
Harmonic Self-Conditioned Flow Matching for joint Multi-Ligand Docking and Binding Site Design

Hannes Stark¹ Bowen Jing¹ Regina Barzilay¹ Tommi Jaakkola¹

Abstract

A significant amount of protein function requires binding small molecules, including enzymatic catalysis. As such, designing binding pockets for small molecules has several impactful applications ranging from drug synthesis to energy storage. Towards this goal, we first develop HARMONICFLOW, an improved generative process over 3D protein-ligand binding structures based on our self-conditioned flow matching objective. FLOWSITE extends this flow model to jointly generate a protein pocket’s discrete residue types and the molecule’s binding 3D structure. We show that HARMONICFLOW improves upon state-of-the-art generative processes for docking in simplicity, generality, and average sample quality in pocket-level docking. Enabled by this structure modeling, FLOWSITE designs binding sites substantially better than baseline approaches.

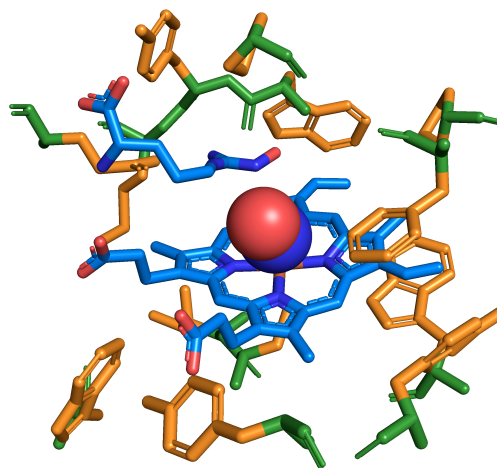


Figure 1. Binding site design. Given the backbone (green) and multi-ligand without structure, FLOWSITE generates residue types and structure (orange) to bind the multi-ligand and its jointly generated structure (blue). The majority of the pocket is omitted for visibility.

1. Introduction

Designing proteins that can bind small molecules has many applications, ranging from drug synthesis to energy storage or gene editing. Indeed, a key part of any protein’s function derives from its ability to bind and interact with other molecular species. For example, we may design proteins that act as antidotes that sequester toxins or design enzymes that enable chemical reactions through catalysis, which plays a major role in most biological processes. We develop FLOWSITE to address this design challenge by building on recent advances in deep learning (DL) based protein design (Dauparas et al., 2022) and protein-molecule docking (Corso et al., 2023).

Specifically, we aim to design protein pockets to bind a certain small molecule (called ligand). We assume that

¹CSAIL, Massachusetts Institute of Technology. Correspondence to: Hannes Stark <hstark@mit.edu>.

we are given a protein pocket via the 3D backbone atom locations of its residues as well as the 2D chemical graph of the ligand. We do not assume any knowledge of the 3D structure or the binding pose of the ligand. Based on this information, our goal is to predict the amino acid identities for the given backbone locations (see Figure 1). We also consider the more challenging task of designing pockets that simultaneously bind multiple molecules and ions (which we call multi-ligand). Such multi-ligand binding proteins are important, for example, in enzyme design, where the ligands correspond to reactants.

This task has not been addressed by deep learning yet. While deep learning has been successful in designing proteins that can bind to other proteins (Watson et al., 2023), designing (multi-)ligand binders is different and arguably harder in various aspects. For example, no evolutionary information is directly available, unlike when modeling interactions between amino acids only. The existing approaches, such as designing 6 drug binding proteins Polizzi & DeGrado (2020) or a single enzyme Yeh et al. (2023), build on expert knowledge and require manual steps. Therefore, we develop

FLowsITE as a more general and automated approach and the first deep learning solution for designing pockets that bind small molecules.

FLowsITE jointly generates discrete (residue identities) and continuous (ligand pose) variables. Our flow matching training criterion guides the model to learn a self-conditioned flow that jointly generates the contact residues and the (multi-)ligand 3D binding pose structures. To achieve this, we first develop HARMONICFLOW as a suitable generative process for 3D poses of (multi-)ligands. FLowsITE extends this process to residue types. Starting from initial residue types and ligand atom locations sampled from a harmonic prior FLowsITE updates them by iteratively following the learned vector field, as illustrated in Figure 2.

The HARMONICFLOW component of FLowsITE performs the task known as docking, i.e., it realizes the 3D binding pose of the multi-ligand. As a method, it is remarkably simple in comparison to existing generative processes for docking, including the state-of-the-art product-space diffusion process of DIFFDOCK (Corso et al., 2023) that operates on ligand’s center of mass, orientation, and torsion angles. HARMONICFLOW simply updates the cartesian coordinates of the atoms, yet manages to produce chemically plausible molecular structures without restricting ligand flexibility to torsions. Moreover, HARMONICFLOW outperforms product-space diffusion in average sample quality on multiple pocket-level docking tasks on PDBBind.

Having established HARMONICFLOW as an improved generative process over ligand positions, we extend it to include discrete residue types to obtain FLowsITE. We also adopt an additional "fake-ligand" data augmentation step where side chains are treated as ligands in order to realize additional training cases. Altogether, FLowsITE is able to recover 47.0% of binding site amino acids compared to 39.4% of a baseline approach. This nearly closes the gap to an oracle method (51.4% recovery) with access to the ground truth 3D structure/pose of the ligand. Next to technical innovations (self-conditioned flow matching or equivariant refinement TFN layers) our main contributions are:

1. FLowsITE as the first deep learning solution to design binding sites for small molecules without prior knowledge of the molecule structure.
2. The FLowsITE framework as a simple approach to jointly generate discrete and continuous data.
3. HARMONICFLOW which improves upon the state-of-the-art generative process for generating 3D ligand binding structures in average sample quality, simplicity, and applicability/generalizability.

2. Related Work

Deep learning for Docking. Designing binding sites with high affinity for a ligand requires reasoning about the binding free energy, which is deeply interlinked with modeling ligand binding 3D structures. This task of molecular docking has recently been tackled with deep-learning approaches (Stärk et al., 2022; Lu et al., 2022; Zhang et al., 2023a) including generative models (Corso et al., 2023; Qiao et al., 2023). These generative methods are based on diffusion models, building on DIFFDOCK (Corso et al., 2023), which combines diffusion processes over the ligand’s torsion angles and position with respect to the protein. For the task of multi-ligand docking, no deep learning solutions exist yet, and we provide the first with HARMONICFLOW. We refer to Appendix D for additional important related work on this and the following topics.

Protein Design. A significant technical challenge for protein design is jointly modeling structure and sequence. Inverse folding approaches (Dauparas et al., 2022; Gao et al., 2023a; Yi et al., 2023; Hsu et al., 2022; Gao et al., 2023b) attempt this by designing new sequences given a protein structure. Existing *ligand aware* inverse folding methods such as Carbonara (Krapp et al., 2023) and LigandMPNN (Dauparas et al., 2023) are limited in their applicability for small molecule binding site design since they assume the bound ligand structure to be provided as input instead of solving the hard problem of docking jointly with designing the backbone residues as achieved by FLowsITE.

The same limitation applies to the classical energy function and search algorithm based POCKETOPTIMIZER (Noske et al., 2023) and the sequence-structure co-design framework FAIR (Zhang et al., 2023b), which both require the bound complex as input. FAIR is further distinct from FLowsITE in that it uses all residue types of a protein as input except for those in contact with the ligand, which simplifies recovering the amino acids of known binders to infilling missing residues based on sequence similarity.

Flow Matching. This recent generative modeling paradigm (Lipman et al., 2022; Albergo & Vanden-Eijnden, 2022; Albergo et al., 2023) generalizes diffusion models (Ho et al., 2020; Song et al., 2021) in a simpler framework. Flow matching admits more design flexibility and multiple works (Tong et al., 2023b; Pooladian et al., 2023) showed how it enables learning flows between arbitrary start and end distributions in a simulation-free manner. It is easily extended to data on manifolds (Chen & Lipman, 2023) and comes with straighter paths that enable faster integration.

Other applications of flow matching to biomolecular problems include generating Boltzmann distributions of small molecules Klein et al. (2023), protein structure generation (Yim et al., 2023a; Bose et al., 2023) and small molecule

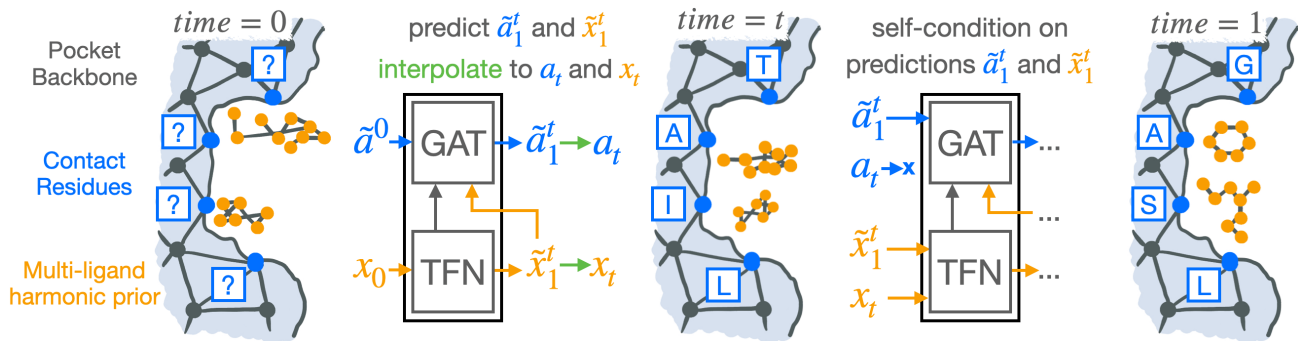


Figure 2. **Overview of FlowSite.** The generative process starts from a protein pocket’s backbone atoms, initial residue types \tilde{a}^0 , and initial ligand positions x_0 . Our joint discrete-continuous self-conditioned flow updates them to a_t, x_t by following its vector field defined by the model outputs $\tilde{a}_1^t, \tilde{x}_1^t$. This integration is repeated until reaching $time = 1$ with the produced sample a_1, x_1 .

generation (Song et al., 2023). We explain flow matching in Section 3.1.

3. Method

Our goal is to design binding pockets for a ligand where we assume the inputs to be the ligand’s 2D chemical graph and the backbone coordinates of the pocket’s residues. In this section, we lay out how FLOWSITE achieves this by first explaining our HARMONICFLOW generative process for docking in 3.1 before covering how FLOWSITE extends it to include discrete residue types in 3.2 and concluding with our model architecture in 3.3.

Overview and definitions. As visualized in Figure 2, FLOWSITE jointly updates discrete residue types and continuous ligand positions. The inputs are a protein pocket’s backbone atoms $y \in \mathbb{R}^{L \times 4 \times 3}$ for L residues with 4 atoms each and the chemical graph of a (multi-)ligand. Based on the ligand connectivity, its initial coordinates $x \in \mathbb{R}^{n \times 3}$ are sampled from a harmonic prior, and we initialize residue types $a \in \{1, \dots, 20\}^L$ with an initial token (we drop the chemical information of the ligands in our notation for brevity).

Given this at time $t = 0$, the flow model v_θ with learned parameters θ iteratively updates residue types and ligand coordinates by integrating the ODE it defines. These integration steps are repeated from time $t = 0$ to time $t = 1$ to obtain the final generated binding pocket designs.

3.1. HarmonicFlow Structure Generation

We first lay out HARMONICFLOW for pure structure generation without residue type estimation. Our notation drops v_θ ’s conditioning on the pocket y and residue estimates a in this subsection (see the Architecture Section 3.3 for how y is included). Simply put, HARMONICFLOW is flow

matching with a harmonic prior, self-conditioning, and x_1 prediction (our refinement TFN layers in Section 3.3 are also important for performance). In more detail:

Conditional Flow Matching. Given the data distribution p_1 of bound ligand structures and any easy-to-sample prior p_0 over $\mathbb{R}^{n \times 3}$, we wish to learn an ODE that pushes the prior forward to the data distribution when integrating it from time 0 to time 1. The ODE will be defined by a time-dependent vector field $v_\theta(\cdot, \cdot) : \mathbb{R}^{n \times 3} \times [0, 1] \mapsto \mathbb{R}^{n \times 3}$. Starting with a sample $x_0 \sim p_0(x_0)$ and following/integrating v through time will produce a sample from the data distribution p_1 .

To see how to train v_θ , let us first assume access to a time-dependent vector field $u_t(\cdot)$ that would lead to an ODE that pushes from the prior p_0 to the data p_1 (it is not straightforward how to construct this u_t). This gives rise to a probability path p_t by integrating u_t until time t . If we could sample $x \sim p_t(x)$ we could train v_θ with the unconditional flow matching objective (Lipman et al., 2022)

$$\mathcal{L}_{FM} = \mathbb{E}_{t \sim \mathcal{U}[0,1], x \sim p_t(x)} \|v_\theta(x, t) - u(x, t)\|^2. \quad (1)$$

Among others, Tong et al. (2023b) show that to construct such a u_t (that transports from prior p_0 to p_1), we can use samples from the data $x_1 \sim p_1(x_1)$ and prior $x_0 \sim p_0(x_0)$ and define u_t via

$$u_t(x) = \mathbb{E}_{x_1 \sim p_1(x_1), x_0 \sim p_0(x_0)} \frac{u_t(x|x_0, x_1)p_t(x|x_0, x_1)}{p_t(x)} \quad (2)$$

where we can choose easy-to-sample conditional flows $p_t(\cdot|\cdot, \cdot)$ that give rise to simple conditional vector fields $u_t(\cdot|\cdot, \cdot)$. We still cannot efficiently compute this $u_t(x)$ and use it in \mathcal{L}_{FM} because we do not know $p_t(x)$, but there is no

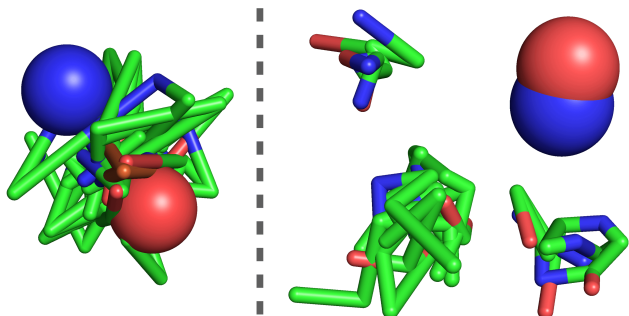


Figure 3. **Harmonic Prior.** Initial positions for the same single multi-ligand from an isotropic Gaussian (*left*) and from a harmonic prior (*right*). (Bound structure for this multi-ligand is in Figure 1).

need to: it is equivalent to instead train with the following conditional flow matching loss since $\nabla_{\theta} \mathcal{L}_{FM} = \nabla_{\theta} \mathcal{L}_{CFM}$.

$$\mathcal{L}_{CFM} = \mathbb{E}_{t \sim \mathcal{U}[0,1], \mathbf{x}_1 \sim p_1(\mathbf{x}_1), \mathbf{x}_0 \sim p_0(\mathbf{x}_0), x \sim p_t(x|\mathbf{x}_0, \mathbf{x}_1)} \|v_{\theta}(\mathbf{x}, t) - u_t(\mathbf{x}|\mathbf{x}_0, \mathbf{x}_1)\|^2. \quad (3)$$

Our simple choice of conditional probability path is $p_t(\mathbf{x}|\mathbf{x}_0, \mathbf{x}_1) = \mathcal{N}(\mathbf{x}|t\mathbf{x}_1 + (1-t)\mathbf{x}_0, \sigma^2)$, which gives rise to the conditional vector field $u_t(\mathbf{x}|\mathbf{x}_0, \mathbf{x}_1) = \mathbf{x}_1 - \mathbf{x}_0$. Notably, we find it helpful to parameterize v_{θ} to predict \mathbf{x}_1 instead of $(\mathbf{x}_1 - \mathbf{x}_0)$.

