protein backbones or sidechains relative to a reference structure. The fluctuation is associated with the B-factor of protein structure determined from crystallography. In an analysis of catalytic residues in 178 enzyme active sites [59], Bartlett *et al* showed that the active-site residues of efficient enzymes generally have a lower B-factor. As such, a lower $\text{RMSD}_{\text{active_site}}$ should be expected for efficient mutant enzymes and designer enzymes in catalysis, albeit the catalytic efficiency may drop under very low $\text{RMSD}_{\text{active_site}}$ range [61]. Besides $\text{RMSD}_{\text{active_site}}$, our group identified a new descriptor to evaluate the overall impact of protein dynamics on substrate positioning [58]. The descriptor, defined as solvent accessible surface area ratio of substrate to active-site residues ($\text{SASA}_{\text{ratio}}$), can be obtained from molecular dynamics (MD) simulations. Using lactonase as a model system, our previous work shows that $\text{SASA}_{\text{ratio}}$ can guide the search of optimal enzyme mutants with enhanced specificity for non-native substrates. Besides protein dynamics, the role of electrostatic environments was reported as a critical factor in mediating enzyme catalysis [62]. Linear correlation was observed between the magnitude of EF in the reaction center and the free energy barrier in ketosteroid isomerase and serine protease [50].

Despite the broad applications of simulation-derived descriptors in guiding enzyme engineering, converging the computation of descriptors in QM and MM simulations is a non-trivial task. Failure of achieving convergence hampers reproducibility of computational outcomes and may misguide experimental designs. This issue is particularly significant for QM-based calculations due to their high computational cost. Benchmarks have been performed to investigate the selection of QM regions that converge the computation of electronic structure descriptors (e.g. partial charge [63–67], charge transfer [67], charge density [66], bond valence [67], and electrostatic potential), energetic properties (e.g. energy barrier, reaction energy, and free energy) [64, 65, 67–75], geometries [64, 73, 74], and nuclear magnetic resonance (NMR) shielding [76, 77] in various model enzymes (peroxidase, methyltransferases, cytochrome P450, and deacetylase) [67]. Rational QM region selection approaches have also been developed, including charge shift analysis [78], Fukui shift analysis [78], and point charge variation analysis [79].

The benchmark studies on descriptors have been mostly performed on wild-type enzymes [72]. However, to understand or predict mutation effects, it is essential to perform convergence tests over multiple enzyme variants. Ideally, the selected active-site region for computing QM or MM properties should be large enough so that further expanding the region size does not substantially change the order of descriptor values across different enzyme variants. In this work, using 18 variants of Kemp eliminase [80], we investigated whether the ranking of descriptor values across enzyme variants approaches convergence as the increase of active-site region sizes used in descriptor computation. We first sampled conformational ensembles for 18 variants using classical molecular dynamics. Based on the sampled conformers, protein dynamics-derived descriptors (i.e. $\text{RMSD}_{\text{active_site}}$ and $\text{SASA}_{\text{ratio}}$) and electronic structure-derived descriptors (i.e. EF along the breaking C–H bond) were evaluated using different sizes of active-site region based on MM or QM methods. For each descriptor, the Spearman correlation matrix was computed to examine the trend of convergence. The study informs the conditions under which different descriptors can be calculated with high fidelity for predicting the impact of mutations on catalytic functions. In addition, as the interplay between protein dynamics and electronic structures emerges as a new direction of study [81, 82], the convergence trend investigated in the current study might inspire the development of new strategies to predict computationally demanding QM properties using MM-derived properties.

2. Computational methods

Protein Structure and Preparation

The crystal structure of KE07-R7-2 was obtained from the Protein Data Bank (PDB ID: 5D38) [56]. All the crystallizing water molecules were removed. To make the amino acid sequence consistent with the original KE07 design [30], the N-terminal alanine was changed to methionine and the residues following Leu253 on the C-terminal were removed. The crystal structure [30] of KE07 in complex with the substrate 5-nitrobenzisoxazole was aligned relative to the KE07-R7-2 crystal structure using PyMol [83]. The coordinates of the substrate were used to construct the KE07-R7-2-substrate complex. The complex was then prepared with the AMBER 18 *tleap* [84] utility for MD simulations. AMBER ff14SB

force field was used for the protein [85]. Parameters for the substrate were obtained using the generalized AMBER force field [86, 87]. The substrate structure was downloaded from PDB under the H5J entry (5-nitro-1,2-benzoxazole). The atomic charges were determined by the AM1-BCC model [88]. The missing atoms were also complemented with *tLeap*.

Molecular Dynamics Simulations

MD simulations for each of the 18 variant-substrate complexes were conducted with a high throughput enzyme modeling platform, EnzyHTP [89]. The 18 variants include one KE07-R7-2 as the ‘wild-type’ and 17 of its mutants, including S48N, H201A, H201K, K222A, R16Q, N25S, I52A, M62A, H84Y, K132N, I199S, I199F, I199A, K132M, K162A, L170A, E185A (supporting information, table S1 and .zip file). Specifically, EnzyHTP automatically generates the structures of enzyme mutants based on the original structure and performs MD simulations using AMBER 18 [84]. The SHAKE algorithm was applied to constrain all the hydrogen-containing bonds [90]. To sample the near transition state conformations throughout the simulations, geometric restraints between the substrate and key amino acid residues were applied from minimization to production runs (supporting information, figure S1). The enzyme complexes were then solvated in a periodic octahedron box with a 10 Å buffer of TIP3P water and were neutralized with Na⁺ counterions. For each variant complex, the whole solvent box was first relaxed using steepest descent method for 10 000 steps followed by conjugate gradient method for another 10 000 steps. After minimization, each box was heated from 0 to 293.15 K within 36 ps with constant volume, equilibrated for 4 ps under constant volume at 293.15 K, and further equilibrated at 293.15 K and 1 atm for 1 ns. In addition to the geometric restraints mentioned above, the backbone C_α, C and N of the amide group were also restrained with a 2 kcal mol⁻¹ Å⁻² weight from the minimization to equilibration. After equilibration, we carried out production runs for 110 ns and output the trajectories every 100 ps. The snapshots derived from the last 100 ns of the production run were used for analyses. This yields a total of 1000 snapshots for each production run. All simulations were performed with a time step of 2 fs. The Langevin thermostat [91] and Berendsen barostat [92] were used throughout the simulations. For each of the 18 variant-substrate complexes, five parallel MD runs were conducted with different random seeds, yielding a total sampling time of 500 ns and 5000 snapshots.

QM/MM Calculations

We conducted QM/MM single-point electronic structure calculations for 500 snapshots sampled from MD production runs with a 1 ns interval. Each snapshot resulted from MD sampling was converted to an image in which the enzyme-substrate complex occupies the center of the box using *autoimage* function of AMBER *cpptraj* utility [93]. QM/MM single-point energies were calculated using TeraChem [94, 95]. The electrostatic interactions between the QM and MM region were treated with the electrostatic embedding method [62]. The QM/MM boundaries cut the backbone C–N bond of the amide group. To cap the unbonded atoms in the QM region, explicit H atoms were placed along the bond vector connecting the QM and MM atoms, and the resulting N–H and C–H bond lengths were set to be 1.09 Å. At the same time, the point charges originally belonging to the QM-region-bonded amide C and N atoms in the MM region were removed, and their charges were redistributed evenly on the remaining MM atoms except for those covalently

bonded to the deleted MM amide C and N atoms. The electronic structures were described using the range-separated exchange-correlation functional ω PBEh [96] ($\omega = 0.2 \text{ bohr}^{-1}$) with 6–31 G(d) [97]. This combination of method and basis set has been validated in the study of large-scale electronic structure effects in catechol *O*-methyltransferase, cytochrome P450cam, lysozyme, and DNA methyltransferase [65, 67]. The restrained electrostatic potential (RESP) point charges [98] of each snapshot were calculated for QM residue EF analyses.

Descriptor Calculations and Analyses

We selected six active-site regions whose boundary's distance to the substrate surface ranges from 3 to 8 Å with a 1 Å interval. For all 18 variants, the active-site regions were classified based on the averaged MD structure of KE07-R7-2. A residue is selected in the region if any one of its heavy atoms is within the distance cutoff from its nearest substrate heavy atom. Based on each of the active-site region, we calculated the enzyme functional descriptors, including mass weighted root-mean-square deviation of an active-site region ($\text{RMSD}_{\text{active_site}}$, in Å), solvent accessible surface area ratio between substrate and active-site residues ($\text{SASA}_{\text{ratio}}$, in Å²), and EF along the breaking C–H bond ($\text{EF}_{\text{C-H}}$, in MV cm⁻¹). The values of each descriptor were first evaluated on individual conformational snapshots, and then averaged over sampled classical MD or QM/MM snapshots (supporting information, tables S2–S5).

