

---

# ACER: Towards Generalizable Protein-ligand Co-folding

---

Anonymous Authors<sup>1</sup>

## Abstract

Predicting protein–ligand complex structures is a central challenge in drug discovery. While recent co-folding models such as AlphaFold-3 achieve accurate structure prediction, they fail to generalize to underexplored binding interfaces – systematically misplacing ligands, particularly for allosteric or structurally novel targets. To address this gap, we present **ACER** (**A**daptive **C**o-folding via pocket **E**xploration and pose **R**anking), a training-free framework that (a) enables co-folding models to systematically explore alternative binding pockets, and (b) leverages the discovered pockets to increase pose accuracy. Our method enables the efficient discovery of non-prevalent pockets without prior expert knowledge. ACER improves pocket discovery and pose accuracy on allosteric targets and structurally novel complexes, successfully modeling binding interfaces that are under-represented or absent from the training set. Our results demonstrate how improved sampling dynamics enhance the generalisability of co-folding models without retraining.

## 1. Introduction

Characterizing protein-ligand interactions at the structural level provides a fundamental understanding of how a small molecule modulates the biological function of one or more macromolecular targets (Zhao et al., 2022; Ciulli, 2013). Traditional structural biology techniques, *e.g.* X-ray crystallography and cryo-electron microscopy, can accurately resolve the structure of protein-ligand complexes, but they remain time-consuming, costly, and difficult to scale for high-throughput applications (Shi, 2014). Co-folding models have recently emerged as scalable computational tools to predict the structures of protein-ligand complexes and complement traditional experimental techniques (Abram-

son et al., 2024; team et al., 2024; OpenFold3 Team, 2024; Zhang et al., 2026; Passaro et al., 2025). Their competitive performance in both structure prediction and molecular docking (Durairaj et al., 2024) highlights their potential to enable high-throughput virtual screening and guide molecule design.

Despite these advances, co-folding models exhibit a critical failure mode (Škrinjar et al., 2025b): their tendency to strongly memorize the training data. The model’s predictive power highly correlates with the similarity between training and test samples, as further evidenced by the fact that models misplace allosteric ligands into orthosteric binding sites for protein targets whose allosteric sites are underrepresented (Nittinger et al., 2025; Parikh et al., 2026). Similarly, protein conformations are strongly affected by the ratio of holo- and apo-structures in the training set (Yu et al., 2026). These findings suggest that co-folding models have learned the template (‘memorization’) rather than the underlying physics of protein-ligand interactions.

One of the main bottlenecks of existing co-folding models is the incorrect ligand placement in the protein pocket, which is estimated to occur in allosteric complexes (Parikh et al., 2026; Nittinger et al., 2025) and in proteins that are underrepresented in the training data (Yu et al., 2026; Škrinjar et al., 2025b). Correct ligand placement is relevant, as most of the targets in drug discovery campaigns are structurally novel. This is particularly consequential for allosteric targets, where the functional pocket is often transient and only forms upon ligand-induced conformational change. This pocket memorization behavior therefore represents a critical gap in co-folding models.

To overcome this bottleneck, explicit pocket constraints have been proposed as guidance for co-folding (Passaro et al., 2025; Zhang et al., 2026). However, this technique requires prior knowledge about the correct pocket, which is usually unavailable for novel or out-of-distribution targets. Existing pocket identification tools (Le Guilloux et al., 2009; Krivák & Hoksza, 2018) predict binding-site residue probabilities in the absence of a ligand, which might miss ligand-induced conformational rearrangements (Meller et al., 2023). In this context, improving pocket discovery for novel targets holds untapped potential on the path to generalizable co-folding models.

---

<sup>1</sup>Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

In this work, we present **ACER** (Adaptive Co-folding for pocket Exploration and pose Ranking), a framework aimed at improving the capacity of co-folding models to discover and prioritize binding pockets – ultimately to retrieve more accurate protein-ligand structures. Our main contributions are:

1. **Inference-time pocket exploration**, which expands the search for binding pockets without retraining or major architectural modification. We implement two complementary strategies to enable the exploration the possible pocket landscape: (1) decoy ligand conditioning, which augments the input representation with auxiliary ligands that occlude default binding pockets; (2) iterative pocket repulsion, which applies negative gradients via Feynmann-Kac (FK) steering (Singhal et al., 2025). These strategies were implemented and tested on Boltz-2 (Passaro et al., 2025) and Protenix-v1 (Team et al., 2026).
2. **Ensemble-based ranking**, aimed at improving the prioritization of predicted ligand poses. This is achieved via normalized scores that derive from protein-ligand conformational ensembles, predicted under each pocket constraint.
3. **Improved generalization on hard co-folding benchmarks**. ACER successfully places ligands into experimentally observed binding pockets for allosteric compounds in those cases when state-of-the-art co-folding mistakenly places ligands into orthosteric pockets. When tested on the challenging subset of the Runs N’ Poses benchmark (Škrinjar et al., 2025a) – where binding interfaces are highly dissimilar to those in training set – this approach substantially outperforms the co-folding baselines.

## 2. Background and Related Works

**Generative models for biomolecular interactions.** Co-folding is built on conditional diffusion governed by two core components (Passaro et al., 2025; Abramson et al., 2024): (1) *trunk representations* that encode geometric features, and (2) a *denoiser* that iteratively recovers the clean structure from noise. Concretely, the model tokenizes protein sequences at the residue level and ligands at the atom level, yielding  $N = N_p + N_\ell$  tokens (protein residues and ligand atom tokens, respectively). The trunk produces a single representation  $\mathbf{s} \in \mathbb{R}^{N \times d}$  and a pair representation  $\mathbf{z} \in \mathbb{R}^{N \times N \times d}$ , which together form the conditioning variable  $\mathbf{Z} = (\mathbf{s}, \mathbf{z})$  for the denoiser. Recent models (Passaro et al., 2025) augment the sampling procedure with training-free Feynmann-Kac (FK) steering via Sequential Monte-Carlo (SMC), applying gradient updates to impose

physico-chemical constraints. Yet, FK steering has been used to drive sampling toward a known pocket, which renders it unsuitable for novel targets or allosteric complexes. While recent works applied Twisted Diffusion (Richman et al., 2026) and latent perturbation (Lee et al., 2026) to explore protein conformational landscapes, its potential for *exploring* alternative distinct ligand binding pockets remains untapped.

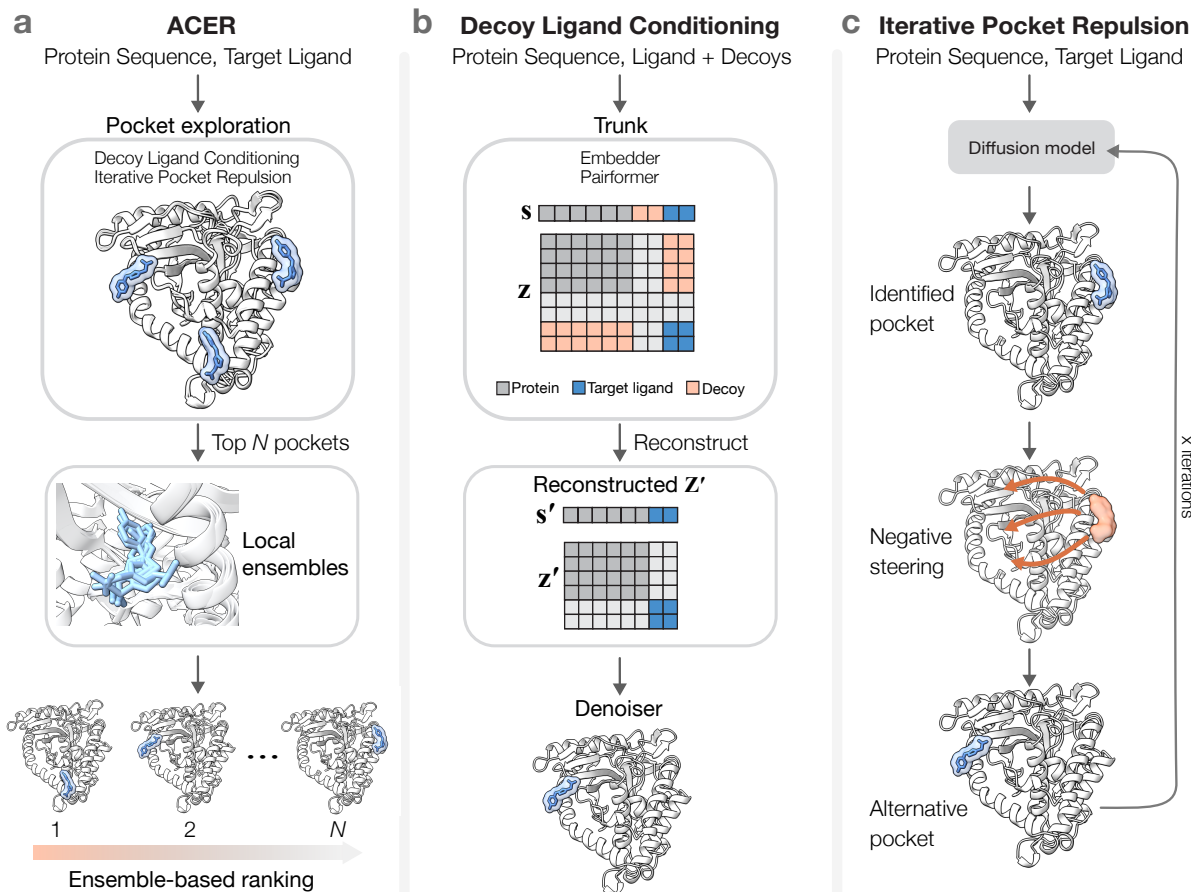
**Memorization in co-folding.** Protein-ligand co-folding models (Abramson et al., 2024; Passaro et al., 2025; OpenFold3 Team, 2024; Zhang et al., 2026) achieve state-of-the-art performance on structure-prediction and blind-docking benchmarks (Škrinjar et al., 2025b; Durairaj et al., 2024). Yet, accuracy drops as the targets diverge from ‘familiar’ biochemical space, and co-folded structures are sensitive to mutations at key binding-site residues (Masters et al., 2025). These findings point to strong memorization of the training data rather than genuine learning of biomolecular interaction physics – a fundamental issue for which the field lacks a principled solution.

**Pocket identification.** Conventional methods (Le Guilloux et al., 2009; Krivák & Hoksza, 2018) detect the pockets directly from protein structures in the absence of ligands, while more recent tools (Wang & Dokholyan, 2025) incorporate chemical probes and learned molecular fingerprints. All share the limitations of co-folding models: their accuracy drops over distribution shifts, e.g., for allosteric sites or unseen pockets where co-folding fails (Parikh et al., 2026). Naively combining the two, therefore, inherits rather than solves the generalization problem. Advances in pocket identification have been reported (Isomorphic Labs Team, 2026); however, the closed-source nature precludes any deeper analysis of gains in co-folding generalizability.

**Pose ranking in molecular docking.** Search-based methods such as AutoDock Vina (Trott & Olson, 2010) rely on physics-based scoring functions, treating pocket side-chains as rigid and ignoring ligand-induced fit effects. Co-folding instead ranks structures via confidence scores, which do not necessarily correlate with the success rate (Team et al., 2025). Recently, Molecular Dynamics (MD) has been combined with Markov state models to surface cryptic and allosteric pockets that are only visible in multi-conformational states (Biswas et al., 2025), but at the cost of expensive simulations.

## 3. Method

ACER addresses ligand misplacement in co-folding models through a two-stage framework (Fig. 1a): (1) *binding pocket exploration* and (2) *pose ranking* via an ensemble approach.



**Figure 1. Overview of ACER** (Adaptive Co-folding for pocket Exploration and pose Ranking). **(a)** ACER co-folds complexes across several candidate binding pockets. Each pocket then serves as a constraint for generating a local ensemble, from which top-ranked poses are selected using ensemble-based scores. **(b)** Given a protein  $\mathcal{P}$  and target ligand  $\ell$ , one or more decoy ligands  $\tilde{\ell}$  are introduced as additional trunk inputs. After processing by the Pairformer, the decoy tokens are discarded to reconstruct conditioning variable that co-folds  $\mathcal{P}$  and  $\ell$  into alternative pockets. **(c)** Inference-time guidance to direct the ligand placement away from the previously explored pockets. At each round, the previously explored pocket is added to a repulsion set  $\mathcal{B}$ , steering the denoiser away from previously explored sites.

### 3.1. Binding pocket exploration

We propose two strategies: (1) *decoy ligand conditioning* (Fig. 1b), which augments trunk inputs with auxiliary ligands to redirect the model’s distributional bias away from dominant pockets, and (2) *pocket repulsion* (Fig. 1c), which steers the denoising trajectory away from previously sampled binding sites. The former leverages the model’s implicit learning of interactions to explore alternative pockets in a parameter-free manner; the latter provides specific, user-defined steering. Together, they expand the effective search space of binding sites without retraining or architectural modifications.

#### 3.1.1. DECOY LIGAND CONDITIONING

Motivated by the evidence that co-folding with two identical allosteric ligands increases the likelihood of placing one into the correct pocket (Nittinger et al., 2025), we introduce

*decoy ligand conditioning* (Fig. 1b). Given a protein  $\mathcal{P}$  and a ligand  $\ell$ , we augment the trunk inputs with one or more decoy ligands  $\tilde{\ell}$  (comprising  $N_{\tilde{\ell}}$  atom tokens). These decoys can be molecules identical to  $\ell$  or generic chemical probes. The augmented inputs are processed by the Pairformer module to compute cross-token interactions between all inputs, producing the following conditioning representations:

$$(\mathbf{s}', \mathbf{z}') = \text{Trunk}(\mathcal{P}, \ell, \tilde{\ell}) \quad (1)$$

where  $\mathbf{s}' \in \mathbb{R}^{(N+N_{\tilde{\ell}}) \times d}$  contains token-level representations enriched by cross-token interactions, and  $\mathbf{z}' \in \mathbb{R}^{(N+N_{\tilde{\ell}}) \times (N+N_{\tilde{\ell}}) \times d}$  encodes pairwise interactions across all token pairs, including those involving  $\tilde{\ell}$ . We then reconstruct the conditioning variable by discarding all rows and columns corresponding to  $\tilde{\ell}$  tokens, yielding:

$$\mathbf{Z}' = (\mathbf{s}'_{\mathcal{P}, \ell}, \mathbf{z}'_{\mathcal{P}, \ell}) \quad (2)$$

where  $s'_{p,\ell} \in \mathbb{R}^{N \times d}$  and  $\mathbf{z}'_{p,\ell} \in \mathbb{R}^{N \times N \times d}$  have the same dimensions as  $\mathbf{Z}$  but are enriched by decoy-induced cross-token interactions. The denoiser generates structures conditioned on this reconstructed variable,  $\mathbf{x}_0 \sim p_\theta(\mathbf{x} \mid \mathbf{Z}')$ , such that  $\tilde{\ell}$  influences only the conditioning representations but does not participate in generation during denoising. Intuitively,  $\tilde{\ell}$  occludes the dominant pocket in  $\mathbf{z}'$ , reducing its dominant signals in  $\mathbf{Z}'$ , and freeing  $\ell$  to explore alternative binding sites.

We use four decoys that are identical to the considered ligand – native ligands. To assess whether non-native binders could improve pocket discovery, we additionally ablated the choice of decoy ligand type using three ‘frequent hitters’ chemical probes (Fig. A7), binding molecules that bind to diverse proteins across the Protein Data Bank (Wang & Dokholyan, 2025). Non-native binders offered no clear advantage over native ligands for pocket exploration.

### 3.1.2. ITERATIVE POCKET REPULSION

While decoy ligand conditioning influences representation learning, inference-time steering is a plug-in, architecture-agnostic strategy for controlling generation to respect specific constraints. To complement the FK-steering framework of Boltz-2 – which accounts for physico-chemical constraints (Passaro et al., 2025) – we introduce a *pocket repulsion potential*, which biases the denoising trajectory away from a specified pocket (via negative gradients), toward alternative binding contacts.

We penalise each ligand atom for being closer than a minimum distance ( $d_{\min}$ ) to the surface of a virtual sphere of radius  $r_g$  centred on the pocket anchor  $\mathbf{c}_g$ , where  $g$  denotes the pockets to steer away from. The pocket anchor  $\mathbf{c}_g$  is computed as the mean centroid of its constituent residue tokens  $t$ :

$$\boldsymbol{\mu}_t := \frac{1}{|\mathcal{A}_t|} \sum_{a \in \mathcal{A}_t} \mathbf{x}_a, \quad \mathbf{c}_g := \frac{1}{|\mathcal{T}_g|} \sum_{t \in \mathcal{T}_g} \boldsymbol{\mu}_t. \quad (3)$$

where  $\mathcal{A}_t$  is the set of atoms in  $t$ ,  $\mathbf{x}_a \in \mathbb{R}^3$  is the position of atom  $a$ , and  $\mathcal{T}_g$  denotes residues within 4 Å from the ligand atoms comprising pocket  $g$ .

The repulsion potential is then defined over all  $G$  pockets and all ligand atoms  $\mathcal{A}_{\text{lig}}$ , penalising any ligand atom that lies within a distance  $d_{\min}$  of the surface of the sphere centred on  $\mathbf{c}_g$  with the hyperparameter of virtual pocket radius  $r_g$ :

$$U_{\text{rep}} = \sum_{g=1}^G \sum_{a \in \mathcal{A}_{\text{lig}}} \max(0, d_{\min} - (\|\mathbf{x}_a - \mathbf{c}_g\|_2 - r_g))^2, \quad (4)$$

We apply these potentials iteratively. Each inference round

collects the occupied pocket as a new anchor, which is then added to the repulsion set  $\mathcal{B}$  for all subsequent rounds (Algorithm 1; full parameter settings in Appendix A.1).

## 3.2. Ensemble-based ranking

Once putative pockets have been identified, it is crucial to distinguish correctly docked poses from geometrically plausible but non-native ones (*i.e.*, those not reflecting the experimentally observed binding mode). Co-folding models generate a single static pose per run, which may miss the correct binding mode due to the intrinsic flexibility of protein-ligand complexes and the dynamics of ligand binding. We therefore generate local conformation ensembles around each candidate pocket — leveraging the fact that Boltz-2 was trained on MD simulations — and rank poses by their consistency with the local ensemble. For each pocket  $p$ , we collect  $n$  pose scores  $\{s_1, \dots, s_n\}$  to compute the z-score *i.e.*  $z_i = (s_i - \mu_p) / \sigma_p$  where  $i = 1, \dots, n$  indexes the poses within pocket  $p$ ,  $\mu_p$  and  $\sigma_p$  are the mean and standard deviation of pose scores within  $p$ .

We then weight each score by its proximity to the mean, either via a Gaussian or Lorentzian normalization, both of which peak at the ensemble mean and decay for outliers. The *Gaussian* weight penalises deviations sharply, whereas the *Lorentzian* weight down-weights outliers more gradually:

$$G_i = \exp\left(-\frac{z_i^2}{2}\right), \quad L_i = \frac{1}{1 + z_i^2}. \quad (5)$$

This ensemble-based weighting scheme derives from two intuitions: (a) a physically reasonable pose docked into the correct pocket should be relatively stable under small conformational perturbations – an outlier that disagrees with its local ensemble is unlikely to reflect the experimentally observed binding mode; and (2) ensemble consistency alone is insufficient – an absolute confidence score (*e.g.*, the interface predicted TM-score, ipTM (Jumper et al., 2021), which measures the predicted quality of the protein-ligand interface) is needed to assess the structural quality of the binding interface. Therefore, multiplying the confidence scores by  $G_i$  or  $L_i$  simultaneously accounts for local ensemble agreement and absolute interface quality.

## 4. Experiments

### 4.1. Study Setup

We evaluate ACER on two tasks: (1) pocket exploration, and (2) ensemble-based ranking. Each stage is assessed on two datasets that directly challenge co-folding models on generalizability (see Appendix A.1 for additional experimental details). Uncertainty is estimated from 1000 bootstrap iterations. All the results are reported in the mean and standard deviation for each metric and dataset.

## 4.1.1. DATASETS

We evaluate ACER on 2 datasets:

1. **Allosteric binders** (Nittinger et al., 2025), comprising proteins present in the training set, which is dominated by orthosteric ligands, i.e., ligands binding to the primary active site rather than to the secondary allosteric sites that modulate protein activity indirectly. This dataset consists of 20 allosteric ligands across 17 proteins.
2. **Runs N’ Poses** (Škrinjar et al., 2025b), consisting of structures entirely unseen by the models. We filter systems released after Boltz-2’s training cutoff (1 June 2023) and focus on failure cases where top-ranked predictions of default co-folding models misplace the ligand (center-of-mass distance  $> 2 \text{ \AA}$  from the X-ray structure) – yielding 88 and 55 protein-ligand pairs for Boltz-2 and Protenix-v1, respectively. ACER, applied to each base model, is evaluated on these two subsets separately to assess model-agnostic benefit. Results are reported under three evaluation settings: (a) the full failure-case set, (b) a filtered set where only the cluster representatives of ligand systems (based on pocket and ligand shape similarity, Škrinjar et al. (2025b)) are considered, and 124 prevalent ligands are excluded (Škrinjar et al., 2025b), and (c) the full scope without restricting to wrong-pocket cases ( $n = 217$ ) for benchmarking.

## 4.1.2. POCKET EXPLORATION

**Experimental design.** We apply the two pocket exploration strategies sequentially: (1) decoy ligand conditioning using five molecules (ligand plus four identical decoys) and five diffusion samples, followed by (2) pocket repulsion for five rounds. This setting guarantees the same sampling budget of 30 samples for each system. We then compare ACER extending Boltz-2 and Protenix-v1 against five baselines: (a) two ligand-free baselines based on volumetric or physicochemical properties, i.e., P2Rank (Krivák & Hoksza, 2018) and FPocket (Le Guilloux et al., 2009), and (b) three co-folding baselines: Boltz-2, OpenFold-3, and Protenix-v1.

**Metrics.** Success is measured by the Distance Center-to-Center (DCC) (Stepniewska-Dziubinska et al., 2020), i.e., the distance between the centers of mass (CoM) of the predicted and the ground-truth ligand (Appendix, Eq. A6). For P2Rank and FPocket – which do not co-fold complexes – DCC is measured between the predicted pocket center and the CoM of the ground-truth ligand. A pocket is considered correctly retrieved if the DCC falls below a given threshold ( $2 \text{ \AA}$  or  $3 \text{ \AA}$  in this work).

## 4.1.3. ENSEMBLE-BASED RANKING

**Experimental design.** We apply ACER on Boltz-2<sup>1</sup> and use the pair interface predicted TM-score (ipTM) (Jumper et al., 2021) as the baseline for ACER (see Appendix, Eq.A10) as a measure of the quality of the protein-ligand interface. We then follow a three-step ranking protocol: (1) We narrow the pocket candidate to  $P$  pockets to discard clearly unreliable interfaces. Using  $P = 10$  or  $20$  can capture almost all the recoverable pockets (Appendix A.1, Table A6). (2) Within each pocket constraint, we generate  $S$  local conformational ensembles, resulting in  $P \times S$  poses in total. (3) We rank the resulting poses by different criteria: (a) pair ipTM (used as a baseline), and (b) ensemble-normalized pair ipTM, using the Gaussian and Lorentzian weights (Eq. 5). Unless explicitly specified, we use  $P = 10$  and  $S = 5$ .

**Metrics.** We measure ranking success by the pose accuracy of the top- $N$  candidates sorted by each scoring metric, using ligand RMSD (L-RMSD) as the accuracy (Appendix, Eq. A7). The success rate is evaluated at  $2 \text{ \AA}$  and  $3 \text{ \AA}$  thresholds (as percentage of systems having an L-RMSD lower than the chosen threshold), combined with pose quality check by PoseBusters (Buttenschoen et al., 2024), ensuring that generated poses are not only geometrically close to the experimental structure but also physically plausible (see Appendix, A.1.3).

## 4.2. Results

## 4.2.1. POCKET EXPLORATION

**Allosteric binders.** Increasing the number of diffusion samples from the baseline models increases the likelihood of correct ligand placement into the right pocket (Table 1), suggesting that the underrepresented allosteric pockets are accessible but undersampled under default settings. ACER surpasses all baselines in terms of success rate (DCC lower than  $2 \text{ \AA}$  or  $3 \text{ \AA}$ , Table 1). Notably, ACER produces gains ranging from +11% to +13% on Boltz-2, and from +32% to +35% on Protenix-v1.

**Failed-pocket subsets.** ACER outperforms its respective base models (Table 2). ACER–Boltz-2 improves DCC success rate compared to Boltz-2 by +13.8% ( $2 \text{ \AA}$ ) and +14.4% ( $3 \text{ \AA}$ ) on the filtered subset ( $n = 46$ ), and by +7.3% and +7.9% on the full set ( $n = 88$ ), though all gains fall within the propagated uncertainty reflecting the limited sample size. Similarly, ACER–Protenix-v1 improves by +14.7% and +14.6% on the filtered subset ( $n = 30$ ) and by +8.8% and +9.7% on the full set ( $n = 55$ ) While not reaching statistical significance, the gains are consis-

<sup>1</sup>As Protenix-v1 does not include MD simulations in the training set, local ensemble generation is inherently unsupported, and, therefore, it was excluded from pose ranking evaluation.

Table 1. **Allosteric binding pocket identification** ( $n = 20$ ). DCC success rate (%): fraction of systems where at least one predicted pocket falls within 2 Å or 3 Å of X-ray ligand center of mass.

Method	DCC < 2 Å (%) <sup>†</sup>	DCC < 3 Å (%) <sup>†</sup>
P2Rank	25.0 ± 9.3	35.0 ± 10.9
FPocket	25.0 ± 9.9	55.0 ± 10.7
OpenFold-3 (preview-2)	35.0 ± 10.8	40.0 ± 10.9
Boltz-2 (× 5 samples)	40.0 ± 10.9	50.0 ± 10.9
Boltz-2 (× 30 samples)	55.0 ± 11.3	60.0 ± 10.9
<b>ACER – Boltz-2</b>	<b>66.6 ± 9.6</b>	<b>73.5 ± 9.1</b>
Protenix-v1 (× 5 samples)	25.0 ± 9.6	50.0 ± 11.4
Protenix-v1 (× 30 samples)	40.0 ± 10.8	50.0 ± 11.3
<b>ACER – Protenix-v1</b>	<b>72.9 ± 9.5</b>	<b>75.8 ± 9.3</b>

tent with the larger improvements observed on the allosteric dataset, suggesting a real but harder-to-detect effect on this more challenging subset. Finally, in the low-similarity bins (similarity  $\leq 40\%$ ), co-folding models show pocket rescue rates near zero (Appendix A.2, Fig. A4), underperforming the ligand-free methods. With similarity up to 40%, ACER – Boltz-2 improves the DCC success rate of the baseline from 2.9% to 35.3% compared to the baseline (30 samples, Fig. A3).

#### 4.2.2. ENSEMBLE-BASED RANKING

**Allosteric binders.** ACER – L-weight consistently performs on par with or outperforms all baselines. At Top-5, ACER Boltz-2 – L-weight gains +20.0% and +25.0% at 2 Å and +25.0% at 3 Å over the strongest baseline (Boltz-2). At Top-1, ACER matches Boltz-2. Unlike pair ipTM from static snapshots (Table 3), the L-weighted and G-weighted variants leverage local conformational ensembles, thereby recovering accurate poses that static-pose ranking would otherwise miss.

**Runs N’ Poses.** On the failed-pocket subset of the *Runs N’ Poses* benchmark (Table 4), ACER consistently outperforms the Boltz-2 baseline across both thresholds (gains at +8.5% and +4.6% on the filtered and full sets, respectively). ACER achieved the largest gains at Top-5 on the filtered subset (23.9%  $\rightarrow$  28.3% at 2 Å and 32.6%  $\rightarrow$  39.1% at 3 Å).

To benchmark ACER against all co-folding baselines, we extend the analysis to the full *Runs N’ Poses* set ( $n = 217$ , clustered and distinct ligands, Fig. 2), without restricting it to wrong-pocket cases. On this set, ACER achieves the highest Top-5 success rate (52.5%) among all co-folding models, outperforming Boltz-2 (49.8%), Protenix-v1 (51.2%), and OpenFold-3 (46.1%). At Top-1, ACER outperforms Protenix-v1 and OpenFold-3 but falls behind Boltz-2 (47.0%  $\rightarrow$  40.0%), suggesting that ensemble-based scoring improves candidate diversity but might struggle to promote

the best pose to the very top position. Notably, ACER is the only method to achieve a non-zero top-1 success rate (7%) in the lowest-similarity bin (up to 20% similarity) (Fig. 2a). When increasing the sampling budget to 200 samples ( $P = 20$  and  $S = 10$ ), ACER successfully samples accurate poses, reaching 21% at Top-20 and 28% at 200-oracle, demonstrating that its generalization can be improved with additional pocket candidates and higher coverage from local ensembles (Fig. 2b,c).

**Failure modes.** To better understand ACER’s benefits, we further analyze all the poses per system from local ensembles from ACER Boltz-2 – L-weighted for the Boltz-2 failed-pocket subset. We identify three main failure modes (Fig. A6): (1) *Correct pocket, low rank* (14.7% of cases). ACER successfully places ligands into the correct pockets but fails to rank them highly to top 5 candidates, indicating the limitation of the ranking procedure (Fig. A6a and Fig. A3). (2) *Correct ligand placement but wrong pose* (26.5% of cases). The ligand placement is correct, but an accurate pose is never generated; this is likely due to ineffective modeling of difficult protein–ligand interfaces by the baseline co-folding model (Fig. A6b). (3) *Correct pocket not rescued* (58.8% of cases), which highlights a persisting limitation of our pocket exploration (Fig. A6c).

#### 4.2.3. RUNTIME PERFORMANCE

ACER incurs additional compute overhead compared to the baselines: each decoy ligand in decoy ligand conditioning adds  $\sim 2\times$  wall-clock time (trajectories are conditioned on reconstructed cached representation) (Table A8 in Appendix), and iterative pocket repulsion introduces  $\sim 3\times$  overhead per round of guided sampling due to SMC resampling. Together, both approaches accumulate to  $\sim 5\times$  per diffusion sample in total. Ranking overhead is modest and comparable to the baseline. Despite this cost, ACER successfully rescues binding pockets and poses that are undersampled under default co-folding. ACER is hence positioned as a *precision-oriented tool* for challenging targets – such as allosteric or structurally novel systems – where accurate binding site characterization justifies the overhead.

## 5. Limitations

While ACER demonstrates consistent improvements on challenging targets, several avenues for improvement remain. The  $\sim 5\times$  runtime overhead currently limits applicability in high-throughput settings, though this could be reduced through more efficient sampling strategies or caching mechanisms. Pocket recovery on *Runs N’ Poses* reaches  $\sim 50\%$  within a budget of 30 samples – a promising result on hard targets, but one that leaves clear room for improvement as sampling budgets and exploration strategies mature.

Table 2. **Pocket-failure recovery.** ACER is evaluated in comparison with the wrongly identified pockets by each base model. ACER consistently correct places ligands into pockets by the base model, evaluated on the *Runs N' Poses* wrong-pocket subset (Škrinjar et al., 2025b). DCC success rate (%) at 2 Å and 3 Å thresholds is reported for the full set and a filtered subset.

Method	Boltz-2 failed-pocket subset				Protenix-v1 failed-pocket subset			
	Filtered ( $n = 46$ )		Total ( $n = 88$ )		Filtered ( $n = 30$ )		Total ( $n = 55$ )	
	DCC < 2 Å (%)↑	DCC < 3 Å (%)↑	DCC < 2 Å (%)↑	DCC < 3 Å (%)↑	DCC < 2 Å (%)↑	DCC < 3 Å (%)↑	DCC < 2 Å (%)↑	DCC < 3 Å (%)↑
P2Rank	24.9 ± 8.9	45.8 ± 10.2	22.9 ± 6.2	43.1 ± 7.3	20.9 ± 9.1	42.9 ± 11.2	23.8 ± 7.6	<b>49.7 ± 9.4</b>
FPocket	19.5 ± 5.9	37.0 ± 7.0	27.2 ± 4.8	42.1 ± 5.3	10.2 ± 5.6	36.7 ± 8.7	22.1 ± 5.8	40.8 ± 6.5
Boltz-2 (× 30 samples)	30.7 ± 6.6	37.3 ± 7.1	36.1 ± 5.2	42.9 ± 5.3	–	–	–	–
<b>ACER – Boltz-2</b>	<b>44.5 ± 7.1</b>	<b>51.7 ± 7.0</b>	<b>43.4 ± 5.0</b>	<b>50.8 ± 5.0</b>	–	–	–	–
Protenix-v1 (× 30 samples)	–	–	–	–	19.4 ± 6.9	30.1 ± 8.0	19.2 ± 5.1	27.7 ± 5.7
<b>ACER – Protenix-v1</b>	–	–	–	–	<b>34.1 ± 7.9</b>	<b>44.7 ± 7.8</b>	<b>28.0 ± 5.8</b>	37.4 ± 5.9

Table 3. **Pose ranking on the allosteric set.** Top- $k$  L-RMSD (%) at 2 Å and 3 Å thresholds is reported, considering only physically valid poses (Buttenschoen et al., 2024). *ACER Boltz-2 – baseline* uses raw pair ipTM as scoring function; *ACER Boltz-2 – L-weight* and *ACER Boltz-2 – G-weight* weight pair ipTM by the Lorentzian and Gaussian (Eq. 5) ensemble normalization, respectively.

Ranking Approach	% L-RMSD < 2 Å & PB Valid ↑			% L-RMSD < 3 Å & PB Valid ↑		
	Top-1	Top-5	Top-10	Top-1	Top-5	Top-10
Protenix-v1 (× 50 samples)	20.2 ± 9.02	24.9 ± 9.6	24.9 ± 9.6	30.3 ± 10.2	35.3 ± 10.6	40.5 ± 10.7
OpenFold-3 (× 50 samples)	20.2 ± 9.02	24.9 ± 9.6	24.9 ± 9.6	24.9 ± 9.6	24.9 ± 9.6	24.9 ± 9.6
Boltz-2 (× 50 samples)	<b>25.0 ± 9.6</b>	30.0 ± 10.1	30.0 ± 10.1	30.0 ± 10.2	40.0 ± 11.0	40.0 ± 11.0
<b>ACER Boltz-2 – baseline</b>	20.0 ± 9.0	25.0 ± 9.6	35.0 ± 10.6	30.0 ± 10.1	40.0 ± 10.8	45.0 ± 10.9
<b>ACER Boltz-2 – L-weight</b>	<b>25.0 ± 9.3</b>	<b>50.0 ± 11.1</b>	<b>50.0 ± 11.1</b>	30.0 ± 9.9	<b>65.0 ± 10.6</b>	<b>70.0 ± 10.3</b>
<b>ACER Boltz-2 – G-weight</b>	<b>25.0 ± 9.3</b>	<b>50.0 ± 11.1</b>	<b>50.0 ± 11.1</b>	<b>35.0 ± 10.3</b>	<b>65.0 ± 10.6</b>	<b>70.0 ± 10.3</b>

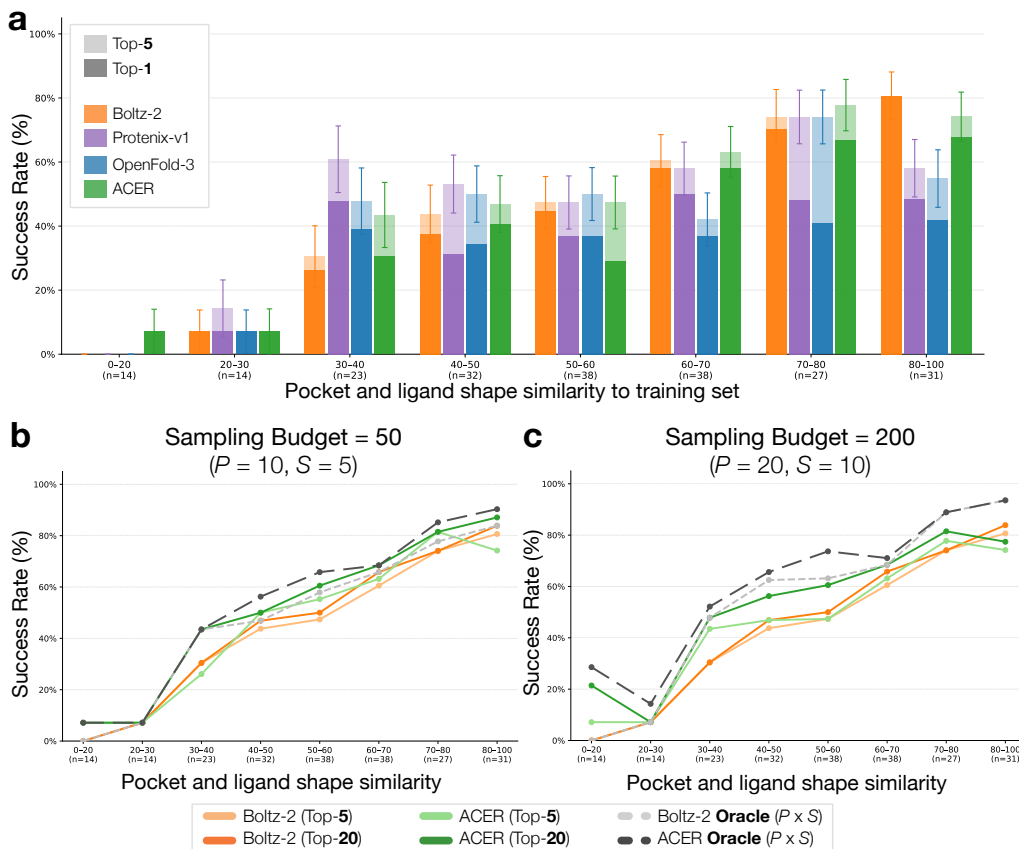
Table 4. **Pose ranking on the *Runs N' Poses* wrong-pocket subset of Boltz-2.** Top- $k$  L-RMSD (%) at 2 Å and 3 Å thresholds is reported, considering only physically valid poses (Buttenschoen et al., 2024). *ACER – baseline* uses raw pair ipTM as scoring function; *ACER – L-weight* and *ACER – G-weight* weight pair ipTM by the Lorentzian and Gaussian (Eq. 5) ensemble normalization.

Ranking Approach	Clustered & Distinct Ligands ( $n = 46$ )				Total ( $n = 88$ )			
	% L-RMSD < 2 Å & PB Valid ↑		% L-RMSD < 3 Å & PB Valid ↑		% L-RMSD < 2 Å & PB Valid ↑		% L-RMSD < 3 Å & PB Valid ↑	
	Top-1	Top-5	Top-1	Top-5	Top-1	Top-5	Top-1	Top-5
Boltz-2 (× 50 samples)	15.2 ± 5.5	15.2 ± 5.5	19.6 ± 6.0	19.6 ± 6.0	26.1 ± 4.8	27.3 ± 4.8	28.4 ± 4.9	29.5 ± 5.0
<b>ACER Boltz-2 – baseline</b>	<b>23.9 ± 6.4</b>	23.9 ± 6.4	<b>32.6 ± 7.1</b>	34.8 ± 7.2	<b>30.7 ± 5.0</b>	31.8 ± 5.1	<b>35.2 ± 5.2</b>	37.5 ± 5.3
<b>ACER Boltz-2 – L-weight</b>	15.2 ± 5.5	26.1 ± 6.7	21.7 ± 6.3	37.0 ± 7.2	21.6 ± 4.4	29.5 ± 4.9	26.1 ± 4.8	36.4 ± 5.3
<b>ACER Boltz-2 – G-weight</b>	15.2 ± 5.5	<b>28.3 ± 6.8</b>	21.7 ± 6.3	<b>39.1 ± 7.2</b>	21.6 ± 4.4	<b>31.8 ± 5.1</b>	26.1 ± 4.8	<b>38.6 ± 5.4</b>

Ensemble-based ranking consistently improves top-5 and beyond, and closing the gap to top-1 represents a natural next step, likely addressable through better scoring functions or learned ranking models. Finally, statistical validation remains challenging given the scarcity of allosteric and out-of-distribution targets – a limitation shared across the field, and one that underscores the need for larger, more diverse benchmarks for generalizable co-folding. Further evaluation on the recent dataset (?) would be a valuable next step.

## 6. Discussion

Co-folding models have emerged as a tool for predicting protein-ligand structures, yet their tendency to memorize the training set distributions can lead to systematic ligand misplacement on allosteric binding and out-of-distribution targets – precisely the most relevant cases for drug discovery. ACER addresses this through a training-free framework combining decoy ligand conditioning and iterative pocket repulsion to redirect the model’s sampling toward



**Figure 2. Ranking on Runs  $N'$  Poses**, measured as success rate (percentage of L-RMSD  $< 2 \text{ \AA}$  and PB valid). **(a)** ACER vs co-folding baselines on the clustered and distinct ligand subset ( $n = 217$ ). **(b, c)** Success rate of Boltz-2 and ACER at sampling budgets of **(b)** 50 and **(c)** 200 samples. Increasing the budget rescues accurate poses in the hardest similarity bin (0–20] particularly at top-20. ( $P$  = no. of pocket candidates to consider;  $S$  = no. of local conformation ensembles).

underrepresented binding sites. Its ensemble-based ranking further enables a principled comparison of binding poses across candidate pockets, providing an informative starting point for downstream refinement, such as energy minimization, specialized docking (Prat et al., 2026), or extensive MD simulations. Taken together, these results suggest that the generalizability gap in co-folding is not a fundamental limitation of the underlying models, but rather a sampling problem amenable to inference-time intervention. ACER demonstrates that redirecting the denoising trajectory – without touching model weights – can unlock binding sites that are systematically missed under default sampling. This reframing opens a broader research direction: as co-folding models continue to improve, inference-time exploration strategies like ACER may prove essential for extending their reach to the underrepresented pocketome. Looking beyond protein-ligand interactions, extending ACER’s rationale to protein-protein or protein-nucleic acid interfaces (where co-folding models are similarly susceptible to memorization) represents a natural and important future direction. ACER establishes a blueprint for training-free generalization in structural biology: rather than collecting more data

or retraining larger models, one can instead design smarter sampling dynamics that push existing models beyond their training distribution.

#### IMPACT STATEMENT

This work advances computational methods for protein-ligand structure prediction, with direct applications to structure-based drug discovery. By enabling co-folding models to explore the underrepresented pocketome, ACER accelerates the identification of novel therapeutic targets and pharmaceutical drugs that are currently inaccessible to standard computational pipelines.

#### REPRODUCIBILITY

To facilitate the reproducibility of our research, upon acceptance of this paper, we commit to making our code publicly available on GitHub (and on Zenodo as a static archive with a persistent DOI) to ensure that the community can fully reproduce our results.

## References

- Abramson, J., Adler, J., Dunger, J., Evans, R., Green, T., Pritzel, A., Ronneberger, O., Willmore, K., Ballard, A. J., Bambrick, J., et al. Accurate structure prediction of biomolecular interactions with AlphaFold 3. *Nature*, 630(8016):493–500, 2024. doi: 10.1038/s41586-024-07487-w. URL <https://doi.org/10.1038/s41586-024-07487-w>.
- Biasini, M., Schmidt, T., Bienert, S., Mariani, V., Studer, G., Haas, J., Johner, N., Schenk, A. D., Philippsen, A., and Schwede, T. Openstructure: an integrated software framework for computational structural biology. *Acta Crystallographica Section D: Biological Crystallography*, 69(5):701–709, May 2013. doi: 10.1107/S0907444913007051.
- Biswas, A. D., Sabato, E., Vittorio, S., Aletayeb, P., Pedretti, A., Mazzolari, A., Gratteri, C., Beccari, A. R., Talarico, C., and Vistoli, G. Novel method for prioritizing protein binding sites using pocket analysis and md simulations. *Heliyon*, 11(10):e43084, 2025. ISSN 2405-8440. doi: <https://doi.org/10.1016/j.heliyon.2025.e43084>. URL <https://www.sciencedirect.com/science/article/pii/S2405844025014653>.
- 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. *Chemical Science*, 15(9):3130–3139, 2024. ISSN 2041-6539. doi: 10.1039/d3sc04185a. URL <http://dx.doi.org/10.1039/D3SC04185A>.
- Ciulli, A. *Biophysical Screening for the Discovery of Small-Molecule Ligands*, pp. 357–388. Humana Press, Totowa, NJ, 2013. ISBN 978-1-62703-398-5. doi: 10.1007/978-1-62703-398-5\_13. URL [https://doi.org/10.1007/978-1-62703-398-5\\_13](https://doi.org/10.1007/978-1-62703-398-5_13).
- Durairaj, J., Adeshina, Y., Cao, Z., Zhang, X., Oleinikovas, V., Duignan, T., McClure, Z., Robin, X., Studer, G., Kovtun, D., Rossi, E., Zhou, G., Veccham, S., Isert, C., Peng, Y., Sundareson, P., Akdel, M., Corso, G., Stärk, H., Tauriello, G., Carpenter, Z., Bronstein, M., Kucukbenli, E., Schwede, T., and Naef, L. Plinder: The protein-ligand interactions dataset and evaluation resource. *bioRxiv*, 2024. doi: 10.1101/2024.07.17.603955. URL <https://www.biorxiv.org/content/early/2024/07/19/2024.07.17.603955.1>.
- Isomorphic Labs Team. Accurate predictions of novel biomolecular interactions with isodde. Technical report, Isomorphic Labs, February 2026. URL [https://storage.googleapis.com/isomorphiclabs-website-public-artifacts/isodde\\_technical\\_report.pdf](https://storage.googleapis.com/isomorphiclabs-website-public-artifacts/isodde_technical_report.pdf).
- Jumper, J., Evans, R., Pritzel, A., Kohli, P., Jaderberg, M., Chen, Z., Kirk, R., Dhar, A., Ahir, A., Wu, D., et al. Highly accurate protein structure prediction with AlphaFold. *Nature*, 596(7873):583–589, 2021. doi: 10.1038/s41586-021-03819-2.
- Krivák, R. and Hoksza, D. P2Rank: machine learning based tool for rapid and accurate prediction of ligand binding sites from protein structure. *Journal of Cheminformatics*, 10(1):39, 2018. doi: 10.1186/s13321-018-0285-8. URL <https://doi.org/10.1186/s13321-018-0285-8>.
- Le Guilloux, V., Schmidtke, P., and Tuffery, P. Fpocket: An open source platform for ligand pocket detection. *BMC Bioinformatics*, 10:168, 2009. doi: 10.1186/1471-2105-10-168. URL <https://doi.org/10.1186/1471-2105-10-168>.
- Lee, M., Kalicki, C., Jeon, M., Qabel, A., Fadini, A., and AlQuraishi, M. Conforntets: Latents-based conformational control in openfold3, 2026. URL <https://arxiv.org/abs/2604.18559>.
- Ma, W., Liu, Z., Yang, J., Lu, C., Zhang, H., and Xiao, W. From dataset curation to unified evaluation: Revisiting structure prediction benchmarks with pxmeter. *bioRxiv*, 2025. doi: 10.1101/2025.07.17.664878. URL <https://www.biorxiv.org/content/early/2025/07/22/2025.07.17.664878>.
- Masters, M. R., Mahmoud, A. H., and Lill, M. A. Investigating whether deep learning models for co-folding learn the physics of protein-ligand interactions. *Nature Communications*, 16(1):8854, 2025. doi: 10.1038/s41467-025-63947-5. URL <https://doi.org/10.1038/s41467-025-63947-5>.
- Meller, A., Ward, M. D., Borowsky, J. H., Lotthammer, J. M., Kshirsagar, M., Oviedo, F., Ferrer, J. L., and Bowman, G. Predicting the locations of cryptic pockets from single protein structures using the pocketminer graph neural network. *Biophysical journal*, 122(3):445a, 2023.
- Nittinger, E., Özge Yoluk, Tibo, A., Olanders, G., and Tyrchan, C. Co-folding, the future of docking – prediction of allosteric and orthosteric ligands. *Artificial Intelligence in the Life Sciences*, 8:100136, 2025. ISSN 2667-3185. doi: <https://doi.org/10.1016/j.aills.2025.100136>. URL <https://www.sciencedirect.com/science/article/pii/S2667318525000121>.
- Olanders, G., Testa, G., Tibo, A., Nittinger, E., and Tyrchan, C. Challenge for deep learning: Protein structure prediction of ligand-induced conformational changes at allosteric and orthosteric sites. *Journal of Chemical Information and Modeling*, 2024.

- 495 OpenFold3 Team. Openfold3-preview2 technical report.  
496 Technical report, Open Molecular Software Foundation,  
497 2024. URL [https://portal.openfold.omsf.io/reports/of3p2\\_technical\\_report.pdf](https://portal.openfold.omsf.io/reports/of3p2_technical_report.pdf).  
498 Accessed: May 2026.  
499
- 500 Parikh, V., Foley, B., Gatlin, W., Ludwick, M., Tu-  
501 rano, L., and Verkhivker, G. M. Decoding the  
502 allosteric paradox: A dual framework integrating  
503 ai cofolding models with landscape-guided inter-  
504 pretable ai framework of ligand-protein binding.  
505 *bioRxiv*, 2026. doi: 10.64898/2026.02.24.707829.  
506 URL [https://www.biorxiv.org/content/  
507 early/2026/02/26/2026.02.24.707829](https://www.biorxiv.org/content/early/2026/02/26/2026.02.24.707829).
- 509 Passaro, S., Corso, G., Wohlwend, J., Reveiz, M.,  
510 Thaler, S., Somnath, V. R., Getz, N., Portnoi, T.,  
511 Roy, J., Stark, H., Kwabi-Addo, D., Beaini, D.,  
512 Jaakkola, T., and Barzilay, R. Boltz-2: Towards  
513 accurate and efficient binding affinity prediction.  
514 *bioRxiv*, 2025. doi: 10.1101/2025.06.14.659707.  
515 URL [https://www.biorxiv.org/content/  
516 early/2025/06/18/2025.06.14.659707](https://www.biorxiv.org/content/early/2025/06/18/2025.06.14.659707).
- 518 Prat, A., Zhang, L., Deane, C. M., Teh, Y. W., and Morris,  
519 G. M. Sigmadock: Untwisting molecular docking with  
520 fragment-based se(3) diffusion, 2026. URL [https://  
521 arxiv.org/abs/2511.04854](https://arxiv.org/abs/2511.04854).
- 522 Richman, D. D., Karaguesian, J., Suomivuori, C.-M.,  
523 and Dror, R. O. Unlocking hidden biomolecular con-  
524 formational landscapes in diffusion models at infer-  
525 ence time, 2026. URL [https://arxiv.org/abs/  
526 2512.03312](https://arxiv.org/abs/2512.03312).
- 528 Shi, Y. A glimpse of structural biology through x-ray  
529 crystallography. *Cell*, 159(5):995–1014, 2014. ISSN  
530 0092-8674. doi: [https://doi.org/10.1016/j.cell.2014.10.  
531 051](https://doi.org/10.1016/j.cell.2014.10.051). URL [https://www.sciencedirect.com/  
532 science/article/pii/S0092867414014238](https://www.sciencedirect.com/science/article/pii/S0092867414014238).
- 534 Singhal, R., Horvitz, Z., Teehan, R., Ren, M., Yu, Z., McK-  
535 eown, K., and Ranganath, R. A general framework for  
536 inference-time scaling and steering of diffusion mod-  
537 els, 2025. URL [https://arxiv.org/abs/2501.  
538 06848](https://arxiv.org/abs/2501.06848).
- 539 Škrinjar, P., Eberhardt, J., Tauriello, G., Schwede,  
540 T., and Durairaj, J. Have protein-ligand cofolding  
541 methods moved beyond memorisation? *bioRxiv*,  
542 2025a. doi: 10.1101/2025.02.03.636309. URL  
543 [https://www.biorxiv.org/content/  
544 early/2025/08/04/2025.02.03.636309](https://www.biorxiv.org/content/early/2025/08/04/2025.02.03.636309).
- 546 Škrinjar, P., Eberhardt, J., Tauriello, G., Schwede,  
547 T., and Durairaj, J. Have protein-ligand cofolding  
548 methods moved beyond memorisation? *bioRxiv*,  
549 2025b. doi: 10.1101/2025.02.03.636309. URL  
[https://www.biorxiv.org/content/  
early/2025/08/04/2025.02.03.636309](https://www.biorxiv.org/content/early/2025/08/04/2025.02.03.636309).
- Stepniewska-Dziubinska, M. M., Zielenkiewicz, P.,  
and Siedlecki, P. Improving detection of protein-  
ligand binding sites with 3d segmentation. *Sci-  
entific Reports*, 10(1):5035, 2020. doi: 10.1038/  
s41598-020-61860-z. URL [https://doi.org/10.  
1038/s41598-020-61860-z](https://doi.org/10.1038/s41598-020-61860-z).
- team, C. D., Boitreaud, J., Dent, J., McPartlon, M.,  
Meier, J., Reis, V., Rogozhonikov, A., and Wu, K.  
Chai-1: Decoding the molecular interactions of life.  
*bioRxiv*, 2024. doi: 10.1101/2024.10.10.615955.  
URL [https://www.biorxiv.org/content/  
early/2024/10/11/2024.10.10.615955](https://www.biorxiv.org/content/early/2024/10/11/2024.10.10.615955).
- Team, G. R., Dobles, A., Jovic, N., Leidal, K., Murugan,  
P., Williams, D. C., Wulsin, D., Gruver, N., Ji, C. X.,  
Pruegsanusak, K., Scarpellini, G., Sharma, A., Swider-  
ski, W., Bootsma, A., Bowen, R. S., Chen, C., Chen,  
J., Dämgen, M. A., DiFrancesco, B., Fishman, J. D.,  
Ivanova, A., Kagin, Z., Li-Bland, D., Liu, Z., Morozov,  
I., Ouyang-Zhang, J., IV, F. C. P., Shah, K. S., Shor, B.,  
da Silva, G. M., Tal, R., Tessmer, M., Tilbury, C., Vetcher,  
C., Zeng, D., Al-Shedivat, M., Faust, A., Feinberg, E. N.,  
LeVine, M. V., and Pan, M. Pearl: A foundation model  
for placing every atom in the right location, 2025. URL  
<https://arxiv.org/abs/2510.24670>.
- Team, P., Zhang, Y., Gong, C., Zhang, H., Ma, W.,  
Liu, Z., Chen, X., Guan, J., Wang, L., Yang, Y.,  
Xia, Y., and Xiao, W. Protenix-v1: Toward high-  
accuracy open-source biomolecular structure prediction.  
*bioRxiv*, 2026. doi: 10.64898/2026.02.05.703733.  
URL [https://www.biorxiv.org/content/  
early/2026/02/22/2026.02.05.703733.1](https://www.biorxiv.org/content/early/2026/02/22/2026.02.05.703733.1).
- Trott, O. and Olson, A. J. AutoDock Vina: improving the  
speed and accuracy of docking with a new scoring func-  
tion, efficient optimization, and multithreading. *Journal  
of Computational Chemistry*, 31(2):455–461, 2010. doi:  
10.1002/jcc.21334.
- Wang, J. and Dokholyan, N. V. Unified protein-small  
molecule graph neural networks for binding site predic-  
tion. *bioRxiv*, 2025. doi: 10.1101/2025.09.03.674017.  
URL [https://www.biorxiv.org/content/  
early/2025/09/08/2025.09.03.674017](https://www.biorxiv.org/content/early/2025/09/08/2025.09.03.674017).
- Xie, J., Wang, S., Xu, Y., Deng, M., and Lai, L. Uncover-  
ing the dominant motion modes of allosteric regulation  
improves allosteric site prediction. *Journal of Chemical  
Information and Modeling*, 62(1):187–195, 2022. doi:  
10.1021/acs.jcim.1c01267.

550 Yu, H., Bekar-Cesaretli, A. A., Lazou, M., Kozakov, D.,  
551 Joseph-McCarthy, D., and Vajda, S. Bias in the Al-  
552 phaFold3 prediction of ligand-induced domain motion in  
553 enzymes. *Proceedings of the National Academy of Sci-*  
554 *ences*, 123(10):e2530709123, 2026. doi: 10.1073/pnas.  
555 2530709123. URL [https://doi.org/10.1073/](https://doi.org/10.1073/pnas.2530709123)  
556 [pnas.2530709123](https://doi.org/10.1073/pnas.2530709123).

557 Zhang, Y., Gong, C., Sun, J., Guan, J., Ren, M.,  
558 Xue, S., Zhang, H., Ma, W., Liu, Z., Chen, X.,  
559 and Xiao, W. Protenix-v2: Broadening the reach  
560 of structure prediction and biomolecular design.  
561 *bioRxiv*, 2026. doi: 10.64898/2026.04.10.717613.  
562 URL [https://www.biorxiv.org/content/](https://www.biorxiv.org/content/early/2026/04/11/2026.04.10.717613)  
563 [early/2026/04/11/2026.04.10.717613](https://www.biorxiv.org/content/early/2026/04/11/2026.04.10.717613).

564  
565 Zhao, L., Zhu, Y., Wang, J., Wen, N., Wang, C., and Cheng,  
566 L. A brief review of protein–ligand interaction prediction.  
567 *Computational and Structural Biotechnology Journal*, 20:  
568 2831–2838, 2022. doi: 10.1016/j.csbj.2022.06.004.  
569

570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593  
594  
595  
596  
597  
598  
599  
600  
601  
602  
603  
604

## A. Appendix

### A.1. Experimental Details

#### A.1.1. DATASET DETAILS

**Allosteric systems.** We adopt the same set of allosteric complexes from (Nittinger et al., 2025) where the details are provided in their supplementary materials. The dataset was originally published by (Xie et al., 2022) and (Olanders et al., 2024). To enhance the quality of the structures, the structures are removed if they are (1) close analogues binding to both orthosteric and allosteric binding site, (2) with different UniProt IDs, (3) where ligands bind to different proteins, and (4) with covalent ligands or peptides. The details of individual system is provided in Table A5.

Table A5. Protein targets with UniProt IDs, gene names, PDB IDs, and ligand CCDs for 20 allosteric compounds.

Protein Name	UniProt ID	Gene Name	AS Source PDBid	Ligand CCD
Serine/threonine protein kinase Chk1	O14757	CHEK1	3JVR	38M
Tyrosine-protein kinase ABL1	P00519	ABL1	3PYY	3YY, AY7
Tyrosine protein kinase ABL1	P00520	Abl1	3K5V	STJ
Insulin-like growth factor 1 receptor	P08069	IGF1R	3LW0	CCX
Cytochrome P450 3A4	P08684	CP3A4	1W0F	STR
Myosin II heavy chain	P08799	mhcA	2JHR	BIT
Androgen receptor	P10275	AR	2YHD	AV6, YLO
Tyrosine-protein phosphatase nonreceptor type 1	P18031	PTPN1	1T49	892
RAC-alpha serine/threonine-protein kinase	P31749	AKT1	3O96	0R4
Glucokinase	P35557	GCK	1V4S	MRK
Mitogen-activated protein kinase 8	P45983	MAPK8	3O2M	46A
Nuclear receptor ROR-gamma	P51449	RORC	4ypq	4F1
Kinesin-like protein KIF11	P52732	KIF11	3ZCW	2AZ, B4S
Beta-lactamase TEM	P62593	bla	1PZO	FTA
Casein kinase 2a	P68400	CSNK2A1	3H30	503
Focal adhesion kinase 1	Q05397	PTK2	4EBW	007
Mitogen-activated protein kinase 14	Q16539	MAPK14	3NEW	3NE

**Runs N’ Poses.** Runs N’ Poses is a well-established benchmark for protein–ligand generalization that groups samples by pocket similarity, pocket sequence identity, ligand similarity, and shape overlay (Škrinjar et al., 2025b). Because we apply ACER on top of Boltz-2, we restrict evaluation to systems deposited after its training cutoff of 1 June 2023. To isolate the wrong-pocket failure mode that ACER is designed to address, we further filter to systems where the baseline prediction misplaces the ligand, defined as a center-of-mass distance greater than 2 Å from the X-ray reference structure. The protein-ion pairs are excluded from the evaluation. The resulting test set is stratified by closest similarity to pre-cutoff training samples, allowing us to analyze performance as a function of target difficulty. All evaluations are reported on two subsets: (i) the full set of post-cutoff systems, and (ii) a clustered and filtered variant that removes 124 prevalent ligands as in (Škrinjar et al., 2025b), retaining only non-overlapping systems to reflect genuine gains in generalization.

To establish the protein-ligand generalizability benchmark, we extend our evaluation to the full scope of Runs N’ Poses. We remain focused on the clustered and distinct ligands subset, as it represents the most challenging subset for generalization and avoids double-counting successes on similar structures. Following the protocol prescribed by (Škrinjar et al., 2025b) and running the code from their official repository, we obtain a final set of 217 systems released after 1 June 2023.

#### A.1.2. PARAMETER SETTINGS

**Pocket exploration.** ACER is implemented on top of Boltz-2 and Protenix-v1, which we select because it is open-source and supports both inference-time steering and local ensemble generation. Exploration proceeds in two phases applied in

sequence. In the first phase, we run the inference of 5 diffusion samples with different number of decoy ligands (1 through 5), where each decoy is a duplicate of the target ligand. In the second phase, we perform 5 additional inference iterations with pocket repulsion, using  $c_g$  computed from pocket residues within 6 Å of ligand atoms,  $d_{\min} = 5$  Å, and  $r_g = 8$  Å. We report the mean and standard deviation of DCC success rate at which uncertainty estimates are computed over 5 seeds (seeds = 231234, 89567, 412903, 451027, 738291). We ablate the parameters associated with each method in detail in A.3

**Ensemble-based ranking.** For ranking tasks, we take the output from a single seed of the exploration stage and pocket constraints from residues whose atoms lie within 6 Å of any ligand atom from top  $P$  pockets ranked by pair ipTM. Leveraging those pocket constraints, We run the inference with  $S$  diffusion samples with MD method conditioning to generate local conformational ensembles. The total sampling budget is therefore  $P \times S$  samples. Unless explicitly specified, we use  $P = 10$  and  $S = 5$  in the main experiment. Unlike pocket exploration, we only report the results and corresponding uncertainty estimates on a single seed (seed = 738291). The ablation study on the total sampling budget is provided in Section A.3.

### A.1.3. METRICS

**Pocket recovery.** For multi-chain systems that contain duplicated protein-ligand pairs, we first superpose the protein chain of predicted structure onto the ground-truth chain by backbone alignment. Similarly to (Ma et al., 2025), we handle homodimers or copies of equivalent protein chains in the reference structure by permuting the alignment between predicted and ground-truth chains to obtain the best results, as formulated below.

$$\text{DCC} = \min_{c \in \text{chains}} \min_{\ell \in \text{ligands}} \|\mathbf{R}_c \bar{\mathbf{r}}_\ell + \mathbf{t}_c - \bar{\mathbf{r}}_{\text{lig}}\|_2, \quad (\text{A6})$$

where  $(\mathbf{R}_c, \mathbf{t}_c)$  is the rigid-body transform from the prediction frame to ground-truth chain  $c$  and  $\bar{\mathbf{r}}_\ell$  is the heavy-atom centroid of predicted ligand chain  $\ell$  (waters and hydrogens excluded). Minimum over  $c$  and  $\ell$  is the permutation-invariant best alignment; a system is counted as *rescued* at threshold  $\tau$  if at least one predicted ligand has  $\text{DCC} < \tau$ .

**Pose accuracy.** Per-pose ligand RMSD (L-RMSD) is computed as:

$$\text{L-RMSD} = \sqrt{\frac{1}{|\mathcal{A}_\ell|} \sum_{a \in \mathcal{A}_\ell} \|\mathbf{x}_a - \mathbf{x}_a^*\|_2^2} \quad (\text{A7})$$

where  $\mathcal{A}_\ell$  denotes the set of ligand heavy atoms,  $\mathbf{x}_a \in \mathbb{R}^3$  is the predicted position of atom  $a$ , and  $\mathbf{x}_a^* \in \mathbb{R}^3$  is its ground-truth position. L-RMSD is computed with OpenStructure (Biasini et al., 2013). The software first establishes a chain correspondence between prediction and reference using its IDDT-based chain-mapping algorithm. Symmetry-aware atom correspondences are enumerated via molecular-graph automorphism on the heavy-atom skeleton and the reported *symmetry-corrected ligand RMSD* is the minimum over all valid atom mappings after superposing the backbone carbon of the assigned protein chain. This yields a permutation-invariant pose-quality score, robust both to protein-chain renaming and to ligand internal symmetry.

**Pose quality.** Besides structural accuracy, we assess whether poses are physically plausible. We use PoseBusters plausibility checks using the PoseBusters package version 0.6.5. The pose is considered valid only if it passes all the following checks.

mol_pred_loaded	internal_steric_clash
mol_true_loaded	aromatic_ring_flatness
mol_cond_loaded	double_bond_flatness
sanitization	internal_energy
molecular_formula	minimum_distance_to_protein
molecular_bonds	minimum_distance_to_organic_cofactors
double_bond_stereochemistry	minimum_distance_to_inorganic_cofactors
tetrahedral_chirality	volume_overlap_with_protein
bond_lengths	volume_overlap_with_organic_cofactors
bond_angles	volume_overlap_with_inorganic_cofactors

**Pair ipTM.** ipTM score, introduced in (Jumper et al., 2021), quantifies the alignment quality between two molecular entities  $A$  and  $B$ . It is defined asymmetrically as:

$$\text{ipTM}(A \rightarrow B) = \max_{i \in A} [\text{mean}_{j \in B, \text{pTM}_{ij}}] \quad (\text{A8})$$

where the predicted TM score (pTM) computes the TM-score from predicted aligned errors (PAE):

$$\text{pTM} = \max_i \left[ \frac{1}{L} \sum_{j=1}^L \frac{1}{1 + \left( \frac{\text{PAE}_{ij}}{d_0} \right)^2} \right] \quad (\text{A9})$$

(Jumper et al., 2021; Abramson et al., 2024) provides ipTM for a pair of molecular entities. In protein-ligand complexes, it provides 2 values of directional ipTM. Therefore, we define the pair ipTM as:

$$\text{pair ipTM} = \max \left( \text{ipTM}_{\text{protein} \rightarrow \text{ligand}}, \text{ipTM}_{\text{ligand} \rightarrow \text{protein}} \right) \quad (\text{A10})$$

#### A.1.4. BASELINE DETAILS

**Pocket exploration.** To establish baselines for pocket exploration, we consider both traditional pocket detection methods and modern co-folding approaches. For the former, P2Rank (Krivák & Hoksza, 2018) predicts ligand binding sites using a random forest classifier trained on local physicochemical features of the protein surface. FPocket (Le Guilloux et al., 2009) detects cavities on the protein surface using Voronoi tessellation, without requiring any training data or ligand information. Both produce ranked lists of candidate pocket centers without explicit ligand placement. For these methods, DCC is measured between the predicted pocket center and the ground-truth ligand center of mass. For co-folding baselines, we use the same inference protocols as ACER for Boltz-2, Protenix-v1 (protenix\_base\_default.v1.0.0), and OpenFold-3 (preview-2) to generate 30 samples but without decoy ligand conditioning or pocket repulsion potential. For Boltz-2 and Protenix-v1, we also enable the default potentials to ensure physical plausibility of the poses.

**Ensemble-based ranking.** To benchmark the pose ranking scheme, we generate  $P \times S$  samples ( $P = 10, S = 5$ ) using Boltz-2, Protenix-v1, and OpenFold-3, ranked by each model’s confidence scores, maintaining equal sampling budget across all methods for a fair comparison with ACER. When reporting best-of-k or top-k results in tables and figures, we consider the best pose among top k candidates, which is both physically valid and has the lowest L-RMSD within the threshold.

## A.2. Extended Results

**Top pocket candidate selection.** We show how %DCC success rate varies by the top  $P$  pockets according to pair ipTM to discard clearly unreliable predicted interfaces. Table A6 shows how selected top  $P$  pockets captures recoverable pockets across subsets, achieving 100% rescue rates at Top-20 under both DCC criteria, while avoiding running local ensemble generation for all the pocket candidates, which can be computationally expensive in the subsequent ranking stage.

Table A6. Pocket rescue rate for the top- $N$  pockets ranked by pocket exploration strategies. DCC  $< 2 \text{ \AA}$  and  $< 3 \text{ \AA}$  denote the distance criteria for pocket recovery. Percentages are computed relative to the ceiling of recoverable pockets.

$N$	Allosteric ( $n = 20$ )		Clustered & Distinct Ligands ( $n = 46$ )		Total ( $n = 88$ )	
	DCC $< 2 \text{ \AA}$ (%) $\uparrow$	DCC $< 3 \text{ \AA}$ (%) $\uparrow$	DCC $< 2 \text{ \AA}$ (%) $\uparrow$	DCC $< 3 \text{ \AA}$ (%) $\uparrow$	DCC $< 2 \text{ \AA}$ (%) $\uparrow$	DCC $< 3 \text{ \AA}$ (%) $\uparrow$
1	35.7	33.3	52.4	52.2	57.9	56.5
5	71.4	80.0	66.7	73.9	73.7	82.6
10	92.9	100.0	85.7	91.3	86.8	93.5
20	100.0	100.0	100.0	100.0	100.0	100.0

**Multi-seed results.** To assess whether wrong-pocket predictions are consistent across random seeds, we repeat each inference 5 times and compute the union of successful placements. Concretely, a system is counted as successful if at

770 least one of the 5 predictions achieves a DCC below the threshold, regardless of the other seeds. As shown in Figure A5,  
771 increasing the number of samples modestly reduces the fraction of wrong-pocket systems compared to the single top-1  
772 prediction.

773 **Pocket and pose recovery by similarity bins.** Figure A3 shows the best L-RMSD per system stratified by pocket-ligand  
774 shape similarity to the training set, restricted to systems where the Boltz-2 baseline fails to identify the correct pocket (DCC  
775  $\geq 2$  Å). ACER rescues are most prevalent in the low-similarity regime (0–30), where Boltz-2 consistently misplaces ligands,  
776 and in several of these rescued cases ACER can achieve an accurate pose (L-RMSD  $< 2$  Å), despite the ranking bottleneck,  
777 as described in 4.2.2.

778 **Protein-ligand generalizability benchmark.** Table A7 details the success rate (% L-RMSD  $\leq 2$  Å or 3 Å thresholds  
779 and PB-Valid). While ACER’s top-1 performance is inferior to co-folding baselines, it offers clear advantages at Top-5  
780 and Top-10 success rate at which. it can sample accurate poses in the challenging similarity bins (0–30]. On the overall  
781 benchmark, ACER achieves 52.5% and 54.8% at Top-5 and Top-10 (L-RMSD  $\leq 2$  Å) with a similar trend at the 3 Å threshold  
782 (64.1% and 65.9%). This pattern indicates that ACER’s ranking effectively high-quality poses within a small candidate pool,  
783 even though its top-ranked pose does not lead.

784 **Runtime, compute and memory overhead** We run ACER primarily on NVIDIA A100-40GB GPUs. For decoy ligand  
785 conditioning, the additional cost arises from the extra tokens that must be processed by the trunk module, since triangle  
786 attention scales quadratically with token count, inflating both runtime and temporary activation memory. Using 1T49 as  
787 a representative structure, Table A8 shows the incremental runtime due to longer trunk processing. Additional diffusion  
788 forward passes conditioned on each reconstructed representation compounds runtime overhead accordingly (Table A9).  
789

790 For iterative pocket repulsion, memory footprint remains unchanged relative to the default, as the input token count is  
791 unaffected. However, FK steering incurs extra runtime due to a particle resampling step. Each generated sample under  
792 pocket repulsion requires  $3\times$  diffusion run time, compared to default sampling (Table A10). Both approaches therefore  
793 increase wall-clock time, with the cost scaling jointly with the number of diffusion forward passes and the resampling  
794 frequency. As a result, a greater number of exploration iterations leads to proportionally longer runtimes. We also report the  
795 overall runtime over Boltz-2 wrong-pocket dataset in Table A11

### 798 A.3. Ablation study

#### 799 A.3.1. POCKET EXPLORATION

800 We ablate each component of ACER on the Runs N’ Poses benchmark across 5 seeds (seeds = ) to isolate the effect of decoy  
801 ligand conditioning and iterative pocket repulsion. All ablations are performed exclusively on the Boltz-2 wrong-pocket  
802 subset, allowing us to directly measure the contribution of each ACER component to DCC success rate. Due to the  
803 computational costs, we omit the detailed ablation study on ACER – Protenix-v1 expected to yield the similar insights to  
804 Boltz-2 case. The results are reported with uncertainty estimates from 95% confidence intervals over 1,000 bootstraps.  
805

806 As shown in Table A12, decoy ligand conditioning alone yields substantially more correct pocket predictions than iterative  
807 pocket repulsion. Notably, pocket repulsion has negligible effect at the strict DCC $<2$ Å threshold, suggesting that repulsion  
808 guidance in the large search space is insufficient for precise ligand placement. However, by combining both components,  
809 ACER attains complementary gains, implying that pocket repulsion becomes effective when the search space is more  
810 constrained by occlusion effect introduced by decoy ligand conditioning. For a single seed (seed = 738291) due to  
811 computational costs, we therefore ablate the two components with higher iterations of pocket repulsion, demonstrating  
812 greater improvements due to increased iterations to substantiate its benefit on top of decoy ligand conditioning (Table A13)  
813

814 Table A14 shows the effect of the number of decoy ligands on the conditioning representation. Performance improves  
815 consistently with more decoys across both thresholds and both subsets, suggesting that a larger ensemble of decoy poses  
816 occludes more of the surface around the incorrect pocket, steering the model toward alternative binding interfaces that are  
817 underexplored in default sampling mode.  
818

819 Table A15 shows DCC success rates from decoy ligand conditioning, accumulated across decoy ligand types for 1-4 decoy  
820 ligands and 1 target ligand. Using decoy ligands identical to the target ligand performs comparably to small and medium  
821 chemical probes (KI2 and FER). In contrast, the bulkier PTY probe underperforms, likely because its steric bulk constrains  
822 the range of plausible binding sites. Since chemical probes offer no clear advantage in exploration, we use native ligands as  
823  
824

Table A7. Performance comparison of co-folding methods across sequence identity bins on the Runs N’ Poses benchmark ( $n = 217$  when top  $P = 10$  pocket candidates and  $S = 5$  samples of local ensembles are considered).

Similarity bin (n)	Ranking Approach	% L-RMSD < 2 Å & PB Valid ↑			% L-RMSD < 3 Å & PB Valid ↑		
		Top-1	Top-5	Top-10	Top-1	Top-5	Top-10
0–20 (14)	Boltz-2	0.0 ± 0.0	0.0 ± 0.0	0.0 ± 0.0	0.0 ± 0.0	0.0 ± 0.0	0.0 ± 0.0
	Protenix-v1	0.0 ± 0.0	0.0 ± 0.0	0.0 ± 0.0	0.0 ± 0.0	0.0 ± 0.0	0.0 ± 0.0
	OpenFold-3	0.0 ± 0.0	0.0 ± 0.0	7.1 ± 6.8	0.0 ± 0.0	7.1 ± 6.8	7.1 ± 6.8
	ACER Boltz-2 L-weight	<b>7.1 ± 7.1</b>	<b>7.1 ± 7.1</b>	<b>7.1 ± 7.1</b>	<b>7.1 ± 7.1</b>	<b>7.1 ± 7.1</b>	<b>14.3 ± 9.3</b>
20–30 (14)	Boltz-2	7.1 ± 6.8	7.1 ± 6.8	7.1 ± 6.8	7.1 ± 6.8	7.1 ± 6.8	7.1 ± 6.8
	Protenix-v1	7.1 ± 6.8	<b>14.3 ± 9.4</b>	<b>14.3 ± 9.4</b>	7.1 ± 6.8	<b>14.3 ± 9.4</b>	<b>14.3 ± 9.4</b>
	OpenFold-3	7.1 ± 6.8	7.1 ± 6.8	7.1 ± 6.8	7.1 ± 6.8	7.1 ± 6.8	7.1 ± 6.8
	ACER Boltz-2 L-weight	7.1 ± 6.8	7.1 ± 6.8	7.1 ± 6.8	<b>14.3 ± 9.4</b>	<b>14.3 ± 9.4</b>	<b>14.3 ± 9.4</b>
30–40 (23)	Boltz-2	26.1 ± 9.1	30.4 ± 9.5	30.4 ± 9.5	30.4 ± 9.4	34.8 ± 9.7	43.5 ± 10.2
	Protenix-v1	<b>47.8 ± 10.8</b>	<b>60.9 ± 10.1</b>	<b>69.6 ± 9.4</b>	<b>47.8 ± 10.8</b>	<b>60.9 ± 10.1</b>	<b>69.6 ± 9.4</b>
	OpenFold-3	39.1 ± 10.5	47.8 ± 10.4	52.2 ± 10.5	<b>47.8 ± 10.8</b>	<b>60.9 ± 10.2</b>	65.2 ± 10.0
	ACER Boltz-2 L-weight	17.4 ± 7.9	26.1 ± 8.9	39.1 ± 10.2	30.4 ± 9.5	47.8 ± 10.5	52.2 ± 10.4
40–50 (32)	Boltz-2	<b>37.5 ± 8.9</b>	43.8 ± 9.2	46.9 ± 9.2	<b>53.1 ± 9.2</b>	56.2 ± 9.1	59.4 ± 9.1
	Protenix-v1	31.2 ± 8.3	<b>53.1 ± 8.9</b>	<b>53.1 ± 8.9</b>	46.9 ± 8.7	62.5 ± 8.6	62.5 ± 8.6
	OpenFold-3	34.4 ± 8.9	50.0 ± 9.2	<b>53.1 ± 9.2</b>	46.9 ± 8.8	59.4 ± 8.8	59.4 ± 8.8
	ACER Boltz-2 L-weight	<b>37.5 ± 8.5</b>	50.0 ± 9.0	50.0 ± 9.0	<b>53.1 ± 8.7</b>	<b>71.9 ± 8.1</b>	<b>71.9 ± 8.1</b>
50–60 (38)	Boltz-2	<b>44.7 ± 8.0</b>	47.4 ± 8.1	50.0 ± 8.2	<b>68.4 ± 7.7</b>	68.4 ± 7.7	<b>71.1 ± 7.5</b>
	Protenix-v1	36.8 ± 8.0	47.4 ± 8.4	50.0 ± 8.4	36.8 ± 8.0	47.4 ± 8.4	52.6 ± 8.4
	OpenFold-3	36.8 ± 7.8	50.0 ± 7.7	55.3 ± 7.8	47.4 ± 8.0	55.3 ± 7.8	57.9 ± 7.7
	ACER Boltz-2 L-weight	31.6 ± 7.6	<b>55.3 ± 8.1</b>	<b>57.9 ± 8.1</b>	52.6 ± 8.2	<b>71.1 ± 7.5</b>	<b>71.1 ± 7.5</b>
60–70 (38)	Boltz-2	<b>57.9 ± 8.0</b>	60.5 ± 7.9	60.5 ± 7.9	<b>63.2 ± 7.7</b>	65.8 ± 7.6	65.8 ± 7.6
	Protenix-v1	50.0 ± 8.0	57.9 ± 8.1	<b>63.2 ± 8.0</b>	52.6 ± 8.0	63.2 ± 7.9	65.8 ± 7.9
	OpenFold-3	36.8 ± 7.6	42.1 ± 7.8	44.7 ± 7.9	52.6 ± 8.1	57.9 ± 8.1	57.9 ± 8.1
	ACER Boltz-2 L-weight	55.3 ± 8.2	<b>63.2 ± 7.8</b>	<b>63.2 ± 7.8</b>	<b>63.2 ± 7.8</b>	<b>71.1 ± 7.4</b>	<b>71.1 ± 7.4</b>
70–80 (27)	Boltz-2	<b>70.4 ± 8.7</b>	74.1 ± 8.3	74.1 ± 8.3	<b>85.2 ± 6.8</b>	<b>85.2 ± 6.8</b>	<b>85.2 ± 6.8</b>
	Protenix-v1	48.1 ± 9.9	74.1 ± 8.7	74.1 ± 8.7	51.9 ± 9.9	77.8 ± 8.3	77.8 ± 8.3
	OpenFold-3	40.7 ± 9.4	74.1 ± 8.4	74.1 ± 8.4	44.4 ± 9.5	77.8 ± 8.1	77.8 ± 8.1
	ACER Boltz-2 L-weight	63.0 ± 9.2	<b>81.5 ± 7.4</b>	<b>81.5 ± 7.4</b>	70.4 ± 8.7	<b>85.2 ± 6.8</b>	<b>85.2 ± 6.8</b>
80–100 (31)	Boltz-2	<b>80.6 ± 6.9</b>	<b>80.6 ± 6.9</b>	<b>83.9 ± 6.4</b>	<b>83.9 ± 6.5</b>	<b>83.9 ± 6.5</b>	<b>90.3 ± 5.3</b>
	Protenix-v1	48.4 ± 8.7	58.1 ± 8.7	61.3 ± 8.8	54.8 ± 9.0	67.7 ± 8.4	67.7 ± 8.4
	OpenFold-3	41.9 ± 8.8	54.8 ± 8.6	58.1 ± 8.6	51.6 ± 8.8	64.5 ± 8.5	67.7 ± 8.4
	ACER Boltz-2 L-weight	61.3 ± 8.9	74.2 ± 7.9	77.4 ± 7.7	64.5 ± 8.8	80.6 ± 7.2	87.1 ± 5.9
<b>Overall (217)</b>	Boltz-2	<b>47.0 ± 3.4</b>	49.8 ± 3.4	51.2 ± 3.4	<b>57.1 ± 3.4</b>	58.5 ± 3.4	61.3 ± 3.3
	Protenix-v1	38.2 ± 3.3	51.2 ± 3.4	53.9 ± 3.4	42.4 ± 3.3	55.3 ± 3.3	57.6 ± 3.3
	OpenFold-3	33.6 ± 3.2	46.1 ± 3.3	49.3 ± 3.3	42.9 ± 3.3	54.8 ± 3.3	56.2 ± 3.3
	ACER Boltz-2 L-weight	40.1 ± 3.4	<b>52.5 ± 3.5</b>	<b>54.8 ± 3.5</b>	50.7 ± 3.5	<b>64.1 ± 3.4</b>	<b>65.9 ± 3.4</b>

Table A8. Incremental trunk runtime and triangle attention memory footprint as a function of the number of decoy ligands processed by the trunk module, using 1T49 as a representative structure. Trunk cost is estimated as the difference between total forward pass time and diffusion time.

Number of ligands	Tokens	Tri. Attn., bf16 (GB)	Trunk overhead (s)
1	286	0.17	14.0 ± 0.4
5	378	0.29	31.5 ± 0.1
10	493	0.50	71.2 ± 0.2
20	608	0.76	199.3 ± 0.4

decoy ligands in this work, as their duplicate representations for conditioning can be cached once and reused across multiple diffusion runs.

We ablate the pocket repulsion algorithm parameters without decoy ligand conditioning to assess how DCC success rate is

Table A9. Runtime comparison between 2 diffusion samples from default sampling and decoy ligand conditioning where each sample was generated in a separate forward pass conditioned on 2 distinct reconstructed representations. Standard deviation is shown over 3 runs.

Mode	# Samples	Diffusion (s)	Wall-clock runtime (s)
Boltz-2 default: 1 forward $\times$ samples=2	2	21.2 $\pm$ 0.6	92.3 $\pm$ 8.5
ACER Boltz-2 – Decoy ligand conditioning: 2 forwards $\times$ samples=1	2	40.8 $\pm$ 0.5	154.0 $\pm$ 1.7

Table A10. Runtime comparison between Boltz-2 default sampling and ACER with pocket repulsion. Standard deviation is shown over 3 runs.

Protein	# Res	Number of diffusion samples	Diffusion Time (s)	
			Default Boltz-2	ACER Boltz-2
4YPQ	243	1	20.9 $\pm$ 0.4	64.2 $\pm$ 0.8
		2	21.4 $\pm$ 0.0	85.1 $\pm$ 1.0
		4	20.9 $\pm$ 0.2	137.4 $\pm$ 1.1
		16	49.0 $\pm$ 0.0	823.4 $\pm$ 14.3
3F9N	322	1	20.5 $\pm$ 0.4	66.6 $\pm$ 2.1
		4	21.0 $\pm$ 0.3	140.1 $\pm$ 1.3
5MO4	495	1	20.3 $\pm$ 0.3	65.1 $\pm$ 0.4
		4	31.6 $\pm$ 0.0	151.6 $\pm$ 1.8

Table A11. Wall-clock runtime breakdown on Boltz-2 wrong-pocket subset for decoy ligand conditioning and iterative pocket repulsion where each diffusion forward pass uses 200 sampling steps. Runtime is reported as mean  $\pm$  std over benchmark tasks.

Stage	Config	Sampling steps	Runtime (s)
Decoy ligand	$n_{\text{ligand}} = 1$	200	37 $\pm$ 9
Decoy ligand	$n_{\text{ligand}} = 2$	400	74 $\pm$ 18
Decoy ligand	$n_{\text{ligand}} = 3$	600	111 $\pm$ 27
Pocket repulsion	$n_{\text{rounds}} = 1$	200	128 $\pm$ 25
Pocket repulsion	$n_{\text{rounds}} = 4$	800	410 $\pm$ 94
Pocket repulsion	$n_{\text{rounds}} = 5$	1000	524 $\pm$ 123
Pocket repulsion	$n_{\text{rounds}} = 8$	1600	828 $\pm$ 194

Table A12. Ablation of ACER components on the wrong-pocket subset of Runs N’ Poses benchmark for Boltz-2 over 5 seeds (seeds = 231234, 89567, 412903, 451027, 738291)

Method	Clustered & Distinct Ligands ( $n = 46$ )		Total ( $n = 88$ )	
	DCC < 2 Å (%) $\uparrow$	DCC < 3 Å (%) $\uparrow$	DCC < 2 Å (%) $\uparrow$	DCC < 3 Å (%) $\uparrow$
Boltz-2 ( $\times 30$ samples)	30.4 $\pm$ 7.0	37.0 $\pm$ 7.2	36.4 $\pm$ 5.1	43.2 $\pm$ 5.3
Only decoy ligand conditioning (1 – 4 decoy ligands)	43.9 $\pm$ 6.8	50.4 $\pm$ 6.6	42.5 $\pm$ 5.0	48.4 $\pm$ 4.8
Only 5-round pocket repulsion	34.8 $\pm$ 6.9	47.8 $\pm$ 7.6	38.6 $\pm$ 4.9	47.7 $\pm$ 5.2
<b>Decoy ligand conditioning + 5-round pocket repulsion</b>	<b>44.3 <math>\pm</math> 6.5</b>	<b>51.7 <math>\pm</math> 6.7</b>	<b>43.6 <math>\pm</math> 5.0</b>	<b>51.1 <math>\pm</math> 5.1</b>

sensitive to the parameters in guidance sampling, including the virtual pocket radius  $r_g$  (Table A16), the minimum distance threshold  $d_{\text{min}}$  (Table A17), and optional protein-ligand contact guidance that preserves protein-ligand contact (Table A18). The results are derived from a single seed (seed = 738291). Uncertainty is estimated by 1000 bootstrap resampling.

Table A16 reports DCC success rate by varying the virtual pocket radius  $r_g$  used to define the repulsion region. Larger  $r_g$  improve the likelihood of correct ligand placement into alternative pockets, with  $r = 8$  achieving the best results. This suggests that larger repulsion against the already explored regions is more effective to discover the new binding interface.

Table A13. Ablation of ACER components on the wrong-pocket subset of Runs N' Poses benchmark for Boltz-2 on a single seed = 738291

Method	Clustered & Distinct Ligands ( $n = 46$ )		Total ( $n = 88$ )	
	DCC<2 Å (%)↑	DCC<3 Å (%)↑	DCC<2 Å (%)↑	DCC<3 Å (%)↑
Boltz-2 ( $\times 30$ samples)	30.4 $\pm$ 6.6	37.0 $\pm$ 7.1	36.4 $\pm$ 5.2	43.2 $\pm$ 5.3
Only decoy ligand conditioning (1 – 4 decoy ligands)	45.7 $\pm$ 7.3	50.0 $\pm$ 7.4	43.2 $\pm$ 5.4	52.3 $\pm$ 5.4
<b>Decoy ligand conditioning + 5-round pocket repulsion</b>	45.7 $\pm$ 7.2	52.2 $\pm$ 7.4	43.2 $\pm$ 5.2	54.5 $\pm$ 5.4
<b>Decoy ligand conditioning + 15-round pocket repulsion</b>	<b>47.8 <math>\pm</math> 7.6</b>	<b>54.3 <math>\pm</math> 7.2</b>	<b>46.6 <math>\pm</math> 5.2</b>	<b>58.0 <math>\pm</math> 5.1</b>

Table A14. Ablation as a function of the number of decoy ligands on the Runs N' Poses and Allosteric Systems benchmarks on a single seed = 738291

# Total ligands	# Decoy ligands	Clustered & Distinct Ligands ( $n = 46$ )		Total ( $n = 88$ )		Allosteric Systems ( $n = 20$ )	
		DCC<2 Å (%)↑	DCC<3 Å (%)↑	DCC<2 Å (%)↑	DCC<3 Å (%)↑	DCC<2 Å (%)↑	DCC<3 Å (%)↑
2	1	41.4 $\pm$ 7.3	47.8 $\pm$ 7.4	36.9 $\pm$ 5.1	46.6 $\pm$ 5.3	44.7 $\pm$ 11.2	50.0 $\pm$ 10.7
3	2	45.6 $\pm$ 7.4	50.5 $\pm$ 7.2	41.2 $\pm$ 5.3	48.8 $\pm$ 5.4	54.4 $\pm$ 10.8	59.6 $\pm$ 10.9
4	3	<b>45.9 <math>\pm</math> 7.5</b>	49.9 $\pm$ 7.3	42.3 $\pm$ 5.1	50.2 $\pm$ 5.3	55.1 $\pm$ 11.5	60.2 $\pm$ 10.9
5	4	45.8 $\pm$ 7.5	<b>50.3 <math>\pm</math> 7.3</b>	<b>43.3 <math>\pm</math> 5.5</b>	<b>51.3 <math>\pm</math> 5.2</b>	<b>55.1 <math>\pm</math> 11.2</b>	<b>65.4 <math>\pm</math> 10.6</b>

Table A15. Ablation of the choice of decoy ligand type on DCC success rate (%) at 2 Å and 3 Å thresholds.

Decoy ligand type	Clustered & Distinct Ligands ( $n = 46$ )		Total ( $n = 88$ )	
	DCC<2 Å (%)↑	DCC<3 Å (%)↑	DCC<2 Å (%)↑	DCC<3 Å (%)↑
Native ( $5 \times$ native)	<b>45.3 <math>\pm</math> 7.2</b>	50.2 $\pm$ 7.3	<b>43.2 <math>\pm</math> 5.1</b>	<b>52.6 <math>\pm</math> 5.2</b>
KI2	26.1 $\pm$ 6.6	<b>35.1 <math>\pm</math> 6.8</b>	29.7 $\pm$ 5.0	38.5 $\pm$ 5.4
PTY	15.2 $\pm$ 5.3	28.3 $\pm$ 6.7	21.7 $\pm$ 4.5	34.0 $\pm$ 5.1
FER	25.9 $\pm$ 6.5	30.3 $\pm$ 6.6	25.1 $\pm$ 4.6	31.8 $\pm$ 4.9

Table A16. Ablation as a function of the pocket radius on the Runs N' Poses benchmark over 5 iterations of pocket repulsion. Boltz-2 baseline runs default sampling without decoy ligand conditioning

Pocket radius	Clustered & Distinct Ligands ( $n = 46$ )		Total ( $n = 88$ )	
	DCC<2 Å (%)↑	DCC<3 Å (%)↑	DCC<2 Å (%)↑	DCC<3 Å (%)↑
Boltz-2 ( $\times 5$ samples without pocket repulsion)	26.1 $\pm$ 6.5	34.9 $\pm$ 7.4	30.6 $\pm$ 5.1	36.6 $\pm$ 5.2
4 Å	28.4 $\pm$ 6.5	32.4 $\pm$ 6.8	34.1 $\pm$ 5.2	38.6 $\pm$ 4.9
6 Å	<b>32.5 <math>\pm</math> 6.8</b>	36.8 $\pm$ 7.1	<b>36.4 <math>\pm</math> 4.8</b>	40.9 $\pm$ 5.2
8 Å	28.3 $\pm$ 6.8	<b>41.4 <math>\pm</math> 7.4</b>	31.7 $\pm$ 5.1	<b>40.9 <math>\pm</math> 5.3</b>

Table A17 shows DCC success rate by the minimum distance threshold  $d_{\min}$  controlling how far decoy ligands must be placed from the restricting binding pocket. Performance degrades at smaller values ( $d_{\min} \in \{4, 6\}$ ), indicating that ligands placed too close to the previous pocket do not effectively explore alternative pockets.

As applying the repulsion potential alone risks pushing the ligand entirely away from the protein surface, we experiment introducing the contact potential  $U_{\text{protein-contact}}$  in Equation A11 as a safeguard to keep the ligand in contact with protein but not so close to cause steric clashes. We therefore ablate the contact weight added to the existing potential where  $\lambda_c$  to quantify its effect. Table A18 reports the effect of varying the contact weight hyperparameter. Disabling contact weighting

Table A17. Ablation as a function of the minimum distance threshold  $d_{\min}$  on the Runs N’ Poses benchmark over 5 iterations of pocket repulsion. Boltz-2 baseline runs default sampling without decoy ligand conditioning

$d_{\min}$	Clustered & Distinct Ligands ( $n = 46$ )		Total ( $n = 88$ )	
	DCC < 2 Å (%)↑	DCC < 3 Å (%)↑	DCC < 2 Å (%)↑	DCC < 3 Å (%)↑
Boltz-2 ( $\times 5$ samples without pocket repulsion)	26.1 $\pm$ 6.5	34.9 $\pm$ 7.4	30.6 $\pm$ 5.1	36.6 $\pm$ 5.2
4 Å	23.8 $\pm$ 6.2	32.4 $\pm$ 7.0	30.7 $\pm$ 4.8	35.5 $\pm$ 5.0
6 Å	28.1 $\pm$ 6.9	34.9 $\pm$ 7.0	<b>34.9 <math>\pm</math> 5.3</b>	38.4 $\pm$ 4.9
8 Å	<b>28.1 <math>\pm</math> 6.8</b>	<b>41.6 <math>\pm</math> 7.5</b>	31.7 $\pm$ 4.9	<b>41.0 <math>\pm</math> 5.3</b>

entirely ( $\lambda_c = 0.0$ ) performs similarly on DCC success rate < 2 Å threshold, suggesting that conditioning representations that direct the denoiser already implicitly enforce structural priors to maintain protein-ligand contact. The protein-ligand contact potential provides limited benefit in this setting.

$$U_{\text{protein-contact}} := \sum_{a \in \mathcal{A}_{\text{lig}}} \max \left( 0, \min_{p \in \mathcal{A}_p} \|\mathbf{x}_a - \mathbf{x}_p\|_2 - d_c \right)^2 \quad (\text{A11})$$

where  $\mathcal{A}_{\text{lig}}$  and  $\mathcal{A}_p$  denote the sets of ligand and protein heavy atoms, respectively,  $\mathbf{x}_p \in \mathbb{R}^3$  is the position of protein atom  $p$ , and  $d_c$  is the maximum allowed distance to the nearest protein atom. The term is added to the existing potentials<sup>2</sup>.

$$U_{\text{total}} := U_{\text{repulsion}}(\mathcal{B}) + \lambda_c U_{\text{contact}} \quad (\text{A12})$$

where  $\lambda_c \geq 0$  is a scalar weight controlling the strength of the contact term. Setting  $\lambda_c = 0$  refers to pure repulsion steering.

Table A18. Ablation as a function of the contact weight on the Runs N’ Poses benchmark over 5 iterations of pocket repulsion. Boltz-2 baseline runs default sampling without decoy ligand conditioning.

Contact weight	Clustered & Distinct Ligands ( $n = 46$ )		Total ( $n = 88$ )	
	DCC < 2 Å (%)↑	DCC < 3 Å (%)↑	DCC < 2 Å (%)↑	DCC < 3 Å (%)↑
Boltz-2 ( $\times 5$ samples without pocket repulsion)	26.1 $\pm$ 6.5	34.9 $\pm$ 7.4	30.6 $\pm$ 5.1	36.6 $\pm$ 5.2
0.0	<b>30.7 <math>\pm</math> 7.0</b>	41.5 $\pm$ 7.1	<b>34.3 <math>\pm</math> 5.3</b>	40.8 $\pm$ 5.3
0.5	28.2 $\pm$ 6.6	38.8 $\pm$ 7.3	33.1 $\pm$ 5.0	<b>40.9 <math>\pm</math> 5.5</b>
1.0	27.8 $\pm$ 6.9	<b>41.7 <math>\pm</math> 7.4</b>	31.9 $\pm$ 5.0	40.6 $\pm$ 5.3

### A.3.2. ENSEMBLE-BASED RANKING

We ablate how the two key sampling-budget parameters affect ensemble-based pose ranking across the allosteric set, the Boltz-2 wrong-pocket subset of *Runs N’ Poses*, and the full clustered-and-distinct ligand subset. Uncertainty is reported as 95% confidence intervals over 1,000 bootstraps.

Given the total sampling budget  $B = P \times S$ ,  $P$  is the number of candidate binding pockets obtained from the exploration phase, while  $S$  represents the local conformational coverage within each pocket. We ablate both parameters across all test sets in Tables A19–A22. A modest budget of  $B = 50$  samples is typically sufficient for ACER to outperform the baselines. However, on the generalizability benchmark, increasing the budget to  $B = 200$  yields substantially higher success rates in the hardest similarity bin ( $\leq 20\%$ ), where pocket discovery is the key contributor and retaining more candidate pockets in the ranking pool becomes critical (Figure 2).

<sup>2</sup>Boltz-steering includes various constraints: chirality, ring planarity, internal geometry, steric clashes, overlapping chains

Table A19. Ablation on selected  $P$  pockets and  $S$  local ensemble conformers within the pocket on allosteric set ( $n = 20$ ).

Ranking Approach	$P$	$S$	% L-RMSD < 2 Å & PB Valid ↑			% L-RMSD < 3 Å & PB Valid ↑		
			Top-1	Top-5	Top-10	Top-1	Top-5	Top-10
<i>Budget = 50</i>								
Boltz-2	10	5	25.0 ± 9.6	30.0 ± 10.1	30.0 ± 10.1	30.0 ± 10.2	40.0 ± 11.0	40.0 ± 11.0
ACER – baseline	10	5	20.0 ± 9.0	25.0 ± 9.6	35.0 ± 10.6	30.0 ± 10.1	40.0 ± 10.8	45.0 ± 10.9
ACER – L-weight	10	5	<b>25.0 ± 9.3</b>	<b>50.0 ± 11.1</b>	<b>50.0 ± 11.1</b>	30.0 ± 9.9	<b>65.0 ± 10.6</b>	<b>70.0 ± 10.3</b>
ACER – G-weight	10	5	<b>25.0 ± 9.3</b>	<b>50.0 ± 11.1</b>	<b>50.0 ± 11.1</b>	<b>35.0 ± 10.3</b>	<b>65.0 ± 10.6</b>	<b>70.0 ± 10.3</b>
<i>Budget = 100</i>								
Boltz-2	10	10	25.0 ± 9.6	30.0 ± 10.1	30.0 ± 10.1	30.0 ± 10.2	40.0 ± 11.0	40.0 ± 11.0
ACER – L-weight	10	10	<b>30.0 ± 10.1</b>	<b>50.0 ± 11.2</b>	<b>55.0 ± 11.1</b>	<b>35.0 ± 10.6</b>	60.0 ± 11.0	<b>70.0 ± 10.3</b>
ACER – G-weight	10	10	<b>30.0 ± 10.1</b>	<b>50.0 ± 11.2</b>	<b>55.0 ± 11.1</b>	<b>35.0 ± 10.6</b>	<b>65.0 ± 11.0</b>	<b>70.0 ± 10.3</b>
<i>Budget = 200</i>								
Boltz-2	20	10	25.0 ± 9.6	30.0 ± 10.1	30.0 ± 10.1	30.0 ± 10.2	40.0 ± 11.0	40.0 ± 11.0
ACER – L-weight	20	10	25.0 ± 9.6	<b>45.0 ± 11.4</b>	<b>55.0 ± 11.1</b>	30.0 ± 10.2	60.0 ± 11.0	<b>70.0 ± 10.3</b>
ACER – G-weight	20	10	25.0 ± 9.6	<b>45.0 ± 11.4</b>	<b>55.0 ± 11.1</b>	30.0 ± 10.2	<b>65.0 ± 11.0</b>	<b>70.0 ± 10.3</b>

Table A20. Ablation on selected  $P$  pockets and  $S$  local ensemble conformers within the pocket on *Runs N' Poses* wrong-pocket subset of Boltz-2 – Clustered & Distinct Ligands ( $n = 46$ ).

Ranking Approach	$P$	$S$	% L-RMSD < 2 Å & PB Valid ↑		% L-RMSD < 3 Å & PB Valid ↑	
			Top-1	Top-5	Top-1	Top-5
<i>Budget = 50</i>						
Boltz-2	10	5	15.2 ± 5.5	15.2 ± 5.5	19.6 ± 6.0	19.6 ± 6.0
ACER – baseline	10	5	<b>23.9 ± 6.4</b>	23.9 ± 6.4	<b>32.6 ± 7.1</b>	34.8 ± 7.2
ACER – L-weight	10	5	15.2 ± 5.5	26.1 ± 6.7	21.7 ± 6.3	37.0 ± 7.2
ACER – G-weight	10	5	15.2 ± 5.5	<b>28.3 ± 6.8</b>	21.7 ± 6.3	<b>39.1 ± 7.2</b>
<i>Budget = 100</i>						
Boltz-2	10	10	15.2 ± 5.5	15.2 ± 5.5	19.6 ± 6.0	19.6 ± 6.0
ACER – baseline	10	10	<b>23.9 ± 6.4</b>	23.9 ± 6.4	<b>32.6 ± 7.1</b>	34.8 ± 7.2
ACER – L-weight	10	10	19.6 ± 6.0	<b>28.3 ± 7.0</b>	28.3 ± 6.8	<b>39.1 ± 7.4</b>
ACER – G-weight	10	10	19.6 ± 6.0	26.1 ± 6.7	28.3 ± 6.8	37.0 ± 7.4
<i>Budget = 200</i>						
Boltz-2	20	10	15.2 ± 5.5	15.2 ± 5.5	19.6 ± 6.0	19.6 ± 6.0
ACER – baseline	20	10	21.7 ± 6.2	23.9 ± 6.4	28.3 ± 6.7	34.8 ± 7.2
ACER – L-weight	20	10	21.7 ± 6.2	<b>26.1 ± 6.8</b>	<b>32.6 ± 7.1</b>	<b>39.1 ± 7.4</b>
ACER – G-weight	20	10	21.7 ± 6.2	23.9 ± 6.5	<b>32.6 ± 7.1</b>	37.0 ± 7.4

Table A21. Ablation on selected  $P$  pockets and  $S$  local ensemble conformers within the pocket on the *Runs N' Poses* wrong-pocket subset of Boltz-2 – Total ( $n = 88$ ).

Method	$P$	$S$	% L-RMSD < 2 Å & PB Valid ↑		% L-RMSD < 3 Å & PB Valid ↑	
			Top-1	Top-5	Top-1	Top-5
<i>Budget = 50</i>						
Boltz-2	10	5	26.1 ± 4.8	27.3 ± 4.8	28.4 ± 4.9	29.5 ± 5.0
ACER – baseline	10	5	<b>30.7 ± 5.0</b>	31.8 ± 5.1	<b>35.2 ± 5.2</b>	37.5 ± 5.3
ACER – L-weight	10	5	21.6 ± 4.4	29.5 ± 4.9	26.1 ± 4.8	36.4 ± 5.3
ACER – G-weight	10	5	21.6 ± 4.4	<b>31.8 ± 5.1</b>	26.1 ± 4.8	<b>38.6 ± 5.4</b>
<i>Budget = 100</i>						
Boltz-2	10	10	26.1 ± 4.8	27.3 ± 4.8	28.4 ± 4.9	29.5 ± 5.0
ACER – baseline	10	10	<b>30.7 ± 5.0</b>	31.8 ± 5.1	<b>35.2 ± 5.2</b>	37.5 ± 5.3
ACER – L-weight	10	10	25.0 ± 4.7	30.7 ± 5.0	29.5 ± 5.0	<b>37.5 ± 5.3</b>
ACER – G-weight	10	10	25.0 ± 4.7	29.5 ± 4.9	29.5 ± 5.0	36.4 ± 5.3
<i>Budget = 200</i>						
Boltz-2	20	10	26.1 ± 4.8	27.3 ± 4.8	28.4 ± 4.9	29.5 ± 5.0
ACER – baseline	20	10	<b>29.5 ± 5.0</b>	<b>33.0 ± 5.2</b>	<b>33.0 ± 5.2</b>	<b>38.6 ± 5.4</b>
ACER – L-weight	20	10	26.1 ± 4.7	29.5 ± 4.9	31.8 ± 5.1	37.5 ± 5.3
ACER – G-weight	20	10	26.1 ± 4.7	28.4 ± 4.8	31.8 ± 5.1	36.4 ± 5.3

Table A22. Ablation on selected  $P$  pockets and  $S$  local ensemble conformers within the pocket on the full scope of *Runs N' Poses* Clustered and Distinct Ligands ( $n = 217$ ).

Ranking Approach	$P$	$S$	% L-RMSD < 2 Å & PB Valid ↑			% L-RMSD < 3 Å & PB Valid ↑		
			Top-1	Top-5	Top-10	Top-1	Top-5	Top-10
<i>Budget = 50</i>								
Boltz-2	10	5	<b>47.0 ± 3.3</b>	49.8 ± 3.4	51.2 ± 3.3	<b>57.1 ± 3.4</b>	58.5 ± 3.4	61.3 ± 3.3
Protenix-v1	10	5	38.2 ± 3.3	51.2 ± 3.5	53.9 ± 3.4	42.4 ± 3.5	55.3 ± 3.4	57.6 ± 3.4
OpenFold-3	10	5	33.6 ± 3.2	46.1 ± 3.4	49.3 ± 3.4	42.9 ± 3.4	54.8 ± 3.4	56.2 ± 3.4
ACER Boltz-2 L-weight	10	5	40.1 ± 3.3	<b>52.5 ± 3.4</b>	<b>54.8 ± 3.4</b>	50.7 ± 3.4	<b>64.1 ± 3.2</b>	<b>65.9 ± 3.2</b>
<i>Budget = 100</i>								
Boltz-2	10	10	<b>47.0 ± 3.3</b>	49.8 ± 3.4	51.2 ± 3.3	<b>57.1 ± 3.4</b>	58.5 ± 3.4	61.3 ± 3.3
Protenix-v1	10	10	38.2 ± 3.3	51.2 ± 3.5	53.9 ± 3.4	42.4 ± 3.5	55.3 ± 3.4	57.6 ± 3.4
OpenFold-3	10	10	33.6 ± 3.2	46.1 ± 3.4	49.3 ± 3.4	42.9 ± 3.4	54.8 ± 3.4	56.2 ± 3.4
ACER Boltz-2 L-weight	10	10	43.8 ± 3.4	<b>53.5 ± 3.4</b>	<b>56.2 ± 3.4</b>	52.5 ± 3.4	<b>63.1 ± 3.2</b>	<b>65.0 ± 3.2</b>
<i>Budget = 200</i>								
Boltz-2	20	5	<b>47.0 ± 3.3</b>	49.8 ± 3.4	51.2 ± 3.3	<b>57.1 ± 3.4</b>	58.5 ± 3.4	61.3 ± 3.3
Protenix-v1	20	5	38.2 ± 3.3	51.2 ± 3.5	53.9 ± 3.4	42.4 ± 3.5	55.3 ± 3.4	57.6 ± 3.4
OpenFold-3	20	5	33.6 ± 3.2	46.1 ± 3.4	49.3 ± 3.4	42.9 ± 3.4	54.8 ± 3.4	56.2 ± 3.4
ACER Boltz-2 L-weight	20	5	41.0 ± 3.3	<b>52.5 ± 3.4</b>	<b>54.8 ± 3.4</b>	50.7 ± 3.4	<b>63.1 ± 3.2</b>	<b>65.0 ± 3.3</b>

Algorithm 1 details the iterative pocket repulsion procedure. Over  $R$  sampling rounds, the model accumulates a set of occupied pockets  $\mathcal{B}$ , from a prior set of pockets to steer away from (e.g., known prevalent or previously sampled pockets). In each round, a pose is sampled under the repulsion potential  $U_{\text{repulsion}}(\mathcal{B})$ , which penalises all pockets in  $\mathcal{B}$  and steers the model toward unexplored regions. The pocket residues  $\mathcal{R}_r$  for the resulting pose are extracted by selecting all residues with any heavy atom within  $6 \text{ \AA}$  of any ligand atom, and added to  $\mathcal{B}$  to repel subsequent rounds. User-specified parameters include the number of rounds  $R$ , the distance threshold for pocket extraction, and the choice of prior pockets in  $\mathcal{B}_0$ .

---

**Algorithm 1** Iterative pocket repulsion
 

---

**Require:** protein  $\mathcal{P}$ , ligand  $\ell$ , number of rounds  $R$ , prior pocket set  $\mathcal{B}_0$

```

0:  $\mathcal{B} \leftarrow \mathcal{B}_0$  {pockets to steer away from}
0: for  $r = 1, \dots, R$  do
0:   Sample  $\mathbf{x}_0^{(r)} \sim p_\theta(\mathbf{x} \mid \mathcal{P}, \ell; U_{\text{repulsion}}(\mathcal{B}))$ 
0:   Extract pocket residues  $\mathcal{R}_r$  occupied by  $\ell$  in  $\mathbf{x}_0^{(r)}$ 
0:    $\mathcal{B} \leftarrow \mathcal{B} \cup \{\mathcal{R}_r\}$ 
0: end for
0: return  $\{\mathbf{x}_0^{(r)}\}_{r=1}^R = 0$ 

```

---

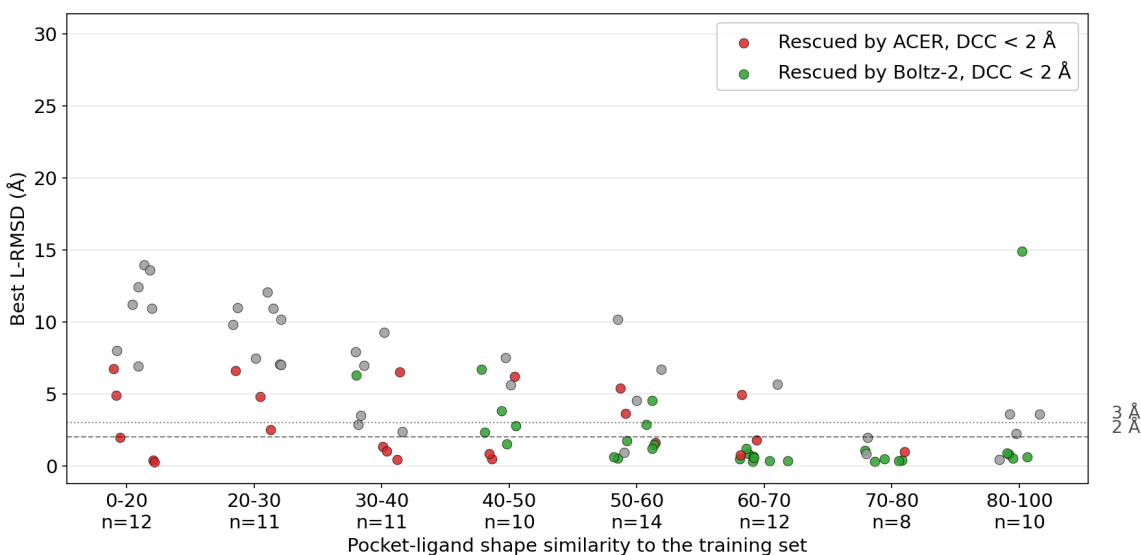


Figure A3. Best L-RMSD per system, stratified by pocket-ligand shape similarity to the training set. Red indicates systems where the Boltz-2 baseline fails ( $\text{DCC} \geq 2 \text{ \AA}$ ) but ACER rescues the correct pocket ( $\text{DCC} < 2 \text{ \AA}$ ); green indicates systems where Boltz-2 alone already succeeds. Gray dots indicate systems where both methods fail. ACER rescue is most prevalent at low similarity (0–30), where Boltz-2 consistently fails to place ligands, and can sample the accurate poses (L-RMSD below  $2 \text{ \AA}$ ) in several cases.

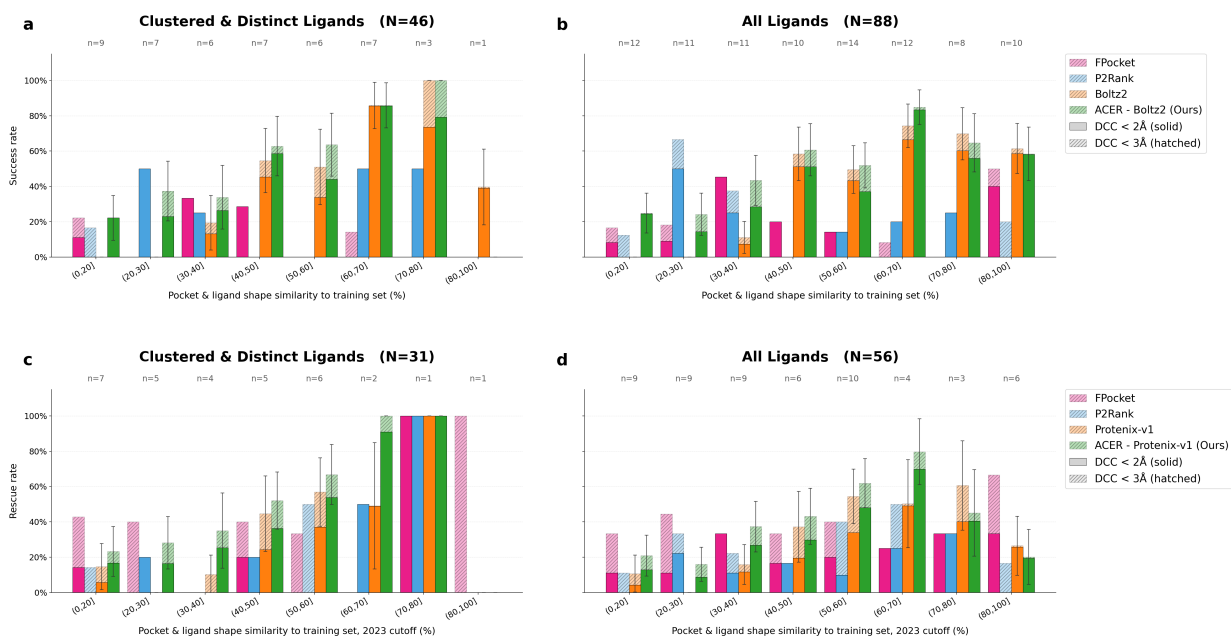


Figure A4. DCC success rate on the wrong-pocket subset of Runs N' Poses (after 2023-06-01), stratified by pocket–ligand shape similarity to the training set. (a, b) Systems where top-ranked Boltz-2 misplaces the ligand (DCC > 2 Å). (c, d) Same, but for Protenix-v1. Filtered ( $n = 46$ ,  $n = 31$ ) and full ( $n = 88$ ,  $n = 55$ ) subsets shown in the left and right columns, respectively.

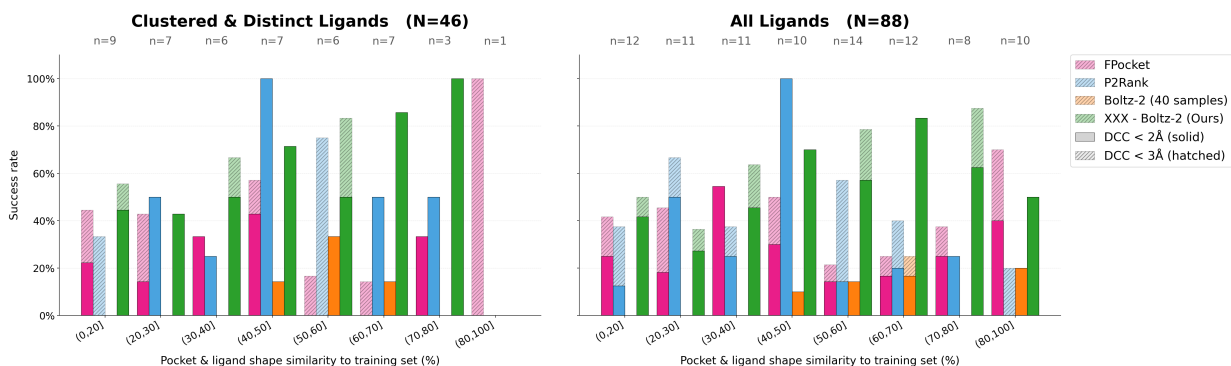


Figure A5. DCC success rate on the wrong-pocket subset of Runs N' Poses when considering oracle across 5 seeds. The system is considered successful if at least one prediction across 5 seeds achieves a DCC below the .

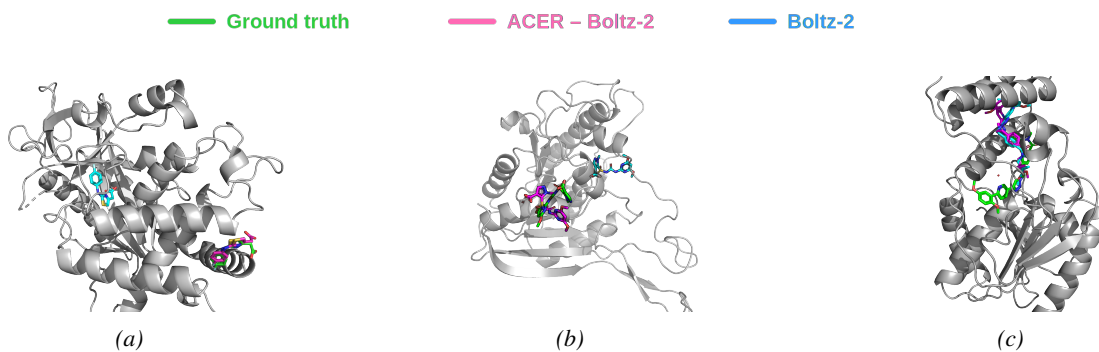
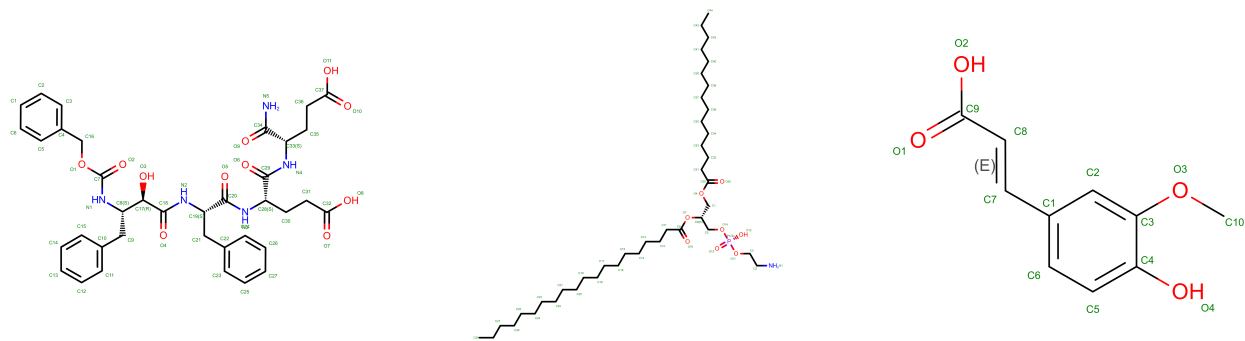


Figure A6. Qualitative examples in the (0, 20] similarity bin of the Runs N' Poses benchmark. (a) 7FTM – where ACER correctly places ligands into the right pocket (DCC < 2 Å) and recovers an accurate ligand pose (L-RMSD < 2 Å). (b) 8GOY – a partial success where ligand placement succeeds (DCC < 2 Å), yet the predicted ligand pose remains inaccurate (L-RMSD > 2 Å). (c) 8TQV – a failure case in which both ACER – Boltz-2 and default Boltz-2 can never rescue the correct pocket.

1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295  
1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319



(a) KI2 from PDB: 1NH0

(b) PTY from PDB: 3AR4

(c) FER from PDB: 3CBG

Figure A7. The structure of chemical probes used in the ablation study of decoy ligand type.