

A new discrete-geometry approach for integrative docking of proteins using chemical crosslinks

Yichi Zhang^{1,#}, Muskaan Jindal^{2,#}, Shruthi Viswanath^{2,*}, and Meera Sitharam^{1,*}

¹CISE Department, University of Florida, Gainesville, Florida 32611-6120, United States

²National Center for Biological Sciences, Tata Institute of Fundamental Research, Bengaluru 560065, India

#Contributed equally

*Corresponding Author: shruthiv@ncbs.res.in (SV), sitharam@cise.ufl.edu (MS)

Abstract

The structures of protein complexes allow us to understand and modulate the biological functions of the proteins. Integrative docking is a computational method to obtain the structures of a protein complex, given the atomic structures of the constituent proteins along with other experimental data on the complex, such as chemical crosslinks or SAXS profiles. Here, we develop a new discrete geometry-based method, wall-EASAL, for integrative rigid docking of protein pairs given the structures of the constituent proteins and chemical crosslinks. The method is an adaptation of EASAL (Efficient Atlasing and Search of Assembly Landscapes), a state-of-the-art discrete geometry method for efficient and exhaustive sampling of macromolecular configurations under pairwise inter-molecular distance constraints. We provide a mathematical proof that the method finds a structure satisfying the crosslink constraints under a natural condition satisfied by energy landscapes. We compare wall-EASAL with IMP (Integrative Modeling Platform), a commonly used integrative modeling method, on a benchmark, varying the numbers, types, and sources of input crosslinks, and sources of monomer structures. The wall-EASAL method performs better than IMP in terms of the average satisfaction of the configurations to the input crosslinks and the average similarity of the configurations to their corresponding native structures. The ensembles from IMP exhibit greater variability in these two measures. Further, wall-EASAL is more efficient than IMP. Although the current study uses crosslinks, the method is general and any source of distance constraints can be used for integrative docking with wall-EASAL. However, the current implementation only supports binary rigid protein docking, *i.e.*, assumes that the monomer structures are known and remain rigid. Additionally, the current implementation is deterministic, *i.e.*, it does not account for uncertainties in the crosslinking data beyond using distance bounds. Neither of these appears to be a theoretical or algorithmic limitation of the EASAL methodology. Structures from wall-EASAL can be incorporated in methods for modeling large macromolecular assemblies, for example by suggesting rigid bodies or restraints for use in these methods. This will facilitate the characterization of assemblies and cellular neighborhoods at increased efficiency, accuracy, and precision. The wall-EASAL method is available at <https://bitbucket.org/geoplexity/easal-dev/src/Crosslink> and the benchmark is available at https://github.com/isblab/Integrative_docking_benchmark.

Introduction

Protein-protein interactions play a crucial role in biological processes, for example, in immune response, metabolism, growth, and development. Characterizing the structures of complexes formed by two or more proteins *via* a single experimental technique can often be challenging. Protein-protein docking methods aim to computationally determine the structure of the complex formed by two proteins, given their three-dimensional structures (Lensink et al., 2023; Wodak et al., 2023). Some protein-protein docking methods employ rigid docking, where it is assumed that the proteins do not undergo significant conformational change upon binding. Rigid docking is computationally efficient since the search space is restricted to rigid translations and rotations of one protein with respect to the other. In integrative docking, additional experimental data, such as residual dipolar couplings from NMR spectroscopy and data from chemical crosslinking mass spectrometry (XLMS) can be used to guide the docking search (Braitbard et al., 2019; Koukos & Bonvin, 2020; Rout & Sali, 2019; Russel et al., 2012; D. Saltzberg et al., 2019; Schneidman-Duhovny et al., 2012). Integrative docking methods aim to compute an ensemble of structures of the complex that are consistent with the experimental data. Here, we develop a new method for integrative rigid docking of protein pairs using inter-protein chemical crosslinks.

Chemical crosslinking involves treating the protein complex of interest with a chemical crosslinker (Rappsilber, 2011; Yu & Huang, 2023). A crosslinker consists of two reactive groups, separated by a spacer that defines the maximum crosslinker length. Common crosslinkers include DSSO (Disuccinimidyl sulfoxide), DMTMM (4-(4,6-Dimethoxy-1,3,5-triazin-2-yl)-4-methylmorpholiniumchloride), ADH (Adipic Dihydrazide), and EDC (1-ethyl-3-(3-dimethylaminopropyl)carbodiimide hydrochloride) (Rappsilber, 2011; Yu & Huang, 2023). The reactive groups in a crosslinker can bond covalently with accessible residues in a complex. Subsequent treatment with trypsin and analysis of the resulting mass spectrometry data provides a list of residue pairs in the complex that are crosslinked. Therefore, XLMS can provide upper bound distances between pairs of residues, which can inform the proximity of these residues in the structure of the complex.

Several methods exist for modeling protein complexes based on atomic structures and crosslinks (Arvindekar et al., 2024; Beck et al., 2024; Graziadei & Rappsilber, 2022; Rappsilber, 2011; Yu & Huang, 2023). Well-known integrative modeling methods such as Haddock, IMP, and Assembline can be used to dock proteins based on crosslinks (Dominguez et al., 2003; Honorato et al., 2024; Rantos et al., 2022; Russel et al., 2012; D. J. Saltzberg et

al., 2021). These methods are broadly applicable, allowing for the use of diverse types of data, including data other than crosslinks and structures. Many of them also allow for determining the structure of larger complexes and assemblies. Some integrative modeling methods were developed specifically for modeling with crosslinking data, such as XLMOD (Ferber et al., 2016) and IMProv (Ziemianowicz et al., 2021). Recently, deep learning-based methods for structure prediction have been extended to use chemical crosslinks as additional inputs. Methods such as Alphalink, Alphalink2, and DistanceAF use the crosslinking data either in the input residue pair representation (Alphalink and Alphalink2) or in the loss function (DistanceAF) (Stahl et al., 2023, 2024; Zhang et al., 2023).

Many of the aforementioned methods for modeling with crosslinks employ randomized sampling, e.g., *via* Markov Chain Monte Carlo (MCMC) methods which is typically not stochastic, *i.e.*, the sampled regions of the landscape depend on the choice of initial random configurations of the sampling trajectories (Arvindkar et al., 2022; Pasani & Viswanath, 2021; D. J. Saltzberg et al., 2021). In particular, exhaustive sampling of complex landscapes is not guaranteed in these methods. Also, the combination of simplified coarse-grained bead representations with hard-sphere excluded volume restraints in some methods may make accurate modeling particularly challenging for protein structures with concavities, such as grooves and pits in the interfaces. These perceived limitations of current approaches motivated us to explore other sampling methods.

EASAL (Efficient Atlasing and Search of Assembly Landscapes) is a state-of-the-art discrete geometry-based methodology for roadmapping, sampling and analyzing the landscape of macromolecular configurations satisfying possible sets of pairwise inter-molecular distance constraints (Ozkan & Sitharam, 2011; Prabhu et al., 2020). EASAL is both a resource-light, stand-alone method and also complements prevailing MC, MD and docking methods (Ozkan et al., 2021), demonstrating superior performance, especially for discontinuous pair-potential energy landscapes. The EASAL methodology (Prabhu et al., 2020) and curated open-source software (Ozkan et al., 2018), <https://bitbucket.org/geoplexity/easal-dev/src/master/> (see also <http://www.cise.ufl.edu/~sitharam/EASALvideo.mpeg>) can efficiently generate an exhaustive ensemble of structures lying within specific pair-potential wells, discretized as a staircase of nested distance-interval constraints. The methodology has been earlier used for effectively predicting virus assembly pathways (Wu et al., 2020), sticky-sphere path integrals (Prabhu et al., 2020) and free energy, configurational entropy, or volume computation (Zhang & Sitharam, 2022, 2024). Here, we leverage the unique features of the EASAL method for integrative docking with crosslinks. Given the structures of two proteins and an input set of crosslinks

between them, the modified method, *wall-EASAL*, produces an ensemble of structures of the complex satisfying the maximum number of input crosslinks.

We compared the performance of *wall-EASAL* with IMP on the problem of integrative docking with crosslinks on thirty protein pairs, varying the number of input crosslinks, the crosslinker length, the source of crosslinks, and the source of monomer structures. Assessing the structure ensembles from these methods based on their average satisfaction of the input crosslinks and their average similarity to the corresponding native structures, we find that *wall-EASAL* performs better than IMP. It is also more efficient. Although the current study uses crosslinks, the method is general and any source of distance constraints can be used for integrative docking with *wall-EASAL*. The limitations are that the current implementation only supports binary protein docking, the monomer structures are assumed to be known and remain rigid, and the uncertainty in the crosslinking experiment is not considered. However, none of these appears to be a theoretical or algorithmic limitation of the EASAL methodology. Structures from *wall-EASAL* can complement methods for modeling large macromolecular assemblies by suggesting rigid bodies or restraints for use in integrative modeling methods (Bryant et al., 2022; Chim & Elofsson, 2024; Dominguez et al., 2003; Rantos et al., 2022; Russel et al., 2012; D. J. Saltzberg et al., 2021; Shor & Schneidman-Duhovny, 2024). This approach is expected to enhance the efficiency, accuracy, and precision at which large assemblies and cellular neighborhoods are characterized.

Notes on terminology. We use (sampled) “structure of a complex” and “configuration” interchangeably. A “restraint” is a probabilistic term with biophysical origin referring to a constraint that may or may not be satisfied. A “constraint” is a geometric condition that is deterministic and Boolean. Either a constraint is satisfied, or it is not. We use “crosslink distance” to refer to the distance between crosslinked residues in a structure of the complex.

Methods

EASAL Background and Crosslink Satisfaction

The unique features of EASAL mitigate the curse of dimensionality in configurational entropy (free energy), pathway, and kinetics computations by achieving the following.

(a) Decoupling exploration from sampling, that is, generating an atlas of the landscape, including a roadmap of basins, barriers, paths, and their neighborhood relationships, with

minimal sampling, using geometric constraints (Sitharam et al., 2019) and rigidity-based roadmap (Fig. 1).

(b) Reducing ambient dimension and convexifying contiguous, constant-potential-energy regions (*macrostates*) using customized, Cayley parametrization (Sitharam & Gao, 2010) which is a distance-based internal coordinate representation of assembly configurations that are constrained by inter-residue distances (Fig. 1).

The Cayley parameterization idea is broadly applicable (Sadjadi et al., 2021; Sitharam & Wang, 2014; Wang & Sitharam, 2015) as it maps landscapes – with complex topologies in high dimensional ambient space defined by distance-interval constraints – into a convex base space of much lower intrinsic dimension. In the context of molecular assembly landscapes, convex Cayley parameterization additionally achieves high sampling efficiency and accuracy, avoiding gradient-descent and repeated or discarded configurations.

The roadmap component of the atlas is a directed acyclic graph (Fig. 1). Each node of the roadmap represents a region of the landscape in the well of a specific set of pair-potentials called the *active constraint graph*. These constraints are imposed by the pair potentials. Each node, or *active constraint region* is a collection of a small number of *macrostates*, or constant-energy, contiguous regions. In case of short-range pair-potentials, each constraint is between an inter-monomer residue-pair whose inter-residue distance lies in a small interval that achieves minimum energy, treated in the limit as a hard-sphere potential, which prevailing methods based on Monte Carlo sampling or molecular dynamics find challenging. The theory extends easily – albeit with an efficiency tradeoff – to longer range potentials, or *crosslink intervals*, discretized as nested distance-interval constraints.

The directed acyclic graph structure represents a stratification of active constraint regions by effective dimension and energy level. The effective co-dimension of an active constraint region (or a node of the roadmap) can be determined directly as a number of constraints of edges of the active constraint graph, using combinatorial rigidity (Sitharam et al., 2019) of the active constraint graph. Consequently, the effective dimension of a macrostate in an active constraint region becomes a proxy for the energy level. Parent regions in the roadmap have one higher dimension or energy level than child regions, which have one more active constraint than the parent, whereby the roadmap facilitates basins and their neighborhoods to be faithfully represented.

Individual sampled configurations, which are traditionally represented using Cartesian parameters, are instead represented using Cayley or distance-based parameters that are customized to the active constraint graph of the macrostates. Each Cayley parameter represents the distance between a residue pair, including active constraint pairs. Furthermore, inverting the Cayley parametrization is cheap: a Cayley representation is mapped to a small number of pre-image Cartesian configurations (representing different chiralities). Crucially, within an active constraint region, using Cayley parameters avoids gradient-descent search used by all prevailing methods to sample constrained regions. Further, under Cayley parameterizations, active constraint regions or macrostates become convex spaces with easily computable bounds (Fig. 1).