For $\text{RMSD}_{\text{active_site}}$, we included all the heavy atoms of the amino acid residues. The reference structure was averaged from sampled MD snapshots. The $\text{SASA}_{\text{ratio}}$ was calculated based on the ratio of SASA_{sub} (substrate's SASA) to $\text{SASA}_{\text{protein}}$ (protein residues' SASA). SASA was quantified using the Shrake and Rupley algorithm [99] embedded in the python library MDTraj [100]. The probe radius was 1.4 Å and the surface of each atom was represented by 5000 grid points. $\text{EF}_{\text{C-H}}$ was calculated to be the projected EF strength at the middle point of the breaking C–H bond of the substrate 5-nitrobenzisoaxazole. The bond vector direction points from C to H. We separately computed $\text{EF}_{\text{C-H}}$ based on RESP charges either derived from molecular mechanics force field or single-point electronic structure calculation. For MM-derived $\text{EF}_{\text{C-H}}$, the $\text{EF}_{\text{C-H}}$ was summed over from all atoms in the selected active-site region based on the RESP charges used in the classical force field. For QM/MM-derived $\text{EF}_{\text{C-H}}$, the $\text{EF}_{\text{C-H}}$ was summed over from all atoms in the QM and MM region. The EF contributions from the capping H atoms were not included.

For each active-site region, the averaged descriptor values were computed and then ranked across 18 enzyme variants from 1 to 18. The Spearman correlation coefficient ρ was calculated as $1 - \frac{6 \sum d_i^2}{n(n^2 - 1)}$ where d_i is the rank difference for the i th enzyme variant and $n = 18$ is the total number of enzyme variants. Spearman correlation matrix for each descriptor was computed that contains the correlation coefficients for each pair of the active-site regions.

3. Results and discussion

3.1. Kemp eliminase variants as the model system

As the first known *de novo*-designed enzyme, Kemp eliminase catalyzes the conversion of benzisoxazole to cyanophenol via C–H deprotonation followed by ring opening (figure 1 top right) [30]. Multiple generations of Kemp eliminase have been reported [30, 42, 51, 52, 56, 80, 101–103]. Some well-known examples include the KE family designed using the ‘inside-out’ protocol by Baker, Houk, Tawfik, and co-workers [30], the HG family using iterative protocol by Hilvert, Houk, Mayo and co-workers [31], and the AlleyCat family using the minimalist approach by Korendovych, Degrado, and coworkers [32, 103]. From their initial reports, the most efficient enzyme variants were identified to be KE07-R7-2 ($k_{\text{cat}}/K_{\text{M}}$ 2590 M⁻¹ s⁻¹), HG-3 ($k_{\text{cat}}/K_{\text{M}}$ 430 M⁻¹ s⁻¹), and AlleyCat (i.e. $k_{\text{cat}}/K_{\text{M}}$ = 128.4 M⁻¹ s⁻¹). All three families of Kemp eliminase involve a general acid-base mechanism, in which the substrate 5-nitrobenzisoxazole is deprotonated by a nearby carboxylate (side chain of Glu or Asp) to form 2-hydroxy-5-nitrobenzonitrile via one single transition state (figure 1, *top-right*). Nonetheless, they involve a different set of active site residues for substrate deprotonation and binding.

We chose the model system to be a member of the KE family, KE07-R7-2 [30, 56], and 17 of its variants with single amino acid substitution reported by Head-Gordon and coworkers [80]. KE07-R7-2 was derived from seven rounds of directed evolution based on a computationally designed enzyme scaffold KE07. In KE07-R7-2 and its variants, the carboxylic sidechain of Glu101 serves as the catalytic base (figure 1, *bottom-right*). These variants were selected in the benchmark for three reasons. First, the mutational spots of the variants span over a wide range of spatial proximity to the substrate (i.e. 3–23 Å, figure 1, *left* and supporting information, table S1). Both close and distal mutations are thus considered in the study. Second, although modeling has been performed for KE07-R7-2 to infer mutational hotspots based on correlated residue motion [80], protein dynamics and electronic structures for the 17 variants of KE07-R7-2 have not been investigated. Third, the kinetic parameters (i.e. k_{cat} or K_{M}) for these variants are known experimentally [80]. This implies that the mutation does not abolish the structural and catalytic integrity of Kemp eliminase. The crystal structure for KE07-R7-2 can be used as a scaffold for mimicking the mutant structures. Notably, although Kemp eliminase is the only model enzyme used in the current study, the general acid-base mechanism of Kemp eliminase represents a general mechanistic scheme shared by hydrolases, isomerases, and many other types of enzymes. As such, the convergence trend identified here can be potentially applied to other cases.

For KE07-R7-2 and its variants, the greater enzyme active site region entails 32 residues, including 6 polar, 20 non-polar, and 6 charged residues (figure 1, *bottom-right* and supporting information, table S6). By design, Glu101, Lys222 and Trp50 directly participate in the reaction or stabilize the transition state [30]. Glu101 is the general base that deprotonates the substrate. Lys222 is the H-bond donor to stabilize the phenoxide intermediate. Trp50 is the π -stacking residue to stabilize the substrate binding and charge-separated transition state. Four polar residues are observed within 5 Å of the substrate, including Tyr128, Ser48, His201, and Arg202. They likely stabilize the substrate binding or

transition state with electrostatic or polar interactions. In addition, a total of eight polar (i.e. Glu101, Lys222, Tyr128, Ser48, His201, Arg202, Asp224, Asn103) and four charged (i.e. Glu101, Lys222, Arg202, Asp224) residues are found within 5.5 Å from the substrate. These residues mediate the EF environment exerted on the breaking C-H bond. Besides dispersion interactions, the nonpolar residues likely contribute to the active site dynamics as described by $\text{RMSD}_{\text{active_site}}$ and $\text{SASA}_{\text{ratio}}$.

3.2. Region selection for descriptor calculation

To calculate simulation-derived descriptors, an active-site region should be defined first. In this study, the calculations of $\text{RMSD}_{\text{active_site}}$, $\text{SASA}_{\text{ratio}}$, and MM-derived $\text{EF}_{\text{C-H}}$ involve only the residues classified within a defined active-site region. The calculation of QM/MM-derived $\text{EF}_{\text{C-H}}$ involves treatment of the active-site region residues using quantum mechanics and the rest of the enzyme residues using molecular mechanics.

To benchmark the region size effect, the active-site regions were defined based on the residues' spatial proximity to the substrate (see section 2, Descriptor Calculations and Analyses). We selected six active-site regions whose boundaries to the substrate range from 3 to 8 Å (with 1 Å interval)—they are named C3–C8, respectively (figure 2). The residues were consistently selected by referencing KE07-R7-2. For C3, only Glu101 is included. Glu101 serves as the catalytic base to deprotonate the residue. Notably, throughout the MD simulations, a distance constraint was applied between Glu101 and the substrate to maintain their favorable catalytic pose. Compared to C3, C4 involves an expansion of seven additional residues. Among these residues, Lys222 and Trp50 appear in the original design of theozyme [30]. These two residues, cooperating with His201, Tyr128, and Ser48, likely facilitate proton transfer needed for the general acid-base mechanism. Unlike C3 which bears a -1 charge, C4 is charge neutral due to the addition of Lys222. In C5, only one additional residue Arg202 is included. This indicates that the catalytic core of KE involves a relatively compact inner cluster of residues surrounding the substrate. The positive charge introduced by Arg202 in C5 is neutralized by Asp224 in the C6 region. Notably, among the newly added residues in C6, three (i.e. Leu10, Phe49, and Val169) out of five are non-polar. This trend is also observed in C7 and C8. For the new additions, only two residues (i.e. Ser144 and Thr78) out of eight are polar in C7; two (i.e. Asp7 and Asp51) out of ten are polar in C8. The excessive number of non-polar residues in the greater active-site region contribute to the stability of Kemp eliminase.

The total number of atoms ranges from 31 in C3 to 495 in C8 (table 1). The region size tested here is comparable to or greater than the optimal region sizes determined from previous benchmark studies, including DNA methyltransferase by Solt *et al* [68]. (300 atoms), histone deacetylase by Morgenstern *et al* [66]. (200 atoms), and catechol O-methyltransferase by Kulik *et al* [64]. (500–600 atoms) and Jindal and Warshel [72]. (60 atoms). Notably, the size effect mentioned above was determined to be optimal for QM-derived quantities, including charge transfer, electron density, reaction enthalpy, and so on. In contrast, the size effects investigated in this study emphasize both MM-derived and QM/MM-derived features.

For each of the six active-site regions, we computed the average descriptor values for the 18 KE07-R7-2 variants based on their conformational ensembles (supporting information, tables S3–S6). We investigated how the ranking of descriptor values across the 18 variants varies with the increase of region size. Instead of benchmarking a certain molecular property against its reference value, this study intends to identify a condition of region size under which further expanding the region boundary minimally changes the ranking of descriptor values across enzyme variants. Notably, the region size condition for a converged trend does not guarantee the convergence of individual property values. Nonetheless, the mutation effect can be reasonably inferred under this condition to guide enzyme engineering.

3.3. Descriptor of protein dynamics: $\text{RMSD}_{\text{active_site}}$ and $\text{SASA}_{\text{ratio}}$

We first investigated the dynamics-derived descriptors, $\text{RMSD}_{\text{active_site}}$ and $\text{SASA}_{\text{ratio}}$. They represent different aspects of protein dynamics. $\text{RMSD}_{\text{active_site}}$ informs the conformational fluctuation of active site residues, while $\text{SASA}_{\text{ratio}}$ informs the dynamic positioning and fitness of substrate in the active site.