Training with the conditional flow matching loss then boils down to 1) Sample data $\mathbf{x}_1 \sim p_1(\mathbf{x}_1)$ and prior $\mathbf{x}_0 \sim p_0(\mathbf{x}_0)$. 2) Interpolate between between the points. 3) Add noise to the interpolation to obtain x . 4) Evaluate and minimize $\mathcal{L}_{CFM} = \|v_{\theta}(\mathbf{x}, t) - \mathbf{x}_1\|^2$ with it. **Inference** is just as straightforward. We sample from the prior $\mathbf{x}_0 \sim p_0(\mathbf{x}_0)$ and integrate from $t = 0$ to $t = 1$ with an arbitrary ODE solver. We use an Euler solver, i.e., we iteratively predict \mathbf{x}_1 as $\tilde{\mathbf{x}}_1 = v_{\theta}(\mathbf{x}_t, t)$, and then calculate the step size scaled velocity estimate from it and add it to the current point $\mathbf{x}_{t+\Delta t} = \mathbf{x}_t + \Delta t(\tilde{\mathbf{x}}_1 - \mathbf{x}_0)$. Training and inference algorithms are in Appendix A.3.

Harmonic Prior. Any prior can be used for p_0 in the flow matching framework. We choose a harmonic prior as in EigenFold (Jing et al., 2023) that samples atoms to be close to each other if they are connected by a bond. Potentially, this inductive bias is especially helpful when dealing with multiple molecules and ions since atoms of different molecules are already spatially separated at $t = 0$ as visualized in Figure 3.

This prior is constructed based on the ligand’s covalent bonds that define a graph with adjacency matrix \mathbf{A} from which we can construct the graph Laplacian $\mathbf{L} = \mathbf{D} - \mathbf{A}$ where \mathbf{D} is the degree matrix. The harmonic prior is then

$p_0(\mathbf{x}_0) \propto \exp(-\frac{1}{2}\mathbf{x}_0^T \mathbf{L} \mathbf{x}_0)$ which can be sampled as a transformed gaussian.

Structure Self-conditioning. With this, we aim to bring AlphaFold2’s (Jumper et al., 2021) successful recycling strategy to flow models for structure generation. Recycling enables training a deeper structure predictor without additional memory cost by performing multiple forward passes while only computing gradients for the last. For flow matching, we achieve the same by adapting the discrete diffusion model self-conditioning approach of Chen et al. (2023), similar to self-conditioning in protein structure generation diffusion models (Yim et al., 2023b; Watson et al., 2023).

Instead of defining the vector field $v_{\theta}(\mathbf{x}_t, t)$ as a function of \mathbf{x}_t and t alone, we additionally condition it on the prediction $\tilde{\mathbf{x}}_1^t$ of the previous integration step and use $v_{\theta}(\mathbf{x}_t, \tilde{\mathbf{x}}_1^t, t)$. At the beginning of **inference** the self-conditioning input is a sample from the harmonic prior $\tilde{\mathbf{x}}_1^0 \sim p_0(\tilde{\mathbf{x}}_1^0)$. In all following steps, it is the flow model’s output (its prediction of \mathbf{x}_1) of the previous step $\tilde{\mathbf{x}}_1^t = v_{\theta}(\mathbf{x}_{t-\Delta t}, \tilde{\mathbf{x}}_1^{t-\Delta t}, t - \Delta t)$. To **train** this, in a random 50% of the training steps, the self-conditioning input is a sample from the prior $\tilde{\mathbf{x}}_1^0$. In the other 50%, we first generate a self-conditioning input $\tilde{\mathbf{x}}_1^{t+\Delta t} = v_{\theta}(\mathbf{x}_t, \tilde{\mathbf{x}}_1^0, t)$, detach it from the gradient computation graph, and then use $v_{\theta}(\mathbf{x}_t, \tilde{\mathbf{x}}_1^{t+\Delta t}, t)$ for the loss computation. Algorithms 3 and 4 show these training and inference procedures.

3.2. FlowSite Binding Site Design

In the FLOWSITE binding site design framework, HARMONICFLOW $\tilde{\mathbf{x}}_1^{t+\Delta t} = v_{\theta}(\mathbf{x}_t, \tilde{\mathbf{x}}_1^t, t)$ is augmented with an additional self-conditioned flow over the residue types to obtain $(\tilde{\mathbf{x}}_1^{t+\Delta t}, \tilde{\mathbf{a}}_1^{t+\Delta t}) = v_{\theta}(\mathbf{x}_t, \tilde{\mathbf{x}}_1^t, \mathbf{a}_t, \tilde{\mathbf{a}}_1^t, t)$. The flow no longer produces $\tilde{\mathbf{x}}_1^{t+\Delta t}$ as an estimate of \mathbf{x}_1 and then interpolates to $\mathbf{x}_{t+\Delta t}$ but instead produces $(\tilde{\mathbf{x}}_1^{t+\Delta t}, \tilde{\mathbf{a}}_1^{t+\Delta t})$ from which we obtain the interpolation $(\mathbf{x}_{t+\Delta t}, \mathbf{a}_{t+\Delta t})$ and use it for the next integration step (see Figure 4). The start $\mathbf{a}_0, \tilde{\mathbf{a}}_1^0$ are initialized as a mask token while the structures $\mathbf{x}_0, \tilde{\mathbf{x}}_1^0$ are drawn from a harmonic prior.

This joint discrete-continuous data process is trained with the same self-conditioning strategy as in **structure self-conditioning**, but with the additional discrete self-conditioning input $\tilde{\mathbf{a}}_1^t$ that is either a model output or a mask token. To the training loss we add the cross-entropy \mathcal{L}_{type} between \mathbf{a} and $\tilde{\mathbf{a}}_1^t$. In practice, we find that the \mathbf{a}_1 prediction $\tilde{\mathbf{a}}_1^t$ already carries most information that is useful for predicting \mathbf{a}_1 and we omit the interpolation \mathbf{a}_t as model input to obtain the simpler $(\tilde{\mathbf{x}}_1^{t+\Delta t}, \tilde{\mathbf{a}}_1^{t+\Delta t}) = v_{\theta}(\mathbf{x}_t, \tilde{\mathbf{x}}_1^t, \tilde{\mathbf{a}}_1^t, t)$. This formulation admits a direct interpretation as recycling (Jumper et al., 2021) and a clean joint discrete-continuous process without defining a discrete data interpolation.

Fake Ligand Data Augmentation. This strategy is based

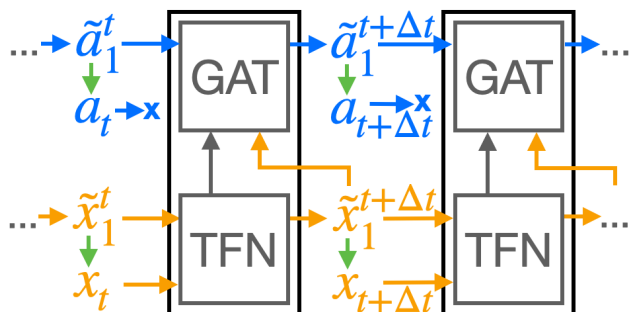


Figure 4. **FlowSite self-conditioned updates.** Residue type predictions \tilde{a}_1^t from invariant GAT layers and position predictions \tilde{x}_1^t from equivariant TFN layers are used as self-conditioning inputs and to *interpolate* to the updates a_t, x_t .

on the evidence of Polizzi & DeGrado (2020) that a protein’s sidechain-sidechain interactions are similar to sidechain-ligand interactions for tight binding. In our optional data augmentation, we train with 20% of the samples having a “fake ligand”. Given a protein, we construct a fake ligand as the atoms of a randomly selected residue that has at least 4 other residues within 4Å heavy atom distance. Additionally, we modify the protein by removing the residue that was chosen as the fake ligand and the residues that are within 7 positions next to that residue in the protein chain (Figure 5). This data augmentation was also employed in concurrent work (Dauparas et al., 2023; Corso et al., 2024).

3.3. Architecture

Here, we provide an overview of the FLOWSITE architecture (visualized in Appendix Figure 6) that outputs ligand positions \tilde{x}_1 and uses them for a residue type prediction \tilde{a}_1 . The structure prediction is produced by a stack of our SE(3)-equivariant refinement TFN layers that are crucial for the performance of HARMONICFLOW’s structure generation. This is followed by invariant layers to predict the invariant residue types, which we found to perform better than the more expressive equivariant layers that seem well suited for structure prediction but not for recovering residue types. The precise architecture definition is in Appendix A.5 and an architecture visualization in Figure 6.

Radius Graph Representation. We represent the (multi-)ligand and the protein as graphs where nodes are connected based on their distances. Each protein residue and each ligand atom is a node. These are connected by protein-to-protein edges, ligand-to-ligand edges, and edges between ligand and protein. While only a single node is assigned to each residue, they contain information about all backbone atom positions (N, Ca, C, O).

Equivariant refinement TFN layers. Based on Tensor

Field Networks (TFN) (Geiger et al., 2020), these layers are a simple yet effective tweak from previous use cases of message passing TFNs (Jing et al., 2022; Corso et al., 2023), where we instead update and refine ligand coordinates with each layer akin to EGNNs (Hoogeboom et al., 2022). See Appendix A.5 for more details.

The k -th refinement TFN layer takes as input the protein positions y , current ligand positions x_t , and features h^{k-1} (with h^0 being zeros for the ligand and vectors between N, Ca, C, O for the protein). We construct equivariant messages for each edge via a tensor-product of neighboring nodes’ invariant and equivariant features. The messages include the *structure self-conditioning* information by using the interatomic distances of the self-conditioning input x_1^t to parameterize the tensor products. We sum the messages to obtain new node features h^{k+1} and use them as input to an O(3) equivariant linear layer to predict intermediate refined ligand coordinates \hat{x}_1^k . Before passing \hat{x}_1^k to the next refinement TFN layer, we detach them from the gradient computation graph for the non-differentiable radius graph building of the next layer.

After a stack of K TFN refinement layers, the positions \hat{x}_1^K are used as final prediction $\tilde{x}_1^{t+\Delta t}$. While $\tilde{x}_1^{t+\Delta t}$ is supervised with the conditional flow matching loss $\mathcal{L}_{CFM} = \|\tilde{x}_1^{t+\Delta t} - x_1\|^2$ the intermediate positions \hat{x}_1^k contribute to an additional refinement loss $\mathcal{L}_{refine} = \sum_{k=1}^{K-1} \|\hat{x}_1^k - x_1\|^2$.

Invariant Network. The inputs to this part of FLOWSITE are the TFN’s ligand structure prediction \tilde{x}_1 , the protein structure y , the invariant scalar features of the refinement TFN layers, and the self-conditioning input a_1^t . From the protein structure, we construct on PiFold’s (Gao et al., 2023a) distance-based invariant edge features and node features that encode the geometry of the backbone. For the edges between protein and ligand, we construct features that encode the distances from a ligand atom to all 4 backbone atoms of a connected residue.

These are processed by a stack of graph attention layers that update ligand and protein node features as well as edge features for each type of edge (ligand-to-ligand, protein-to-protein, and between the molecules). For each edge, the convolutional layers first predict attention weights from the edge features and the features of the nodes they connect. We then update a node’s features by summing messages from each incoming edge weighted by the attention weights. Then, we update an edge’s features based on its nodes’ new features. A precise definition is in Appendix A.5. From the residue features after a stack of these convolutions, we predict new residue types $a_{t+\Delta t}$ together with side chain torsion angles α . We use those in an auxiliary loss $\mathcal{L}_{torsion}$ defined as in AlphaFold2’s Appendix 1.9.1 (Jumper et al., 2021). Thus, the complete loss for FLOWSITE is a weighted sum of $\mathcal{L}_{CFM}, \mathcal{L}_{refine}, \mathcal{L}_{type}$, and $\mathcal{L}_{torsion}$, while HAR-

MONICFLOW only uses \mathcal{L}_{CFM} and \mathcal{L}_{refine} .

4. Experiments

We evaluate FLOWSITE with multiple datasets, splits, and settings. Every reported number is averaged over 10 generated samples for each ligand. Precise experimental details are in Appendix E and code to reproduce each experiment is at <https://github.com/HannesStark/FlowSite>. The main questions we seek to answer with the experiments are:

- **1. Structure Generation:** How does HARMONICFLOW compare with SOTA binding structure generative models? Does it work for *multi-ligands*?
- **2. Binding Site recovery:** What is the improvement in recovering residue types of known binding sites with FLOWSITE compared with baseline approaches?
- **3. Ablations and Flow Matching Investigations:** How does flow matching behave for structure generation, and how much does each component help? (self-conditioning, ...)

4.1. Datasets

We use **PDBBind** version 2020 with 19k complexes to evaluate the structure generation capability of flow matching and the ability of FLOWSITE to design binders for a single connected ligand. We employ two dataset splits. The first is based on time, which has been heavily used in the DL community (Stärk et al., 2022; Corso et al., 2023). The second is sequence-based with a maximum of 30% chain-wise similarity between train, validation, and test data. Buttenschoen et al. (2023) found DL docking methods to be significantly more challenged by sequence similarity splits.

For many binding pocket design tasks, it is required to bind multi-ligands. For example, when designing enzymes for multiple reactants. Such multi-ligands are present in **Binding MOAD**. We use its 41k complexes with a 30% sequence similarity split carried out as described above. We construct our *multi-ligands* as all molecules and ions that have atoms within 4Å of each other. An example of an enzyme with all substrates in the pocket as multi-ligand is in Figure 1.

4.2. Question 1: HarmonicFlow Structure Generation

Here, we consider the HARMONICFLOW component of FLOWSITE and investigate its binding structure generation capability. This is to find out whether HARMONICFLOW is fit for binder design where good structure generation is necessary to take the bound ligand structure into account.

Task Setup. The architecture only contains refinement TFN layers, and there is no sequence prediction. The inputs are the (multi-)ligand’s chemical graph and the protein pocket’s backbone atoms and residue types (see Appendix Table 7 for experiments without residue type inputs). From this, the binding structure of the (multi-)ligand has to be inferred. There is also no *fake ligand augmentation*.