In a Cayley-convexifiable macrostate or active constraint region, there is a collection of residue pairs (the Cayley parameters) satisfying the following property: very roughly speaking, an assembly system can follow straight-line paths when parameterized using Cayley parameters and still avoid breaking energy barriers, *i.e.*, while remaining in the same macrostate (a suitable reaction coordinate basis). Convexification improves sampling efficiency for assembly landscapes, significantly reducing the number of repeated and discarded configurations. Overall, the methodology directly addresses the curse of dimension and complexity of landscapes while giving formal guarantees of efficiency, accuracy, robustness, and trade-offs for the core algorithm.

Since EASAL's roadmapping and sampling algorithm is integrally based on using inter-monomer residue-pair distances as both distance-interval constraints and Cayley parameters, EASAL is almost tailored for sampling configurations satisfying the crosslink (and collision) constraints. Maximal subset of crosslinks can be chosen that guarantee convexification of the corresponding Cayley configuration space. Further crosslinks outside this subset are checked *a posteriori* for each sampled configuration. However, as mentioned earlier, due to the relatively wide crosslink distance intervals, there is significant loss of efficiency in directly using EASAL for sampling configurations that satisfy all crosslinks.

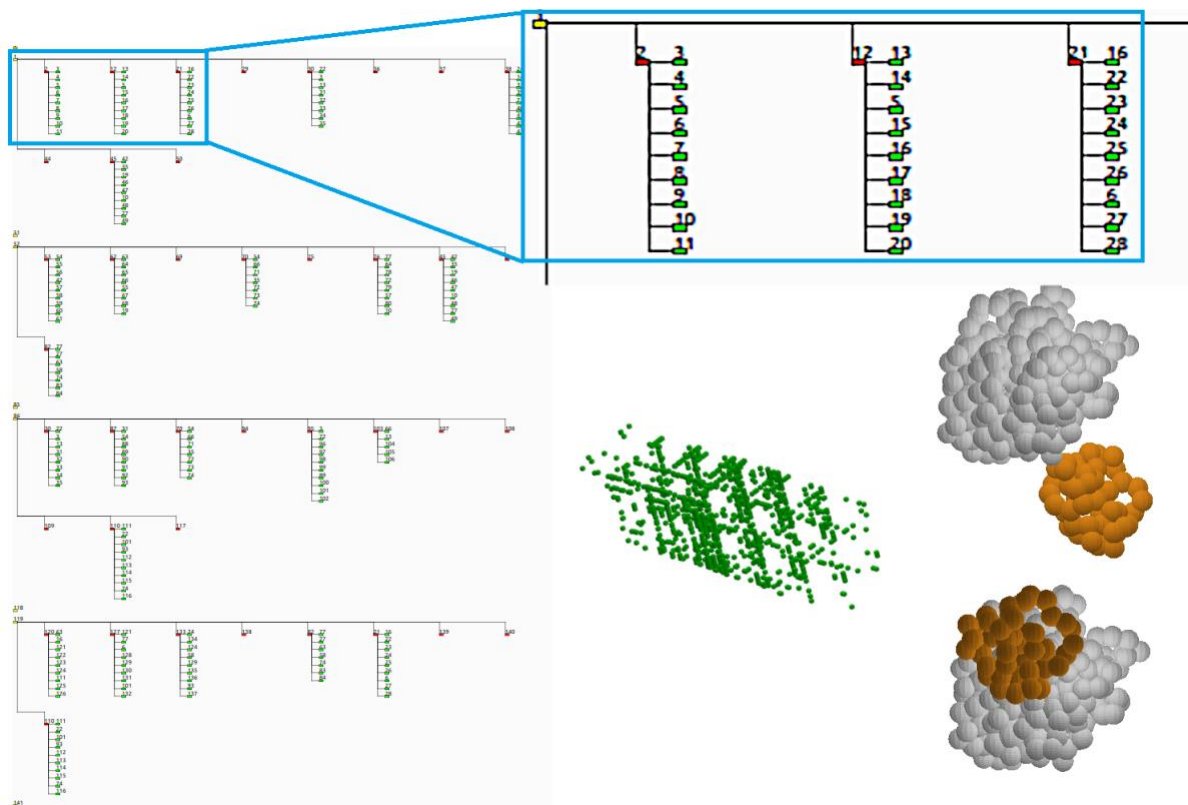


Figure 1: EASAL Background. Illustrative screenshots of EASAL from input case 1r0r/DSSO/3. Top: EASAL roadmap (directed acyclic graph represented as a tree by repeating nodes). Left: enlarged portion of the roadmap for detail; Bottom mid: view of the Cayley configuration space using distances between crosslinked residues as Cayley parameters, each green cube represents a feasible (collision-free and satisfying all crosslinks) configuration. Bottom right: selected configurations of the system satisfying all crosslinks and collision constraints, gray - Monomer A, orange - Monomer B.

Wall-EASAL

To boost the efficiency that deteriorates when directly using EASAL to deal with a 6-dimensional region defined by the large distance intervals arising from the crosslink ranges, we develop a novel adaptation of EASAL, called *wall-EASAL* for reasons that will be clear from the discussion below. Wall-EASAL fully leverages the EASAL methodology's ability to deal with *exact distance or small distance-interval* constraints, by mapping the constrained configuration spaces from their high ambient dimension to a convex Cayley space in their typically much lower intrinsic dimension. The advantage of the mapping is that it completely avoids the gradient descent used by prevailing methods to enforce distance constraints.

The key intuition behind wall-EASAL (mathematically proven in the Supporting Information, see Fig. 2) is that if the collision free configuration space is *path-connected*, and there is a

feasible configuration satisfying all crosslink and collision constraints, then there must also exist a feasible configuration in which some crosslink attains either its maximum or its minimum distance *exactly*, *i.e.*, one of the extremes or *walls* of the crosslink distance interval. Then the strategy is to replace each crosslink's distance interval constraint $c \in (l, u)$ with two distance constraints, $c = l$ and $c = u$. This effectively splits each k -dimensional active constraint region in the direct EASAL approach into 2^{6-k} subregions, each being a "wall" of the original region of all feasible configurations.

If necessary, more walls can be introduced distributed in the interior of the crosslink distance-interval, *i.e.*, between the two extremes of the interval, at the expense of reducing efficiency (see Discussion section). In this paper, however, wall-EASAL only uses the two extreme walls, which are shown to nevertheless provide a representative set of configurations.

Using this process, wall-EASAL returns all collision-free configurations *in the wall subset*, *i.e.*, in which at least 1 crosslink has its distance at its wall, and all (or the maximum number of) crosslinks are satisfied.

The path-connectivity condition is important. Otherwise, disconnectivity in the collision-free configuration space could result in the feasible configuration space (intersection of collision-free and crosslink-satisfying configuration spaces) lying in the interior of the crosslink-satisfying configuration space as in Fig. 2 (right). In this case, no crosslink distance in such configurations lies at the extremes of the crosslink distance interval, and wall-EASAL will fail to find a feasible configuration, although one exists.

Next, we discuss two potential issues that could impact the performance of wall-EASAL. (1) How realistic is the path-connectivity assumption on the collision-free configuration space, which is crucial to guarantee that wall-EASAL finds a collision-free configuration that satisfies all crosslinks if one exists (or one that satisfies the maximum number of crosslinks)? (2) How representative are the wall configurations in the space of all (interior) feasible configurations?

Walls and Pockets

Path-connectivity of the collision-free configuration space is guaranteed unless there is a collision-free configuration from which the two monomers are unable to untangle and break free while following a configurational path that avoids collisions. This scenario is avoided in most situations as long as (a) the monomers cannot be arranged into a *knot* (formally a *link* in

topology terminology) configuration, and (b) there is no *pocket* in one monomer into which some part of the other monomer is jammed with no collision-free exit from the jammed configuration (Fig. 3).

Although true links and pockets are rare, such artifacts can arise from coarse sampling. *i.e.*, although the collision-free configuration space may be path-connected, there could be a narrow bottleneck through which all paths must pass, effectively disconnecting coarsely sampled regions. Thus, the notion of path-connectivity that guarantees wall-EASAL's accuracy is in fact a relative notion that depends on the coarseness of sampling (Fig. 3).

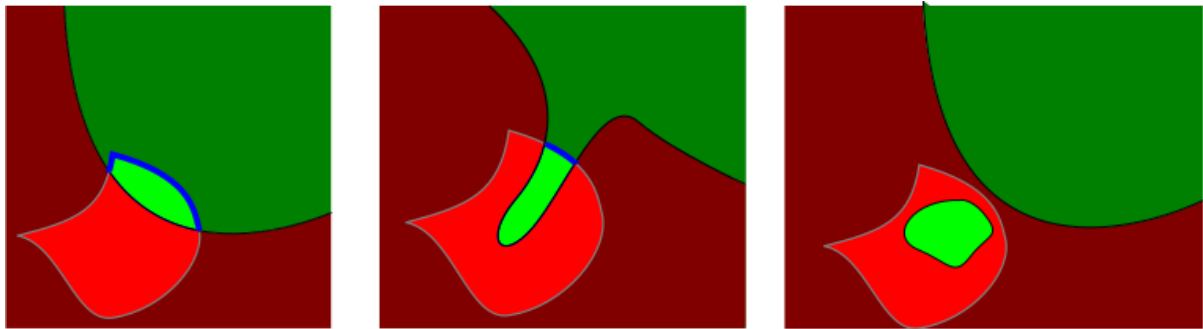


Figure 2: A schematic illustration of configuration spaces relevant to wall-EASAL. Shades of green: space of collision-free configurations, shades of red: space of colliding configurations. Lighter color: all crosslinks satisfied, darker color: not all crosslinks satisfied. The boundary between light and dark is a wall to which wall-EASAL sampling is restricted. Regular EASAL returns configurations in the light green region, wall-EASAL returns configurations in the blue curve only. Left: common/standard input cases; mid: input cases with a narrow bottleneck in collision-free configuration space; right: input cases disconnected with a pocket.

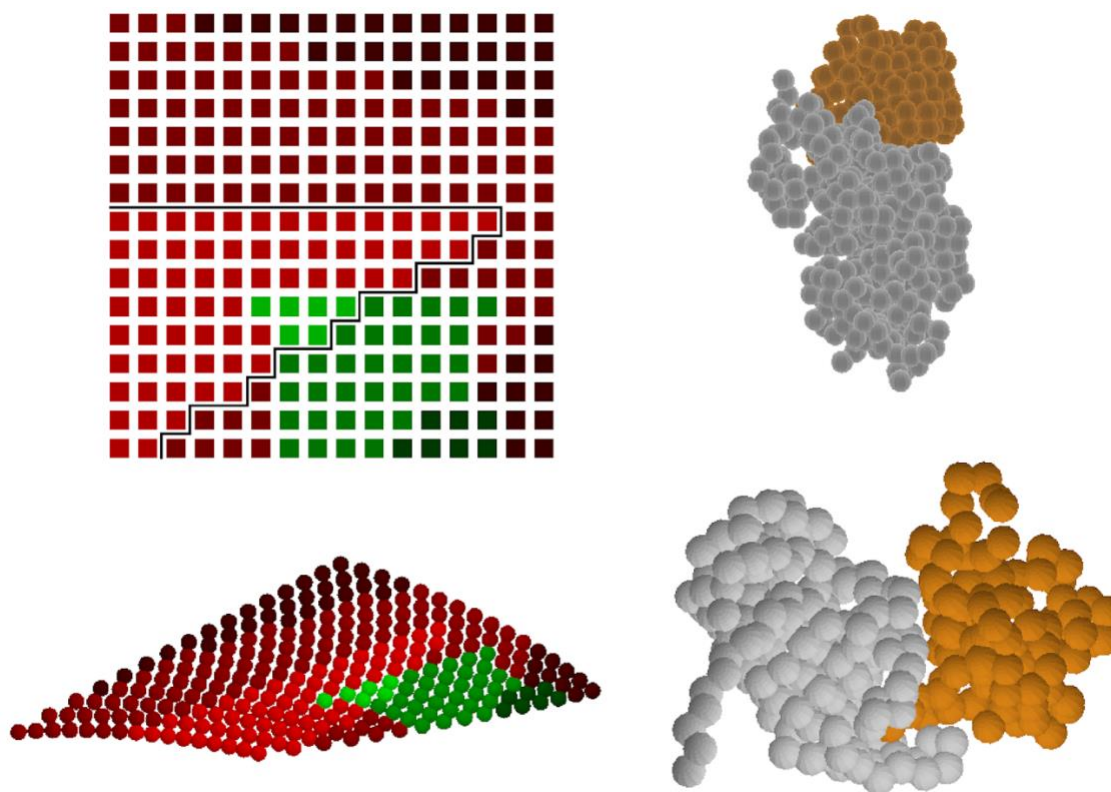


Figure 3: Pocket Artifact. Left: Cayley (above) and Cartesian (below) representations of a typical 2-dimensional slice of the configuration space for 2b42/DMTMM/10 in the neighborhood of a collision-free configuration (Top Right) satisfying all 10 crosslinks found by IMP. Wall-EASAL failed to find such a configuration due to pocket artifacts from *coarse* sampling. The slices were *finely* sampled for purposes of analysis/diagnosis. Each dot represents a sampled configuration, with Red = collision, Green = collision-free. Lighter shade of color means more crosslinks satisfied with the lightest being 10 (all crosslinks satisfied). Grey denotes the boundary of the region in which all 10 crosslinks are satisfied. Although the boundary between the lightest and slightly darker shade of green in fact consists of feasible wall configurations, since the entire feasible region is narrow, *coarse* sampling created a pocket artifact and caused wall-EASAL to miss this wall. However, for the input 2hle/DMTMM/9 with nearly the same number of crosslinks as 2b42, a variety of feasible configurations satisfying the maximum number of crosslinks (matching IMP) were found by wall-EASAL (Bottom Right).