Figure 3 shows the Spearman correlation matrix for $\text{RMSD}_{\text{active_site}}$ (*left*) and $\text{SASA}_{\text{ratio}}$ (*right*). Each element of the matrix represents a Spearman correlation coefficient (i.e. ρ) between descriptor values derived from two regions with a distinct size. For $\text{RMSD}_{\text{active_site}}$, a high ρ value (i.e. ≥ 0.70) is observed for almost all pairs of regions except those that involve C3. The moderate ρ values between C3 and C5–C8 (i.e. 0.4–0.6) are caused by the small size of C3 that involves only one residue in the region (i.e. Glu101). The correlation coefficients tend to be higher for regions that are close in size (e.g. $\rho > 0.9$ for C5–C6, C6–C7, and C7–C8) and lower for regions with a larger size gap (e.g. $\rho = 0.40$ for C3–C8 and 0.73 for C4–C8). Notably, it is unexpected that the correlation coefficient is still as high as 0.40 between C3 and C8 because their numbers of residues differ by 31 and of atoms by 464. This indicates that the $\text{RMSD}_{\text{active_site}}$ ranking calculated from C3 can still partially inform the ranking of dynamic fluctuation exhibited by larger-sized regions. This is likely caused by the collective motions of residues in the enzyme active site, where all residues are somewhat interconnected in a complex, dynamic network.

Unlike $\text{RMSD}_{\text{active_site}}$, which emphasizes protein dynamics, $\text{SASA}_{\text{ratio}}$ represents the interplay of dynamic motion between substrate and its surrounding active-site residues. The Spearman correlation matrix of $\text{SASA}_{\text{ratio}}$ shows a similar trend to that of $\text{RMSD}_{\text{active_site}}$ (figure 3, *right*). For each pair of regions, the Spearman correlation coefficient of $\text{SASA}_{\text{ratio}}$ is generally greater than that of $\text{RMSD}_{\text{active_site}}$. The ρ values are greater for correlations between larger regions that are closer in size (e.g. $\rho = 0.96$, 0.96, and 0.93 for C5–C6, C6–C7, and C7–C8, respectively). Notably, the $\text{SASA}_{\text{ratio}}$ is computed by the SASA ratio of substrate to active-site residues. For different active-site regions, the SASA value of the substrate always remains constant. This helps dampen the perturbation of expanding region size on the ranking of descriptor values across variants.

To determine a convergence cutoff for computing dynamics-derived descriptors, we investigated the change of Spearman correlation coefficients between adjacent active-site regions versus the increase of region size (figure 4). For both $\text{RMSD}_{\text{active_site}}$ and $\text{SASA}_{\text{ratio}}$, the correlation coefficient appears greater than 0.90 after C5 (i.e. 5.0 Å from the substrate).

With a Spearman ρ value greater than 0.90, the ranking of descriptor values computed from one active-site region is largely preserved in another. As such, the convergence cutoff for dynamics-derived descriptors is determined to be 5.0 Å. Notably, from C5 to C8, the atomic charge varies from 1 (C5), to 0 (C6 and C7), then to -2 (i.e. C8). The correlation coefficients remain high even between regions of different charges. This observation confirms that the dynamics-derived descriptors used here are approximately independent from electrostatic effects—they are insensitive to electrostatic perturbation in the protein environment.

3.4. Descriptor of electrostatic environment: EF along breaking C–H bond

Next, we investigated the descriptor for enzyme electrostatics, EF_{C-H} , the EF along the breaking C–H bond. The interior EF in Kemp eliminase has been proposed as a factor to stabilize the developing dipole moment along the C–H bond [52]. Optimizing the EF through mutagenesis has also been demonstrated as an effective strategy to improve enzyme catalytic efficiency [52, 104].

Figure 5 shows the Spearman correlation matrix for EF_{C-H} that were separately computed using MM (*left*) and QM/MM (*right*) method. MM-derived EF_{C-H} involves only the local residues that are classified in the active-site region. This approach is similar to the distance cutoff method used in Rosetta score functions for computing electrostatic interactions [105]. QM/MM-derived EF_{C-H} employs QM to treat residues in the active-site region and MM for residues in the rest of the enzyme. This approach incorporates the effects of long-range electrostatics. Notably, to ensure a consistent comparison, the same set of MD-derived snapshots was used in the calculations for both MM- and QM/MM-derived EF_{C-H} values. In our test cases, the QM/MM optimization consistently increases the resulting EF_{C-H} values by a few tens of $MV\ cm^{-1}$ (supporting information, table S7).

Unlike dynamics-derived descriptors, low correlation coefficients are more frequently observed between active-site regions of different sizes, especially between regions with a larger size gap. For example, the Spearman ρ values for C3–C6 (i.e. differ by 13 residues), C3–C7 (i.e. differ by 21 residues), and C3–C8 (i.e. differ by 31 residues) are 0.07, 0.07, and 0.05, respectively, for MM-derived EF_{C-H} (figure 5, *left*). The low correlation strength indicates that the ranking of EF_{C-H} values derived from a smaller active-site region cannot be used to infer the ranking from a larger active-site region. Different from dynamics-derived descriptors, the EF depends more sensitively on the active-site regions used in the calculation. From C3 to larger active-site regions, individual residues added to the active site region, especially polar and charged residues, can significantly affect the representation of mutation effects on interior enzyme electrostatics.

Similar to dynamics-derived descriptors, the Spearman ρ values are greater for correlations of EF_{C-H} rankings between larger regions that are closer in size. The ρ values for C4–C5, C6–C7, and C7–C8 are 0.99, 0.99, and 0.98, respectively, for MM-derived EF_{C-H} (figure 5, *left*); and are 0.99, 0.89, and 0.94, respectively, for QM/MM-derived EF_{C-H} (figure 5, *right*). Interestingly, the rankings derived from C4 to C5 are highly consistent, albeit their difference in the total charge of active-site residues by -1. The charge difference is caused by the addition of Arg202 in C5. Despite having a +1 charge, Arg202 has a trivial influence on EF_{C-H} due to it being perpendicular to the breaking C–H bond vector. For C6–C7 and

C7–C8, the newly added residues are mostly nonpolar and are distant from the breaking C–H bond in the substrate (i.e. >6.3 Å). As EF strength is inversely proportional to the square of the distance, the impact of remote residues dies off quickly. As such, a consistent ranking of EF_{C-H} values is observed between regions beyond C6. For MM-derived EF_{C-H} , the Spearman ρ value for C5–C6 (i.e. 0.54) is significantly lower than that for C4–C5 or C6–C7 (figure 5, *left*). This is because the newly added charged residue in C6, Asp224, is positioned along the direction of the breaking C–H bond vector. As such, the impact of Asp224 on the ranking of EF_{C-H} values is substantial. Notably, the drop of Spearman ρ value for C5–C6 was also observed for Mulliken charge-based QM/MM-derived EF_{C-H} values (supporting information, table S8 and figure S2), albeit the correlation coefficients remain consistently high (i.e. >0.90) for larger QM regions.

By comparing the Spearman ρ values for C4–C5, C5–C6, and C7–C8, the results show that for both MM- and QM/MM-derived EF_{C-H} values, the ranking is dependent more on the spatial distribution of charged residues relative to the breaking C–H bond than on the total atomic charge in the active-site regions. This finding can potentially help rational identification of residues for tuning interior enzyme EFs for selective bond activation. To determine a convergence cutoff for computing electrostatic descriptors (i.e. for MM- and QM/MM-derived EF_{C-H}), we investigated the change of Spearman correlation coefficients between adjacent active-site regions versus the increase of region size (figure 6). Consistent with the dynamics-derived descriptors, we adopted a Spearman ρ value of 0.90 as the criterion for determining the convergence cutoff. For MM-derived EF_{C-H} with truncated enzyme models, the convergence occurs at 6 Å. For QM/MM-derived EF_{C-H} with whole enzyme models, the correlation coefficients are consistently high even at minimal QM region C3/C4 (figure 6, *right*). This indicates that MM charges are sufficiently accurate and can even replace QM-derived RESP charges in describing the ranking of EF values across different mutants (supporting information, figure S3). The results are likely caused by the fact that the reactant state of Kemp eliminase does not involve longer-range charge transfer with residues in the greater protein environment. The trend of correlation likely changes when the deprotonation transition state is bound to the active site.

However, we should note that the computed EF values are still dependent on the QM region selection. As such, we investigated the trend of convergence for absolute QM/MM-derived EF_{C-H} values under different QM region sizes (supporting information, figure S4). The distribution of EF values systematically drops as the QM region enlarges from C3 to C4 and becomes steady until the QM region size hits C7 where the distribution bumps up due to reorganization of active site residues (i.e. Ala9, Thr78, Gly80, and Asp224). Combining the region size benchmark of correlation matrix and absolute EF values, we determine C4 as the convergence cutoff. To examine the influence of MM region, we calculated the C4-based QM/MM-derived EF_{C-H} values under different MM region sizes (supporting information, table S9 and figure S5). The ρ values for C4–C5, C6–C7, C7–C8, and C8-whole enzyme are 0.93, 0.98, 0.99, and 0.95, respectively. This result shows that the ranking of EF values across mutants is not sensitive to MM region sizes.