We test docking on the pocket level since that is the structure modeling capability required for the binding site design task (in Appendix C, we show preliminary results for docking to the whole protein). We define the binding pocket in two ways. In the *Distance-Pocket* definition, we calculate the distances of all ligand heavy atoms to the protein’s alpha carbons, add Gaussian noise with $\sigma = 0.5$ to them, and include a residue in the pocket if it has any noisy distance smaller than 14 Å. The center of the pocket (where the prior p_0 is centered) is the center of mass of all residues within 8 Å of noisy ligand distance. We additionally add Gaussian noise with $\sigma = 0.2$ to the pocket center. The motivation for the noisy distance cutoffs and pocket center is to alleviate distribution shifts during inference time and to prevent the models from inferring the ligand positions from the cutoff with which the pocket was constructed.

In *Radius-Pockets*, we first obtain the center of mass of residues within 8 Å of any ligand heavy atom. The pocket includes all residues with a noisy distance to the center of mass that is less than a specific radius. This radius is 7Å plus the minimum of 5Å and half of the ligand’s diameter. The pocket center is obtained in the same way as for *Distance-Pockets*.

Baseline. We compare with the state-of-the-art product-space diffusion process of DIFFDOCK (Corso et al., 2023) which has recently also been proven successful for pocket level docking (Plainer et al., 2023). *Note that this is not the full DIFFDOCK docking pipeline:* Both HARMONICFLOW and DIFFDOCK’s diffusion can generate multiple samples and, for the task of docking, a further discriminator (called confidence model in DIFFDOCK) could be used to select the most likely poses. We only compare the 3D structure generative models and neither use language model residue embeddings. We train product-space diffusion with our pocket definitions using 5 of its default TFN layers followed by its pseudotorque and center-convolution (the TFN layers are identical to ours apart from our position updates). This uses their training parameters and the same 32 scalar and 8 vector features of our model. We use this deep learning method for comparison since (like HARMONICFLOW) it is able to dock to protein structures without an all-atom representation (unlike traditional docking methods, e.g., Autodock VINA (Trott & Olson, 2010)). This is required for the pocket design task where the sidechain atom locations are unknown.

Table 1. **HARMONICFLOW vs. PRODUCT SPACE DIFFUSION.** Comparison on PDBBind splits for docking into *Distance-Pockets* (residues close to ligand) and *Radius-Pockets* (residues within a radius of the pocket center). The columns "%<2" show the fraction of predictions with an RMSD to the ground truth that is less than 2Å (higher is better). "Med." is the median RMSD (lower is better).

Method	Sequence Similarity Split				Time Split			
	Distance-Pocket		Radius-Pocket		Distance-Pocket		Radius-Pocket	
	%<2	Med.	%<2	Med.	%<2	Med.	%<2	Med.
PRODUCT SPACE DIFFUSION	27.2	3.2	16.1	4.0	20.8	3.8	15.2	4.3
HARMONICFLOW	30.1	3.1	20.5	3.4	42.8	2.5	28.3	3.2

Table 2. **Multi-Ligand Docking.** Structure generation performance on Binding MOAD’s *multi-ligands*. "%<2" means the fraction of predictions with an RMSD to the ground truth less than 2Å (higher better). "Med." is the median RMSD (lower better).

Method	%<2	%<5	Med.
EIGENFOLD DIFFUSION	39.7	73.5	2.4
HARMONICFLOW	44.4	75.0	2.2

PDBBind docking results. In Table 1, we find that our flow matching based HARMONICFLOW outperforms product-space diffusion in average sample quality. The sampled conformations in Figure 12 show that HARMONICFLOW produces chemically plausible structures and well captures the physical constraints of interatomic interactions without the need to restrict conformational flexibility to torsion angles. No separate losses, rotation updates, or expensive torsion angle updates are required - HARMONICFLOW is arguably a cleaner and simpler solution. Thus, it is a promising future direction to further explore HARMONICFLOW for docking and other biomolecular structure generation tasks.

Binding MOAD multi-ligand docking results. For binding site design, it is often necessary to model multiple ligands and ions (e.g., reactants for an enzyme). We test this with Binding MOAD, which contains such multi-ligands. Since no deep learning solutions for multi-ligands exist yet and traditional docking methods would require side-chain atom locations, we compare with EIGENFOLD’s (Jing et al., 2023) Diffusion and provide qualitative evaluation in Appendix Figure 12. For EIGENFOLD DIFFUSION, we use the same model as HARMONICFLOW and predict x_0 (in what corresponds to x_0 prediction in diffusion models), which we found to work better. Table 2 shows HARMONICFLOW as viable for docking multi-ligands - thus, the first ML method for this task with important applications besides binding site design.

4.3. Question 2: FlowSite Binding Site Recovery

Setup. The input to FLOWSITE is the binding pocket/site specified by its backbone and the chemical identity of the ligand (without its 3D structure). With the pocket, in practice, chosen by the user and well known, we use the *Distance-*

Pocket definition here.

Metrics. We use two metrics, sequence recovery and our *BLOSUM score*. Sequence recovery is the percentage of generated residue types of the contact residues (those with a heavy atom within 4Å of the ligand) that are the same as in the original binding site. This metric only rewards exact matches, it cannot account for multiple correct solutions, and there is no notion of amino acid similarity encoded in it. Having lower penalties for mismatches of very similar amino acids would be more meaningful. To address this, we propose *BLOSUM score*, which takes evolutionary and physicochemical similarity into account. A precise definition is in Appendix A.1.

We note that metrics that correlate with binding affinity, such as MM-PBSA or docking scores of traditional docking software, usually require the atomic structure of the side chains as input, which is not available. Furthermore, they were only validated for discriminating different ligands, while in our design task, we desire to discriminate different binding sites.

Baselines. PIFOLD (*no ligand*) is the architecture of Gao et al. (2023a) and does not use any ligand information and PROTEINMPNN (*no ligand*) is the analog with the architecture of (Dauparas et al., 2022). In PIFOLD (*2D ligand*), we first process the ligand with PNA (Corso et al., 2020) message passing and pass its features as additional input to the PIFOLD architecture. GROUND TRUTH POS and RANDOM LIGAND POS use the architecture of FLOWSITE without the ligand structure prediction layers. Instead, the ligand positions are either the ground truth bound structure or sampled from a standard Normal at the pocket’s alpha carbon center of mass. Similarly, DIFFDOCK-POCKET POS uses fixed positions and the same architecture, but the positions are given by the pocket-level docking tool DiffDock-Pocket (Plainer et al., 2023) (we only have results for single ligands since it is not implemented for multi-ligands). The oracle GROUND TRUTH POS method also uses fake ligand data augmentation.

Pocket Recovery Results. Table 3 shows that FLOWSITE consistently is able to recover the original pocket better than simpler treatments of the (multi-)ligand, closing the gap

Table 3. **Binding Site Recovery.** Comparison on PDBBind and Binding MOAD sequence similarity splits for recovering residues of binding sites. *Recovery* is the percentage of correctly predicted residues, and *BLOSUM score* takes residue similarity into account. 2D ligand refers to a simple GNN encoding of the ligand’s chemical graph as additional input. The GROUND TRUTH POS row has access to the, in practice, unknown ground truth 3D crystal structure of the ligand and protein.

Method	Binding MOAD		PDBBind	
	<i>BLOSUM score</i>	<i>Recovery</i>	<i>BLOSUM score</i>	<i>Recovery</i>
PROTEINMPNN (no ligand)	–	–	40.3	36.3
PIFOLD (no ligand)	35.2	39.4	40.7	43.5
PIFOLD (2D ligand)	35.7	40.4	42.2	44.5
RANDOM LIGAND POS	38.2	41.8	41.5	43.7
DIFFDOCK-POCKET POS	–	–	42.6	45.0
FLowsITE	44.3	47.0	47.6	49.5
GROUND TRUTH POS	48.4	51.4	51.3	51.2

Table 4. **Flow matching investigation.** Ablations of HARMONICFLOW. Column %<2* indicates performance when selecting the best of 5 generated samples.

Method	%<2	%<2*	Med.
GAUSSIAN PRIOR	17.0	29.2	3.8
VELOCITY PREDICTION	11.9	28.8	3.8
STANDARD TFN LAYERS	13.7	25.4	3.6
NO REFINEMENT LOSS	9.8	22.1	3.7
NO SELF-CONDITIONING	14.3	29.8	3.7
HARMONICFLOW $\sigma = 0$	18.3	31.3	3.5
HARMONICFLOW $\sigma = 0.5$	20.5	34.5	3.4

to the oracle method that has access to the ground truth ligand structure. The joint structure generation helps in determining the original residue types (keeping in mind that these are not necessarily the only or best). RANDOM LIGAND POS further confirms that inferring approximate ligand coordinates, like HARMONICFLOW in FLOWSITE, is crucial for recovering the binding pocket.

4.4. Question 3: Ablations and Flow Matching

Here, we attempt to understand better the behavior of flow-matching generative models for biomolecular structure generation via experiments with HARMONIC FLOW. We use *Radius-Pockets* on the sequence similarity split of PDBBind.

Investigations. GAUSSIAN PRIOR uses an isotropic Gaussian as prior instead of our harmonic prior. In VELOCITY PREDICTION, the TFN model predicts $(\mathbf{x}_1 - \mathbf{x}_0)$ instead of \mathbf{x}_1 meaning that $\mathcal{L}_{CFM} = \|v_\theta - (\mathbf{x}_1 - \mathbf{x}_0)\|^2$. In STANDARD TFN LAYERS, our [refinement TFN layers](#) are replaced, meaning that there are no intermediate position updates - only the last layer produces an update. In NO REFINEMENT LOSS, the loss \mathcal{L}_{refine} is dropped from the final weighted sum of losses. NO SELF-CONDITIONING does not

use our [structure self-conditioning](#). SIGMA=0 uses $\sigma = 0$ for the conditional flow, corresponding to a deterministic interpolant for training.

Results. Table 4 shows the importance of our self-conditioned flow matching objective, which enables refinement of the binding structure prediction $\tilde{\mathbf{x}}_1^t$ next to updates of \mathbf{x}_t at little additional training time - a 12.8% increase in this experiment. Furthermore, the refinement TFN layers improve structure prediction substantially. Lastly, parameterizing the vector field to predict \mathbf{x}_1 instead of $(\mathbf{x}_1 - \mathbf{x}_0)$ appears more suitable for flow matching applications in molecular structure generation.

5. Conclusion

We proposed the HARMONICFLOW generative process for binding structure generation and FLOWSITE for binding site design. Our HARMONICFLOW improves upon the state-of-the-art generative process for docking in simplicity, applicability, and performance in various docking settings. We investigated how flow matching contributes to this, together with our technical innovations such as self-conditioned flow matching, harmonic prior ligands, or equivariant refinement TFN layers.

With FLOWSITE, we leverage our superior binding structure generative process and extend it to discrete residue types, resulting in a joint discrete-continuous process for designing ligand binding pockets—an important task for which no general and no deep learning solutions exist yet. FLOWSITE improves upon various baselines in recovering native binding sites without requiring prior knowledge of the bound protein-ligand complex. Thus, FLOWSITE is a useful step toward generally applicable binding site design, which has important applications in fields ranging from drug discovery to enzyme design.

Acknowledgments

We thank Jason Yim, Gabriele Corso, Rachel Wu, Felix Faltings, Jeremy Wohlwend, MinGyu Choi, Juno Nam, Soojung Yang, Nick Polizzi, Jody Mou, Michael Plainier, Julia Balla, Shangyuan Tong, and Peter Holderrieth for insightful discussions and feedback.

This work was supported by the NSF Expeditions grant (award 1918839: Collaborative Research: Understanding the World Through Code), the Machine Learning for Pharmaceutical Discovery and Synthesis (MLPDS) consortium, the Abdul Latif Jameel Clinic for Machine Learning in Health, the DTRA Discovery of Medical Countermeasures Against New and Emerging (DOMANE) threats program, the DARPA Accelerated Molecular Discovery program, the NSF AI Institute CCF-2112665, the NSF Award 2134795, and the GIST-MIT Research Collaboration grant.

Impact Statement

We develop a general framework for generating 3D point clouds in HARMONICFLOW and for jointly generating discrete and continuous data in FLOWSITE. As generative modeling frameworks, these are general methods with many possible societal impacts that do not require specific mention. We apply these methods to the tasks of docking and binding site design. In the past and present, the societal impacts of these applications have been overwhelmingly positive, with advances possible in drug discovery, enzyme design, or biological sensing. Future applications with negative outcomes can also be imagined, such as decreasing the barrier of entry for developing biological weapons.