Walls versus Interiors

Although the wall subset of configurations is non-empty provided the entire configuration space is either non-empty or path-connected, are walls representative of the entire feasible region including the interiors?

We answer the question affirmatively by providing quantitative results in the next section. The geometric intuition for the wall being representative of the whole is that the volume of a high-dimensional object (such as the feasible region in discussion here) mostly lies close to its boundary, or “most of a high-dimensional orange’s volume is at its peel”. Fig. 4 pictorially illustrates for the input case 1r0r/DSSO/3 how the wall subset is a good representation of the entire, much larger set of feasible configurations.

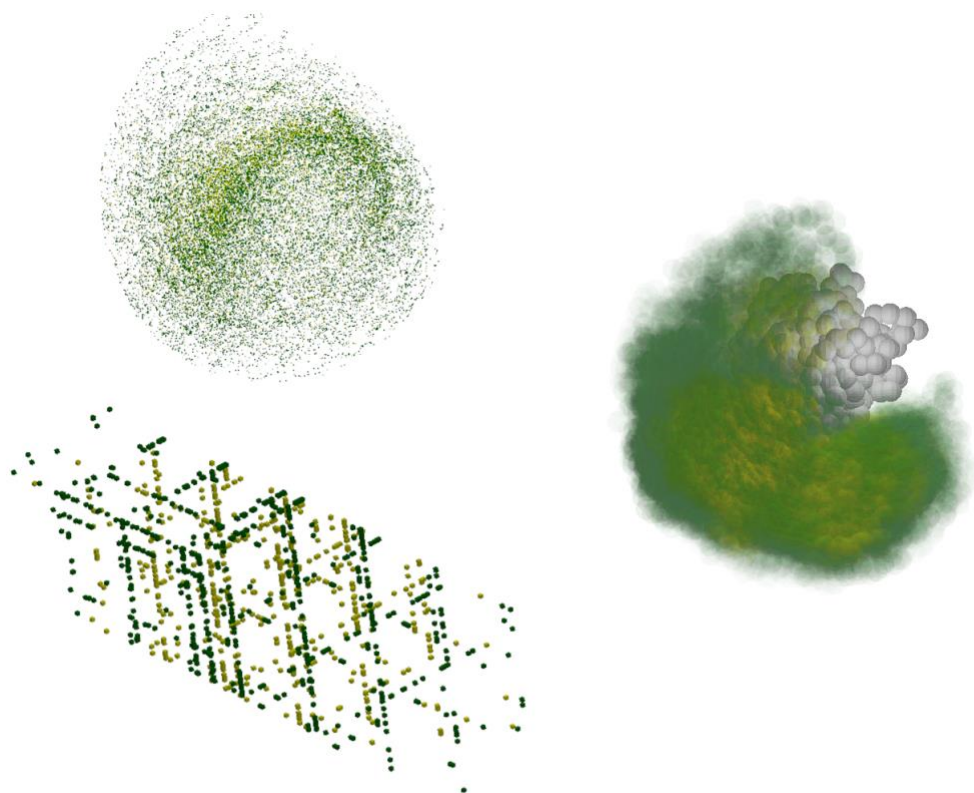


Figure 4. Representativeness of wall-EASAL sampling. Three different views of input case 1r0r/DSSO/3 showing all configurations on crosslink constraint walls (darker green) and not on walls (yellowish lighter green). Top Left: sampled configurations projected to 3 Cartesian dimensions (x, y, z), Bottom: projected to 3 crosslink distances used as Cayley parameters (the same as Fig. 1 Cayley configuration space). Right: sweep view of wall and interior feasible configurations of Monomer B (dark green - on wall, light green - off wall) with respect to Monomer A (gray) held fixed.

Benchmark creation

Structures

We constructed a benchmark consisting of thirty integrative modeling cases of binary complexes (Table S1). There are twelve hetero-dimers, comprising seven complexes with

experimentally solved structures from the Zlab 5.5 protein-protein docking benchmark (Guest et al., 2021), and five AlphaFold-multimer predicted hetero-dimers from a crosslinking study (O'Reilly et al., 2023). Protein pairs with concave interfaces, containing grooves and pits were selected for docking, as these are perceived to be more difficult to dock than those with flat interfaces (Fig. S1). The monomer structures from the bound structure of the complex were used as input to the docking methods; the bound and unbound structures of the monomers are very similar for these cases.

Crosslinks

We generated two kinds of crosslinks: a longer crosslinker between lysines (DSSO) and a shorter crosslinker between aspartic acid and glutamic acid residues (DMTMM). DSSO and DMTMM crosslinks were simulated using Jwalk (Bullock et al., 2016). The maximum distances between crosslinked residues in Jwalk were set to 32 Å (20 Å) for DSSO (DMTMM) (Bullock et al., 2018). A false positive rate of 20% was used (default in Jwalk). Random subsets of inter-protein crosslinks from Jwalk were used for the benchmark cases.

In all, the benchmark consisted of twenty-five input cases with simulated crosslinks on experimentally determined structures of complexes from the Zlab 5.5 benchmark; the number and type of input crosslinks were varied across these cases. Additionally, there were five input cases with crosslinks from experimental studies alongside AlphaFold-multimer predictions of the complexes (Table S1).

Running IMP on the benchmark

We used the Integrative Modeling Platform's Python Modeling Interface (PMI) (IMP 2.17.0; <https://integrativemodeling.org>) for integrative docking. The modeling protocol was adapted from previous studies (Arvindkar et al., 2022; Liu et al., 2024; Pasani et al., 2023; Russel et al., 2012; D. J. Saltzberg et al., 2021).

The monomers were represented as independent rigid bodies based on their structures and coarse-grained at one residue per bead centered at the C α atom. Bayesian crosslinking restraints were used, along with excluded volume and sequence connectivity restraints (Arvindkar et al., 2022; D. J. Saltzberg et al., 2021; Shi et al., 2014). The Gibbs sampling Replica Exchange Markov Chain Monte Carlo (MCMC) algorithm was used for structural sampling. We started with initial random configurations for each protein in each pair. A configuration, *i.e.*, structure of the complex, was saved after every ten Gibbs sampling steps,

each of which consisted of a cycle of Monte Carlo moves comprising random translations and rotations of the monomer rigid bodies. We performed twenty independent runs with four replicas and ten thousand MCMC steps per run, resulting in eight million configurations for each benchmark input. Following this, we applied the analysis and validation protocols used in previous studies (Arvindekar et al., 2022; D. J. Saltzberg et al., 2021; Viswanath et al., 2017).

Running wall-EASAL on the benchmark

In wall-EASAL, each monomer was similarly represented as a rigid body, coarse-grained at one residue per bead, where each bead was centered at the C α atom and had a radius of 2.5 Å. Each crosslink was encoded by two distance constraints ('walls'), one corresponding to the upper bound and another corresponding to the lower bound on the distance between the centers of C α atoms of the crosslinked residues. The lower bound distance for both DMTMM and DSSO crosslinks was set to 10 Å based on the lengths of the crosslinked lysine and glutamic/aspartic acid side chains and the rigid parts of the crosslinker. The upper bound distance was set to 32 Å (20 Å) for DSSO (DMTMM) crosslinks based on the linker and side chain lengths and the backbone flexibility of the monomers (Bullock et al., 2018; Gong et al., 2020; Merkley et al., 2014; Shi et al., 2014). For each crosslink, two active constraints were defined in wall-EASAL, corresponding to the lower and upper bounds of the crosslink. An active constraint in wall-EASAL is of the form $\lambda * (r_i + r_j) + \delta$, where λ and δ are constants and r_i and r_j correspond to the radii of the crosslinked beads (Ozkan et al., 2018; Prabhu et al., 2020; Sitharam et al., 2019). Based on the above crosslink bounds, the values of λ and δ were set to 2 and 0 respectively for the lower bound constraint, and 0 and 32 (or 20) for DSSO (DMTMM) for the upper bound constraint. For all crosslink constraints, $r_i = r_j = 2.5$ Å was used. The threshold for the number of crosslinks to be satisfied by a configuration ('*crossLinkSatisfyThres*') was set to $n - 2$ where n is the number of crosslinks. The collision constraint, similar to the excluded volume restraint in IMP, was used to avoid overlapping of beads. The parameters for the collision constraints, λ and δ , were set to 1 and 0, respectively. Lastly, we used a step size of 5, which defines the resolution of sampling.

Analyzing wall-EASAL and IMP configurations

The ensemble of configurations from wall-EASAL and IMP were compared based on their crosslink satisfaction and similarity to the corresponding native structures on the benchmark.

The crosslink satisfaction for an ensemble was determined based on two measures: the maximum percentage of crosslinks satisfied by any configuration in an ensemble and the average distance between crosslinked residues across all crosslinks and configurations in the ensemble. A crosslink was satisfied by a configuration if the distance between the bead centers, corresponding to the crosslinked residues, was within an upper bound of 32 Å (20 Å) for DSSO (DMTMM) crosslinks.

The native structure for each complex was the experimentally determined PDB structure of the complex (first twenty-five cases in Table S1) or the AF-multimer prediction (last five cases in Table S1). The similarity of the configurations to the corresponding native structure was determined by two measures: the distance between crosslinked residues difference and the ligand RMSD (root-mean-square deviation). The former is the difference between the distance between crosslinked residues in a sampled configuration to the corresponding distance in the native structure across all crosslinks and across all configurations in an ensemble. The ligand RMSD between a configuration and the native structure was calculated as the C α RMSD of the second protein (ligand) after superposing the first protein (receptor) in the complex (Lensink et al., 2023; Wodak et al., 2023).

Results

Here, we compare and contrast the performance of two methods for integrative docking of protein pairs using chemical crosslinks. We first demonstrate that the new method, wall-EASAL, which samples the configurations at the wall, is representative, by comparing it with the vanilla EASAL which performs exhaustive sampling in the entire search space under distance constraints which we call *interior-EASAL*, in order to clearly differentiate the methods. Next, we compare and contrast the performance of wall-EASAL with that of IMP on a benchmark set of input cases. This comparison is based on the crosslink satisfaction of the respective configurations, the similarity of the configurations to the native structures, and the efficiency of the methods. Finally, we visualize the structures of the complex produced for a few input cases.

Coverage of the Interior using Wall-EASAL

To show that the ensemble generated by wall-EASAL sampling is indeed representative of the entire region of configurations satisfying constraints, we ran a coverage test similar to

(Ozkan & Sitharam, 2014; Zhang & Sitharam, 2024) by sampling both the entire region (using interior-EASAL) and the walls only (using wall-EASAL) and comparing these results to see if the smaller sample set “covers” the larger one. The coverage experiment is designed as follows.