To pinpoint the origin of the difference between MM- and QM/MM-derived EF_{C-H} values at the QM region of C4, we calculated the atomic RESP charges derived from both MM

and QM methods based on the KE07-R7-2 wild-type. The atoms with the top ten largest deviations are backbone atoms residing on the QM/MM boundary (supporting information, table S10). Although the QM/MM boundary effects do not appear to influence the ranking of EF values across mutants, larger QM region should be used when the actual values of physical properties are of interest to the research.

Considering the low computational cost of MM-derived EF_{C-H} values, we would recommend a hybrid approach for future practice of computational enzyme engineering. This hybrid approach involves using MM-derived EF_{C-H} values (of the whole enzyme) for pre-screening of a large number of mutants, followed by assessment of QM/MM-derived EF_{C-H} values to identify mutants for experimental tests. For reactions with stronger polarization and charge transfer, a larger QM region may be used, but the region size could potentially be reduced by using rational QM determination approaches such as charge shift analysis [78], Fukui shift analysis [78], and point charge variation analysis [79].

4. Conclusions

In this work, we investigated how large an active-site region should be to converge the description of mutation effects on enzyme dynamics and electrostatics. For 18 KE07-R7-2 variants, dynamics-derived descriptors ($RMSD_{active_site}$ and $SASA_{ratio}$, both derived from classical MD) and electrostatic descriptors (MM- and QM/MM-derived EF_{C-H} were computed across six active-site regions with various boundary distances (i.e. 3–8 Å) to the substrate. For each descriptor, we employed the Spearman correlation matrix to determine the region size condition under which further expansion of the region boundary does not substantially change the ranking of descriptor values.

Using a Spearman ρ value of 0.9 as a criterion for convergence, we observed that the ranking for $RMSD_{active_site}$ and $SASA_{ratio}$ converges at 5 Å; MM- and QM/MM-derived EF_{C-H} converge at 6.0 and 4.0 Å, respectively. The ranking of EF_{C-H} values derived from MM charges (i.e. including all atoms in the enzyme) is predictive to that from QM/MM charges, albeit the absolute EF values still exhibit dependence on the QM region sizes. As such, we recommend a hybrid approach for future practice of computational enzyme engineering, which involves a pre-screening of a large number of mutants based on MM-derived EF_{C-H} values, followed by an assessment of QM/MM-derived EF_{C-H} values on a smaller number of pre-screened mutants. Notably, the convergence of rankings does not ensure the convergence of measured descriptor values. Nonetheless, the ranking is most useful to guide experimental selection of function-enhancing enzyme mutants. Additionally, the current study emphasizes a designer enzyme, Kemp eliminase. Future studies should entail more types of enzymes with various catalytic actions.

Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

Acknowledgments

This research was supported by the startup grant from Vanderbilt University. Z J Yang, Y Jiang, and Q Shao are supported by the National Institute of General Medical Sciences of the National Institutes of Health under Award Number R35GM146982. This work was carried out in part using computational resources from the Extreme Science and Engineering Discovery Environment (XSEDE), which is supported by National Science Foundation Grant No. TG-BIO200057 [106].

Data availability statement

The data that support the findings of this study are openly available at the following URL/DOI: [10.5281/zenodo.7140835](https://doi.org/10.5281/zenodo.7140835).

References

- [1]. Koeller KM and Wong C-H 2001 Enzymes for chemical synthesis Nature 409 232–40 [PubMed: 11196651]
- [2]. Strohmeier GA, Pichler H, May O and Gruber-Khadjawi M 2011 Application of designed enzymes in organic synthesis Chem. Rev 111 4141–64 [PubMed: 21553913]
- [3]. Petchey MR and Grogan G 2019 Enzyme-catalysed synthesis of secondary and tertiary amides Adv. Synth. Catal 361 3895–914
- [4]. Chundawat SP, Beckham GT, Himmel ME and Dale BE 2011 Deconstruction of lignocellulosic biomass to fuels and chemicals Annu. Rev. Chem. Biomol. Eng 2 121–45 [PubMed: 22432613]
- [5]. Yang B, Dai Z, Ding S-Y and Wyman CE 2011 Enzymatic hydrolysis of cellulosic biomass Biofuels 2 421–49
- [6]. Sweeney MD and Xu F 2012 Biomass converting enzymes as industrial biocatalysts for fuels and chemicals: recent developments Catalysts 2 244–63
- [7]. Horn SJ, Vaaje-Kolstad G, Westereng B and Eijsink VG 2012 Novel enzymes for the degradation of cellulose Biotechnol. Biofuels 5 45 [PubMed: 22747961]
- [8]. Austin HP et al. 2018 Characterization and engineering of a plastic-degrading aromatic polyesterase Proc. Natl Acad. Sci. USA 115 E4350–7 [PubMed: 29666242]
- [9]. Knott BC et al. 2020 Characterization and engineering of a two-enzyme system for plastics depolymerization Proc. Natl Acad. Sci. USA 117 25476–85 [PubMed: 32989159]
- [10]. Tiso T et al. 2021 Towards bio-upcycling of polyethylene terephthalate Metab. Eng 66 167–78 [PubMed: 33865980]
- [11]. Ellis LD> Rorrer NA, Sullivan KP, Otto M, McGeehan JE> Román-Leshkov Y, Wierckx N and Beckham GT 2021 Chemical and biological catalysis for plastics recycling and upcycling Nat. Catal 4 539–56
- [12]. Fessner ND 2019 P450 monooxygenases enable rapid late-stage diversification of natural products via C–H bond activation ChemCatChem 11 2226–42 [PubMed: 31423290]
- [13]. Hong B, Luo T and Lei LX 2020 Late-stage diversification of natural products ACS Cent. Sci. 6 622–35 [PubMed: 32490181]
- [14]. Romero E, Jones BS, Hogg BN, Rué Casamajo A, Hayes MA, Flitsch SL, Turner NJ and Schnepel C 2021 Enzymatic late-stage modifications: better late than never Angew. Chem., Int. Ed 60 16824–55
- [15]. Craven EJ, Latham J, Shepherd SA, Khan I, Diaz-Rodriguez A, Greaney MF and Micklefield J 2021 Programmable late-stage C–H bond functionalization enabled by integration of enzymes with chemocatalysis Nat. Catal 4 385–94
- [16]. Gordon SR, Stanley EJ, Wolf S, Toland A, Wu SJ, Hadidi D, Mills JH, Baker D, Pultz IS and Siegel JB 2012 Computational design of an alpha-gliadin peptidase J. Am. Chem. Soc 134 20513–20 [PubMed: 23153249]
- [17]. Sun S, Jiang D, Fan M, Li H, Jin C and Liu W 2020 Selection of a versatile Lactobacillus plantarum for wine production and identification and preliminary characterisation of a novel histamine-degrading enzyme Int. J. Food Sci. Technol 55 2608–18

- [18]. Samadi N, Heiden D, Klems M, Salzmann M, Rohrhofer J, Weidmann E, Koidl L, Jensen-Jarolim E and Untersmayr E 2021 Gastric enzyme supplementation inhibits food allergy in a BALB/c mouse model *Nutrients* 13 738 [PubMed: 33652629]
- [19]. Bloom JD, Meyer M, Meinhold P, Otey C, Macmillan D and ARNOLD F 2005 Evolving strategies for enzyme engineering *Curr. Opin. Struct. Biol* 15 447–52 [PubMed: 16006119]
- [20]. Alejaldre L, Pelletier JN and Quaglia D 2021 Methods for enzyme library creation: which one will you choose? *BioEssays* 43 2100052
- [21]. Ravikumar Y, Nadarajan SP, Yoo TH, Lee C-S and Yun H 2015 Unnatural amino acid mutagenesis-based enzyme engineering *Trends Biotechnol.* 33 462–70 [PubMed: 26088007]
- [22]. Kolkman JA and Stemmer WPC 2001 Directed evolution of proteins by exon shuffling *Nat. Biotechnol* 19 423–8 [PubMed: 11329010]
- [23]. Akbulut N, Tuzlakoglu Ozturk M, Pijning T, Issever Ozturk S and Gumusel F 2013 Improved activity and thermostability of *Bacillus pumilus* lipase by directed evolution *J. Biotechnol* 164 123–9 [PubMed: 23313890]
- [24]. Reetz MT, Bocola M, Carballeira JD, Zha D and Vogel A 2005 Expanding the range of substrate acceptance of enzymes: combinatorial active-site saturation test *Angew. Chem., Int. Ed. Engl* 44 4192–6 [PubMed: 15929154]
- [25]. Reetz MT, Carballeira JD, Peyralans J, Höbenreich H, Maichele A and Vogel A 2006 Expanding the substrate scope of enzymes: combining mutations obtained by CASTing *Chemistry* 12 6031–8 [PubMed: 16789057]
- [26]. Arnold FH and Volkov AA 1999 Directed evolution of biocatalysts *Curr. Opin. Chem. Biol* 3 54–59 [PubMed: 10021399]
- [27]. Packer MS and Liu DR 2015 Methods for the directed evolution of proteins *Nat. Rev. Genet* 16 379–94 [PubMed: 26055155]
- [28]. Arnold FH 2018 Directed evolution: bringing new chemistry to life *Angew. Chem., Int. Ed. Engl* 57 4143–8 [PubMed: 29064156]
- [29]. Wang Y, Xue P, Cao M, Yu T, Lane ST and Zhao H 2021 Directed evolution: methodologies and applications *Chem. Rev* 121 12384–444 [PubMed: 34297541]
- [30]. Rothlisberger D et al. 2008 Kemp elimination catalysts by computational enzyme design *Nature* 453 190–5 [PubMed: 18354394]
- [31]. Privett HK, Kiss G, Lee TM, Blomberg R, Chica RA, Thomas LM, Hilvert D, Houk KN and Mayo SL 2012 Iterative approach to computational enzyme design *Proc. Natl Acad. Sci. USA* 109 3790–5 [PubMed: 22357762]
- [32]. Korendovych IV, Kulp DW, Wu Y, Cheng H, Roder H and DeGrado WF 2011 Design of a switchable eliminase *Proc. Natl Acad. Sci. USA* 108 6823–7 [PubMed: 21482808]
- [33]. Korendovych IV 2018 *Protein Engineering: Methods and Protocols* ed Bornscheuer UT and Höhne M (New York: Springer) pp 15–23
- [34]. Alonso-Cotchico L, Rodríguez-Guerra J, Lledos A and Marechal J-D 2020 Molecular modeling for artificial metalloenzyme design and optimization *Acc. Chem. Res* 53 896–905 [PubMed: 32233391]
- [35]. Khersonsky O et al. 2018 Automated design of efficient and functionally diverse enzyme repertoires *Mol. Cell* 72 178–86.e5 [PubMed: 30270109]
- [36]. Le P, Zhao J and Franzen S 2014 Correlation of heme binding affinity and enzyme kinetics of dehaloperoxidase *Biochemistry* 53 6863–77 [PubMed: 25330337]
- [37]. Luo Q, Chen D, Boom RM and Janssen AEM 2018 Revisiting the enzymatic kinetics of pepsin using isothermal titration calorimetry *Food Chem.* 268 94–100 [PubMed: 30064809]
- [38]. Richard JP 2019 Protein flexibility and stiffness enable efficient enzymatic catalysis *J. Am. Chem. Soc* 141 3320–31 [PubMed: 30703322]
- [39]. Amrein BA, Steffen-Munsberg F, Szeler I, Purg M, Kulkarni Y and Kamerlin SCL 2017 CADEE: computer-aided directed evolution of enzymes *IUCrJ* 4 50–64
- [40]. Yao J, Chen X, Zheng F and Zhan C-G 2018 Catalytic reaction mechanism for drug metabolism in human carboxylesterase-1: cocaine hydrolysis pathway *Mol. Pharm* 15 3871–80 [PubMed: 30095924]