References

- Albergo, M. S. and Vanden-Eijnden, E. Building normalizing flows with stochastic interpolants. In *The Eleventh International Conference on Learning Representations*, 2022.
- Albergo, M. S., Boffi, N. M., and Vanden-Eijnden, E. Stochastic interpolants: A unifying framework for flows and diffusions. *arXiv preprint arXiv:2303.08797*, 2023.
- Berman, H., Henrick, K., and Nakamura, H. Announcing the worldwide Protein Data Bank. *Nat Struct Biol*, 10 (12):980, Dec 2003.
- Bortoli, V. D., Thornton, J., Heng, J., and Doucet, A. Diffusion schrödinger bridge with applications to score-based generative modeling, 2023.
- Bose, A. J., Akhound-Sadegh, T., Fatras, K., Huguet, G., Rector-Brooks, J., Liu, C.-H., Nica, A. C., Korablyov, M., Bronstein, M., and Tong, A. Se(3)-stochastic flow matching for protein backbone generation, 2023.
- Buttenschoen, M., Morris, G. M., and Deane, C. M. Posebusters: Ai-based docking methods fail to generate physically valid poses or generalise to novel sequences, 2023.
- Campbell, A., Harvey, W., Weilbach, C., Bortoli, V. D., Rainforth, T., and Doucet, A. Trans-dimensional generative modeling via jump diffusion models, 2023.
- Chen, R. T. and Lipman, Y. Riemannian flow matching on general geometries. *arXiv preprint arXiv:2302.03660*, 2023.
- Chen, T., Liu, G.-H., and Theodorou, E. Likelihood training of schrödinger bridge using forward-backward SDEs theory. In *International Conference on Learning Representations*, 2022a. URL <https://openreview.net/forum?id=nioAdKCEdXB>.
- Chen, T., Zhang, R., and Hinton, G. Analog bits: Generating discrete data using diffusion models with self-conditioning, 2023.
- Chen, Y., Chen, Q., and Liu, H. Depact and pacmatch: A workflow of designing de novo protein pockets to bind small molecules. *Journal of Chemical Information and Modeling*, 62(4):971–985, Feb 2022b. ISSN 1549-9596. doi: 10.1021/acs.jcim.1c01398. URL <https://doi.org/10.1021/acs.jcim.1c01398>.
- Corso, G., Cavalleri, L., Beaini, D., Liò, P., and Veličković, P. Principal neighbourhood aggregation for graph nets. *Advances in Neural Information Processing Systems*, 33: 13260–13271, 2020.
- Corso, G., Stärk, H., Jing, B., Barzilay, R., and Jaakkola, T. Diffdock: Diffusion steps, twists, and turns for molecular docking, 2023.
- Corso, G., Deng, A., Fry, B., Polizzi, N., Barzilay, R., and Jaakkola, T. Deep confident steps to new pockets: Strategies for docking generalization, 2024.
- Dauparas, J., Anishchenko, I., Bennett, N., Bai, H., Ragotte, R. J., Milles, L. F., Wicky, B. I. M., Courbet, A., de Haas, R. J., Bethel, N., Leung, P. J. Y., Huddy, T. F., Pellock, S., Tischer, D., Chan, F., Koepnick, B., Nguyen, H., Kang, A., Sankaran, B., Bera, A. K., King, N. P., and Baker, D. Robust deep learning-based protein sequence design using proteinmpnn. *Science*, 378(6615):49–56, 2022. doi: 10.1126/science.add2187.
- Dauparas, J., Lee, G. R., Pecoraro, R., An, L., Anishchenko, I., Glasscock, C., and Baker, D. Atomic context-conditioned protein sequence design using ligandmpnn. *bioRxiv*, 2023.
- Dou, J., Doyle, L., Greisen, Jr, P., Schena, A., Park, H., Johnsson, K., Stoddard, B. L., and Baker, D. Sampling

- and energy evaluation challenges in ligand binding protein design. *Protein Sci.*, 26(12):2426–2437, December 2017.
- Gao, Z., Tan, C., Chacón, P., and Li, S. Z. Pifold: Toward effective and efficient protein inverse folding, 2023a.
- Gao, Z., Tan, C., and Li, S. Z. Knowledge-design: Pushing the limit of protein design via knowledge refinement, 2023b.
- Geiger, M., Smidt, T., M., A., Miller, B. K., Boomsma, W., Dice, B., Lapchevskiy, K., Weiler, M., Tyszkiewicz, M., Batzner, S., Uhrin, M., Frellsen, J., Jung, N., Sanborn, S., Rackers, J., and Bailey, M. Euclidean neural networks: e3nn, 2020.
- Halgren, T. A., Murphy, R. B., Friesner, R. A., Beard, H. S., Frye, L. L., Pollard, W. T., and Banks, J. L. Glide: a new approach for rapid, accurate docking and scoring. 2. enrichment factors in database screening. *Journal of medicinal chemistry*, 2004.
- Ho, J., Jain, A., and Abbeel, P. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- Hoogetboom, E., Satorras, V. G., Vignac, C., and Welling, M. Equivariant diffusion for molecule generation in 3d, 2022.
- Hsu, C., Verkuil, R., Liu, J., Lin, Z., Hie, B., Sercu, T., Lerer, A., and Rives, A. Learning inverse folding from millions of predicted structures. *bioRxiv*, 2022. doi: 10.1101/2022.04.10.487779. URL <https://www.biorxiv.org/content/early/2022/09/06/2022.04.10.487779>.
- Jin, W., Wohlwend, J., Barzilay, R., and Jaakkola, T. Iterative refinement graph neural network for antibody sequence-structure co-design, 2022.
- Jing, B., Corso, G., Chang, J., Barzilay, R., and Jaakkola, T. Torsional diffusion for molecular conformer generation. *arXiv preprint arXiv:2206.01729*, 2022.
- Jing, B., Erives, E., Pao-Huang, P., Corso, G., Berger, B., and Jaakkola, T. S. Eigenfold: Generative protein structure prediction with diffusion models. In *ICLR 2023-Machine Learning for Drug Discovery workshop*, 2023.
- Jumper, J., Evans, R., Pritzel, A., Green, T., Figurnov, M., Ronneberger, O., Tunyasuvunakool, K., Bates, R., Žídek, A., Potapenko, A., et al. Highly accurate protein structure prediction with alphafold. *Nature*, 596(7873):583–589, 2021.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Klein, L., Krämer, A., and Noé, F. Equivariant flow matching, 2023.
- Krapp, L. F., Meireles, F. A., Abriata, L. A., and Peraro, M. D. Context-aware geometric deep learning for protein sequence design. *bioRxiv*, 2023. doi: 10.1101/2023.06.19.545381. URL <https://www.biorxiv.org/content/early/2023/06/19/2023.06.19.545381>.
- Lin, H., Huang, Y., Liu, M., Li, X., Ji, S., and Li, S. Z. Diffbp: Generative diffusion of 3d molecules for target protein binding, 2022.
- Lipman, Y., Chen, R. T., Ben-Hamu, H., Nickel, M., and Le, M. Flow matching for generative modeling. In *The Eleventh International Conference on Learning Representations*, 2022.
- Liu, X., Gong, C., and Liu, Q. Flow straight and fast: Learning to generate and transfer data with rectified flow, 2022.
- Liu, Z., Su, M., Han, L., Liu, J., Yang, Q., Li, Y., and Wang, R. Forging the basis for developing protein–ligand interaction scoring functions. *Accounts of Chemical Research*, 50(2):302–309, 2017.
- Lu, W., Wu, Q., Zhang, J., Rao, J., Li, C., and Zheng, S. Tankbind: Trigonometry-aware neural networks for drug-protein binding structure prediction. *Advances in neural information processing systems*, 2022.
- Malisi, C., Schumann, M., Toussaint, N. C., Kageyama, J., Kohlbacher, O., and Höcker, B. Binding pocket optimization by computational protein design. *PLoS One*, 7(12): e52505, December 2012.
- Martinkus, K., Ludwiczak, J., Cho, K., Liang, W.-C., Lafrance-Vanasse, J., Hotzel, I., Rajpal, A., Wu, Y., Bonneau, R., Gligoričević, V., and Loukas, A. Abdifuser: Full-atom generation of in-vitro functioning antibodies, 2023.
- McNutt, A. T., Francoeur, P., Aggarwal, R., Masuda, T., Meli, R., Ragoza, M., Sunseri, J., and Koes, D. R. Gnina 1.0: molecular docking with deep learning. *Journal of cheminformatics*, 13(1):1–20, 2021.
- Méndez-Lucio, O., Ahmad, M., del Rio-Chanona, E. A., and Wegner, J. K. A geometric deep learning approach to predict binding conformations of bioactive molecules. *Nature Machine Intelligence*, 3(12):1033–1039, 2021.
- Neklyudov, K., Brekelmans, R., Severo, D., and Makhzani, A. Action matching: Learning stochastic dynamics from samples, 2023.

- Nelson, D. L. and Cox, M. M. *Lehninger Principles of Biochemistry, Fourth Edition*. Cold Spring Harbor Laboratory, fourth edition edition, 2004.
- Noske, J., Kynast, J. P., Lemm, D., Schmidt, S., and Höcker, B. PocketOptimizer 2.0: A modular framework for computer-aided ligand-binding design. *Protein Sci.*, 32(1):e4516, January 2023.
- Plainer, M., Toth, M., Dobers, S., Stark, H., Corso, G., Marquet, C., and Barzilay, R. Diffdock-pocket: Diffusion for pocket-level docking with sidechain flexibility. In *NeurIPS 2023 Workshop on New Frontiers of AI for Drug Discovery and Development*, 2023.
- Polizzi, N. F. and DeGrado, W. F. A defined structural unit enables de novo design of small-molecule-binding proteins. *Science*, 369(6508):1227–1233, 2020. doi: 10.1126/science.abb8330. URL <https://www.science.org/doi/abs/10.1126/science.abb8330>.
- Pooladian, A.-A., Ben-Hamu, H., Domingo-Enrich, C., Amos, B., Lipman, Y., and Chen, R. T. Q. Multisample flow matching: Straightening flows with minibatch couplings, 2023.
- Qiao, Z., Nie, W., Vahdat, A., au2, T. F. M. I., and Anandkumar, A. State-specific protein-ligand complex structure prediction with a multi-scale deep generative model, 2023.
- Schneuing, A., Du, Y., Harris, C., Jamasb, A., Igashov, I., Du, W., Blundell, T., Lió, P., Gomes, C., Welling, M., Bronstein, M., and Correia, B. Structure-based drug design with equivariant diffusion models, 2023.
- Shi, Y., Bortoli, V. D., Campbell, A., and Doucet, A. Diffusion schrödinger bridge matching, 2023.
- Somnath, V. R., Pariset, M., Hsieh, Y.-P., Martinez, M. R., Krause, A., and Bunne, C. Aligned diffusion schrödinger bridges, 2023.
- Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Ermon, S., and Poole, B. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*, 2021.
- Song, Y., Gong, J., Xu, M., Cao, Z., Lan, Y., Ermon, S., Zhou, H., and Ma, W.-Y. Equivariant flow matching with hybrid probability transport for 3d molecule generation. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=hHUZ5V9XFu>.
- Stiel, A. C., Nellen, M., and Höcker, B. PocketOptimizer and the design of ligand binding sites. *Methods Mol. Biol.*, 1414:63–75, 2016.
- Stärk, H., Ganea, O.-E., Pattanaik, L., Barzilay, R., and Jaakkola, T. Equibind: Geometric deep learning for drug binding structure prediction, 2022.
- Thomas, N., Smidt, T., Kearnes, S., Yang, L., Li, L., Kohlhoff, K., and Riley, P. Tensor field networks: Rotation-and translation-equivariant neural networks for 3d point clouds. *arXiv preprint*, 2018.
- Thomsen, R. and Christensen, M. H. Moldock: a new technique for high-accuracy molecular docking. *Journal of medicinal chemistry*, 49(11):3315–3321, 2006.
- Tong, A., Malkin, N., Fatras, K., Atanackovic, L., Zhang, Y., Huguet, G., Wolf, G., and Bengio, Y. Simulation-free schrödinger bridges via score and flow matching, 2023a.
- Tong, A., Malkin, N., Huguet, G., Zhang, Y., Rector-Brooks, J., Fatras, K., Wolf, G., and Bengio, Y. Improving and generalizing flow-based generative models with mini-batch optimal transport, 2023b.
- Trott, O. and Olson, A. J. Autodock vina: improving the speed and accuracy of docking with a new scoring function, efficient optimization, and multithreading. *Journal of computational chemistry*, 31(2):455–461, 2010.
- Verma, Y., Heinonen, M., and Garg, V. Abode: Ab initio antibody design using conjoined odes, 2023.
- Vignac, C., Krawczuk, I., Siraudin, A., Wang, B., Cevher, V., and Frossard, P. Digress: Discrete denoising diffusion for graph generation, 2023a.
- Vignac, C., Osman, N., Toni, L., and Frossard, P. Midi: Mixed graph and 3d denoising diffusion for molecule generation, 2023b.
- Watson, J. L., Juergens, D., Bennett, N. R., Trippe, B. L., Yim, J., Eisenach, H. E., Ahern, W., Borst, A. J., Ragotte, R. J., Milles, L. F., Wicky, B. I. M., Hanikel, N., Pellock, S. J., Courbet, A., Sheffler, W., Wang, J., Venkatesh, P., Sappington, I., Torres, S. V., Lauko, A., De Bortoli, V., Mathieu, E., Ovchinnikov, S., Barzilay, R., Jaakkola, T. S., DiMaio, F., Baek, M., and Baker, D. De novo design of protein structure and function with rfdiffusion. *Nature*, 620(7976):1089–1100, Aug 2023.
- Yeh, A. H.-W., Norn, C., Kipnis, Y., Tischer, D., Pellock, S. J., Evans, D., Ma, P., Lee, G. R., Zhang, J. Z., Anishchenko, I., Coventry, B., Cao, L., Dauparas, J., Halabiya, S., DeWitt, M., Carter, L., Houk, K. N., and Baker, D. De novo design of luciferases using deep learning. *Nature*, 614(7949):774–780, Feb 2023. ISSN 1476-4687. doi: 10.1038/s41586-023-05696-3. URL <https://doi.org/10.1038/s41586-023-05696-3>.

Yi, K., Zhou, B., Shen, Y., Liò, P., and Wang, Y. G. Graph denoising diffusion for inverse protein folding, 2023.

Yim, J., Campbell, A., Foong, A. Y., Gastegger, M., Jiménez-Luna, J., Lewis, S., Satorras, V. G., Veeling, B. S., Barzilay, R., Jaakkola, T., et al. Fast protein backbone generation with se (3) flow matching. *arXiv preprint arXiv:2310.05297*, 2023a.

Yim, J., Trippe, B. L., De Bortoli, V., Mathieu, E., Doucet, A., Barzilay, R., and Jaakkola, T. Se(3) diffusion model with application to protein backbone generation. *arXiv preprint*, 2023b.

Zhang, Y., Cai, H., Shi, C., Zhong, B., and Tang, J. E3bind: An end-to-end equivariant network for protein-ligand docking, 2023a.

Zhang, Z., Lu, Z., Hao, Z., Zitnik, M., and Liu, Q. Full-atom protein pocket design via iterative refinement, 2023b.

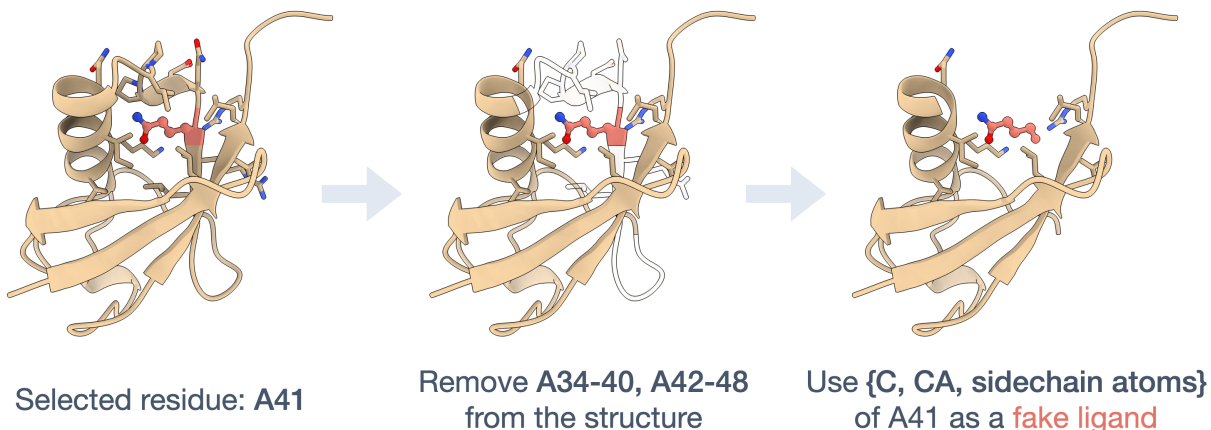


Figure 5. **Visualization of Fake Ligand creation.** Depicted is a fake ligand created for the Ubiquitin protein. Out of all residues that have at least 4 contacts with other residues (apart from those that are within 7 locations in the chain) a residue is randomly selected as the fake ligand. Then we remove the residue itself from the protein and all residues that are within 7 locations in the chain.