For each input case, both sampling methods were executed yielding 2 sets of feasible configurations. Then the entire configuration space was partitioned into a 6-dimensional hypercube grid, and points generated by interior-EASAL were mapped into the hypercubes. We iteratively coarsened the grid (making each grid cube larger) until at least 90% of the occupied cubes had at least γ sampled points in them. Here the coefficient is defined as $\gamma = (\text{interior}/\text{wall})^{1/6}$ to avoid bias caused by difference in number of points between sampling methods. After grid size was determined, points sampled with wall-EASAL were mapped into such a grid, and all cubes with at least γ interior samples were checked to determine if they have at least 1 wall sample in it. The ratio between numbers of those cubes was then tabulated.

Since wall volume is expected to be a lower percentage of the larger interior volume when there are fewer crosslinks, we ran four representative cases with few crosslinks (2, 3, 4, and 5, respectively) to demonstrate the representativeness. In all these cases, wall-EASAL provides an ensemble as good as EASAL in terms of coverage. Specifically, the coverage rate of wall-EASAL is constantly over 80%, corroborating our conclusion that sampling only the walls of the feasible region provides a good coverage of the entire feasible region including the interior.

Table 1: Coverage result of wall-EASAL. Number of sampled configurations in each method and grid cubes they are mapped to are tabulated here.

Crossli nker	PDB	Number of crosslinks	Coverage ratio	Interior sample count	Wall sample count	Interior volume	Wall volume
DSSO	1clv	2	0.84	13604	2941	192	163
	1r0r	3	0.97	36105	5305	81	79
	2ayo	4	0.97	108092	10444	172	168
	2hle	5	1	128	128	51	51

Next, we compared the performance of wall-EASAL with that of IMP on a benchmark of thirty input cases.

Percentage of crosslinks satisfied

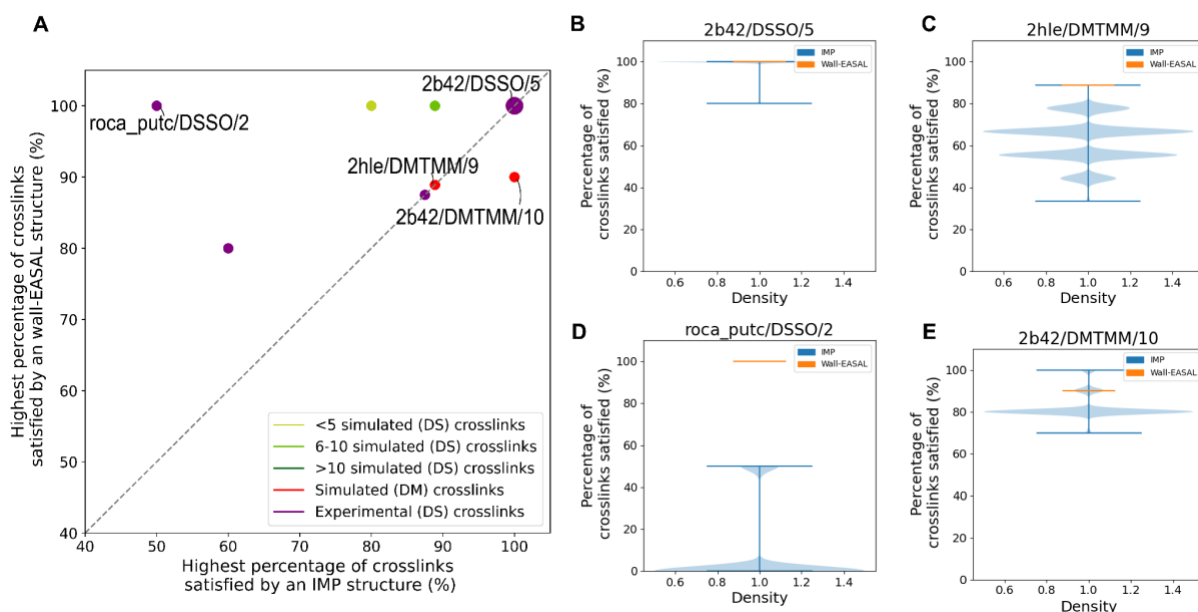


Figure 5. Percentage of crosslinks satisfied in wall-EASAL and IMP structures. For 30 benchmark inputs, the wall-EASAL (orange) and IMP (blue) ensembles are compared. **(A)** The highest percentage of crosslinks satisfied in any configuration in both the ensembles, where the larger point (top right) depicts the majority of the cases in which at least one configuration satisfies all the crosslinks. **(B-E)** Distribution of the percentage of crosslinks satisfied per configuration in the two ensembles for four cases. DS and DM refer to DSSO and DMTMM crosslinks, respectively.

We first examined the number of crosslinks satisfied by the IMP and wall-EASAL ensembles in the thirty input cases in the benchmark. An integrative structure satisfies an input crosslink if the corresponding Ca-Ca distance between the crosslinked residues is less than the upper bound; the upper bounds depend on the linker lengths and were set to 32 Å (20 Å) for DSSO (DMTMM) linkers (Arvindekar et al., 2022; D. Saltzberg et al., 2019; D. J. Saltzberg et al., 2021).

The wall-EASAL and IMP configurations satisfy the crosslinks similarly well in terms of the highest percentage of crosslinks satisfied by a single configuration in the ensemble (Fig. 5A, Fig. S2). However, the distributions of crosslink percentages in the ensembles suggest that the wall-EASAL configurations satisfy a greater percentage of crosslinks on average (Fig. 5B-5E, Fig. S2). The IMP ensemble is more diverse in terms of crosslinks satisfaction.

We provide a few examples. In most input cases, both the wall-EASAL and IMP configurations satisfy the crosslinks equally well. For example, in 2b42/DSSO/5 all the wall-EASAL and IMP configurations satisfy all the crosslinks (Fig. 5B, Fig. S2), and in 2hle/DMTMM/9, the highest percentage of crosslinks satisfied by a configuration is 88% in both ensembles (Fig. 5C, Fig. S2). However, in a few input cases, wall-EASAL ensembles satisfy more crosslinks. For example, in roca_putc/DSSO/2, the highest percentage of crosslinks satisfied by an IMP configuration is only 50%, whereas the wall-EASAL configurations satisfy all the crosslinks (Fig. 5D, Fig. S2).

A sole exception is 2b42/DMTMM/10, in which an IMP configuration satisfies more crosslinks (100%) than a wall-EASAL configuration (90%) (Fig. 5E, Fig. S2). In this case, wall-EASAL performed slightly inferior to IMP primarily because of Wall-EASAL experiments' coarse sampling. As pointed out earlier, the resulting pocket or disconnectivity artifacts in the collision-free configuration space cause wall-EASAL to miss wall regions whose projection on the chosen Cayley parameters (*i.e.*, inter-monomer residue-pair distances) is narrower than sampling step size.

For IMP, we observe that increasing the number of crosslinks improves the performance of integrative docking. The more input crosslinks, the higher the percentage of crosslinks satisfied per configuration, as shown by the shift in the blue distributions to higher crosslink percentages (Fig. S2A-S2C). This could indicate that, in randomized sampling guided by restraints, increasing the quantity of input data facilitates the sampling of more good-scoring configurations (configurations consistent with the input data). In contrast, wall-EASAL's performance is largely independent of the number of crosslinks.

Average crosslink distance

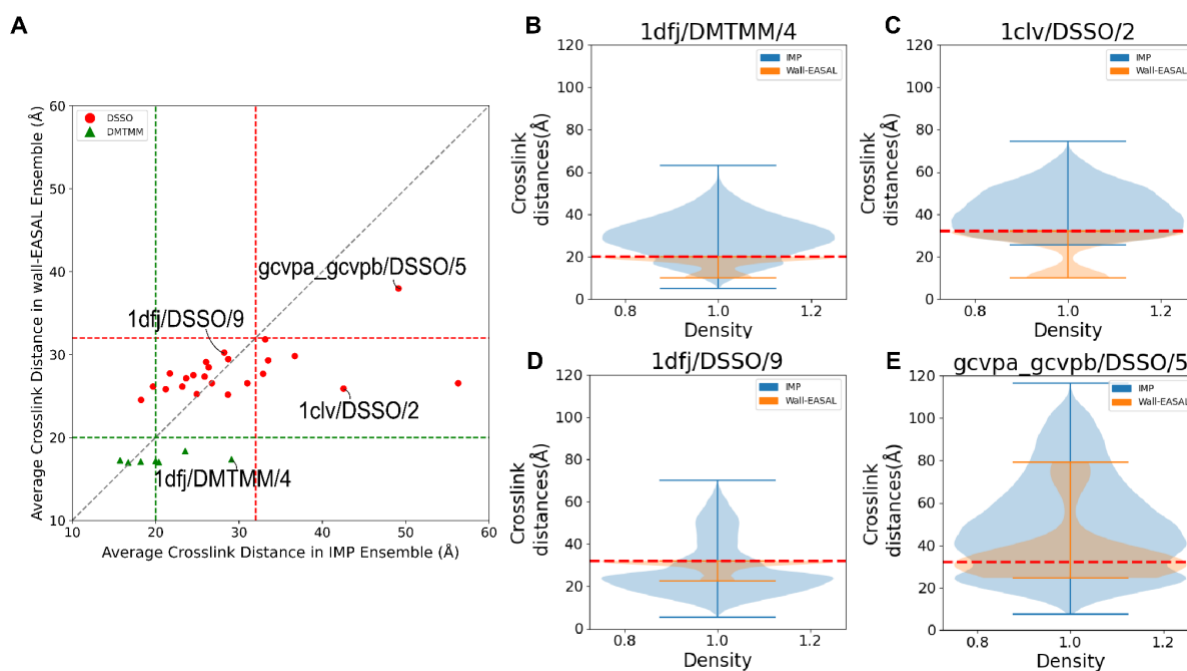


Figure 6. Distances between crosslinked residues in wall-EASAL and IMP structures. (A) For the 23 cases with DSSO crosslinker (red circles) and 7 cases with DMTMM crosslinker (green triangles), the average distances between the crosslinked residues were computed across all crosslinks and all configurations in the wall-EASAL and IMP ensembles. The crosslink upper bound (red and green dashed line) was set to 32 Å (20 Å) for DSSO (DMTMM) crosslinks. **(B-E)** The distribution of the distance between crosslinked residues in the two ensembles is shown for four cases.

Next, we computed the distances between the crosslinked residues (crosslink distances) in the two ensembles. The wall-EASAL configurations satisfy crosslinks better in terms of the crosslink distances (Fig. 6, Fig. S3). In 29 (19) cases, the average crosslink distances were within the upper bounds in the wall-EASAL (IMP) ensembles (Fig. 6A). All wall-EASAL crosslink distances were usually within the upper bounds; in contrast, a fraction of the IMP configurations violated the crosslinks in all cases (Fig. S3). The range of crosslink distances is smaller for wall-EASAL and larger for IMP (height of violin plots, Fig. S3).

In many input cases, the average crosslink distance was within the upper bound for both ensembles, e.g., 1dfj/DSSO/9 (Fig. 6B). However, for some cases, such as 1clv/DSSO/2 and 1dfj/DMTMM/4, the average crosslink distance exceeded the upper bound in the IMP, but not in the wall-EASAL ensemble (Fig. 6C-6D, Fig. S3). Finally, there were a small number of input cases, such as gcvpa_gcvpb/DSSO/5, where the average crosslink distance was much higher than the upper bound in both ensembles (Fig. 6E).

In most wall-EASAL configurations, the crosslink distances were close to the upper bound, especially as the number of crosslinks increased (Fig. S3A-S3C). This is unsurprising as wall-EASAL explicitly samples and finds only configurations where at least one crosslink distance takes an extreme value at one of the endpoints of its allowed interval.

Finally, for IMP, the crosslink distances reduce with an increase in the number of crosslinks, consistent with the decrease in docking difficulty with an increase in the number of crosslinks as shown earlier (Fig. S3A-S3C).