- [41]. Hur S and Bruice TC 2003 The near attack conformation approach to the study of the chorismate to prephenate reaction Proc. Natl Acad. Sci. USA 100 12015–20 [PubMed: 14523243]
- [42]. Khersonsky O, Kiss G, RÖthlisberger D, Dym O, Albeck S, Houk KN, Baker D and Tawfik DS 2012 Bridging the gaps in design methodologies by evolutionary optimization of the stability and proficiency of designed Kemp eliminase KE59 Proc. Natl Acad. Sci. USA 109 10358–63 [PubMed: 22685214]
- [43]. Vaissier Welborn V and Head-Gordon T 2019 Computational design of synthetic enzymes Chem. Rev 119 6613–30 [PubMed: 30277066]
- [44]. Mehmood R, Vennelakanti V and Kulik HJ 2021 Spectroscopically guided simulations reveal distinct strategies for positioning substrates to achieve selectivity in nonheme Fe(II)/ α -ketoglutarate-dependent halogenases ACS Catal. 11 12394–408
- [45]. Liu CT, Francis K, Layfield JP, Huang X, Hammes-Schiffer S, Kohen A and Benkovic SJ 2014 Escherichia colidihydrofolate reductase catalyzed proton and hydride transfers: temporal order and the roles of Asp27 and Tyr100 Proc. Natl Acad. Sci. USA 111 18231–6 [PubMed: 25453098]
- [46]. Masterson JE and Schwartz SD 2015 Evolution alters the enzymatic reaction coordinate of dihydrofolate reductase J. Phys. Chem. B 119 989–96 [PubMed: 25369552]
- [47]. Liao Q, Kulkarni Y, Sengupta U, Petrovi D, Mulholland AJ, van der Kamp MW, Strodel B and Kamerlin SCL 2018 Loop motion in triosephosphate isomerase is not a simple open and shut case J. Am. Chem. Soc 140 15889–903 [PubMed: 30362343]
- [48]. Gao S, Thompson EJ, Barrow SL, Zhang W, Iavarone AT and Klinman JP 2020 Hydrogen–deuterium exchange within adenosine deaminase, a TIM barrel hydrolase, identifies networks for thermal activation of catalysis J. Am. Chem. Soc 142 19936–49 [PubMed: 33181018]
- [49]. Bunzel HA, Kries H, Marchetti L, Zeymer C, Mittl PRE, Mulholland AJ and Hilvert D 2019 Emergence of a negative activation heat capacity during evolution of a designed enzyme J. Am. Chem. Soc 141 11745–8 [PubMed: 31282667]
- [50]. Fried SD and Boxer SG 2017 Electric fields and enzyme catalysis Annu. Rev. Biochem 86 387–415 [PubMed: 28375745]
- [51]. Bhowmick A, Sharma SC and Head-Gordon T 2017 The importance of the scaffold for de novo enzymes: a case study with Kemp eliminase J. Am. Chem. Soc 139 5793–800 [PubMed: 28383910]
- [52]. Vaissier V, Sharma SC, Schaettle K, Zhang T and Head-Gordon T 2018 Computational optimization of electric fields for improving catalysis of a designed Kemp eliminase ACS Catal. 8 219–27
- [53]. Welborn VV and Head-Gordon T 2019 Fluctuations of electric fields in the active site of the enzyme ketosteroid isomerase J. Am. Chem. Soc 141 12487–92 [PubMed: 31368302]
- [54]. Yang Z, Liu F, Steeves AH and Kulik HJ 2019 Quantum mechanical description of electrostatics provides a unified picture of catalytic action across methyltransferases J. Phys. Chem. Lett 10 3779–87 [PubMed: 31244268]
- [55]. Bím D and Alexandrova AN 2021 Local electric fields as a natural switch of heme-iron protein reactivity ACS Catal. 11 6534–46 [PubMed: 34413991]
- [56]. Hong N-S et al. 2018 The evolution of multiple active site configurations in a designed enzyme Nat. Commun 9 3900 [PubMed: 30254369]
- [57]. Carlin DA et al. 2016 Kinetic characterization of 100 glycoside hydrolase mutants enables the discovery of structural features correlated with kinetic constants PLoS One 11 e0147596 [PubMed: 26815142]
- [58]. Jiang Y, Yan B, Chen Y, Juarez RJ and Yang ZJ 2022 Molecular dynamics-derived descriptor informs the impact of mutation on the catalytic turnover number in lactonase across substrates J. Phys. Chem. B 126 2486–95 [PubMed: 35324218]
- [59]. Bartlett GJ, Porter CT, Borkakoti N and Thornton JM 2002 Analysis of catalytic residues in enzyme active sites J. Mol. Biol 324 105–21 [PubMed: 12421562]
- [60]. Lodola A et al. 2010 Structural fluctuations in enzyme-catalyzed reactions: determinants of reactivity in fatty acid amide hydrolase from multivariate statistical analysis of quantum mechanics/molecular mechanics paths J. Chem. Theory Comput 6 2948–60 [PubMed: 26616091]