A. Method Details and Explanations

A.1. BLOSUM Score

Next to sequence recovery, we also evaluate with our *BLOSUM Score* in an attempt to penalize amino acid predictions less if the predicted residue type is similar yet different from the original residue. With $\mathbf{A} \in \mathbb{R}^{20 \times 20}$ being the BLOSUM62 matrix, $\mathbf{X} \in \mathbb{R}^{n \times 20}$ the one hot encoded ground truth residues types and $\hat{\mathbf{X}} \in \mathbb{R}^{n \times 20}$ the predicted residues types the *BLOSUM Score* is:

$$\text{Score}(\mathbf{X}, \hat{\mathbf{X}}) = \frac{\mathbf{1}^T \text{diag}(\mathbf{X} \mathbf{A} \hat{\mathbf{X}}^T)}{\mathbf{1}^T \text{diag}(\mathbf{X} \mathbf{A} \mathbf{X}^T)} \quad (4)$$

A.2. Fake Ligand Data Augmentation Visualization

In Figure 5, we visualize the construction of our fake ligands as described in Section 3.2. When constructing the fake ligand from a residue, we drop the backbone oxygen and nitrogen of the amino acid and keep the carbon, alpha carbon, and the side chain as the ligand’s atoms.

A.3. Flow Matching Training and Inference

In Section 3.1, we lay out the conditional flow matching objective as introduced by Lipman et al. (2022) and extended to arbitrary start and end distributions by multiple works concurrently (Albergo & Vanden-Eijnden, 2022; Albergo et al., 2023; Pooladian et al., 2023; Tong et al., 2023b). We presented conditional flow matching in this more general scenario where the prior p_0 and the data p_1 can be arbitrary distributions, as long as we can sample from the prior.

Many choices of conditional flows and conditional vector fields are possible. For different applications and scenarios, some choices perform better than others. We find it to already work well to use a very simple choice of conditional probability path $p_t(\mathbf{x}|\mathbf{x}_0, \mathbf{x}_1) = \mathcal{N}(\mathbf{x}|t\mathbf{x}_1 + (1-t)\mathbf{x}_0, \sigma^2)$, which gives rise to the conditional vector field $u_t(\mathbf{x}|\mathbf{x}_0, \mathbf{x}_1) = \mathbf{x}_1 - \mathbf{x}_0$. With this conditional flow and with parameterizing v_θ to predict \mathbf{x}_1 , the optimization and inference is remarkably straightforward as algorithms 1 and 2 show.

A.4. Self-conditioned Flow Matching Training and Inference

In Section 3.1, we also explain the self-conditioning training and inference procedure. When additionally using self-conditioning, the training and inference algorithms are only slightly modified and still very simple as presented in algorithms

Algorithm 1 Conditional Flow Matching training with x_1 prediction and simple constant width gaussian conditional path.

Input: Training data distribution p_1 , prior p_0 , σ , and initialized vector field v_θ

while Training **do**

$x_0 \sim p_0(x_0)$; $x_1 \sim p_1(x_1)$; $t \sim \mathcal{U}(0, 1)$

$\mu_t \leftarrow tx_1 + (1-t)x_0$

$x \sim \mathcal{N}(\mu_t, \sigma^2 I)$

$\mathcal{L}_{CFM} \leftarrow \|v_\theta(x, t) - x_1\|^2$

$\theta \leftarrow \text{Update}(\theta, \nabla_\theta \mathcal{L}_{CFM})$

end while

return v_θ

Algorithm 2 Conditional Flow Matching inference with x_1 prediction and simple constant width gaussian conditional path.

Input: Prior p_0 , number of integration steps T , and trained vector field v_θ

$steps \leftarrow 1$

$\Delta t \leftarrow 1/T$

$t \leftarrow 0$

$x_0 \sim p_0(x_0)$

$x_t \leftarrow x_0$

while $steps \leq T - 1$ **do**

$\tilde{x}_1 \leftarrow v_\theta(x_t, t)$

$x_t \leftarrow x_t + \Delta t(\tilde{x}_1 - x_t)/(1-t)$

$t \leftarrow t + \Delta t$

end while

return x_t

3 and 4.

A.5. FLOWSITE Architecture

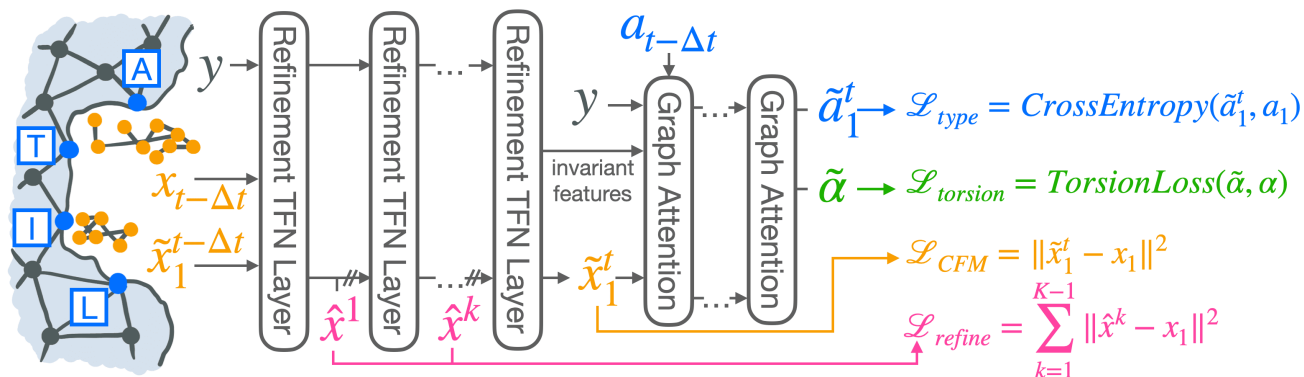


Figure 6. FLOWSITE architecture. The refinement TFN layers (also in HARMONICFLOW) first update the ligand coordinates $x_{t-\Delta t}$ multiple times to produce the structure prediction \hat{x}_1^t (from which \tilde{x}_t is computed during inference). The TFN’s invariant features and \hat{x}_1^t are fed to invariant layers to produce side chain angles $\tilde{\alpha}$ and the new residue estimate \tilde{a}_1^t .

Here, we detail the FLOWSITE architecture as visualized in Figure 6 in more detail. The first half of the architecture is an equivariant Tensor Field Network (Thomas et al., 2018) while the second part is an invariant architecture with graph attention layers similar to the architecture of PiFOLD (Gao et al., 2023a) where edge features are also initialized and updated.

Radius Graph. The protein and (multi-)ligand are represented as graphs: each residue corresponds to a node, and each ligand atom is a node. Edges are drawn between residue nodes if they are within 50 Å, between ligand nodes if they are within 50 Å, and between the two molecules’ nodes if they are within 30 Å. The locations of the residue nodes are given by

Algorithm 3 Conditional Flow Matching training with \mathbf{x}_1 prediction and simple constant width gaussian conditional path.

Input: Training data distribution p_1 , prior p_0 , σ , and initialized vector field v_θ

while Training **do**

$\mathbf{x}_0 \sim p_0(\mathbf{x}_0)$; $\mathbf{x}_1 \sim p(\mathbf{x}_1)$; $t \sim \mathcal{U}(0, 1)$; $s \sim \mathcal{U}(0, 1)$

$\mu_t \leftarrow t\mathbf{x}_1 + (1 - t)\mathbf{x}_0$

$\mathbf{x} \sim \mathcal{N}(\mu_t, \sigma^2 I)$

$\tilde{\mathbf{x}}_1 \sim p_0(\tilde{\mathbf{x}}_1)$

if $s > 0.5$ **then**

$\tilde{\mathbf{x}}_1 \leftarrow v_\theta(\mathbf{x}, \tilde{\mathbf{x}}_1, t)$

end if

$\mathcal{L}_{CFM} \leftarrow \|v_\theta(\mathbf{x}, \tilde{\mathbf{x}}_1, t) - \mathbf{x}_1\|^2$

$\theta \leftarrow \text{Update}(\theta, \nabla_\theta \mathcal{L}_{CFM})$

end while

return v_θ

Algorithm 4 Conditional Flow Matching inference with \mathbf{x}_1 prediction and simple constant width gaussian conditional path.

Input: Prior p_0 , number of integration steps T , and trained vector field v_θ

$steps \leftarrow 1$

$\Delta t \leftarrow 1/T$

$t \leftarrow 0$

$\tilde{\mathbf{x}}_1 \sim p(\mathbf{x}_0)$

$\mathbf{x}_0 \sim p(\mathbf{x}_0)$

$\mathbf{x}_t \leftarrow \mathbf{x}_0$

while $steps \leq T - 1$ **do**

$\tilde{\mathbf{x}}_1 \leftarrow v_\theta(\mathbf{x}, \tilde{\mathbf{x}}_1, t)$

$\mathbf{x}_t \leftarrow \mathbf{x}_t + \Delta t(\tilde{\mathbf{x}}_1 - \mathbf{x}_t)/(1 - t)$

$t \leftarrow t + \Delta t$

end while

return \mathbf{x}_t

their alpha carbons, while the atom locations provide the node positions for the ligand nodes.

Node Features. The ligand features as input to the TNF and to the invariant part of the architecture are atomic number; chirality; degree; formal charge; implicit valence; the number of connected hydrogens; hybridization type; whether or not it is in an aromatic ring; in how many rings it is; and finally, 6 features for whether or not it is in a ring of size 5 or 6.

The initial receptor features for the TFN are scalar feature encodings of the invariant residue types together with vector features, which are three vectors from the alpha carbon to N, C, and O.

For the invariant graph attention layer stack, the residue inputs are the invariant geometric encodings of PiFOLD (Gao et al., 2023a). Additionally, they contain the residue type self-conditioning information via embeddings of the previously predicted features $\tilde{\mathbf{a}}_1^t$ and the invariant scalar node features of the last refinement TFN layer.

Additionally, radial basis encodings of the sampling time t of the conditional flow are added to all initial node features.

Edge Features. For the Tensor Field Network, the edge features are a radial basis embedding of the alpha carbon distances for the protein-to-protein edges, atom distances for the ligand-to-ligand edges, and alpha carbon to ligand atom distances for the edges between the protein and the ligand. Additionally, the ligand-to-ligand edges features obtain information of the structure self-conditioning by also adding the radial basis interatomic distance embeddings of the previously predicted ligand coordinates $\tilde{\mathbf{x}}_1^t$ to them.

Meanwhile, for the invariant graph attention part of the architecture, the ligand-to-ligand edge features are only radial basis embeddings of the interatomic distances. The protein-to-protein edge features are given by radial basis encodings of all pairwise distances between the backbone atoms N, C, Ca, O, and an additional virtual atom (as introduced by PiFOLD) associated with each residue. The edges between the protein and ligand are featurized as the embeddings of the four possible

distances between a single ligand atom and the four backbone atoms of a residue.

Tensor Field Network. The equivariant part of FLOWSITE uses our equivariant refinement TFN layers based on tensorfield networks (Thomas et al., 2018) and implemented using the e3nn library (Geiger et al., 2020). These rely on tensor products between invariant and equivariant features. We denote the tensor products as \otimes_w where w are the path weights. Further, we write the i -th node features after the k -th layer as \mathbf{h}_i^k for the equivariant Tensorfield network layers. \mathbf{h}_i^0 is initialized as described above in the Node Features paragraph. Lastly, \mathcal{N}_i denotes the neighbors of the i -th node in the radius graph.

Equivariant TFN Refinement Layer. Each layer has a different set of weights for all four types of edges: ligand-to-ligand, protein-to-protein, ligand-to-protein, and protein-to-protein. The layers first update node features before updating ligand coordinates based on them. For every edge in the graph, a message is constructed based on the invariant and equivariant features of the nodes it connects. This is done in an equivariant fashion via tensor products. The tensor product is parameterized by the edge embeddings and the invariant scalar features of nodes that are connected by the edge. To obtain a new node embedding, the messages are summed:

$$\mathbf{h}_i^{k+1} \leftarrow \mathbf{h}_i^k + \text{BN} \left(\frac{1}{|\mathcal{N}_i|} \sum_{j \in \mathcal{N}_i} Y(\hat{r}_{ij}) \otimes_{\psi_{ij}} \mathbf{h}_j^k \right) \quad (5)$$

with $\psi_{ij} = \Psi(e_{ij}, \mathbf{h}_i^k, \mathbf{h}_j^k)$

Here, BN is the (equivariant) batch normalization of the e3nn library. The orders of all features are always restricted to a maximum of 1. The neural networks Ψ have separate sets of weights for all 4 kinds of edges. Using these new node features and the previous layer’s ligand position update $\hat{\mathbf{x}}^k$ (or the input positions $\hat{\mathbf{x}}^0 = \mathbf{x}_t$ for the first layer), the next ligand position update $\hat{\mathbf{x}}^{k+1}$ is produced via an O(3) equivariant linear layer Φ of the e3nn library:

$$\hat{\mathbf{x}}^{k+1} \leftarrow \hat{\mathbf{x}}^{k+1} + \Phi(\mathbf{h}^{k+1}) \quad (6)$$

Invariant Graph Attention Layers. These layers are based on PiFOLD and update both node and edge features. The initial features are described in the paragraphs above. We denote these as \mathbf{h}_i^l and \mathbf{e}_{ji}^l for the l -th graph attention layer to disambiguate with the features \mathbf{h}_i^k of the equivariant refinement TFN layers. When aggregating the features for the i -th node, attention weights are first created and then used to weight messages from each neighboring node. With \parallel denoting concatenation and Ω , Ξ , and Π being feed-forward neural networks, the update is defined as:

$$\begin{aligned} w_{ji} &\leftarrow \Pi(\mathbf{h}_j^l \parallel \mathbf{e}_{ji}^l \parallel \mathbf{h}_i^l) \\ a_{ji} &\leftarrow \frac{\exp w_{ji}}{\sum_{a \in \mathcal{N}_i} \exp w_{ai}} \\ \mathbf{v}_j &= \Xi(\mathbf{e}_{ji}^l \parallel \mathbf{h}_j^l) \\ \mathbf{h}_i^{l+1} &= \sum_{j \in \mathcal{N}_i} a_{ji} \mathbf{v}_j. \end{aligned} \quad (7)$$

We drop the *global context attention* used in PiFOLD as we did not find them to be helpful for sequence recovery in any of our experiments. This was with and without ligands.