Comparison of crosslink distances: sampled configuration vs native structure

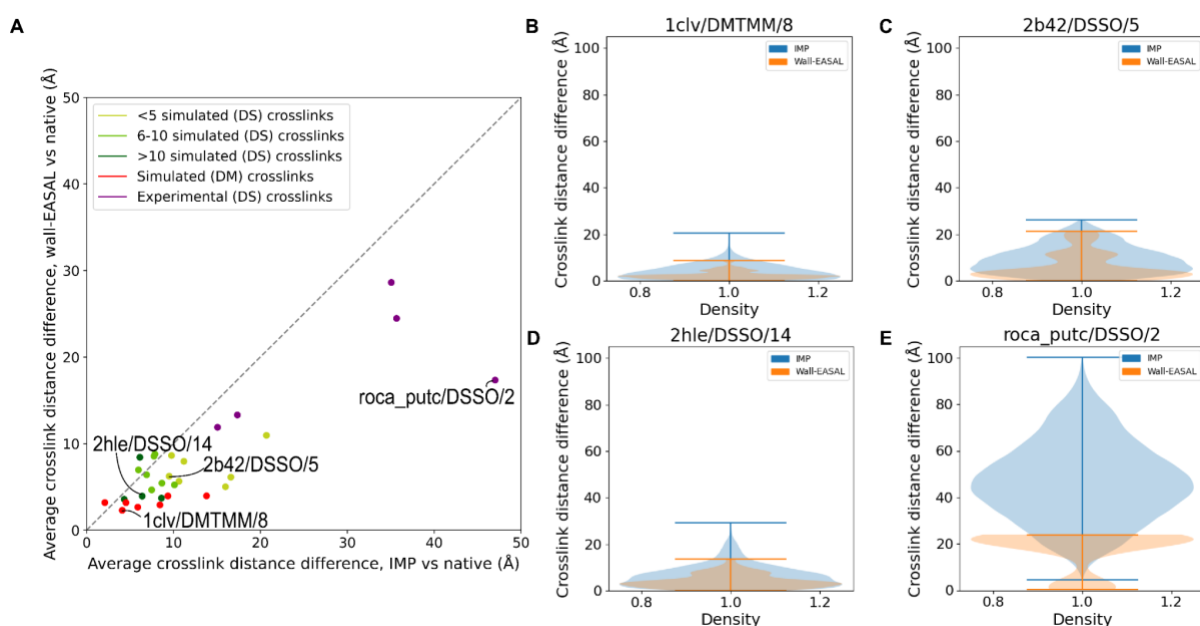


Figure 7. Comparison between crosslink distances in the sampled configurations and the native structure. (A) The average difference between the crosslink distance in a sampled configuration and the crosslink distance in the native structure was compared for the IMP and wall-EASAL ensembles. (B-E) The difference in crosslink distance between the native structure and each configuration in the two ensembles for four cases. DS and DM refer to DSSO and DMTMM crosslinks, respectively.

Further, we compared the crosslink distances in the wall-EASAL and IMP configurations with the corresponding distances in the native structure. The crosslink distances in the wall-EASAL configurations were closer to those in the native structure, implying that the wall-EASAL ensemble contained more near-native configurations (Fig. 7A, Fig. S4). The differences between crosslink distances were more variable for the IMP ensemble (Fig. 7B-7E, Fig. S4).

For most input cases with simulated DMTMM crosslinks, the crosslink distances from both the wall-EASAL and IMP configurations were close to those in the native structure, e.g., 1clv/DMTMM/8 (Fig. 7B). For input cases with simulated DSSO crosslinks, the distances from the wall-EASAL configurations were closer to those in the native structure, e.g., 2b42/DSSO/5 and 2hle/DSSO/14 (Fig. 7C-7D). The differences in crosslink distances in both sets of configurations were larger for the input cases with crosslinks from experiments and AF2-predicted monomer structures, e.g., roca_putc/DSSO/2, implying that these cases were more difficult for both methods (Fig. 7E). An increase in the number of crosslinks was associated with smaller distance differences in the IMP configurations, consistent with the earlier trends (Fig. S4A-S4C).

Accuracy of the IMP and wall-EASAL configurations

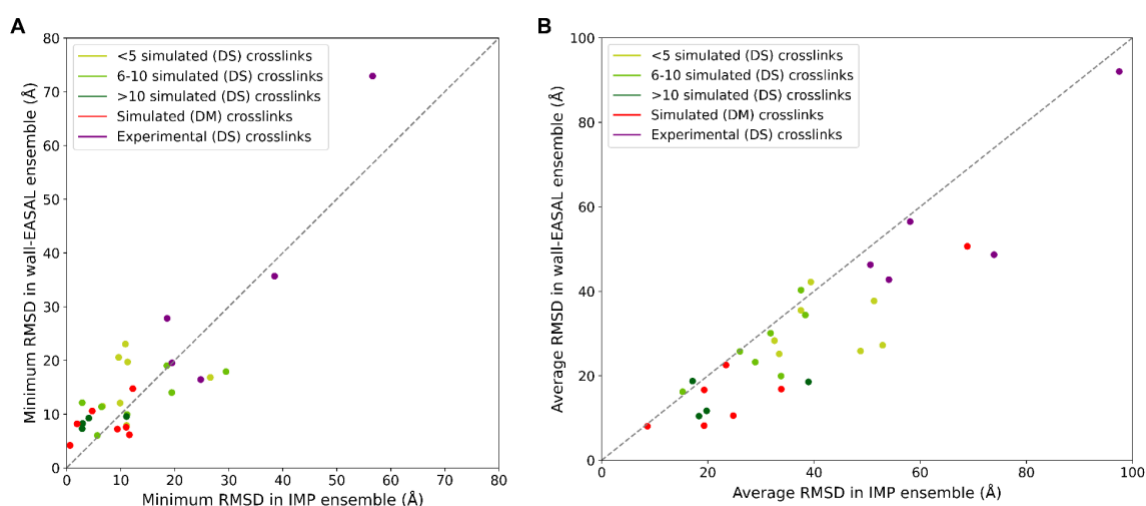


Figure 8. RMSD of the wall-EASAL and IMP configurations. (A) Minimum ligand RMSD and **(B)** Average ligand RMSD of a configuration in the IMP and wall-EASAL ensemble to the native structure. DS and DM refer to DSSO and DMTMM crosslinks, respectively.

We also computed the ligand RMSDs (Root-Mean-Square Deviation) of the configurations in the IMP and wall-EASAL ensembles with respect to the corresponding native structures (Lensink et al., 2023; Wodak et al., 2023). Both methods performed similarly in recovering a single near-native configuration in the ensemble; however, the configurations in the wall-EASAL ensemble were closer to the native structure, on average (Fig. 8, Fig. S5). For 13 (12) input cases for IMP (wall-EASAL), the ligand RMSD of the best configuration was within 10 Å of the native structure (acceptable by CAPRI standards) (Fig. 8A). The average ligand RMSD

to the native structure was 30 Å (38 Å) for wall-EASAL (IMP) (Fig. 8B). Configurations from both ensembles had higher RMSDs for the input cases with crosslinks from experiments and AF2-predicted monomer structures, e.g., phes_phet/DSSO/5, consistent with the higher docking difficulty for these cases noted earlier (Fig. S5E).

Efficiency of IMP and wall-EASAL

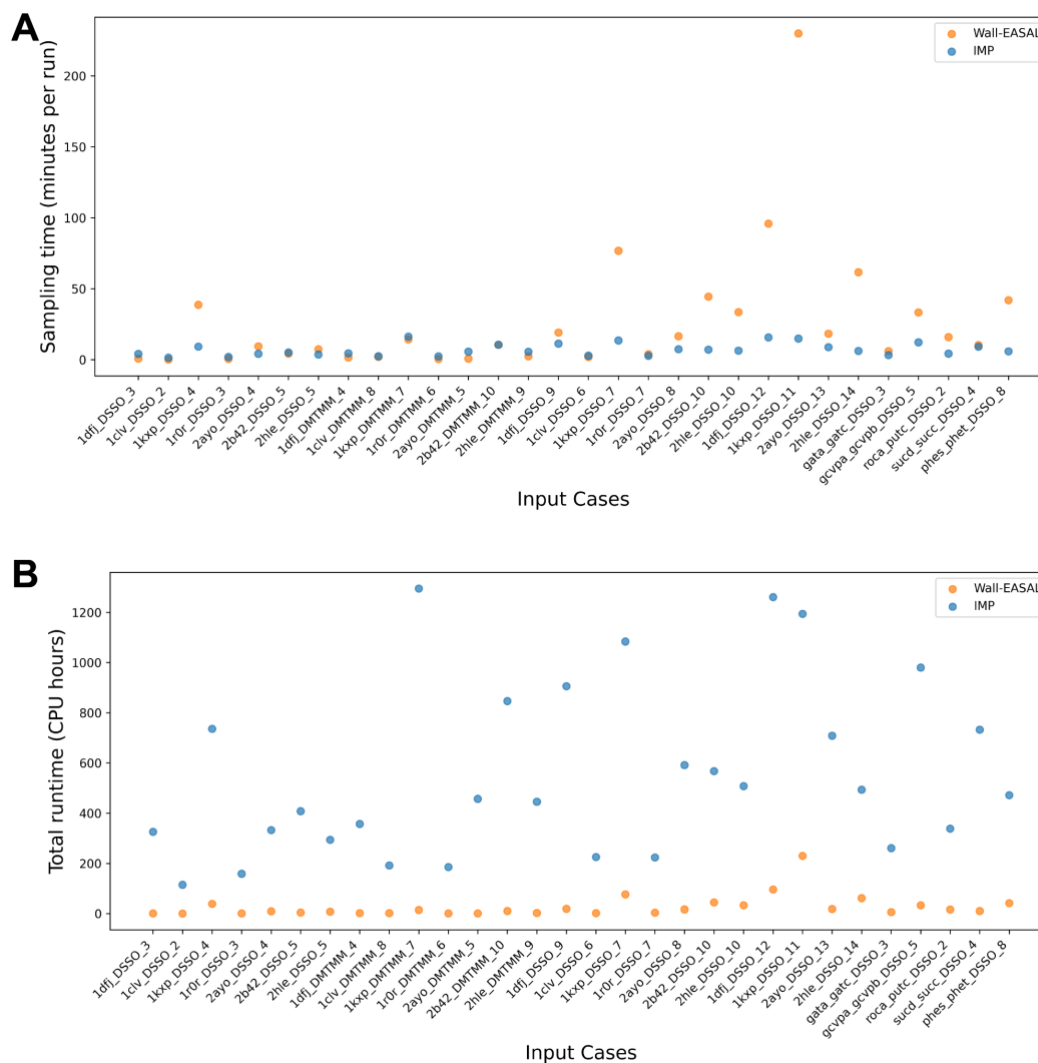


Figure 9. Time efficiency of IMP and wall-EASAL. The distribution of sampling times in terms of the **(A)** average time in CPU minutes per run and **(B)** average total number of CPU hours across the thirty benchmark cases for both methods. All times are on an AMD Ryzen Threadripper 3990x 64-core processor with 256 GB RAM and 2.2 GHz clock speed.

IMP relies on randomized sampling that requires multiple independent runs starting from random initial configurations (Pasani & Viswanath, 2021; Russel et al., 2012; D. Saltzberg et al., 2019; D. J. Saltzberg et al., 2021). In contrast, wall-EASAL is a deterministic method that

requires a single run per input. The sampling time per independent run for IMP is lower than the time for a wall-EASAL run (Fig. 9A). However, the total sampling time for IMP, in terms of the number of CPU hours, is much higher (Fig. 9B). Moreover, the total runtime for IMP will include the time for analysis, which will add another 25% to the sampling time. Therefore, wall-EASAL is more efficient than IMP.

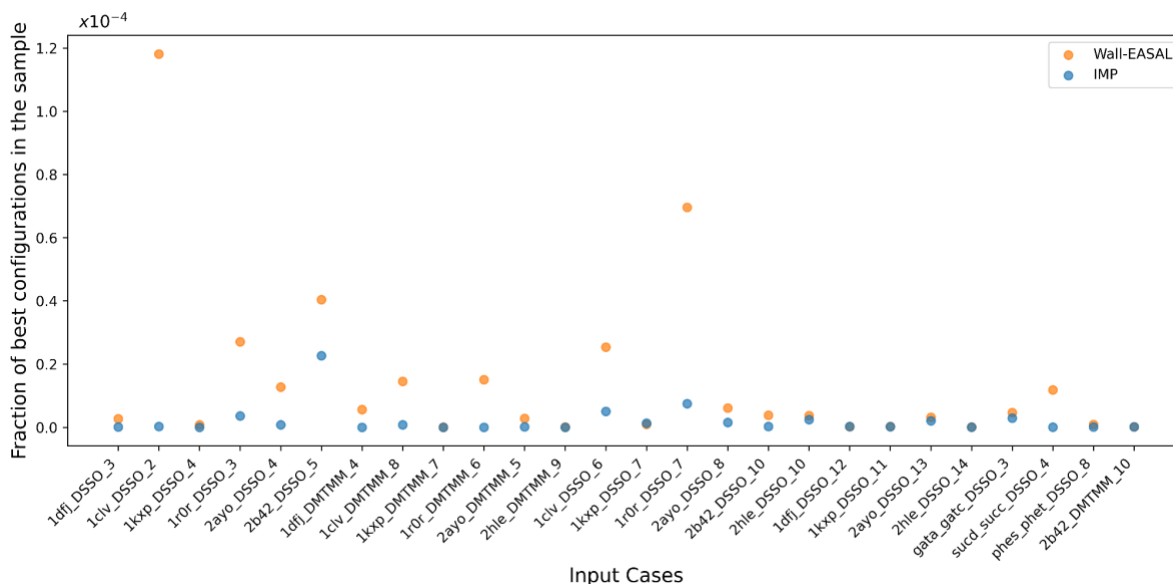


Figure 10. Sampling efficiency of IMP and wall-EASAL. Fraction of configurations in an ensemble with the maximum number of crosslinks satisfied among the total configurations sampled by IMP and EASAL.