- [61]. Yabukarski F, Biel JT, Pinney MM, Doukov T, Powers AS, Fraser JS and Herschlag D 2020 Assessment of enzyme active site positioning and tests of catalytic mechanisms through x-ray-derived conformational ensembles Proc. Natl Acad. Set. USA 117 33204–15
- [62]. Warshel A and Levitt M 1976 Theoretical studies of enzymic reactions: dielectric, electrostatic and steric stabilization of the carbonium ion in the reaction of lysozyme J. Mol. Biol 103 227–49 [PubMed: 985660]
- [63]. Vanpoucke DE, Olah J, De Proft F, Van Speybroeck V and Roos G 2015 Convergence of atomic charges with the size of the enzymatic environment J. Chem. Inf. Model 55 564–71 [PubMed: 25668288]
- [64]. Kulik HJ, Zhang J, Klinman JP and Martínez TJ 2016 How large should the QM region be in QM/MM calculations? The case of catechol O-Methyltransferase J. Phys. Chem. B 120 11381–94 [PubMed: 27704827]
- [65]. Karelina M and Kulik HJ 2017 Systematic quantum mechanical region determination in QM/MM simulation J. Chem. Theory Comput 13 563–76 [PubMed: 28068092]
- [66]. Morgenstern A, Jaszai M, Eberhart ME and Alexandrova AN 2017 Quantified electrostatic preorganization in enzymes using the geometry of the electron charge density Chem. Sci 8 5010–8 [PubMed: 28970888]
- [67]. Mehmood R and Kulik HJ 2020 Both configuration and QM region size matter: zinc stability in QM/MM models of DNA methyltransferase J. Chem. Theory Comput 16 3121–34 [PubMed: 32243149]
- [68]. Solt I, Kulhánek P, Simon I, Winfield S, Payne MC, Csányi G and Fuxreiter M 2009 Evaluating boundary dependent errors in QM/MM simulations J. Phys. Chem. B 113 5728–35 [PubMed: 19341253]
- [69]. Sumowski CV and Ochsenfeld C 2009 A convergence study of QM/MM isomerization energies with the selected size of the QM region for peptidic systems J. Phys. Chem. A 113 11734–41 [PubMed: 19585981]
- [70]. Liao R-Z and Thiel W 2013 Convergence in the QM-only and QM/MM modeling of enzymatic reactions: a case study for acetylene hydratase J. Comput. Chem 34 2389–97 [PubMed: 23913757]
- [71]. Sadeghian K, Flaig D, Blank ID, Schneider S, Strasser R, Stathis D, Winnacker M, Carell T and Ochsenfeld C 2014 Ribose-protonated DNA base excision repair: a combined theoretical and experimental study Angew. Chem., Int. Ed. Engl 53 10044–8 [PubMed: 25065673]
- [72]. Jindal G and Warshel A 2016 Exploring the dependence of QM/MM calculations of enzyme catalysis on the size of the QM region J. Phys. Chem. B 120 9913–21 [PubMed: 27552257]
- [73]. Benediktsson B and Bjornsson R 2017 QM/MM study of the nitrogenase MoFe protein resting state: broken-symmetry states, protonation states, and QM region convergence in the FeMoco active site Inorg. Chem 56 13417–29 [PubMed: 29053260]
- [74]. Kang H and Zheng M 2021 Influence of the quantum mechanical region size in QM/MM modelling: a case study of fluoroacetate dehalogenase catalyzed CF bond cleavage Comput. Theor. Chem 1204 113399
- [75]. Demapan D, Kussmann J, Ochsenfeld C and Cui Q 2022 Factors that determine the variation of equilibrium and kinetic properties of QM/MM enzyme simulations: QM region, conformation, and boundary condition J. Chem. Theory Comput 18 2530–42 [PubMed: 35226489]
- [76]. Flaig D, Beer M and Ochsenfeld C 2012 Convergence of electronic structure with the size of the QM region: example of QM/MM NMR shieldings J. Chem. Theory Comput 8 2260–71 [PubMed: 26588959]
- [77]. Hartman JD, Neubauer TJ, Caulkins BG, Mueller LJ and Beran GJO 2015 Converging nuclear magnetic shielding calculations with respect to basis and system size in protein systems J. Biomol. NMR 62 327–40 [PubMed: 25993979]
- [78]. Qi HW, Karelina M and Kulik HJ 2018 Quantifying electronic effects in QM and QM/MM biomolecular modeling with the Fukui function Acta Phys.-Chim. Sin 34 81–91
- [79]. Brandt F and Jacob CR 2022 Systematic QM region construction in QM/MM calculations based on uncertainty quantification J. Chem. Theory Comput 18 2584–96 [PubMed: 35271768]

- [80]. Bhowmick A, Sharma SC, Honma H and Head-Gordon T 2016 The role of side chain entropy and mutual information for improving the de novo design of Kemp eliminases KE07 and KE70 *Phys. Chem. Chem. Phys.* 18 19386–96 [PubMed: 27374812]
- [81]. Yang Z, Hajlasz N, Steeves AH and Kulik HJ 2021 Quantifying the long-range coupling of electronic properties in proteins with ab initio molecular dynamics** *Chem.–Methods* 1 362–73
- [82]. Steeves AH and Kulik HJ 2022 Insights into the stability of engineered mini-proteins from their dynamic electronic properties *Electron. Struct.* 4 034005
- [83]. Schrödinger, LLC 2015 The PyMOL molecular graphics system, version 2.4
- [84]. Case DA et al. 2018 AMBER 2018 (San Francisco: University of California)
- [85]. Maier JA, Martinez C, Kasavajhala K, Wickstrom L, Hauser KE and Simmerling C 2015 ff14SB: improving the accuracy of protein side chain and backbone parameters from ff99SB *J. Chem. Theory Comput.* 11 3696–713 [PubMed: 26574453]
- [86]. Wang NX and Wilson AK 2004 The behavior of density functionals with respect to basis set. I. The correlation consistent basis sets *J. Chem. Phys.* 121 7632–46 [PubMed: 15485223]
- [87]. Wang J, Wang W, Kollman PA and Case DA 2006 Automatic atom type and bond type perception in molecular mechanical calculations *J. Mol. Graph. Model.* 25 247–60 [PubMed: 16458552]
- [88]. Jakalian A, Jack DB and Bayly CI 2002 Fast, efficient generation of high-quality atomic charges. AM1-BCC model: II. Parameterization and validation *J. Comput. Chem.* 23 1623–41 [PubMed: 12395429]
- [89]. Shao Q, Jiang Y and Yang ZJ 2022 EnzyHTP: a high-throughput computational platform for enzyme modeling/*J. Chem. Inf. Model.* 62 647–55 [PubMed: 35073075]
- [90]. Ryckaert J-P, Ciccotti G and Berendsen HJC 1977 Numerical integration of the Cartesian equations of motion of a system with constraints: molecular dynamics of n-alkanes *J. Comput. Phys.* 23 327–41
- [91]. Loncharich RJ, Brooks BR and Pastor RW 1992 Langevin dynamics of peptides: the frictional dependence of isomerization rates of N-acetylalanyl-N'-methylamide *Biopolymers* 32 523–35 [PubMed: 1515543]
- [92]. Berendsen HJC, Postma JPM, Gunsteren WFV, DiNola A and Haak JR 1984 Molecular dynamics with coupling to an external bath *J. Chem. Phys.* 81 3684–90
- [93]. Roe DR and Cheatham TE 3rd 2013 PTRAJ and CPPTRAJ: software for processing and analysis of molecular dynamics trajectory data *J. Chem. Theory Comput.* 9 3084–95 [PubMed: 26583988]
- [94]. Ufimtsev IS and Martinez TJ 2009 Quantum chemistry on graphical processing units. 3. Analytical energy gradients, geometry optimization, and first principles molecular dynamics *J. Chem. Theory Comput.* 5 2619–28 [PubMed: 26631777]
- [95]. Titov AV, Ufimtsev IS, Luehr N and Martinez TJ 2013 Generating efficient quantum chemistry codes for novel architectures /*J. Chem. Theory Comput.* 9 213–21 [PubMed: 26589024]
- [96]. Rohrdanz MA, Martins KM and Herbert JM 2009 A long-range-corrected density functional that performs well for both ground-state properties and time-dependent density functional theory excitation energies, including charge-transfer excited states *J. Chem. Phys.* 130 054112 [PubMed: 19206963]
- [97]. Hariharan PC and Pople JA 1973 The influence of polarization functions on molecular orbital hydrogenation energies *Theor. Chim. Acta* 28 213–22
- [98]. Cornell WD, Cieplak P, Bayly CI and Kollman PA 1993 Application of RESP charges to calculate conformational energies, hydrogen bond energies, and free energies of solvation *J. Am. Chem. Soc.* 115 9620–31
- [99]. Shrake A and Rupley JA 1973 Environment and exposure to solvent of protein atoms. Lysozyme and insulin *J. Mol. Biol.* 79 351–71 [PubMed: 4760134]
- [100]. McGibbon RT, Beauchamp K, Harrigan M, Klein C, Swails J, Hernández C, Schwantes C, Wang L-P, Lane T and Pande V 2015 MDTraj: a modern open library for the analysis of molecular dynamics trajectories *Biophys. J.* 109 1528–32 [PubMed: 26488642]
- [101]. Alexandrova AN, Rothlisberger D, Baker D and Jorgensen WL 2008 Catalytic mechanism and performance of computationally designed enzymes for Kemp elimination *J. Am. Chem. Soc.* 130 15907–15 [PubMed: 18975945]

- [102]. Khersonsky O, Röthlisberger D, Dym O, Albeck S, Jackson CJ, Baker D and Tawfik DS 2010 Evolutionary optimization of computationally designed enzymes: Kemp eliminases of the KE07 series *J. Mol. Biol.* 396 1025–42 [PubMed: 20036254]
- [103]. Caselle EA et al. 2019 Kemp eliminases of the AlleyCat family possess high substrate promiscuity *ChemCatChem* 11 1425–30 [PubMed: 31788134]
- [104]. Welborn VV, Pestana LR and Head-Gordon T 2018 Computational optimization of electric fields for better catalysis design *Nat. Catal* 1 649–55
- [105]. Fleishman SJ et al. 2011 RosettaScripts: a scripting language interface to the Rosetta macromolecular modeling suite *PLoS One* 6 e20161 [PubMed: 21731610]
- [106]. Towns J et al. 2014 XSEDE: accelerating scientific discovery *Comput. Sci. Eng* 16 62–74