Based on the new node features, the edge features are updated as follows:

$$\mathbf{e}_{ji}^{l+1} = \Omega(\mathbf{h}_j^{l+1} \parallel \mathbf{e}_{ji}^l \parallel \mathbf{h}_i^{l+1}) \quad (8)$$

B. Discussion

HARMONICFLOW has the ability to produce arbitrary bond lengths and bond angles. This distinguishes it from DIFFDOCK (Corso et al., 2023), which only changes torsion angles, translation, and rotation of an initial seed conformer. Thus, unlike DIFFDOCK, HARMONICFLOW would be able to produce unrealistic local structures. That this is not the case, as shown in Figure 12 attests to how HARMONICFLOW learns physical constraints. Still, we argue that the role of deep learning generative models should be to solve the hard problem of finding the correct coarse structure. If one desires a conformer

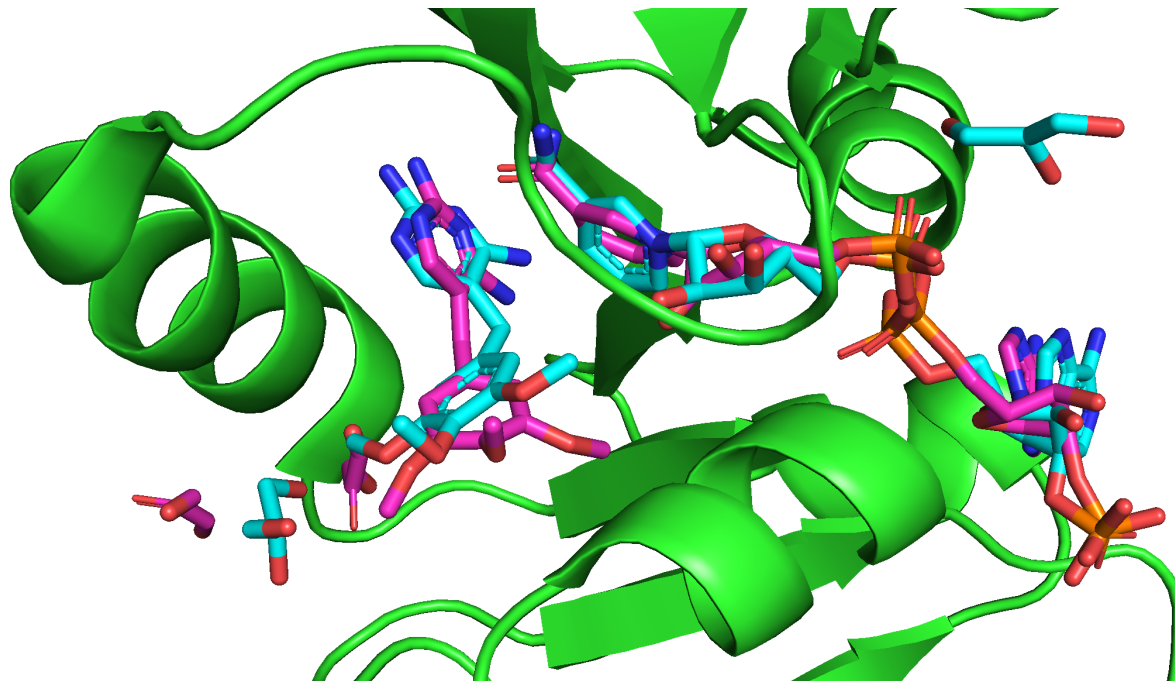


Figure 7. Multi-ligand docking case study. A randomly picked complex with a multi-ligand from the Binding MOAD test set with 4 molecules as multi-ligand. We show the ground truth ligand structure (blue) and a sample of HARMONICFLOW (purple). We find that the predicted structure mostly matches the ground truth for both of the large ligands, but both small ligands are placed on the left while one of them should be on the right.

with low energy with respect to some energy function, this can be easily and quickly obtained by relaxing with that energy function.

C. Additional Results

C.1. Multi-ligand Docking Case Studies

In this subsection we provide 4 case studies of Multi-liand docking to assess the behaviour of HARMONICFLOW for this task. For this purpose, we randomly choose 4 complexes from the Binding MOAD test set under the condition that they contain different ligands.

Overall, we find that while the RMSDs of EIGENFOLD DIFFUSION in Table 2 are similar to those of HARMONICFLOW, the complexes that HARMONICFLOW generates are often more physically plausible. For instance, HARMONICFLOW’s rings have the appropriate shape and planar systems are actually planar which often is not the case for EIGENFOLD DIFFUSION.

C.2. Analysis of joint ODE dynamics

To analyze the behavior of the learned ODE over discrete and continuous data, we provide Figure 11. This figure shows the evolution of the x_1 prediction’s RMSD and the evolution of the output entropy of the residue type probabilities. The results show how a more determined structure prediction correlates with a decrease in uncertainty of the residue type prediction.

C.3. FlowSite Ablations

In Table 5 we provide additional ablations for FLOWSITE. In NO REFINEMENT LOSS, the loss \mathcal{L}_{refine} is dropped from the final weighted sum of losses. in NO SIDE CHAIN TORSION LOSS, the loss $\mathcal{L}_{torsion}$ is dropped from the final weighted sum of losses. In BACKBONE NOISE, we add Gaussian noise to the input protein backbone coordinates with standard deviation 0.2. In ONLY EQUIVARIANT LAYERS, we replace the invariant graph attention layers in the FLOWSITE architecture with equivariant refinement TFN layers, showing how the invariant architecture is crucial for predicting the discrete residue types.

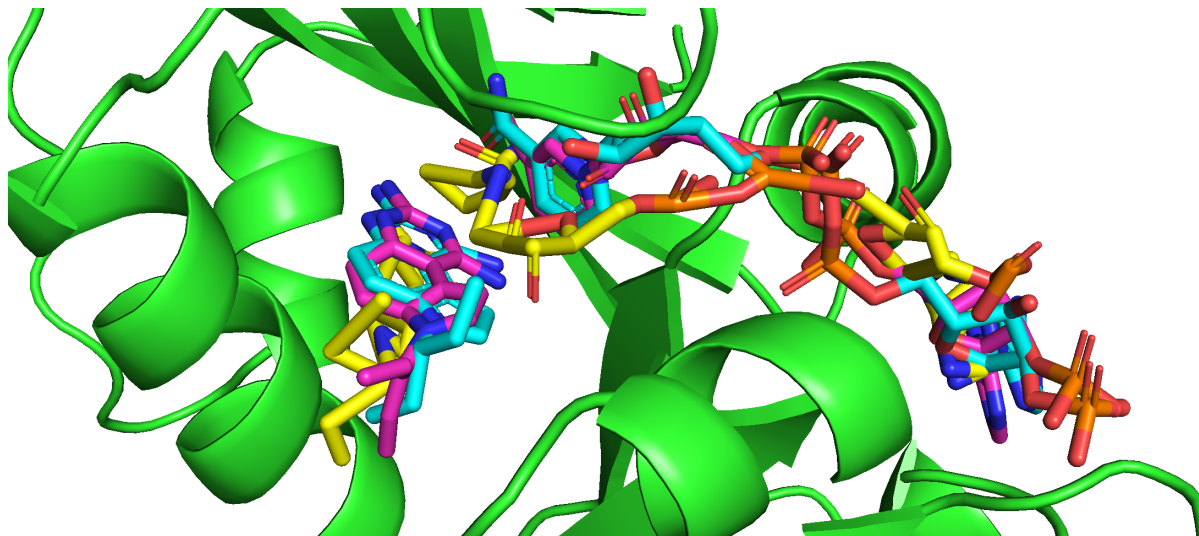


Figure 8. **Multi-ligand docking case study.** A randomly picked complex with a multi-ligand from the Binding MOAD test set for which we show the ground truth ligand structure (blue), a sample of HARMONICFLOW (purple) and a sample of EIGENFOLD DIFFUSION (yellow). We find that the rings of HARMONICFLOW have the correct shapes and that the bond lengths and angles are physically plausible while EIGENFOLD DIFFUSION fails to produce planar rings.

Table 5. **Binding Site Recovery ablations.** Results on the PDBBind sequence similarity split. *Recovery* is the percentage of correctly predicted residues, and *BLOSUM score* takes residue similarity into account.

Method	<i>BLOSUM score</i>	<i>Recovery</i>
NO SIDE CHAIN TORSION LOSS	45.9	47.7
BACKBONE NOISE	39.6	42.4
NO FAKE LIGAND AUGMENTATION	45.5	46.7
NO REFINEMENT LOSS	45.7	47.6
ONLY EQUIVARIANT LAYERS	29.8	35.3
FLowsITE discrete loss weight = 0.8	46.9	47.8
FLowsITE discrete loss weight = 0.2	47.6	49.5

C.4. Docking without residue identities

For our binding site design, it is important that the structure modeling of the ligand is accurate given the evidence that having a good model of the (multi-)ligand structure is important for recovering pockets and given the interlink between 3D structure and binding affinity / binding free energy. In the main text Section 4.2, we investigated HARMONICFLOW’s performance for docking with known residue identities. However, when using HARMONICFLOW for binding site design, the residue identities are not known a priori, and structure reasoning abilities in this scenario are required.

C.5. Blind Docking

In blind docking, the binding site/pocket of the protein is unknown, and the task is to predict the binding structure given the whole protein. While in, e.g., drug discovery efforts and in our binding site design task, the pocket is known, many important applications exist where discovering the binding site is necessary. In these experiments, the runs take longer to converge than in the pocket-level experiments. Thus, the DiffDock runs were trained for 500 epochs while the HARMONICFLOW runs were trained for 250 epochs instead of the 150 epochs in the pocket-level experiments. Table 7, shows these preliminary blind docking results, which are promising for a deeper investigation of HARMONICFLOW for blind docking and for optimizing it toward this task. This could include training a confidence model as in DIFFDOCK (Corso et al., 2023).

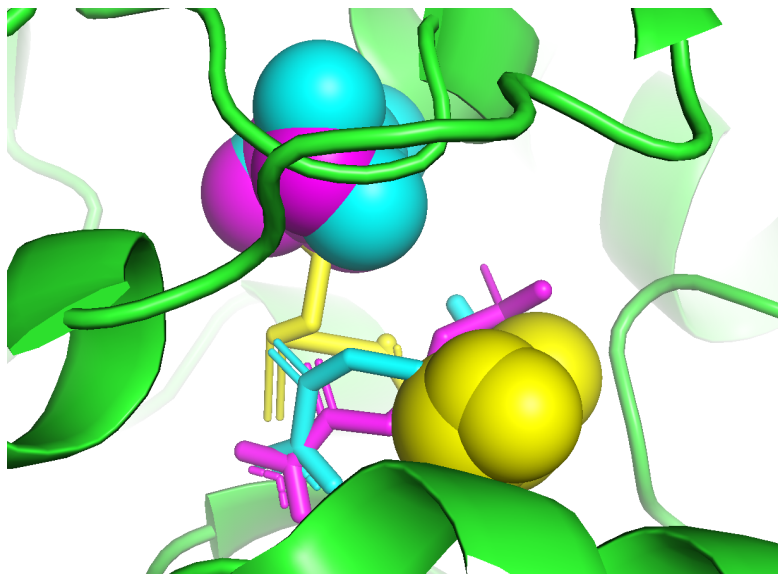


Figure 9. **Multi-ligand docking case study.** A randomly picked complex with a multi-ligand from the Binding MOAD test set for which we show the ground truth ligand structure (blue), a sample of HARMONICFLOW (purple) and a sample of EIGENFOLD DIFFUSION (yellow). We find that EIGENFOLD DIFFUSION incorrectly placed the sulfate and its second ligand prediction is further from the ground truth while HARMONICFLOW almost perfectly places the sulfate.

Table 6. **HARMONICFLOW vs. PRODUCT DIFFUSION without residue identities.** Comparison on PDBBind splits for docking without residue identities into *Distance-Pockets* (residues close to ligand) and *Radius-Pockets* (residues within a radius of the pocket center). The columns "%<2" show the fraction of predictions with an RMSD to the ground truth that is less than 2Å (higher is better). "Med." is the median RMSD (lower is better).

Method	Sequence Similarity Split				Time Split			
	Distance-Pocket		Radius-Pocket		Distance-Pocket		Radius-Pocket	
	%<2	Med.	%<2	Med.	%<2	Med.	%<2	Med.
PRODUCT DIFFUSION	23.2	3.4	14.3	4.0	16.9	4.3	12.3	4.6
HARMONICFLOW	35.5	2.8	20.8	3.5	39.3	2.8	30.7	3.3

C.6. Predicted Complex Visualizations

We visualize generated structures of HARMONICFLOW in Figure 12 from the PDBBind test set under the time-based split of Stärk et al. (2022) in which there are no ligands whose SMILES string was already in the training data. The generated complexes show very chemically plausible ligand structures even though there are no local structure constraints as in DIFFDOCK and HARMONICFLOW has full flexibility in modeling bond angles and bond lengths.

In Table 6, we provide the docking results without residue identities, justifying HARMONICFLOW’s use in FLOWSITE for binding site design.

D. Additional Related Work

D.1. Flow Matching, Stochastic Interpolants, and Schrodinger Bridges

While our exposition of flow matching in the main text focused on the works of Lipman et al. (2022) and Tong et al. (2023b), the innovations in this field were made by multiple papers concurrently. Namely, Action Matching (Neklyudov et al., 2023), stochastic interpolants (Albergo & Vanden-Eijnden, 2022), and rectified flow (Liu et al., 2022) also proposed procedures for learning flows between arbitrary start and end distributions.

An improvement to learning such flows would be if their transport additionally performs the optimal transport between

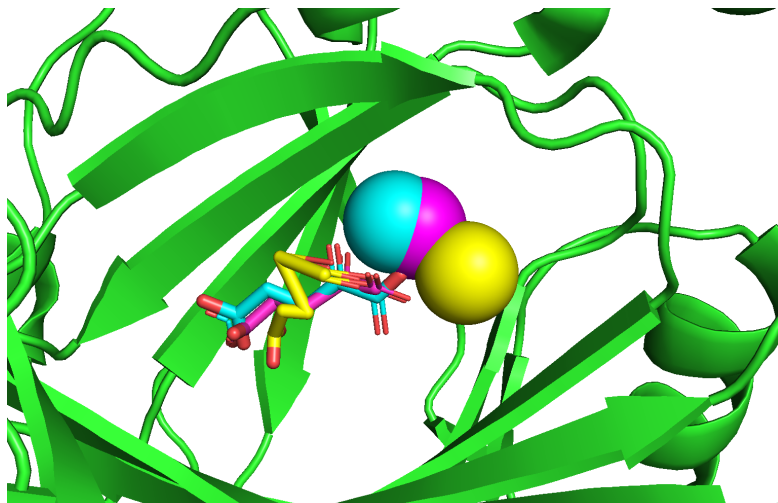


Figure 10. **Multi-ligand docking case study.** A randomly picked complex with a multi-ligand from the Binding MOAD test set for which we show the ground truth ligand structure (blue), a sample of HARMONICFLOW (purple) and a sample of EIGENFOLD DIFFUSION (yellow). We can see unrealistic bond lengths and angles for EIGENFOLD DIFFUSION’s prediction and the predicted ion position is further from the ground truth than for HARMONICFLOW.

Table 7. **HARMONICFLOW vs. PRODUCT DIFFUSION for blind docking.** Comparison on PDBBind splits for blind docking where the binding pocket of the protein is not known, and the whole protein is given as input. The columns “%<2” show the fraction of predictions with an RMSD to the ground truth that is less than 2Å (higher is better). “Med.” is the median RMSD in Å (lower is better). Top 5 and Top 10 refers to the performance when generating 5 or 10 samples and selecting the best among them.

Method	Sequence Split			Time Split		
	%<2	%<5	Med.	%<2	%<5	Med.
PRODUCT DIFFUSION	10.7	40.6	5.9	12.6	44.1	5.6
HARMONICFLOW	11.4	40.5	5.8	25.5	52.2	4.7
	Top 5 %<2	Top 10 %<2		Top 5 %<2	Top 10 %<2	
PRODUCT DIFFUSION	26.6	34.2		28.9	32.5	
HARMONICFLOW	17.1	20.0		30.9	32.6	

the two distributions with respect to some cost. With shorter paths with respect to the cost metric, even fewer integration steps can be performed, and integration errors are smaller. Towards this, Tong et al. (2023b) and Pooladian et al. (2023) concurrently propose mini-batch OT where they train with conditional flow matching but define the conditional paths between the optimal transport solution within a minibatch. They show that in the limit of the batch size, the flow will learn the optimal coupling.

This can be extended to learning Schrodinger bridges in a simulation-free manner via an iterative flow-matching and coupling definition procedure (Shi et al., 2023) akin to rectified flows. Similarly, Tong et al. (2023a) learn a flow and a score simultaneously to reproduce stochastic dynamics as in a Schrodinger bridge and Somnath et al. (2023) learn a Schrodinger bridge between aligned data in a simulation-free manner. Simulation-free here means that the learned vector fields no longer need to be rolled out / simulated during training, which is memory and time-consuming and prohibits learning Schrodinger bridges for larger applications. This was required for previous procedures for learning Schrodinger bridges (Bortoli et al., 2023; Chen et al., 2022a).

D.2. Pocket Design

Yeh et al. (2023) successfully designed a novel Luciferase, which is an enzyme catalyzing the reaction of a Luciferin ligand.

Harmonic Prior Flow Matching

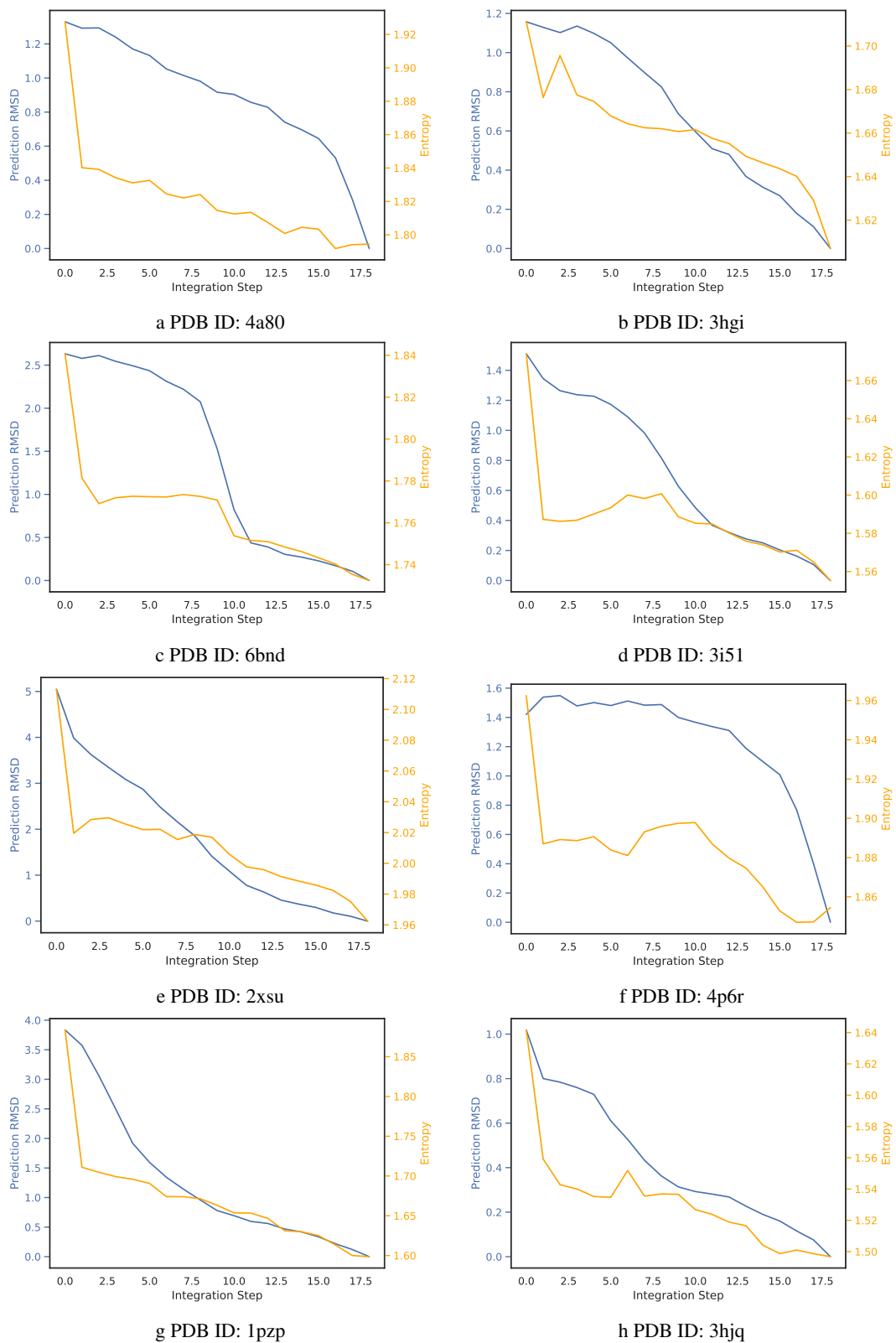


Figure 11. Analysis of joint ODE dynamics. For multiple examples, we visualize two quantities of the trajectory of the ODE. On the x-axis is the integration step of the ODE. On the y-axis on the left in blue, the RMSD between the step's \mathbf{x}_1 prediction and the final \mathbf{x}_1 prediction. On the y-axis on the right in yellow, we show the entropy of the predicted probability distribution averaged over all residues.

In their pipeline, they 1) Choose a protein with a cavity of the right size for the reactants. 2) Keep its backbone 3D structure without amino acid identities and part manually, part computationally decide the amino acids of the cavity to bind the ligands (which we aim to do with FLOWSITE). 3) Design the rest of the protein’s residues with existing tools such as ProteinMPNN (Dauparas et al., 2022).

The POCKETOPTIMIZER line of works (Malisi et al., 2012; Stiel et al., 2016; Noske et al., 2023) develops a pipeline for pocket residue design based on physics-inspired energy functions and search algorithms. Starting from a bound protein-ligand complex, POCKETOPTIMIZER samples different structures and residue types, which are scored with multiple options of energy functions and optimized with different search algorithms for which multiple options are available. Meanwhile, DEPACT (Chen et al., 2022b) and Dou et al. (2017) design pockets by searching databases for bound structures of similar ligands, using their protein’s residues as proposals, and selecting the best combinations based on scoring functions.

D.3. Antibody Design

Another domain where joint sequence and structure design has already been heavily leveraged is antibody design (Jin et al., 2022; Verma et al., 2023; Martinkus et al., 2023). In this task, the goal is to determine the residue types of the complementary determining regions/loops of an antibody to bind an epitope. These epitopes are proteins, and we have the opportunity to leverage evolutionary information. A modeling approach here only has to learn the interactions with the 20 possible amino acids that the epitope is built out of. Meanwhile, in our design task, where we wish to bind arbitrary small molecules, we are faced with a much wider set of possibilities for the ligand.

D.4. Small molecule design

Another frontier where designing structure and "2D" information simultaneously has found application is in molecule generation. For instance, Vignac et al. (2023a) and Vignac et al. (2023b) show how a joint diffusion process over a small molecule’s positions and its atom types can be used to successfully generate novel realistic molecules. This task was initially tackled by EDM (Hoogeboom et al., 2022) and recently was used to benchmark diffusion models with changing numbers of dimensions (Campbell et al., 2023).

Often, it is relevant to generate molecules conditioned on context. In particular, a highly valuable application, if it works well enough, would be generating molecules conditioned on a protein pocket to bind to that pocket (Lin et al., 2022; Schneuing et al., 2023). These applications would be most prominent in the drug discovery industry, where the first step in many drug design campaigns is often to find a molecule that binds to a particular target protein that is known to be relevant for a disease. In our work with FLOWSITE, we consider the opposite task where the small molecule is already given, and we instead want to design a pocket to bind this molecule. Here, the applications range from enzyme design (for which the first step of catalysis is binding the reactants (Nelson & Cox, 2004)) over antidote design to producing new biomedical marker proteins for use in medicinal diagnosis and biology research.

D.5. Protein-Ligand Docking

Historically, docking was performed with search-based methods (Trott & Olson, 2010; Halgren et al., 2004; Thomsen & Christensen, 2006) that have a scoring function and a search algorithm. The search algorithm would start with an initial random conformer and explore the energy landscape defined by the scoring function before returning the best scoring pose as the final prediction. Recently, such scoring functions have been parameterized with machine learning approaches (McNutt et al., 2021; Méndez-Lucio et al., 2021). In these traditional docking methods, to the best of our knowledge, only extensions of Autodock Vina (Trott & Olson, 2010) support multiligand docking. However, this still requires knowledge of the complete sidechains, which is not available in our binding site design scenario.

E. Experimental Setup Details

In this section, we provide additional details on how our experiments were run next to the exact commands and code to reproduce the results available at <https://github.com/HannesStark/FlowSite>. In all of the paper, we only consider heavy atoms (no hydrogens).

Training Details. For optimization, we use the Adam optimizer (Kingma & Ba, 2014) with a learning rate of 0.001 for all experiments. The batch size for pure structure prediction experiments is 4, while that for binding site recovery experiments

is 16. To choose the best model out of all training epochs, we run inference every epoch for experiments that do not involve structure modeling and every 5 epochs for the ones that do. The model that is used for the test set is the one with the best metric in terms of sequence recovery or fraction of predictions with an RMSD below 2 Å. When training for binding site recovery, we limit the number of heavy atoms in the ligand to 60. We note that for the structure prediction experiments for Binding MOAD in Table 2, the dataset construction for both methods had a mistake where ligands were selected based on their residue ID, which is incorrect because a ligand in a different chain could have the same residue ID - we will correct this in the next version of the manuscript. All models were trained on a single A100 GPU. The models that involve structure prediction were trained for 150 epochs, while those without structure modeling and pure sequence prediction converge much faster in terms of their validation metrics and are only trained for 50 epochs. On PDBBind, FlowSite took 58.6 hours to train, and on MOAD 115.6 hours, both on an RTX A600 GPU. The DIFFDOCK models are all trained for 500 epochs.

Runtime. Averaged over 4350 generated samples, the average runtime of HARMONICFLOW is 0.223 seconds per generated structure and that of FLOWSITE is 0.351 seconds per sequence.

Hyperparameters. We tuned hyperparameters on small-scale experiments in the Distance-Pocket setup for HARMONICFLOW and transferred these parameters to FLOWSITE, whose additional parameters we tested separately. The tuning for both methods was light, and we mainly stuck with the initial settings that we already found to work well. By default, our conditional probability path $p_t(\mathbf{x}|\mathbf{x}_0, \mathbf{x}_1) = \mathcal{N}(\mathbf{x}|t\mathbf{x}_1 + (1-t)\mathbf{x}_0, \sigma^2)$ uses $\sigma = 0.5$ for which we also tried 0.1, 0.3, 0.5, 0.8. The number of integration steps we use is 20 for all methods, including EIGENFOLD DIFFUSION and PRODUCT SPACE DIFFUSION. The number of scalar features we use is 32, and we have 8 vector features and 6 of our equivariant refinement TFN layers.

PRODUCT SPACE DIFFUSION baseline. This only uses the score model, the diffusion generative model component of DIFFDOCK (Corso et al., 2023). We do not use the confidence model, which is a significant part of their docking pipeline. Such a discriminator could also be used on top of HARMONICFLOW, and here, we only aim to compare the generative models. For this, we use the code at <https://github.com/gcorso/DiffDock> to train PRODUCT SPACE DIFFUSION with our pocket definitions using the same number of scalar features and vector features using 5 of its default TFN layers followed by its pseudo torque convolution and center-convolution. We train all experiments with PRODUCT SPACE DIFFUSION for 500 epochs.

EIGENFOLD DIFFUSION baseline. Here, we use an identical architecture as for HARMONICFLOW and only replace the flow matching training and inference with the diffusion training and inference approach of EIGENFOLD (Jing et al., 2023). The models were trained in the same settings, and most parameters that we use in HARMONICFLOW were first optimized with EIGENFOLD DIFFUSION since we used it initially.

F. Dataset Details

Here, we lay out the details of the PDBBind and Binding MOAD datasets as we use them to evaluate HARMONICFLOW’s docking abilities and FLOWSITE’s binding site design.

F.1. PDBBind

We use PDBBind dataset (Liu et al., 2017) with protein-ligand complexes of high binding affinity extracted and hand curated from the Protein Data Bank (PDB) (Berman et al., 2003). For this, we use two splits.

Splits. Firstly, the time split proposed by Stärk et al. (2022), which now is commonly used in the machine learning literature when benchmarking docking approaches, although Buttenschoen et al. (2023) among others found many shortcomings of this split, especially for blind docking. Chiefly among them is the fact that of the 363 test complexes, only 144 are not already included in the training data if a protein is counted the same based on UniProtID. The split has 17k complexes from 2018 or earlier for training/validation, and the mentioned 363 test samples are from 2019. Additionally, there is no ligand overlap with the training complexes based on SMILES identity. The data can be downloaded from <https://zenodo.org/record/6408497> as preprocessed by These files were preprocessed by Stärk et al. (2022) with Open Babel before ”correcting” hydrogens and flipping histidines with by running `reduce` <https://github.com/rlabduke/reduce>. For benchmarking traditional docking software, this preprocessed data should not be employed since the hydrogen bond lengths are incorrect. For our deep learning approaches that only consider heavy atoms, this is not relevant.

Secondly, a sequence similarity, which [Buttenschoen et al. \(2023\)](#) found to be a more difficult split than the time split for the blind docking scenario. To create this split, we cluster each chain of every protein with 30% sequence similarity. The clusters for training, validation, and test are then chosen such that each protein’s chains have at least 30% sequence similarity with any other chain in another part of the split. This way, we obtain 17741 train, 688 validation, and 469 test complexes. After filtering for complexes that have at least one contact (a protein residue with a heavy atom within 4Å), 17714 train complexes remain while no validation or test complexes are filtered out.

Dataset Statistics. In Figure 13, we show the number of atoms per ligand in two histograms, while Figure 14 shows the number of contacts (a protein residue with a heavy atom within 4Å) per ligand. These statistics are for the training data.

F.2. Binding MOAD Dataset

Split. The sequence similarity split that we use for BindingMOAD is carried out equivalently as for PDBBind described in Section F.1. This way, we obtain 56649 of Binding MOAD’s biounits for training, 1136 for validation, and 1288 as the test set. We discarded some of the biounits and only ended up with 54575 of them since 2.1k of them did not contain any other atoms besides protein atoms and waters. From these, we only use the complexes denoted as the first biounit to reduce redundancy and have only one biounit per PDB ID after which 38477 training complexes remain. We further filter out all ligands that have only one contact (a protein residue with a heavy atom within 4Å) with their protein to obtain 36203 train, 734 validation, and 756 test proteins with a unique PDB ID for each of them.

Dataset Statistics. Here, we provide statistics for the Binding MOAD training data. In Figure 15, we show the number of ligands per protein that is obtained under our definition of ligands and multi-ligands. Each ligand in the depicted histogram can either be a multi-ligand or a single molecule. Each multi-ligand is only counted once. In Figure 16, we show the number of atoms per ligand in two histograms, while Figure 17 shows the number of contacts per ligand.

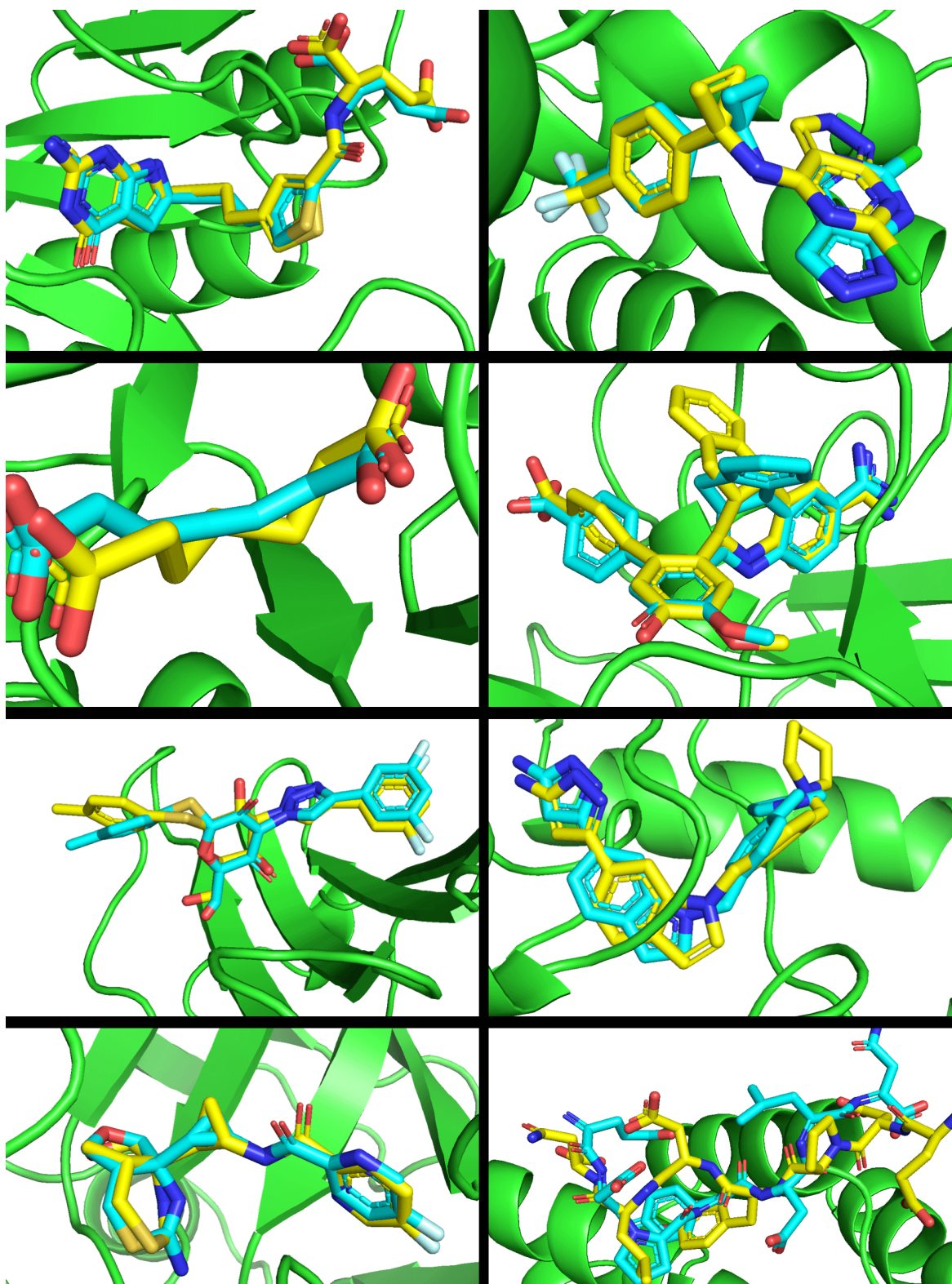


Figure 12. **HARMONICFLOW generated complexes.** Generated complexes of HARMONICFLOW for eight randomly chosen complexes in the PDBBind test set in the Distance-Pocket setup with a time-split where none of the ligands were seen during training.

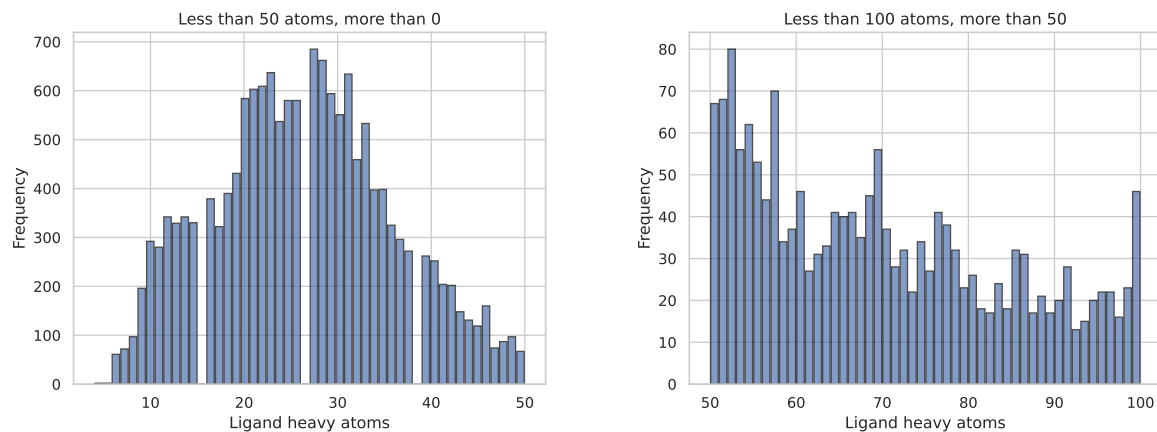


Figure 13. **Number of atoms per ligand: PDBBind.** Histograms showing the number of heavy atoms for all ligands under our ligand definition. This includes many ions, which can be important to filter out if not relevant to the desired application.

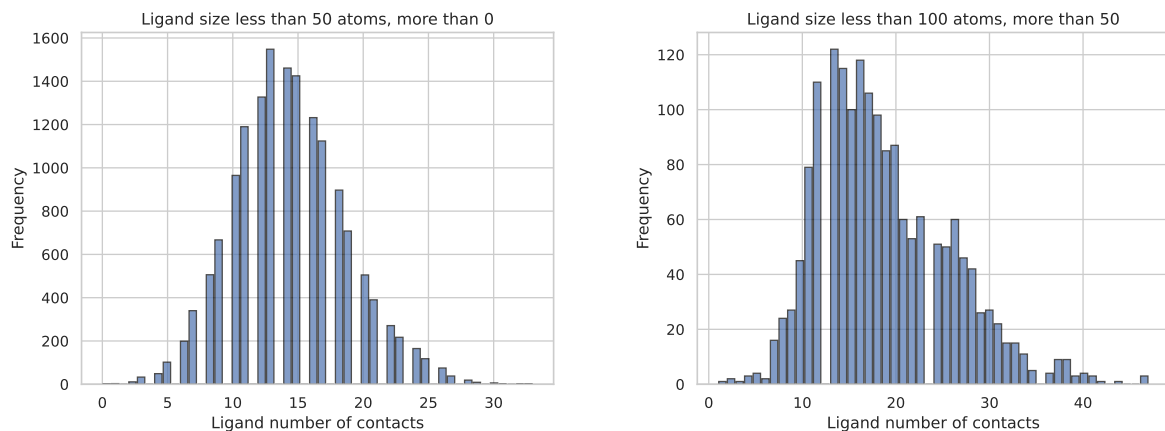


Figure 14. **Number of protein contacts per Ligand: PDBBind.** Histograms showing the number of contacts that each ligand has with its protein. A contact is defined as having a residue with a heavy atom within 4Å of any ligand heavy atom.

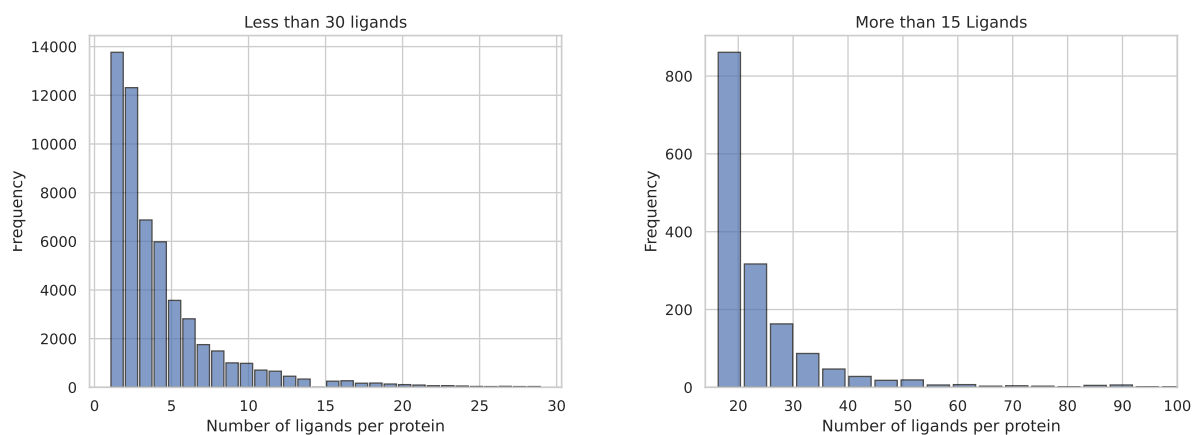


Figure 15. **Number of Ligands per Protein: Binding MOAD.** Histograms showing the number of (multi-)ligands per protein in the Binding MOAD dataset under our ligand definition. Each ligand here can be a multi-ligand. In that case, it is only counted once.

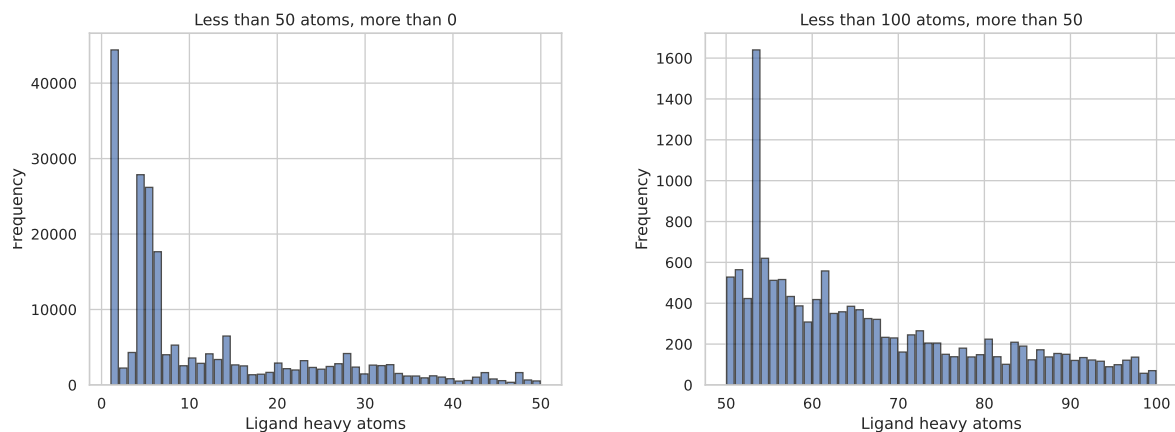


Figure 16. **Number of atoms per ligand: Binding MOAD.** Histograms showing the number of heavy atoms for all ligands under our ligand definition. This includes many ions, which can be important to filter out if not relevant to the desired application.

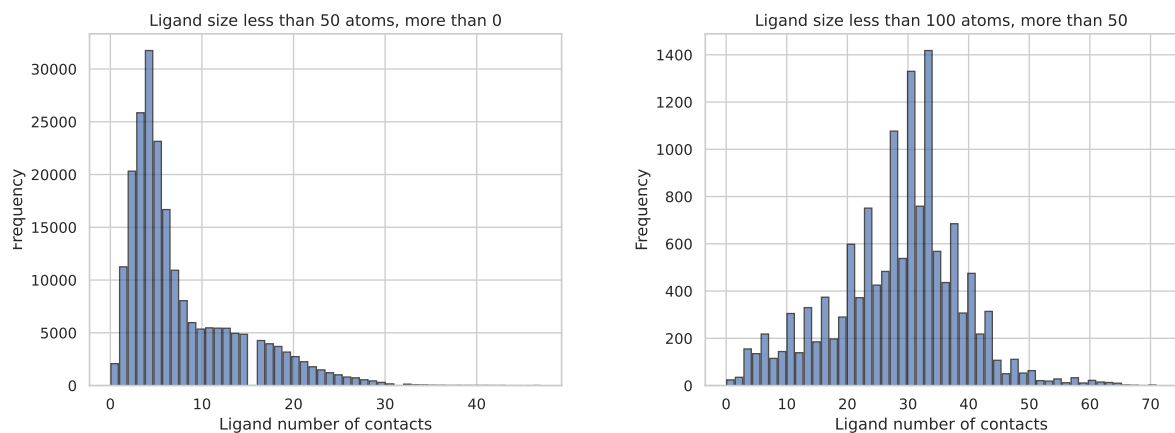


Figure 17. **Number of protein contacts per Ligand: Binding MOAD.** Histograms showing the number of contacts that each ligand has with its protein. A contact is defined as having a residue with a heavy atom within 4Å of any ligand-heavy atom.