We also compared the efficiency of wall-EASAL and IMP in terms of the number of samples required to obtain the structures that satisfy the input crosslinks sufficiently well. Efficiency was defined by the fraction of structures in the respective ensembles that satisfied the most crosslinks. This comparison was performed for all the cases where the highest percentage of crosslinks satisfied by a single structure in the ensemble was the same for IMP and wall-EASAL. As a general rule, IMP requires many more samples than wall-EASAL to obtain the same maximum crosslink satisfaction as the latter (Fig. 10).

For example, in 2hle/DMTMM/9, although the highest percentage of crosslinks satisfied by a configuration is 88% in both ensembles (Fig. 5D, Fig. S2), the fraction of IMP samples that satisfy the maximum number of crosslinks is only 7/100 of the same fraction for wall-EASAL, *i.e.*, the sampling efficiency of wall-EASAL is superior. However, there are exceptions to this rule. For example, in 2b42/DMTMM/10, an IMP configuration satisfies more crosslinks (100%) than a wall-EASAL configuration (90%) (Fig. 5E, Fig. S2), and yet the fraction of wall-EASAL samples that satisfy the maximum number is about the same as for IMP.

Visualization of structures

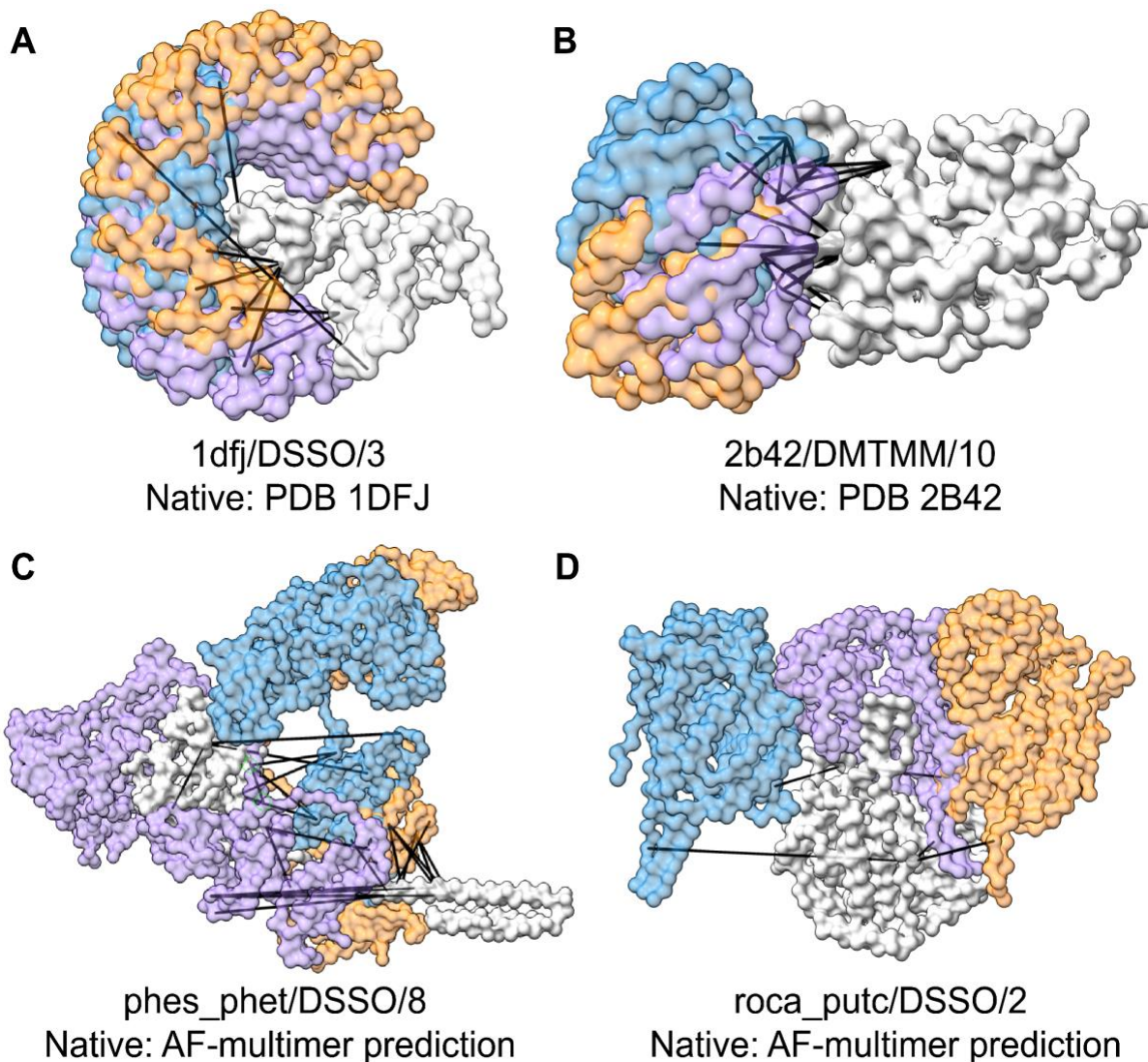


Figure 11. Visualization of wall-EASAL and IMP configurations. The best IMP and wall-EASAL configurations (least ligand RMSD to the native) are superposed on the native structure. **(A-D)** The sampled configuration and native structure are superposed on the receptor (light gray); the ligands in the native structure (purple), wall-EASAL configuration (orange), and IMP configuration (blue) are shown for representative input cases. Crosslinks in the wall-EASAL, IMP, and native configuration are shown by the black lines.

Finally, for four input cases, we visualized the best wall-EASAL and IMP configurations, as defined by the configuration with the least ligand RMSD to the corresponding native structure (Fig. 11). For 1dfj/DSSO/3, the best configurations in both the ensembles are similar to the native structure, consistent with our earlier observations that these configurations satisfy the crosslinks well (Fig. 11A). For 2b42/DMTMM/10, the IMP configuration was closer to the native structure, consistent with the higher crosslink satisfaction observed in the IMP configurations (Fig. 11B, Fig. 5E).

For phes_phet/DSSO/8, both the configurations have large ligand RMSDs from the native structure, *i.e.*, the corresponding AF-multimer prediction. This is intriguing, given that both the ensembles satisfy 87.5% of the crosslinks (Fig. 11C, Fig. S2). This discrepancy might arise because half of the inter-protein crosslinks were violated in the AF-multimer prediction of the phes_phet complex (O'Reilly et al., 2023). In this case, our integrative structures are more consistent with the data. For roca_putc/DSSO/2, the wall-EASAL configuration was closer to the native structure, *i.e.*, the corresponding AF-multimer prediction, compared to the IMP configuration, consistent with the earlier observation that the wall-EASAL configurations satisfy the crosslinks better for this case; all crosslinks are also satisfied in the AF-multimer prediction of roca_putc (Fig. 11D, Fig. 5D).

Discussion

Here, we developed wall-EASAL, a new method for integrative modeling of binary protein-protein complexes given the atomic structures of the constituent proteins and inter-protein chemical crosslinks. The method is based on an efficient discrete geometry algorithm for roadmapping and sampling distance-constrained configurational spaces using distance-based parameterization for dimension reduction and convexification. On a benchmark of thirty input cases, we compared the performance of wall-EASAL with IMP, an integrative modeling method based on randomized sampling. The configurations from wall-EASAL satisfy the crosslinks better as well as resemble the corresponding native structures more closely, on average. Wall-EASAL is also, in general, more efficient than IMP with respect to both measures of efficiency we considered, although there are exceptions to this rule.

Here, we discuss the advantages, uses, limitations, and future directions for integrative docking using wall-EASAL. On the examined benchmark, the method was efficient and produced ensembles that satisfy the input crosslinks well and were close to the native structure. The sampling of configuration space was also demonstrated to be representative, by comparison to the regular EASAL (interior sampling) method for a few input cases. Surprisingly, although at least one crosslink in wall-EASAL's sample configurations is guaranteed to take an extreme value, in contrast to IMP, wall-EASAL's distribution of crosslink distances is usually narrower than the IMP distribution of crosslink distances (Fig. 6, Fig. S3).

The assumptions in the method are that the constituent proteins are docked rigidly, their atomic structures are known, and they can be represented by spherical beads coarse-grained at the residue level. Currently, the method as implemented is applicable only for docking of pairs of proteins. In fact, restriction to a protein pair is not a theoretical or algorithmic requirement for the EASAL methodology. The theory behind EASAL encompasses many monomers and many rigid components, some of which could belong to the same monomer and the corresponding algorithmic extension is given in (Prabhu et al., 2020), awaiting implementation. Finally, the crosslink restraint is implemented by a simple distance constraint with upper and lower distance bounds. It does not account for uncertainties in the crosslinking experiment, such as false positive crosslinks (D. Saltzberg et al., 2019; Schneidman-Duhovny et al., 2014; Shi et al., 2014). However, the deterministic distance interval constraint checks can immediately be made probabilistic according to any given noise distribution. Although this could be viewed as more realistic modeling, there is no theoretical guarantee that the deterministic model (more consistent with the Occam's Razor principle of modeling) would be any less accurate or efficient than the more elaborate probabilistic model.

Here, we discuss parameters that may need to be tuned to improve the performance of wall-EASAL. First, the sampling in wall-EASAL can be made finer (coarser) by reducing (increasing) the step size. We observed that wall-EASAL can get solutions with a step size of 20 (*'stepSize'*), however, for a few input cases using a lower step size of 5 was required to find the best configurations. There is a trade-off between the step size and the sampling time. The sampling time could increase up to 2.5 fold upon reducing the step size by half. Second, one may need to decrease the crosslink satisfaction tolerance (*'crossLinkSatisfyThres'*) value if configurations that satisfy the specified number of crosslinks are not found.

It is conceivable to introduce walls at intermediate distances (*'smartCrosslinkMode'*) in the interior of the crosslink interval to provide a larger configurational space for wall-EASAL to sample configurations from. However, additional walls in the interior would result in a corresponding decrease in efficiency with questionable returns since we have demonstrated that in the absence of pockets or links or their coarse sampling artifacts, wall-EASAL guarantees a feasible solution satisfying all (or the maximum possible number) of the crosslinks, and moreover provides a representative collection of configurations of the entire feasible region including the interior.

Any information on distances between the residues or domains of constituent proteins can be used in wall-EASAL; although the current study uses chemical crosslinks, other types of distances can also be used, for example from NMR or genetic interaction assays (Echeverria

et al., 2023). Structures of constituent proteins can be derived from experiments or AI-based predictions (Abramson et al., 2024). The method may be of particular interest in cases where the structures of constituent proteins have not been experimentally determined, but reliable Alphafold predictions of the monomer are available, along with crosslinks (Bartolec et al., 2023; McCafferty et al., 2023; O'Reilly et al., 2023). Structures of antibody-antigen complexes are also of special interest since AI-based predictions of these complexes are not currently reliable (Ambrosetti et al., 2023; T. Cohen et al., 2023; Giulini et al., 2023).

The structures of binary complexes from wall-EASAL can complement methods for integrative modeling of macromolecular assemblies. For instance, these structures can suggest rigid bodies or restraints on protein interfaces for use in IMP, Assemblin, or Haddock (Dominguez et al., 2003; Honorato et al., 2024; Rantos et al., 2022; Russel et al., 2012; D. J. Saltzberg et al., 2021). Such information on pairs of proteins can then be combined with other information to model a larger complex. The structures from wall-EASAL can also be used as inputs for methods that perform combinatorial searches for structures of large assemblies based on the component binary complexes, such as CombFold and MCTreeSearch (Bryant et al., 2022; Chim & Elofsson, 2024; Shor & Schneidman-Duhovny, 2024).

New geometric deep learning methods that predict the optimal distance between crosslinked residues can be used to further refine the inputs to methods such as wall-EASAL (S. Cohen & Schneidman-Duhovny, 2023). Future planned extensions of the method include parallelizing it for efficiency and modifying the algorithm to scale to larger constituent proteins and protein complexes with multiple subunits. Bayesian formulations of restraints can be used instead of simple distance restraints (Shi et al., 2014). The method can be extended to include restraints other than pairwise distance restraints, such as EM-based shape restraints, by devising ways to convert them to equivalent distance restraints when possible. Incorporating wall-EASAL in integrative modeling methods such as IMP will facilitate the characterization of assemblies and cellular neighborhoods at increased efficiency, accuracy, and precision.

Data and software availability

The implementation of the wall-EASAL method is available at <https://bitbucket.org/geoplexity/easal-dev/src/Crosslink>. The integrative docking benchmark is available at https://github.com/isblab/Integrative_docking_benchmark. The benchmark is also available at Zenodo: <https://doi.org/10.5281/zenodo.13959115>.

Supporting information

Supporting information contains figures showing the input structures (Fig. S1), percentage of crosslink satisfaction (Fig. S2), average crosslink distance (Fig. S3), crosslink distance difference in the sampled configurations and the native structure (Fig. S4), and RMSD of the wall-EASAL and IMP configurations (Fig. S5) in each benchmark case. Table S1 contains the description of the benchmark. The Mathematical proof that wall-EASAL finds a feasible configuration satisfying crosslink constraints if one exists is also given.

Acknowledgments

We thank ISB Lab members Shreyas Arvindkar, Kartik Majila, Omkar Golatkar, and Mubashira KP for their useful comments on the manuscript. Molecular graphics images were produced using the UCSF Chimera and UCSF ChimeraX packages from the Resource for Biocomputing, Visualization, and Informatics at the University of California, San Francisco (supported by NIH P41 RR001081, NIH R01-GM129325, and National Institute of Allergy and Infectious Diseases). The authors acknowledge University of Florida's UFIT Research Computing for providing the Hipergator computational resources that have contributed to the research results reported in this publication.

Funding

This work has been supported by the following grants: Department of Atomic Energy (DAE) TIFR grant RTI 4006, Department of Science and Technology (DST) SERB grant SPG/2020/000475, and Department of Biotechnology (DBT) BT/PR40323/BTIS/137/78/2023 from the Government of India to S.V.

The work has also been partially supported by National Science Foundation grants NSF DMS-1563234 and DMS-1564480 awarded to M.S.

Conflicts of Interest declaration

None declared.

References

- Abramson, J., Adler, J., Dunger, J., Evans, R., Green, T., Pritzel, A., Ronneberger, O., Willmore, L., Ballard, A. J., Bambrick, J., Bodenstein, S. W., Evans, D. A., Hung, C.-C., O'Neill, M., Reiman, D., Tunyasuvunakool, K., Wu, Z., Žemgulytė, A., Arvaniti, E., ... Jumper, J. M. (2024). Accurate structure prediction of biomolecular interactions with AlphaFold 3. *Nature*, 630(8016), 493–500. <https://doi.org/10.1038/s41586-024-07487-w>
- Ambrosetti, F., Jandova, Z., & Bonvin, A. M. J. J. (2023). Information-Driven Antibody–Antigen Modelling with HADDOCK. In K. Tsumoto & D. Kuroda (Eds.), *Computer-Aided Antibody Design* (Vol. 2552, pp. 267–282). Springer US. https://doi.org/10.1007/978-1-0716-2609-2_14
- Arvindekar, S., Jackman, M. J., Low, J. K. K., Landsberg, M. J., Mackay, J. P., & Viswanath, S. (2022). Molecular architecture of nucleosome remodeling and deacetylase sub-complexes by integrative structure determination. *Protein Science*, 31(9), e4387. <https://doi.org/10.1002/pro.4387>
- Arvindekar, S., Majila, K., & Viswanath, S. (2024). *Recent methods from statistical inference and machine learning to improve integrative modeling of macromolecular assemblies* (Version 4). arXiv. <https://doi.org/10.48550/ARXIV.2401.17894>
- Bartolec, T. K., Vázquez-Campos, X., Norman, A., Luong, C., Johnson, M., Payne, R. J., Wilkins, M. R., Mackay, J. P., & Low, J. K. K. (2023). Cross-linking mass spectrometry discovers, evaluates, and corroborates structures and protein–protein interactions in the human cell. *Proceedings of the National Academy of Sciences*, 120(17), e2219418120. <https://doi.org/10.1073/pnas.2219418120>
- Beck, M., Covino, R., Hänel, I., & Müller-McNicoll, M. (2024). Understanding the cell: Future views of structural biology. *Cell*, 187(3), 545–562. <https://doi.org/10.1016/j.cell.2023.12.017>
- Braitbard, M., Schneidman-Duhovny, D., & Kalisman, N. (2019). Integrative Structure Modeling: Overview and Assessment. *Annual Review of Biochemistry*, 88(1), 113–135. <https://doi.org/10.1146/annurev-biochem-013118-111429>

- Bryant, P., Pozzati, G., Zhu, W., Shenoy, A., Kundrotas, P., & Elofsson, A. (2022). Predicting the structure of large protein complexes using AlphaFold and Monte Carlo tree search. *Nature Communications*, *13*(1), 6028. <https://doi.org/10.1038/s41467-022-33729-4>
- Bullock, J. M. A., Schwab, J., Thalassinos, K., & Topf, M. (2016). The Importance of Non-accessible Crosslinks and Solvent Accessible Surface Distance in Modeling Proteins with Restraints From Crosslinking Mass Spectrometry. *Molecular & Cellular Proteomics*, *15*(7), 2491–2500. <https://doi.org/10.1074/mcp.M116.058560>
- Bullock, J. M. A., Sen, N., Thalassinos, K., & Topf, M. (2018). Modeling Protein Complexes Using Restraints from Crosslinking Mass Spectrometry. *Structure*, *26*(7), 1015–1024.e2. <https://doi.org/10.1016/j.str.2018.04.016>
- Chim, H. Y., & Elofsson, A. (2024). MoLPC2: Improved prediction of large protein complex structures and stoichiometry using Monte Carlo Tree Search and AlphaFold2. *Bioinformatics*, *40*(6), btae329. <https://doi.org/10.1093/bioinformatics/btae329>
- Cohen, S., & Schneidman-Duhovny, D. (2023). A deep learning model for predicting optimal distance range in crosslinking mass spectrometry data. *PROTEOMICS*, *23*(17), 2200341. <https://doi.org/10.1002/pmic.202200341>
- Cohen, T., Halfon, M., Carter, L., Sharkey, B., Jain, T., Sivasubramanian, A., & Schneidman-Duhovny, D. (2023). Multi-state modeling of antibody-antigen complexes with SAXS profiles and deep-learning models. In *Methods in Enzymology* (Vol. 678, pp. 237–262). Elsevier. <https://doi.org/10.1016/bs.mie.2022.11.003>
- Dominguez, C., Boelens, R., & Bonvin, A. M. J. J. (2003). HADDOCK: A protein-protein docking approach based on biochemical or biophysical information. *Journal of the American Chemical Society*, *125*(7), 1731–1737. <https://doi.org/10.1021/ja026939x>
- Echeverria, I., Braberg, H., Krogan, N. J., & Sali, A. (2023). Integrative structure determination of histones H3 and H4 using genetic interactions. *The FEBS Journal*, *290*(10), 2565–2575. <https://doi.org/10.1111/febs.16435>
- Ferber, M., Kosinski, J., Ori, A., Rashid, U. J., Moreno-Morcillo, M., Simon, B., Bouvier, G., Batista, P. R., Müller, C. W., Beck, M., & Nilges, M. (2016). Automated structure modeling of large protein assemblies using crosslinks as distance restraints. *Nature Methods*, *13*(6), 515–520. <https://doi.org/10.1038/nmeth.3838>
- Giulini, M., Schneider, C., Cutting, D., Desai, N., Deane, C. M., & Bonvin, A. M. J. J. (2023). *Towards the accurate modelling of antibody-antigen complexes from sequence using machine learning and information-driven docking* (p. 2023.11.17.567543). bioRxiv. <https://doi.org/10.1101/2023.11.17.567543>
- Gong, Z., Ye, S.-X., & Tang, C. (2020). Tightening the Crosslinking Distance Restraints for Better Resolution of Protein Structure and Dynamics. *Structure*, *28*(10), 1160–

- 1167.e3. <https://doi.org/10.1016/j.str.2020.07.010>
- Graziadei, A., & Rappsilber, J. (2022). Leveraging crosslinking mass spectrometry in structural and cell biology. *Structure*, *30*(1), 37–54.
<https://doi.org/10.1016/j.str.2021.11.007>
- Guest, J. D., Vreven, T., Zhou, J., Moal, I., Jeliazkov, J. R., Gray, J. J., Weng, Z., & Pierce, B. G. (2021). An expanded benchmark for antibody-antigen docking and affinity prediction reveals insights into antibody recognition determinants. *Structure*, *29*(6), 606-621.e5. <https://doi.org/10.1016/j.str.2021.01.005>
- Honorato, R. V., Trellet, M. E., Jiménez-García, B., Schaarschmidt, J. J., Giulini, M., Reys, V., Koukos, P. I., Rodrigues, J. P. G. L. M., Karaca, E., Van Zundert, G. C. P., Roel-Touris, J., Van Noort, C. W., Jandová, Z., Melquiond, A. S. J., & Bonvin, A. M. J. J. (2024). The HADDOCK2.4 web server for integrative modeling of biomolecular complexes. *Nature Protocols*. <https://doi.org/10.1038/s41596-024-01011-0>
- Koukos, P. I., & Bonvin, A. M. J. J. (2020). Integrative Modelling of Biomolecular Complexes. *Journal of Molecular Biology*, *432*(9), 2861–2881.
<https://doi.org/10.1016/j.jmb.2019.11.009>
- Lensink, M. F., Brysbaert, G., Raouraoua, N., Bates, P. A., Giulini, M., Honorato, R. V., Van Noort, C., Teixeira, J. M. C., Bonvin, A. M. J. J., Kong, R., Shi, H., Lu, X., Chang, S., Liu, J., Guo, Z., Chen, X., Morehead, A., Roy, R. S., Wu, T., ... Wodak, S. J. (2023). Impact of AlphaFold on structure prediction of protein complexes: The CASP15-CAPRI experiment. *Proteins: Structure, Function, and Bioinformatics*, *91*(12), 1658–1683.
<https://doi.org/10.1002/prot.26609>
- Liu, X., Zhang, Y., Wen, Z., Hao, Y., Banks, C. A. S., Cesare, J., Bhattacharya, S., Arvindekar, S., Lange, J. J., Xie, Y., Garcia, B. A., Slaughter, B. D., Unruh, J. R., Viswanath, S., Florens, L., Workman, J. L., & Washburn, M. P. (2024). An integrated structural model of the DNA damage-responsive H3K4me3 binding WDR76:SPIN1 complex with the nucleosome. *Proceedings of the National Academy of Sciences*, *121*(33), e2318601121. <https://doi.org/10.1073/pnas.2318601121>
- McCafferty, C. L., Pennington, E. L., Papoulas, O., Taylor, D. W., & Marcotte, E. M. (2023). Does AlphaFold2 model proteins' intracellular conformations? An experimental test using cross-linking mass spectrometry of endogenous ciliary proteins. *Communications Biology*, *6*(1), 421. <https://doi.org/10.1038/s42003-023-04773-7>
- Merkley, E. D., Rysavy, S., Kahraman, A., Hafen, R. P., Daggett, V., & Adkins, J. N. (2014). Distance restraints from crosslinking mass spectrometry: Mining a molecular dynamics simulation database to evaluate lysine–lysine distances. *Protein Science*, *23*(6), 747–759. <https://doi.org/10.1002/pro.2458>
- O'Reilly, F. J., Graziadei, A., Forbrig, C., Bremenkamp, R., Charles, K., Lenz, S., Elfmann,

- C., Fischer, L., Stülke, J., & Rappsilber, J. (2023). Protein complexes in cells by AI-assisted structural proteomics. *Molecular Systems Biology*, 19(4), e11544. <https://doi.org/10.15252/msb.202311544>
- Ozkan, A., Prabhu, R., Baker, T., Pence, J., Peters, J., & Sitharam, M. (2018). Algorithm 990: Efficient Atlasing and Search of Configuration Spaces of Point-Sets Constrained by Distance Intervals. *ACM Transactions on Mathematical Software*, 44(4), 1–30. <https://doi.org/10.1145/3204472>
- Ozkan, A., & Sitharam, M. (2011). EASAL (Efficient Atlasing, Analysis and Search of Molecular Assembly Landscapes). *Proceedings of the ISCA 3rd International Conference on Bioinformatics and Computational Biology*, 233–238.
- Ozkan, A., & Sitharam, M. (2014). *Best of Both Worlds: Uniform sampling in Cartesian and Cayley Molecular Assembly Configuration Space* (arXiv:1409.0956). arXiv. <http://arxiv.org/abs/1409.0956>
- Ozkan, A., Sitharam, M., Flores-Canales, J. C., Prabhu, R., & Kurnikova, M. (2021). Baseline Comparisons of Complementary Sampling Methods for Assembly Driven by Short-Ranged Pair Potentials toward Fast and Flexible Hybridization. *Journal of Chemical Theory and Computation*, 17(3), 1967–1987. <https://doi.org/10.1021/acs.jctc.0c00945>
- Pasani, S., Menon, K. S., & Viswanath, S. (2023). *The molecular architecture of the desmosomal outer dense plaque by integrative structural modeling* [Preprint]. Biophysics. <https://doi.org/10.1101/2023.06.13.544884>
- Pasani, S., & Viswanath, S. (2021). A Framework for Stochastic Optimization of Parameters for Integrative Modeling of Macromolecular Assemblies. *Life*, 11(11), 1183. <https://doi.org/10.3390/life11111183>
- Prabhu, R., Sitharam, M., Ozkan, A., & Wu, R. (2020). Atlasing of Assembly Landscapes using Distance Geometry and Graph Rigidity. *Journal of Chemical Information and Modeling*, 60(10), 4924–4957. <https://doi.org/10.1021/acs.jcim.0c00763>
- Rantos, V., Karius, K., & Kosinski, J. (2022). Integrative structural modeling of macromolecular complexes using Assemble. *Nature Protocols*, 17(1), Article 1. <https://doi.org/10.1038/s41596-021-00640-z>
- Rappsilber, J. (2011). The beginning of a beautiful friendship: Cross-linking/mass spectrometry and modelling of proteins and multi-protein complexes. *Journal of Structural Biology*, 173(3), 530–540. <https://doi.org/10.1016/j.jsb.2010.10.014>
- Rout, M. P., & Sali, A. (2019). Principles for Integrative Structural Biology Studies. *Cell*, 177(6), 1384–1403. <https://doi.org/10.1016/j.cell.2019.05.016>
- Russel, D., Lasker, K., Webb, B., Velázquez-Muriel, J., Tjioe, E., Schneidman-Duhovny, D., Peterson, B., & Sali, A. (2012). Putting the Pieces Together: Integrative Modeling

- Platform Software for Structure Determination of Macromolecular Assemblies. *PLoS Biology*, 10(1), e1001244. <https://doi.org/10.1371/journal.pbio.1001244>
- Sadjadi, M., Hagh, V. F., Kang, M., Sitharam, M., Connelly, R., Gortler, S. J., Theran, L., Holmes-Cerfon, M., & Thorpe, M. F. (2021). Realizations of Isostatic Material Frameworks. *Physica Status Solidi (b)*, 258(9), 2000555. <https://doi.org/10.1002/pssb.202000555>
- Saltzberg, D., Greenberg, C. H., Viswanath, S., Chemmama, I., Webb, B., Pellarin, R., Echeverria, I., & Sali, A. (2019). Modeling Biological Complexes Using Integrative Modeling Platform. *Methods in Molecular Biology (Clifton, N.J.)*, 2022, 353–377. https://doi.org/10.1007/978-1-4939-9608-7_15
- Saltzberg, D. J., Viswanath, S., Echeverria, I., Chemmama, I. E., Webb, B., & Sali, A. (2021). Using Integrative Modeling Platform to compute, validate, and archive a model of a protein complex structure. *Protein Science: A Publication of the Protein Society*, 30(1), 250–261. <https://doi.org/10.1002/pro.3995>
- Schneidman-Duhovny, D., Pellarin, R., & Sali, A. (2014). Uncertainty in integrative structural modeling. *Current Opinion in Structural Biology*, 28, 96–104. <https://doi.org/10.1016/j.sbi.2014.08.001>
- Schneidman-Duhovny, D., Rossi, A., Avila-Sakar, A., Kim, S. J., Velázquez-Muriel, J., Strop, P., Liang, H., Krukenberg, K. A., Liao, M., Kim, H. M., Sobhanifar, S., Dötsch, V., Rajpal, A., Pons, J., Agard, D. A., Cheng, Y., & Sali, A. (2012). A method for integrative structure determination of protein-protein complexes. *Bioinformatics*, 28(24), 3282–3289. <https://doi.org/10.1093/bioinformatics/bts628>
- Shi, Y., Fernandez-Martinez, J., Tjioe, E., Pellarin, R., Kim, S. J., Williams, R., Schneidman-Duhovny, D., Sali, A., Rout, M. P., & Chait, B. T. (2014). Structural Characterization by Cross-linking Reveals the Detailed Architecture of a Coatomer-related Heptameric Module from the Nuclear Pore Complex. *Molecular & Cellular Proteomics*, 13(11), 2927–2943. <https://doi.org/10.1074/mcp.M114.041673>
- Shor, B., & Schneidman-Duhovny, D. (2024). CombFold: Predicting structures of large protein assemblies using a combinatorial assembly algorithm and AlphaFold2. *Nature Methods*, 21(3), 477–487. <https://doi.org/10.1038/s41592-024-02174-0>
- Sitharam, M., & Gao, H. (2010). Characterizing Graphs with Convex and Connected Cayley Configuration Spaces. *Discrete & Computational Geometry*, 43(3), 594–625. <https://doi.org/10.1007/s00454-009-9160-8>
- Sitharam, M., St. John, A., & Sidman, J. (2019). *Handbook of geometric constraint systems principles*. CRC press Taylor & Francis group.
- Sitharam, M., & Wang, M. (2014). How the Beast really moves: Cayley analysis of mechanism realization spaces using CayMos. *Computer-Aided Design*, 46, 205–210.

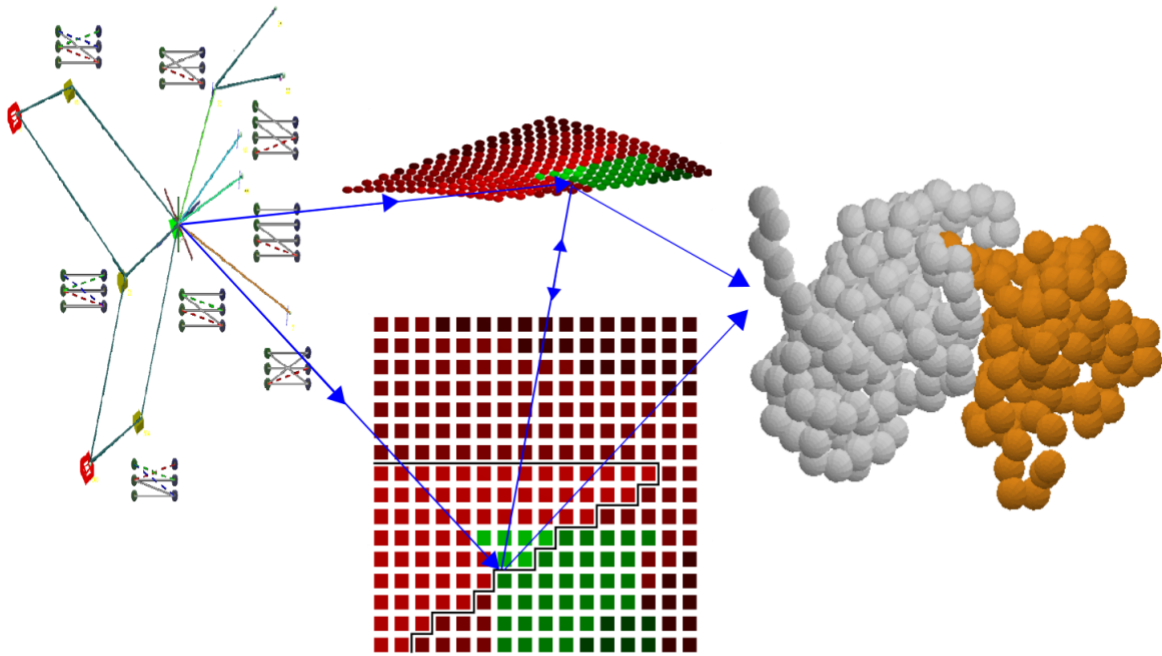
- <https://doi.org/10.1016/j.cad.2013.08.033>
- Stahl, K., Demann, L., Bremenkamp, R., Warneke, R., Hormes, B., Stülke, J., Brock, O., & Rappsilber, J. (2024). *Modelling protein complexes with crosslinking mass spectrometry and deep learning* (p. 2023.06.07.544059). bioRxiv.
<https://doi.org/10.1101/2023.06.07.544059>
- Stahl, K., Graziadei, A., Dau, T., Brock, O., & Rappsilber, J. (2023). Protein structure prediction with in-cell photo-crosslinking mass spectrometry and deep learning. *Nature Biotechnology*, 1–10. <https://doi.org/10.1038/s41587-023-01704-z>
- Viswanath, S., Chemmama, I. E., Cimermancic, P., & Sali, A. (2017). Assessing Exhaustiveness of Stochastic Sampling for Integrative Modeling of Macromolecular Structures. *Biophysical Journal*, 113(11), 2344–2353.
<https://doi.org/10.1016/j.bpj.2017.10.005>
- Wang, M., & Sitharam, M. (2015). Algorithm 951: Cayley Analysis of Mechanism Configuration Spaces using CayMos: Software Functionalities and Architecture. *ACM Transactions on Mathematical Software*, 41(4), 1–8. <https://doi.org/10.1145/2699462>
- Wodak, S. J., Vajda, S., Lensink, M. F., Kozakov, D., & Bates, P. A. (2023). Critical Assessment of Methods for Predicting the 3D Structure of Proteins and Protein Complexes. *Annual Review of Biophysics*, 52(1), 183–206.
<https://doi.org/10.1146/annurev-biophys-102622-084607>
- Wu, R., Prabhu, R., Ozkan, A., & Sitharam, M. (2020). Rapid prediction of crucial hotspot interactions for icosahedral viral capsid self-assembly by energy landscape atlas validated by mutagenesis. *PLOS Computational Biology*, 16(10), e1008357.
<https://doi.org/10.1371/journal.pcbi.1008357>
- Yu, C., & Huang, L. (2023). New advances in cross-linking mass spectrometry toward structural systems biology. *Current Opinion in Chemical Biology*, 76, 102357.
<https://doi.org/10.1016/j.cbpa.2023.102357>
- Zhang, Y., & Sitharam, M. (2022). *Best of two worlds: Cartesian sampling and volume computation for high dimensional configuration spaces using Cayley coordinates*. paper 8417.
- Zhang, Y., & Sitharam, M. (2024). *Best of two worlds: Cartesian sampling and volume computation for distance-constrained configuration spaces using Cayley coordinates* (arXiv:2408.16946). arXiv. <http://arxiv.org/abs/2408.16946>
- Zhang, Y., Zhang, Z., Kagaya, Y., Terashi, G., Zhao, B., Xiong, Y., & Kihara, D. (2023). *Distance-AF: Modifying Predicted Protein Structure Models by Alphafold2 with User-Specified Distance Constraints* (p. 2023.12.01.569498). bioRxiv.
<https://doi.org/10.1101/2023.12.01.569498>
- Ziemianowicz, D. S., Saltzberg, D., Pells, T., Crowder, D. A., Schröder, C., Hepburn, M.,

Sali, A., & Schriemer, D. C. (2021). IMProv: A Resource for Cross-link-Driven Structure Modeling that Accommodates Protein Dynamics. *Molecular & Cellular Proteomics*, 20, 100139. <https://doi.org/10.1016/j.mcpro.2021.100139>

For Table of Contents Use Only

A new discrete-geometry approach for integrative docking of proteins using chemical crosslinks

Yichi Zhang, Muskaan Jindal, Shruthi Viswanath, and Meera Sitharam



Legend: Given input structures of two proteins and chemical crosslinks between them, Wall-EASAL optimizes sampling to guarantee an output set of representative configurations of the complex satisfying the input constraints.