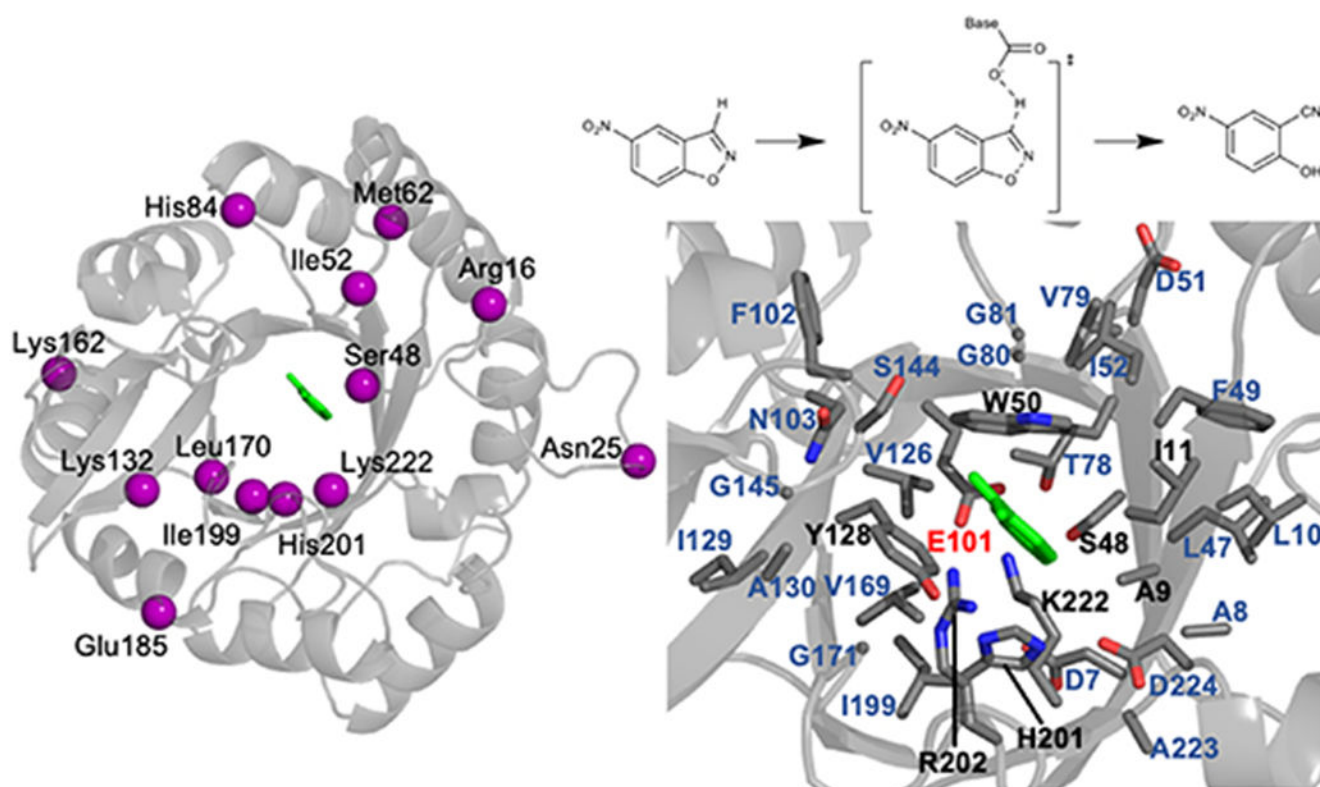


Figure 1. Mutation spots, catalyzed reaction, and active site residues of Kemp eliminase KE07-R7-2. (*Left*) Spatial distribution of mutation spots. The C_{α} atom of each site is shown in purple sphere, and the substrate is shown in green sticks. (*Top-right*) Catalyzed Kemp elimination reaction. The single transition state involves the deprotonation of a carbon atom. The transition state is stabilized by a general base from an amino acid side chain. The partial negative charge on the oxygen atom is stabilized by a hydrogen bond donor, which can be an amino acid side chain or a solvent water molecule. (*Bottom-right*) Active site residues are shown in stick. The substrate is shown in green, and the rest of the residues are shown in gray. The catalytic base is labeled in red. The residues that are within 5 Å from the substrate are labeled in black. The residues that are between 6 and 8 Å from the substrate are labeled in darker blue.

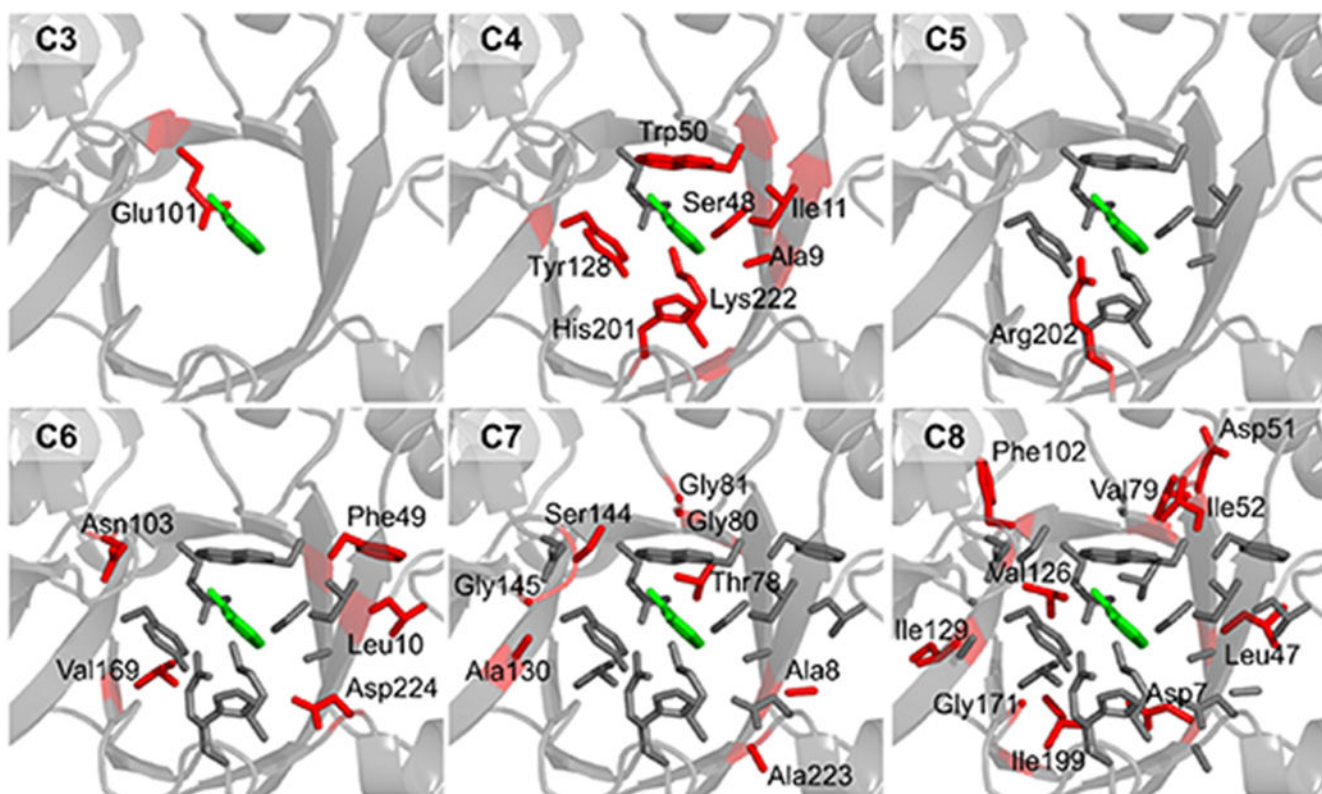


Figure 2. Six active-site regions with various boundary distances to the substrate. The distances used here range from 3 to 8 Å; the regions are named C3–C8, respectively. For each region, the selected residues are shown in stick. Compared to an adjacent region with a smaller size, the newly-added residues shown in red stick and labeled with a residue name; the existing residues are shown in gray stick.

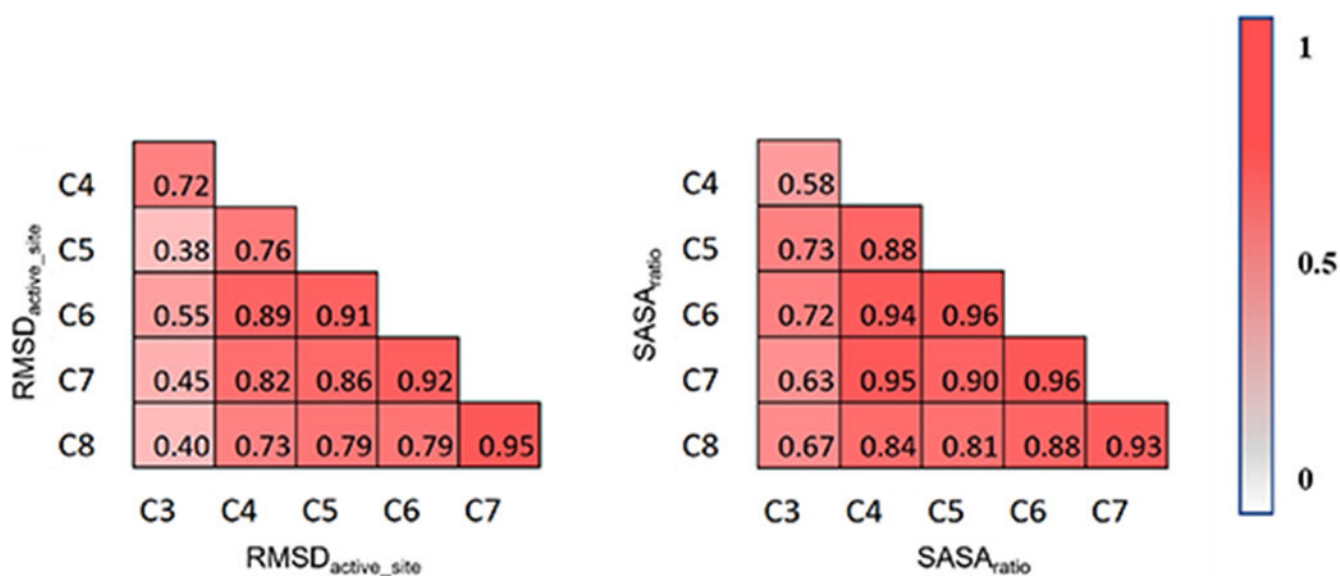


Figure 3. Spearman correlation matrices for protein dynamics-derived descriptors, RMSD_{active_site} (*left*) and SASA_{ratio} (*right*). Each matrix element represents a Spearman correlation coefficient for a pair of active-site regions with a distinct size. The magnitude of the correlation is coded by a gradient color bar that ranges from 0 (white) to 1 (red).

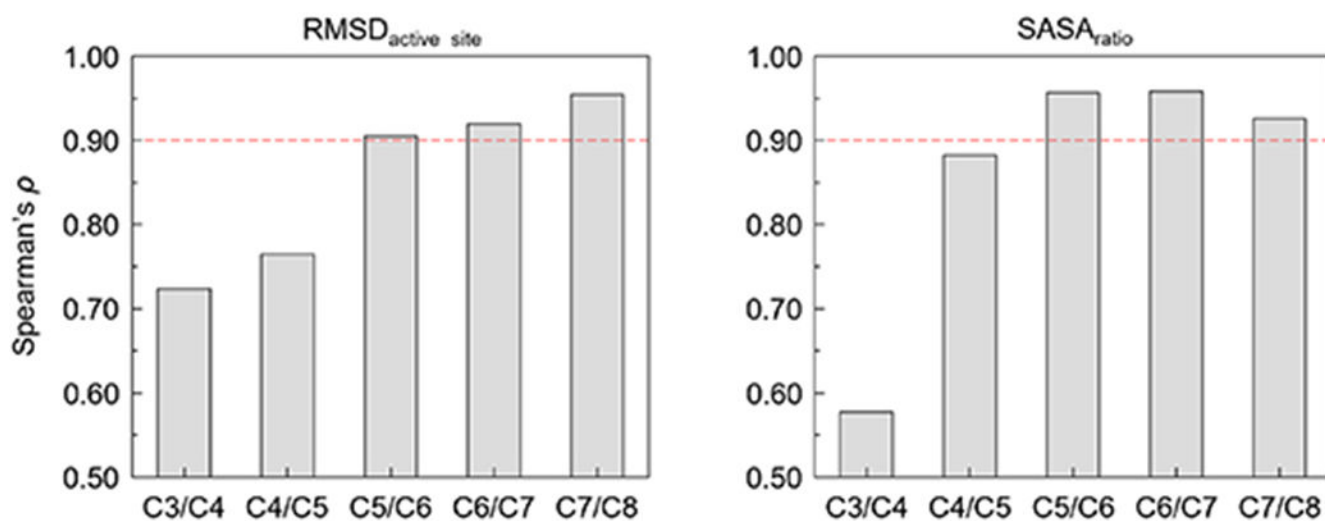


Figure 4. Spearman correlation coefficients for the dynamics-derived descriptors, RMSD_{active_site} (*left*) and SASA_{ratio} (*right*) between regions that are close in size. The red dashed line indicates the convergence cutoff of 0.90.

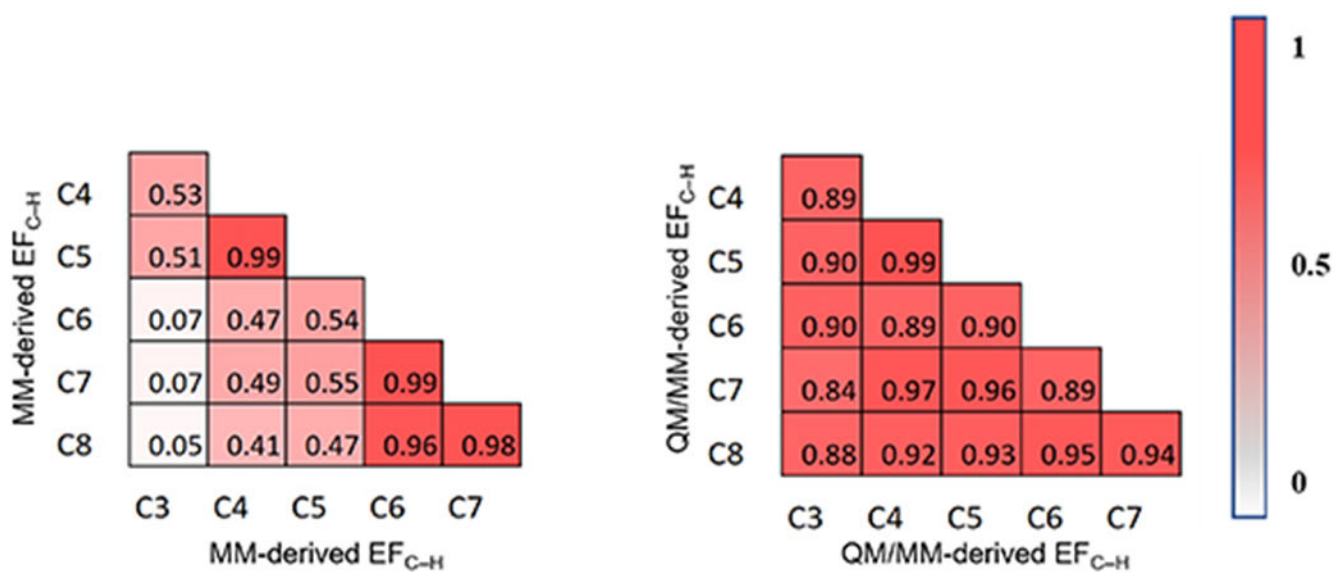


Figure 5. Spearman correlation matrix for MM-derived EF_{C-H} (*left*) and QM/MM-derived EF_{C-H} (*right*). Each matrix element represents a Spearman correlation coefficient for a pair of active-site regions with different region sizes. The magnitude of the correlation is coded by a gradient color bar that ranges from 0 (white) to 1 (red).

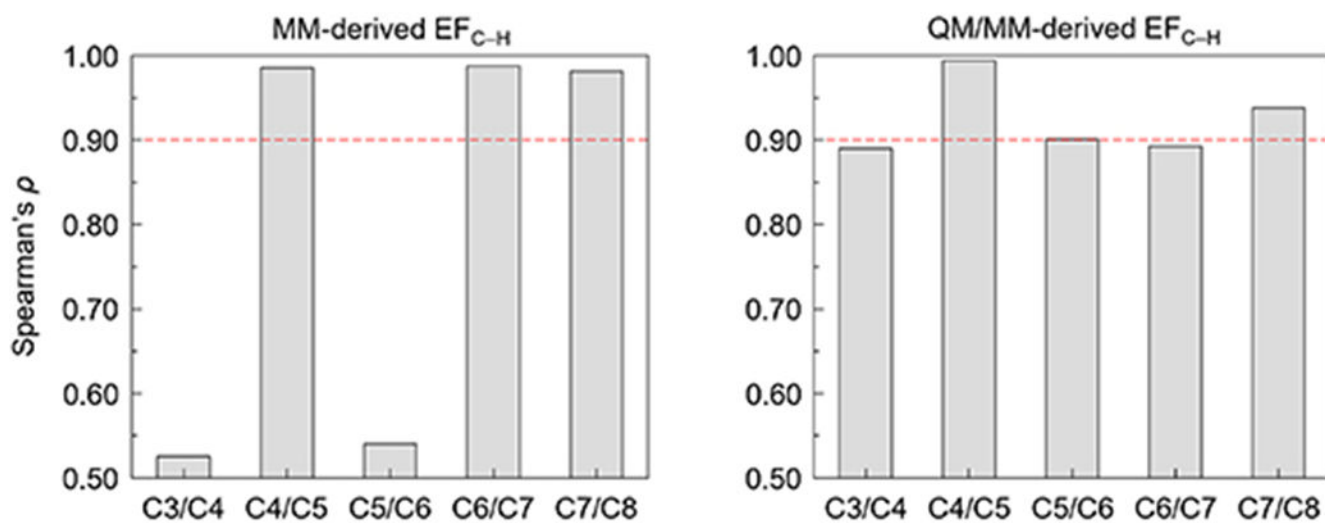


Figure 6. Spearman correlation coefficients for the dynamics-derived descriptors, MM-derived EF_{C-H} (*left*) and QM/MM-derived EF_{C-H} (*right*) between regions that are close in size. The red dashed line indicates the convergence cutoff of 0.90.

Table 1.

Number of residues, number of atoms, and net charge for the six active-site regions of KE07-R7-2 with different region sizes. Substrate atoms are counted in the number of atoms.

Cutoff (Å)	Number of residues	Number of atoms	Net Charge
3	1	31	-1
4	8	155	0
5	9	179	1
6	14	260	0
7	22	336	0
8	32	495	-2

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript