

# STRUCTURAL PLAUSIBILITY WITHOUT BINDING SPECIFICITY: LIMITS OF AI-BASED ANTIBODY-ANTIGEN STRUCTURE PREDICTION SCORES

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Antibodies are central to modern immunotherapy, yet accurately predicting antibody-antigen binding interfaces remains a major challenge for computational modeling. While recent AI-based structure prediction methods can generate plausible antibody-antigen complexes, it remains unclear whether they can reliably discriminate cognate binding partners and identify correct paratope-epitope interfaces in realistic discovery settings. Here, we introduce a controlled benchmarking framework to systematically evaluate publicly available state-of-the-art structure prediction models (AlphaFold3, Boltz-2, and Chai-1) in their ability to distinguish real antibody-antigen complexes-cognate pairs, extracted from experimentally solved structures ( $n=106$ ), from shuffled complexes, which serve as artificial non-cognate negative controls ( $n=11,342$ ). We assessed structural accuracy, interface correctness, and the discriminative power of commonly used confidence metrics to distinguish cognate from non-cognate complexes (including ipTM and DockQ) when applied without access to structural ground truth, reflecting real-world deployment scenarios. We further analyzed sequence-level, structural, and sequence-structure features associated with high or low prediction confidence, independent of pairing correctness, and evaluated the trade-offs between computational cost and performance gains from increased sampling. To support community-driven mining, benchmarking, and method development, we release a large-scale dataset of 1.8 million (561,800 VHH-antigen pairings predicted by 3 different tools) complex structures. Our results show that current confidence scores (ipTM) often fail to discriminate cognate from non-cognate interactions (high false positive rate), even with extensive sampling, highlighting key limitations in current antibody-antigen modeling pipelines. This work provides a biologically grounded benchmark for antibody-antigen interface prediction and outlines critical directions for improving computational screening strategies in antibody discovery.

## 1 INTRODUCTION

Antibodies are key immunotherapeutic biomolecules characterized by antigen-specific binding. This specificity enables antibodies to identify and bind molecular targets such as tumor- or pathogen-associated antigens (Meng et al., 2024), underpinning a wide range of therapeutic applications and positioning antibodies as the largest class of biotherapeutics (Bielska et al., 2025). With the continued growth of monoclonal antibody (mAb) therapeutics, there is increasing interest in developing *in silico* antibody discovery and design tools (Akbar et al., 2022; Bashour et al., 2024; Greiff et al., 2020; Levine et al., 2026; Overath et al., 2026). While experimental discovery pipelines still rely heavily on antibody libraries and screening (Laustsen et al., 2021), improved prediction of paratope-epitope interactions may enable *in silico* discovery strategies (Akbar et al., 2022).

Antigen recognition is primarily mediated by the antibody’s six hypervariable complementarity-determining region (CDR) loops, which together form the paratope (Ruffolo et al., 2023). These loops engage antigen surface regions (the epitope) through a combination of shape complementarity, physicochemical interactions, and conformational adaptability (Smorodina et al., 2025a; Kuroda & Gray, 2016). Although recent advances in deep learning have revolutionized protein structure pre-

054 diction for many molecular systems (Abramson et al., 2024; Passaro et al., 2025; Discovery et al.,  
055 2024), accurately predicting antibody-antigen complexes in general, and identifying the correct  
056 paratope-epitope interface in particular, remains challenging (Smorodina et al., 2022; Yin & Pierce,  
057 2024). Previous benchmarking studies report success rates of approximately 20% for antibody-  
058 antigen docking using AlphaFold-Multimer (v2.3.0) and Rosetta-based protocols (Hitawala & Gray,  
059 2024; Harmalkar et al., 2023; Weitzner et al., 2017; Ambrosetti et al., 2020). More recent evalu-  
060 ations indicate improved performance, with approximately 35% success using a single stochastic  
061 seed and up to around 60% success when extensive sampling (up to 1000 seeds) is combined with  
062 confidence-based ranking (Abramson et al., 2024; Hitawala & Gray, 2025), albeit at substantial  
063 computational cost.

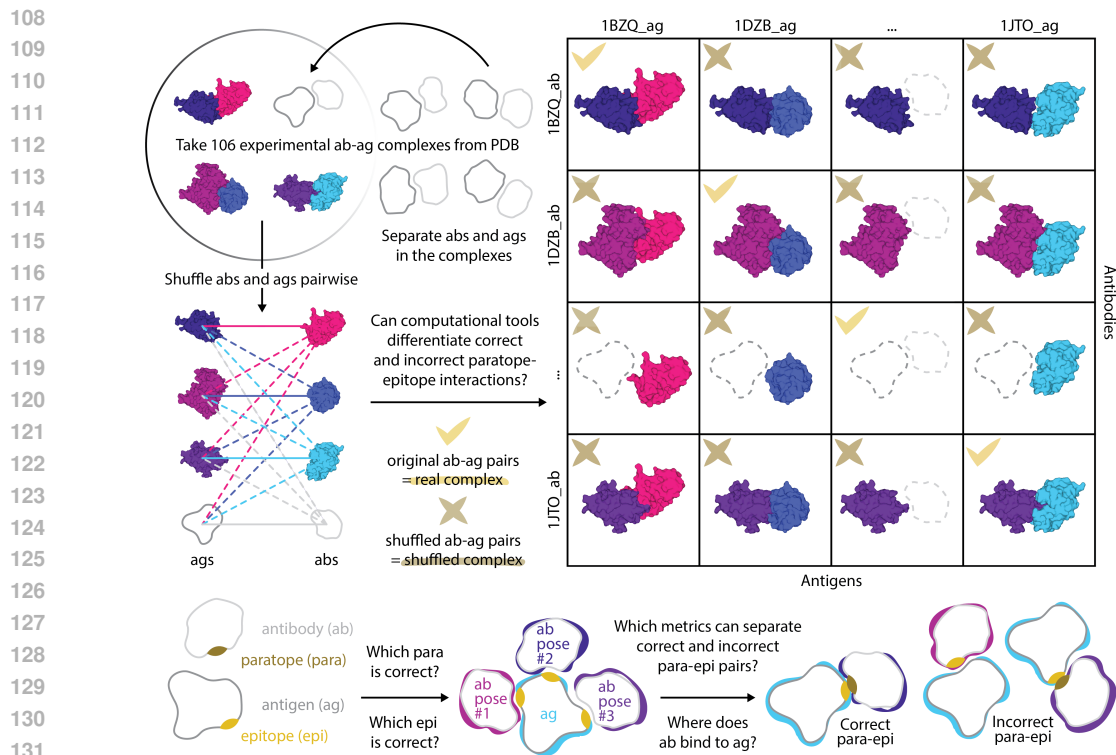
064 Recently, several state-of-the-art molecular structure prediction models have emerged, including Al-  
065 phaFold2/3 (Abramson et al., 2024; Jumper et al., 2021), Boltz-1/2 (Passaro et al., 2025; Wohlwend  
066 et al., 2024), and Chai-1/2 (Boitreau et al., 2024; Team et al., 2025). Performance evaluations  
067 on diverse benchmarks spanning protein monomers, multimers, and small-molecule interactions  
068 have suggested that some of these models outperform earlier generations such as AlphaFold-  
069 Multimer (Fromm et al., 2025; Unsal et al., 2025; Grieswelle et al., 2025; Bradley, 2023; Evans  
070 et al., 2021). However, systematic evaluation of their ability to identify correct antibody-antigen  
071 interfaces-and, critically, their inability to reject incorrect interactomes-remains absent. Evaluat-  
072 ing computational complex-prediction models can highlight key metrics that discriminate cognate  
073 binding partners (“real” complexes) from non-cognate ones (“shuffled” complexes used as negative  
074 controls) (Overath et al., 2025). This distinction is particularly important in discovery and screening  
075 contexts, where large numbers of candidate antibody-antigen pairings must be evaluated to separate  
076 binders from non-binders and incorrect-epitope binders, with false positives remaining a persistent  
077 challenge (Hitawala & Gray, 2025; Holt et al., 2025; Dang et al., 2023).

078 To explicitly address the problem of partner discrimination, we introduce a benchmarking frame-  
079 work that enables ground-truth distinction between real and shuffled antibody (here, VHH) and  
080 antigen complexes (Figure 1). Real complexes are defined as cognate VHH-antigen pairs extracted  
081 from experimentally solved structures, where the VHH and antigen are known biological binding  
082 partners within the same Protein Data Bank (PDB) entry. In contrast, shuffled non-cognate com-  
083 plexes are artificially generated by pairing a VHH sequence from one experimental complex with an  
084 antigen sequence from a different experimental complex. These non-cognate shuffled pairs do not  
085 correspond to known biological interactions and serve as ground-truth decoys. This benchmark de-  
086 sign allows direct comparison of computational predictions across cognate and non-cognate pairings  
087 under identical modeling conditions.

088 Using this framework, we benchmark AlphaFold3, Boltz-2, and Chai-1 for predicting VHH-antigen  
089 paratope-epitope interactions, quantifying their ability to distinguish real cognate from shuffled non-  
090 cognate pairs. This enables assessment of structural and interface accuracy, evaluation of whether  
091 confidence metrics (e.g., ipTM (Evans et al., 2021; Genz et al., 2025), DockQ (Basu & Wallner,  
092 2016; Mirabello & Wallner, 2024)) can discriminate real from shuffled complexes, and identification  
093 of sequence or structural features associated with high or low confidence independently of pairing  
094 correctness. Overall, we present a ground-truth-based benchmarking pipeline that exposes current  
095 tool limitations and highlights opportunities to improve antibody-antigen complex prediction for  
096 practical screening workflows.

## 097 2 COMPUTATIONAL PREDICTORS LARGELY FAIL TO DISCRIMINATE REAL 098 FROM “SHUFFLED” ANTIBODY-ANTIGEN COMPLEXES 099

100  
101 To assess whether state-of-the-art structure prediction tools can distinguish biologically observed  
102 antibody-antigen interactions from shuffled complexes, we evaluated three independent AI-based  
103 complex structure prediction tools (AF3, Boltz-2, and Chai-1) across a panel of 561,800 ( $106^2 \times 50$ )  
104 VHH-antigen pairings (Figure 2): 106 real systems (complexes) from 91 unique PDB IDs and  
105 11,342 shuffled complexes with 50 replicates each, where a replicate is defined as a prediction  
106 sample generated by a model. The selection of 106 real VHH-antigen complexes was based on se-  
107 quence and structural diversity, size, and uniqueness (see Methods). For all complexes, information  
on whether they were present in the training set of the respective AI tool was also available.



Next, we asked to what extent scores differ between pairs that belong to the train or test sets within each dataset split (Figure 2C). Violin plots summarizing ipTM score distributions confirmed substantial similarity of interface scores between real and shuffled complexes for all tools and across training, test, and mixed (*both*) splits. AF3 showed bimodal distributions for real complexes, with a subset of high-scoring interactions, but shuffled complexes still occupied overlapping score ranges. Boltz-2 assigned consistently high ipTM values to both real and shuffled interactions (mean  $\sim 0.85-0.91$ ), with minimal separation. Chai-1 yielded uniformly low scores for both classes, again with nearly identical distributions. Across all tools, mean scores and variances were comparable between real and shuffled complexes, indicating that score magnitude alone does not encode interaction authenticity.

We used TopModel’s clash score (Fernández-Quintero et al., 2023) to evaluate structural quality across real and shuffled antibody-antigen complexes (Figure 2D). The clash score accounts for both the number of steric clashes between all atoms and protein length, with lower values indicating better geometric quality. Across all tools and dataset splits, clash score distributions were broad and overlapped between real and shuffled complexes. For AF3, top-ranked (highest ipTM) models showed mean clash scores of  $25.7 \pm 9.1$  (train),  $26.7 \pm 9.8$  (test), and  $26.8 \pm 9.7$  (*both*) for shuffled complexes, compared with  $22.7 \pm 9.4$ ,  $24.2 \pm 9.9$ , and  $26.8 \pm 9.7$  for real complexes across the same splits. Boltz-2 produced systematically higher clash scores than AF3, with shuffled complexes averaging  $30.0 \pm 9.5$  (train),  $36.1 \pm 12.4$  (test), and  $33.0 \pm 11.3$  (*both*), while real complexes averaged  $27.9 \pm 10.7$ ,  $34.8 \pm 12.2$ , and  $32.9 \pm 11.3$ , respectively. Chai-1 exhibited markedly poorer structural quality by this metric: mean clash scores were 4-5 $\times$  higher and more variable, reaching  $\sim 80-90$  on average, with extreme dispersion (e.g.,  $62.0 \pm 65.0$  and  $90.7 \pm 125.2$  in the test split). These elevated and highly variable scores were observed for both real and shuffled complexes, indicating frequent steric clashes in top-ranked predictions and a lack of effective internal filtering for geometric plausibility. Importantly, real complexes did not systematically achieve lower clash scores than shuffled complexes in any tool or dataset split. In several cases, shuffled complexes displayed comparable or even slightly lower average clash scores than real complexes, despite being biologically incorrect.

To assess whether different prediction tools agree on which antibody-antigen pairs are confident or uncertain, we directly compared ipTM score matrices generated by AF3, Boltz-2, and Chai-1 (Figure 2A). Global correlations between flattened score matrices were uniformly low (Pearson  $r = 0.14$  for Chai-1 vs. Boltz-2,  $r = 0.18$  for Chai-1 vs. AF3, and  $r = 0.13$  for Boltz-2 vs. AF3), indicating low agreement across tools in their assessment of interaction confidence.

Per-system correlation analysis (testing whether tools rank the same systems as similarly confident or unconfident) further revealed substantial heterogeneity across both VHHs and antigens (Figure 2E). For each pair of tools, many systems exhibited weak or even negative correlations, highlighting disagreement in how individual VHHs or antigens are scored. These discrepancies were not confined to a small subset of problematic systems but were broadly distributed across the benchmark, suggesting that tool-specific inductive biases strongly shape confidence assignment. Outlier analysis reinforces this conclusion: high-confidence shuffled outliers identified by one tool rarely overlapped with those from another. Specifically, only four shuffled outliers overlapped between Chai-1 and Boltz-2, fourteen between Chai-1 and AF3, and just one between Boltz-2 and AF3, despite all tools being evaluated on the same set of VHH-antigen pairings. Representative examples illustrate the magnitude of these discrepancies. For instance, the system (4NC2, 7R24) was assigned high confidence by Boltz-2 (ipTM  $\approx 0.81$ ) but substantially lower confidence by Chai-1 (ipTM  $\approx 0.60$ ), while (9EMY, 8UKV) achieved high confidence in Chai-1 (ipTM  $\approx 0.73$ ) but was scored markedly lower by Boltz-2 (ipTM  $\approx 0.44$ ). Similarly, several systems such as (7TGF, 6OBO) and (8K4Q, 8EW6) were consistently high-confidence outliers in AF3 (ipTM  $\geq 0.8$ ) yet did not emerge as outliers in the other tools (Supplementary Table 2).

Together, these results show that both ipTM confidence and structural accuracy remain limited across tools, with AF3 producing the lowest average clash scores, Boltz-2 intermediate values, and Chai-1 the poorest overall structural quality. However, even the best-performing method by this metric fails to reliably discriminate real antibody-antigen complexes from shuffled mismatches. This inability to distinguish real from shuffled complexes persists across training, test, and combined datasets, indicating that the observed overlap is unlikely to result from overfitting or data leakage. Instead, it reflects a limitation of current structure-based confidence scores and simple structural quality measures, which capture generic interface properties but fail to reflect biological correctness in antibody-antigen complex prediction.

216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269

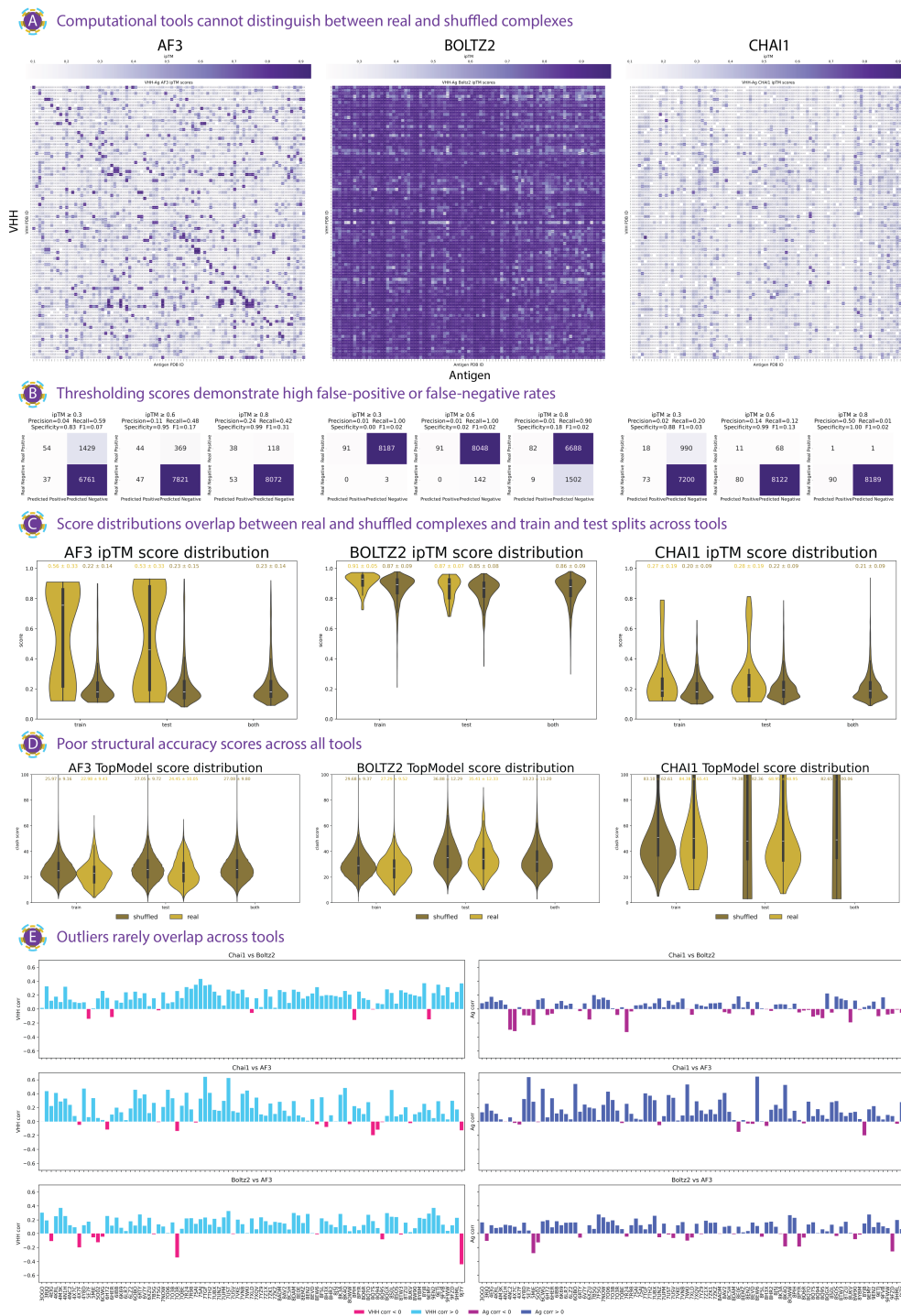


Figure 2: Current structure prediction tools and confidence scores fail to capture antibody-antigen specificity. (A) All-by-all interaction score landscapes. Heatmaps show the best interface confidence scores (ipTM) across 50 stochastic predictions for all VHH-antigen pairings generated by AF3, Boltz-2, and Chai-1. Rows denote VHHs and columns antigens; color scales indicate ipTM values. (B) Overlapping confidence distributions for real and shuffled complexes. Violin plots summarize ipTM scores across training, test, and mixed (*both*) splits. Real and shuffled distributions largely overlap for all models. (C) Threshold-based classification is unreliable. Confusion matrices report precision, specificity, and F1 scores at ipTM thresholds of 0.3, 0.6, and 0.8 for each tool. (D) Structural quality metrics are non-discriminative. Violin plots show clash score distributions for top-ranked predictions across dataset splits; real and shuffled complexes strongly overlap for all tools. (E) Tool-specific confidence outliers. Per-system correlations of ipTM scores between tool pairs reveal weak agreement (Pearson correlation) and inconsistent high-confidence outliers across VHHs and antigens.

## A APPENDIX

### A.1 CROSS-TOOL SCORE AGREEMENT IS WEAK AND SAMPLING IMPROVES STRUCTURE BUT NOT CONFIDENCE ALIGNMENT

We next asked whether models "know" when they fail—that is, whether confidence scores (ipTM) reliably indicate structural prediction quality (DockQ) under different sampling regimes (Figure 3). Confidence calibration—the correspondence between predicted confidence and actual quality (Guo et al., 2017; Kull et al., 2019)—is critical for practical screening workflows where ground-truth structural information is unavailable. Accordingly, confidence calibration analysis was performed only on real complexes, as shuffled pairings lack experimental reference structures against which DockQ can be computed (Mirabello & Wallner, 2024). We found that confidence calibration differs systematically across tools when examining the correlation between DockQ and ipTM for the best, initial (sample<sub>0</sub>), and worst of 50 predictions (Figure 3A).

To interpret calibration patterns, we defined four diagnostic quadrants based on thresholds of DockQ = 0.23 (acceptable quality) and ipTM = 0.5 (confident prediction): Q1 (confident and correct), Q2 (overconfident failures), Q3 (uncertain and incorrect), and Q4 (underconfident success) (Supplementary Table 3). Among the three methods, AF3 exhibits the strongest calibration, with a Pearson correlation of  $r = 0.736$  for best-DockQ structures and relatively few overconfident failures (Q2: 2%; Q4: 28%). AF3 maintained high correlation even at the single-sample level ( $r = 0.888$  for sample<sub>0</sub>), indicating comparatively robust confidence assignment without extensive sampling. In contrast, Boltz-2 was systematically overconfident. Although its best-DockQ calibration remained moderate (Pearson  $r = 0.665$ ), a substantial fraction of predictions fell into the overconfident failure regime (Q2: 18%), where ipTM is high despite poor structural quality. This effect was exacerbated for single-sample predictions, where correlation dropped sharply (Pearson  $r = 0.493$ ), highlighting the unreliability of Boltz-2 confidence scores in low-sampling screening workflows. Conversely, Chai-1 exhibited systematic underconfidence. Although a subset of predictions achieved high structural quality (DockQ  $\approx 0.4$ - $0.9$ ), these models were frequently assigned low ipTM scores, resulting in a substantial fraction of underconfident successes (Q4: 24%). Consequently, the relationship between confidence and best-case structural quality was weaker for Chai-1 (best-DockQ vs. ipTM,  $r = 0.612$ ), implying that confidence-based filtering would exclude many structurally accurate predictions. Extreme quadrant cases further illustrate observed failure modes. Overconfident failures (Q2) correspond to geometrically poor structures assigned high confidence (e.g., AF3 system\_24\_24\_6V7Y), whereas underconfident successes (Q4) correspond to high-quality complexes that would be missed by thresholding (e.g., AF3 system\_27\_27\_7B5G). These cases provide concrete examples of how confidence scores can diverge from structural reality in both directions (see Supplementary Table 3 for the classification of all systems across the four quadrants).

Next, we examined whether sampling improves structure but not confidence alignment. Across all three tools, saturation sampling improved best-case DockQ (Supplementary Figure 1A), confirming that additional sampling can refine geometry. DockQ improvement distributions were right-skewed for all models, with AF3 and Boltz-2 exhibiting the largest absolute gains, while Chai-1 showed smaller but consistent improvements. The proportion of systems achieving acceptable or higher quality (DockQ > 0.23) increased substantially after saturation sampling (Figure 3B; Supplementary Figure 1B). For AF3, incorrect predictions decreased by nearly 20%, with Boltz-2 showing a comparable reduction and a marked increase in medium- and high-quality outcomes; in contrast, Chai-1 remained dominated by incorrect predictions despite sampling. These changes reflect systems "rescued" above the quality threshold through sampling alone (the "improvers," Figure 3B).

However, the relationship between structural improvement and confidence remained weak. Although all three models showed significant improvements in DockQ scores (Wilcoxon  $p < 0.001$ ; Supplementary Figure 1C), paired analyses showed no significant improvement in ipTM for AF3 or Chai-1 ( $p = 0.179$  and  $p = 0.222$ , respectively). While Boltz-2 exhibited a statistically significant change ( $p < 0.001$ ), the ipTMs actually slightly decreased (Supplementary Figure 1D). In particular, correlations between changes in DockQ and changes in ipTM were near zero across all models ( $r = -0.027$ ,  $-0.040$ , and  $-0.019$  for AF3, Boltz-2, and Chai-1, respectively; Figure 3C), indicating that confidence scores were largely insensitive to structural improvements achieved through sampling and remained "locked in" to the initial structural hypothesis.

324 To assess whether confidence scores provide consistent rankings across models, we compared ipTM  
325 values for each system between model pairs (Supplementary Figure 1E). Cross-model correlations  
326 were modest (Pearson  $r = 0.29-0.39$ ), and importantly, higher confidence in one model did not re-  
327 liably correspond to higher structural quality as measured by DockQ. From a screening perspective,  
328 this behavior has important implications. Single-sample predictions from Boltz-2 are frequently  
329 overconfident relative to achieved DockQ, whereas AF3 displays comparatively better calibration  
330 even prior to saturation sampling. Nonetheless, for all tools, confidence does not reliably track  
331 structural refinement, limiting its utility as a post-sampling ranking criterion.

332 Finally, we examined prediction consistency using epitope-paratope variation across 15 antigens,  
333 each crystallized with two distinct VHHs (Supplementary Figure 2A; Methods). This provides a  
334 natural test of whether models that succeed with one VHH can generalize to alternative binders of  
335 the same antigen. AF3 produced acceptable predictions (DockQ > 0.23 for both pairs) in 10 out of  
336 15 cases (67%), compared to 6/15 (40%) for Boltz-2 and 2/15 (13%) for Chai-1 (Supplementary  
337 Figure 2B). Boltz-2 and especially Chai-1 more frequently succeeded with only one VHH or failed  
338 on both partners (Supplementary Figure 2A), suggesting that AF3 better captures the underlying  
339 antigen structure independently of the specific paratope geometry.

340 Together, these results demonstrate that confident failures are largely tool-specific rather than reflect-  
341 ing a shared set of universally difficult or ambiguous systems. The low overlap in outlier systems  
342 and weak agreement in confidence-based rankings across predictors highlight the absence of a com-  
343 mon confidence landscape, underscoring the risk of relying on any single model's confidence scores  
344 for candidate selection-particularly after saturation sampling, where structural quality improves but  
345 confidence remains misaligned.

346  
347  
348  
349  
350  
351  
352  
353  
354  
355  
356  
357  
358  
359  
360  
361  
362  
363  
364  
365  
366  
367  
368  
369  
370  
371  
372  
373  
374  
375  
376  
377

378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431

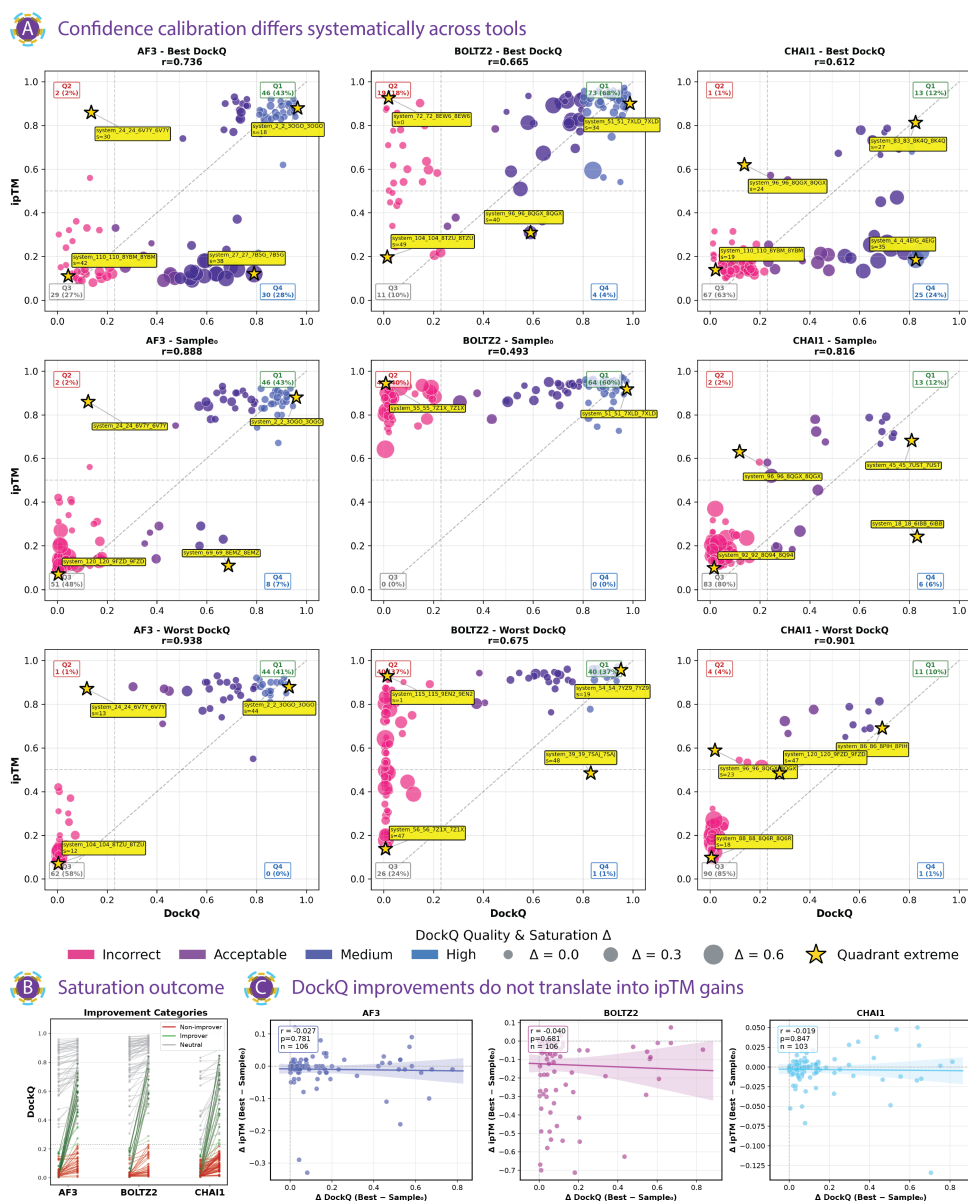


Figure 3: Cross-tool disagreement and confidence calibration of antibody-antigen docking predictions. **(A)** Confidence calibration differs systematically across ML structure prediction tools. Scatter plots show DockQ versus ipTM for AF3, Boltz-2, and Chai-1 under three sampling regimes: best DockQ, initial sample (sample<sub>0</sub>), and worst DockQ. Points are colored by DockQ quality category (incorrect, acceptable, medium, high) and sized by saturation range ( $\Delta$ DockQ from sample<sub>0</sub> to best). Dashed lines indicate DockQ= 0.23 and ipTM= 0.5. Pearson correlation coefficients are shown for each panel. **(B)** Saturation outcome categories reveal when sampling rescues docking quality. Line plots show per-system trajectories from sample<sub>0</sub> to the best sampled prediction, grouped into non-improvers, improvers, and neutral systems. **(C)** DockQ gains from sampling do not translate into confidence-score gains. Scatter plots show the relationship between  $\Delta$ DockQ and  $\Delta$ ipTM for each model, with linear regression fits and 95% confidence intervals. Near-zero correlations indicate that confidence scores fail to track structural improvements achieved through sampling.

## A.2 STRUCTURAL ACCURACY SCORES ARE UNIFORMLY POOR WHILE INTERFACE CONTACTS DOES NOT REFLECT SPECIFICITY

We next evaluated whether predicted antibody-antigen complexes recovered the experimentally defined antigen epitope and whether this information can discriminate real from shuffled pairings (Fig-

ure 4). Epitope recall was quantified as the fraction of experimental epitope residues recovered by model predictions, where a residue was considered recovered if it contacted the VHH in at least 50% of stochastic replicate predictions ( $n = 50$ ). Across all three models, epitope recall values were generally low and sparsely distributed, with substantial overlap between real and shuffled complexes (Figure 4A). Although real complexes lie along the diagonal of the epitope-recall matrices, off-diagonal (shuffled) complexes frequently exhibited comparable levels of epitope recovery. This indicates that models often identify plausible antigen contact regions even in non-cognate pairings, limiting the specificity of epitope-based discrimination. Consistent with this observation, direct comparison of epitope-recall distributions shows that shuffled complexes frequently overlap or exceed the recall observed for real complexes across all models and data splits (Figure 4B). Focusing on cognate pairings and whether the respective models recovered at least one residue from the true epitope (epitope recall  $> 0$ ), AF3 recovered the most true epitopes in the test set of complexes ( $n = 29$ ), followed by Chai-1 ( $n = 26$ ) and Boltz-2 ( $n = 17$ ).

To assess the impact of data leakage, we examined epitope-recall distributions for real complexes split into train-leakage and true test sets (Figure 4B). Across models, train-leakage complexes showed slightly higher median epitope recall than true test complexes, consistent with partial memorization or bias toward known interfaces. However, the overall distributions remained broad, and substantial overlap persisted between train-leakage and true test predictions. Notably, shuffled complexes spanned similar recall ranges, indicating that leakage alone does not explain the limited discriminatory power of epitope recall. Thus, high epitope recall does not reliably indicate cognate binding even in the absence of explicit train-test overlap.

We next tested whether increasing stochastic sampling improves epitope discrimination by assessing epitope-enrichment significance as a function of replicate count (Figure 4C). Enrichment significance was computed using a one-sided binomial test for each VHH-antigen pair. Let  $N$  denote the number of stochastic structure predictions and  $K$  the number of predictions containing at least one epitope contact. Under a null model in which contacts occur uniformly across the antigen surface with probability  $p_{\text{epitope}}$  (estimated from the fraction of solvent-accessible antigen residues belonging to the epitope), the probability of observing  $K$  or more epitope hits is given by the binomial survival function  $P(X \geq K | X \sim \text{Binomial}(N, p_{\text{epitope}}))$ . Resulting  $p$ -values were corrected across all VHH-antigen pairs using the Benjamini-Hochberg procedure, and pairs with  $\text{FDR} < 0.05$  were considered significantly enriched. To assess the effect of ensemble size, we recomputed binomial  $p$ -values for hypothetical replicate counts  $m$  by rescaling the observed hit fraction ( $K/N$ ) and repeating the same multiple-testing correction. For all three models, the fraction of VHH-antigen pairs declared significantly enriched increased rapidly with the number of stochastic predictions and saturated by approximately 10-15 replicates. Importantly, real and shuffled complexes exhibited highly similar saturation behavior, regardless of whether the shuffled pairs involved train leakage of the VHH, the antigen, both, or neither. This indicates that increasing ensemble size primarily reinforces existing contact preferences rather than sharpening discrimination between cognate and non-cognate interactions. AF3 and Boltz-2 retained a modest residual separation between real and shuffled complexes across replicate counts, with real pairs showing slightly higher significant-enrichment rates than any shuffled cohort. In contrast, Chai-1 showed near-complete convergence: shuffled complexes rapidly achieved significance frequencies comparable to, or exceeding, those of real complexes, consistent with a strong epitope-seeking bias that limits the interpretability of enrichment-based metrics for this model. Notably, the leakage-stratified shuffled cohorts largely overlapped within each model, further indicating that these effects are not primarily driven by explicit train-test leakage but instead reflect intrinsic model priors amplified by ensembling.

Finally, we examined the relationship between structural consistency across predictions and epitope recovery within each model (Figure 4D). Structural consistency was quantified as the mean pairwise RMSD across  $n = 50$  stochastic replicates for each complex. Across all models and data splits, epitope recall was negatively correlated with replicate RMSD, indicating that complexes predicted with consistent orientations tend to recover a larger fraction of epitope residues. However, this relationship held for both train-leakage and true test complexes and did not differentiate real from shuffled pairings. Pearson correlation coefficients ranged from  $r = -0.22$  to  $-0.30$  for AF3,  $r = -0.14$  to  $-0.27$  for Boltz-2, and  $r = -0.39$  to  $-0.59$  for Chai-1, demonstrating that while structural consistency improves epitope recovery, it does not guarantee biological correctness. Notably, RMSD distributions of the 50 replicates showed slightly higher values for shuffled complexes compared to

486 real complexes (Figure 4E), indicating a potential-but limited-signal for distinguishing cognate from  
487 non-cognate interfaces.

488 Together, these analyses show that epitope recovery and enrichment metrics are dominated by  
489 generic contact biases and structural consistency rather than interaction specificity. Increasing  
490 stochastic sampling improves apparent epitope enrichment but does so similarly for real and shuffled  
491 complexes, limiting its utility as a discriminative signal. These findings further reinforce the conclu-  
492 sion that current structure prediction tools preferentially identify plausible antigen contact regions  
493 but struggle to distinguish true cognate paratope-epitope interactions from non-cognate alternatives.  
494

### 495 A.3 COMPUTATIONAL COST AND SAMPLING EFFICIENCY DIFFER SUBSTANTIALLY ACROSS 496 STRUCTURE PREDICTION TOOLS 497

498 Using ML structure prediction tools for high-throughput antibody screening requires understanding  
499 not only prediction accuracy but also computational cost. With GPU resources unevenly distributed  
500 across institutions (ref, 2025) and growing concern over the energy footprint of large-scale com-  
501 putation (Morand et al., 2024; Lannelongue & Inouye, 2023), identifying cost-effective sampling  
502 strategies has practical importance. Saturation sampling and seed selection can both improve pre-  
503 diction quality (Abramson et al., 2024; Hitawala & Gray, 2025; Fromm et al., 2025; Johansson-Åkhe  
504 & Wallner, 2022), yet their relative contributions and associated costs remain uncharacterized. We  
505 therefore evaluated each system ( $n = 106$ , 91 unique PDBs; no deposition-year bias observed; Sup-  
506 plementary Figure 4A) across five random independent seeds with increasing numbers of diffusion  
507 samples ( $N = 1, 10, 25, 50$ , and 100), monitoring GPU energy consumption throughout (Figure 5).  
508 For this analysis, we additionally included Boltz-1, the predecessor to Boltz-2, to check consistency  
509 of saturation effects across model versions. Because each saturation level uses an independent seed,  
510 this design enables analysis of combined seed and saturation effects; extracting only the first sample  
511 from each seed isolates seed-dependent variation.

512 To characterize how computational cost scales with sampling depth, we computed the median en-  
513 ergy consumption across all systems at each saturation level ( $N = 1, 10, 25, 50$ , and 100). Energy  
514 consumption varied markedly between tools, revealing distinct cost profiles. At the highest satura-  
515 tion level ( $N = 100$  samples), median energy usage per system ranged from 23.0 Wh for AF3 to  
516 82.9 Wh for Chai-1, corresponding to a 3.6-fold difference (Figure 5A; Supplementary Figure 3A-  
517 C). Boltz-1 and Boltz-2 exhibited intermediate energy-usage profiles, with substantially lower cost  
518 than Chai-1 at high saturation but higher baseline costs than AF3. When aggregated across all 106  
519 real systems and saturation levels, total energy consumption ranged from 13.6 kWh for AF3 to  
520 22.7 kWh for Chai-1.

521 These differences reflect distinct architectural choices. AF3 incurred a high baseline energy cost for  
522 MSA computation (82.4 Wh per system, computed once regardless of sampling depth) but showed  
523 relatively low marginal cost for additional diffusion samples. In contrast, Chai-1 exhibited a low  
524 baseline cost (6.6 Wh) but a steep increase in energy consumption with increasing saturation, con-  
525 sistent with its embedding-based workflow that omits MSA computation (Methods). Boltz-based  
526 models showed intermediate behavior. Across all tools, both energy consumption and runtime  
527 scaled approximately linearly with system length, with Pearson correlation coefficients ranging from  
528  $r = 0.60$  for AF3 to  $r = 0.98$  for Boltz-1 (Supplementary Figure 3D,E).

529 Increasing saturation sampling improved prediction quality for all models with diminishing returns,  
530 as measured by DockQ (Figure 5B). We quantified this by identifying, for each system, the maxi-  
531 mum DockQ achieved at each saturation level and computing the median across systems. Median  
532 maximum DockQ increased from baseline ( $N = 1$ ) to  $N = 100$  by  $\Delta = +0.43$  for AF3 (0.24 to  
533 0.68),  $\Delta = +0.23$  for Boltz-2 (0.57 to 0.80),  $\Delta = +0.20$  for Boltz-1 (0.06 to 0.26), and  $\Delta = +0.15$   
534 for Chai-1 (0.04 to 0.19). While all models benefited from increased sampling, the magnitude of  
535 improvement differed substantially, with AF3 showing the largest gains.

536 Marginal quality gain was computed as the per-system difference in maximum DockQ between con-  
537 secutive saturation levels and summarized as the median across systems. Across models, the largest  
538 marginal improvements occurred at low saturation levels. The transition from baseline to  $N = 10$   
539 samples captured the majority of achievable quality improvement, while additional sampling beyond  
 $N = 25$  yielded progressively smaller gains (Figure 5C). Consistent with this pattern, the propor-

540 tion of incorrect predictions ( $\text{DockQ} < 0.23$ ) decreased sharply from baseline to  $N = 10$  but showed  
541 more modest ( $\sim 5\text{-}15\%$ ) reductions at higher saturation levels ( $N = 25\text{-}100$ ; Figure 5D).

542 Analyzing cumulative DockQ improvement as a function of cumulative energy expenditure revealed  
543 a consistent "efficiency frontier" across tools (Figure 5C). All models exhibited steep initial slopes,  
544 indicating high quality gains per unit energy at low sampling depths, followed by pronounced flat-  
545 tening as energy expenditure increased. Despite differing absolute costs, the efficiency frontiers  
546 converge across tools, indicating that early sampling dominates quality gains regardless of architec-  
547 ture. These trends indicate that  $N \approx 10\text{-}25$  samples capture most attainable quality improvement at  
548 a fraction of the computational cost required for  $N = 100$ .

549 Beyond sampling depth, seed choice had a substantial independent effect on prediction outcomes  
550 (Figure 5E). To isolate seed effects from saturation, we extracted only the first sample (sample<sub>0</sub>)  
551 from each seed and computed the per-system DockQ range (maximum minus minimum across  
552 seeds). The median range was 0.04-0.05, though extreme cases spanned nearly the full quality  
553 spectrum (e.g., system\_70\_70\_8EN2: DockQ 0.01-0.93 across five AF3 seeds). Across models, ap-  
554 proximately 10-15% of systems were strongly seed-sensitive and benefited disproportionately from  
555 deeper saturation sampling (Supplementary Figure 4D; Supplementary Figure 5A).

556 Cross-seed trajectory analysis further showed that deeper sampling can partially rescue poor initial  
557 seed choices. For each seed, we tracked cumulative maximum DockQ (the best quality achieved  
558 up to each sample number) and computed the median across systems (Figure 5F). Seeds producing  
559 high-quality structures early tended to maintain their advantage throughout saturation, whereas seeds  
560 that began in low-quality regions could still achieve substantial improvements with deeper sampling,  
561 though typically plateauing at lower absolute quality than favorable seeds (Figure 5F; Supplemen-  
562 tary Figure 5B). Together, these results indicate that seed selection and saturation sampling act as  
563 orthogonal optimization mechanisms: seed choice determines which region of the solution land-  
564 scape is explored, while saturation sampling refines predictions within that region and can partially  
565 mitigate but rarely fully overcome suboptimal initial trajectories (approximately 85-90% of systems  
566 remained below their best-seed ceiling).

567 Finally, we assessed whether confidence metrics reflect structural improvements achieved through  
568 saturation sampling. Across all models, changes in ipTM from baseline to  $N = 100$  showed near-  
569 zero correlation with corresponding changes in DockQ (Pearson  $r = 0.00\text{-}0.16$ ; Supplementary  
570 Figure 4B,C). This indicates that confidence scores are largely determined by properties of the initial  
571 prediction and remain insensitive to substantial quality gains obtained through additional sampling.  
572 Reliable confidence tracking would enable early stopping without ground-truth structures, but the  
573 observed ipTM-DockQ decoupling precludes such use. This behavior mirrors the DockQ-ipTM  
574 misalignment observed above, reinforcing that current confidence metrics do not reflect convergence  
575 toward higher-quality docking solutions.

576 Together, these results demonstrate that moderate sampling ( $N \approx 10\text{-}25$ ) captures most achievable  
577 quality gains at a fraction of full saturation cost, seed selection and saturation act as orthogonal  
578 optimization mechanisms, and current model confidence metrics fail to track these improvements.

594  
595  
596  
597  
598  
599  
600  
601  
602  
603  
604  
605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647

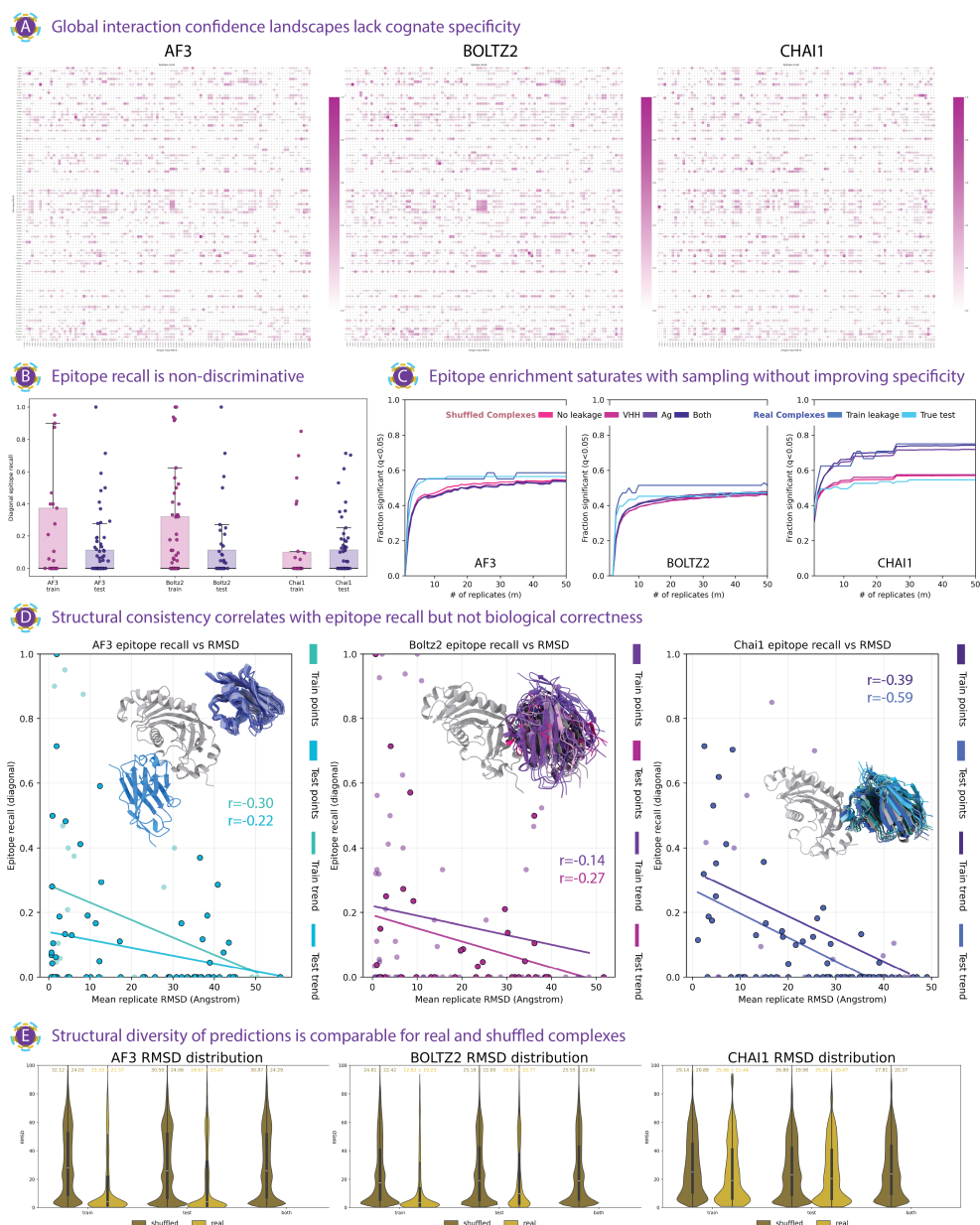


Figure 4: Epitope recovery and enrichment fail to distinguish cognate from non-cognate antibody-antigen interactions. **(A)** Epitope recall across real and shuffled complexes. Heatmaps show epitope recall for predicted antibody-antigen complexes, defined as the fraction of experimentally determined epitope residues recovered by the model. An epitope residue is considered recovered if it contacts the VHH in at least 50% of the  $n = 50$  stochastic replicate predictions. Diagonal entries correspond to real (cognate) pairs, while off-diagonal entries represent shuffled (non-cognate) complexes. **(B)** Distribution of epitope recall across real and shuffled complexes. Boxplots summarize epitope recall values stratified by train-leakage and true test sets. Shuffled complexes frequently overlap or exceed recall observed for real complexes. Numbers of systems with epitope recall  $> 0$  are indicated for each model and split. **(C)** Epitope-enrichment significance saturates rapidly with stochastic sampling. Line plots show the fraction of VHH-antigen pairs declared significantly enriched (FDR  $< 0.05$ ) as a function of the number of stochastic predictions. Real and shuffled complexes exhibit similar saturation behavior across models. **(D)** Structural consistency correlates with epitope recovery but not interaction correctness. Scatter plots show epitope recall versus mean pairwise RMSD across  $n = 50$  replicates, with linear fits shown separately for train-leakage and true test sets. Pearson correlation coefficients are reported. **(E)** Structural diversity of predicted complexes. Violin plots show RMSD distributions for real, shuffled, and mixed ("both") complexes across train-leakage and true test sets. Similar distributions indicate that structural consistency alone is insufficient to discriminate cognate from non-cognate interactions.

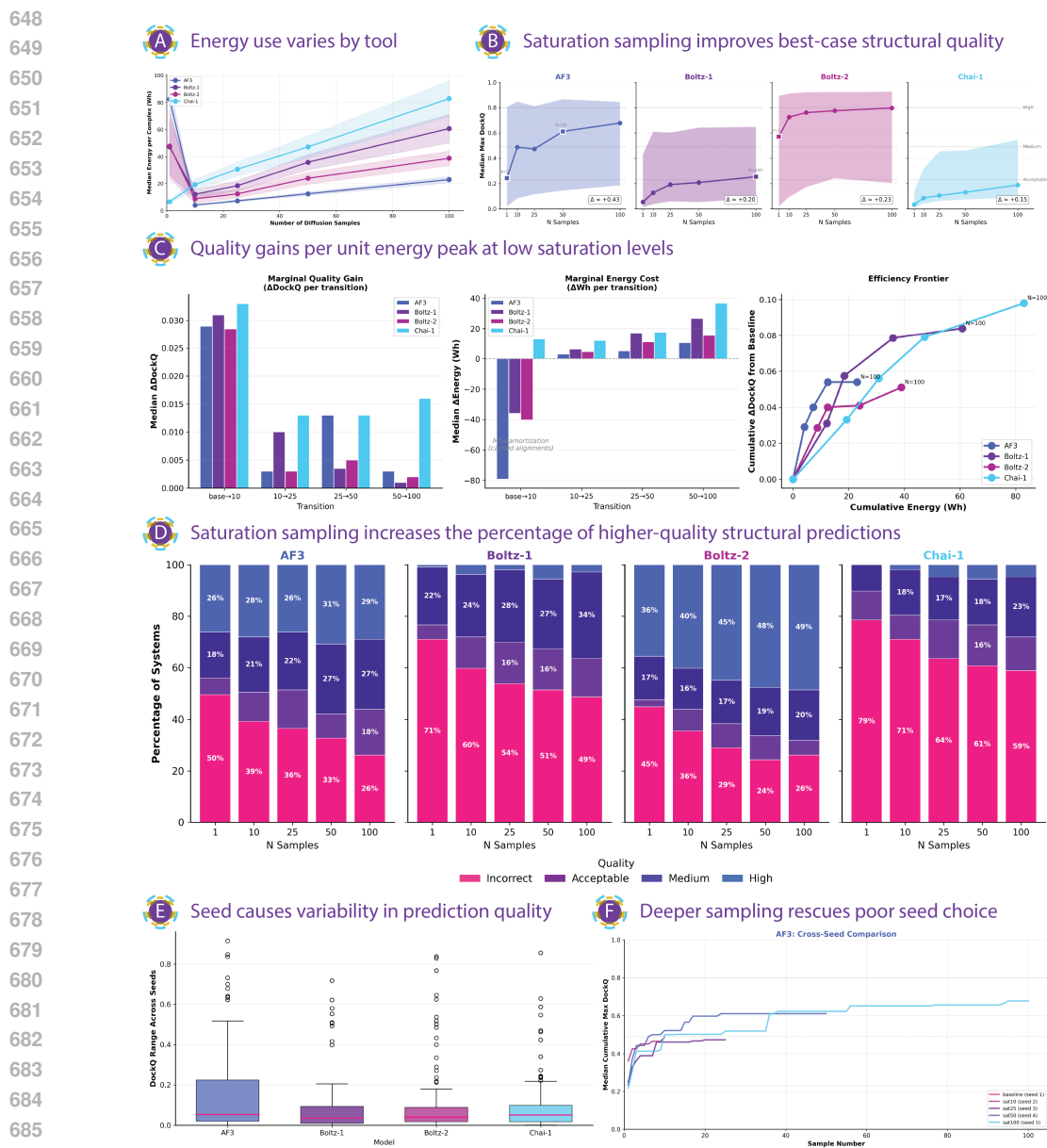
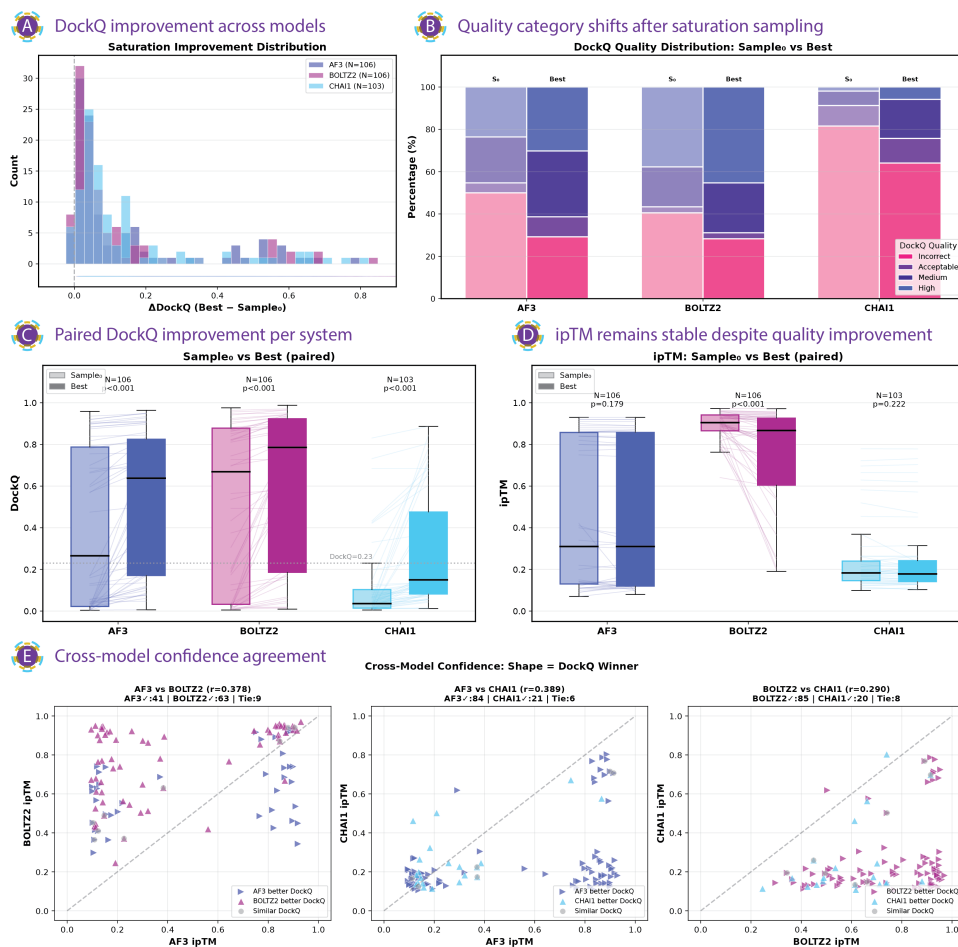


Figure 5: Computational cost, sampling efficiency, and seed-dependent optimization in antibody-antigen structure prediction. **(A)** Energy consumption differs substantially across tools and sampling depth. Median energy usage per predicted complex is shown as a function of diffusion sampling depth ( $N = 1, 10, 25, 50, 100$ ), with each saturation level corresponding to an independent random seed. Shaded regions indicate interquartile ranges across 106 systems. **(B)** Saturation sampling improves best-case structural quality with model-dependent magnitude. Median maximum DockQ is shown as a function of sampling depth. Horizontal dashed lines mark DockQ quality thresholds separating incorrect ( $< 0.23$ ), acceptable ( $0.23-0.49$ ), medium ( $0.49-0.80$ ), and high ( $\geq 0.80$ ) quality regimes. **(C)** Quality gains per unit energy peak at low saturation levels. Bar plots show median gains in maximum DockQ between consecutive saturation levels. Line plots show efficiency frontiers relating cumulative DockQ improvement to cumulative energy expenditure. **(D)** Saturation sampling reduces incorrect predictions but plateaus at higher depth. Stacked bar plots show DockQ quality-category distributions across saturation levels. **(E)** Seed selection introduces substantial variability in prediction quality. Boxplots show the DockQ range obtained from five independent seeds using only the first sample (sample<sub>0</sub>). **(F)** Different seeds explore distinct solution landscapes. Cross-seed saturation trajectories show median cumulative maximum DockQ across systems as a function of sample number. Early advantages persist, while poor initial seeds plateau at lower quality ceilings.

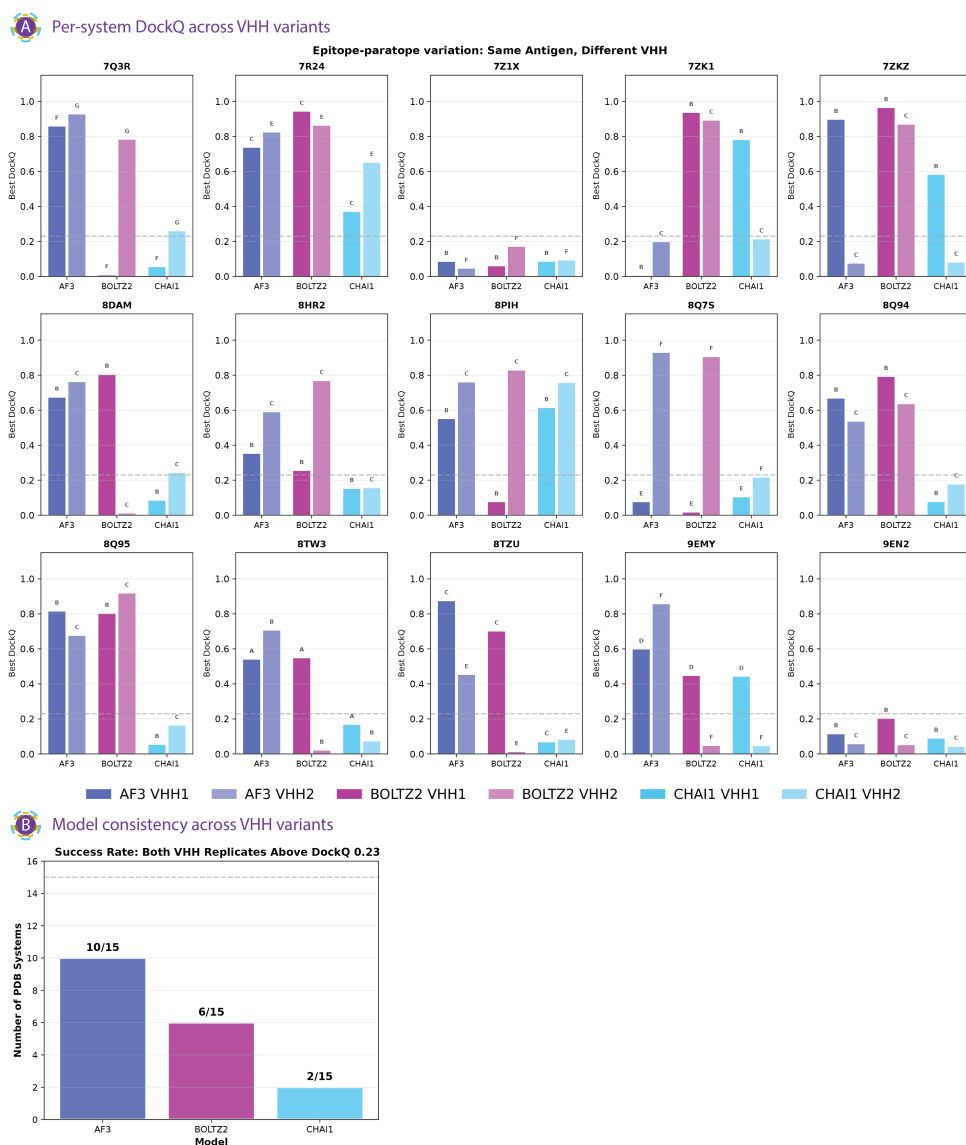
## B SUPPLEMENTARY FIGURES

Supplementary Figures 1–5 provide additional analyses that support and extend the results presented in the main text and Appendix. These figures elaborate on sampling behavior, confidence calibration, seed dependence, and computational cost across models, offering deeper quantitative insight into trends discussed earlier.



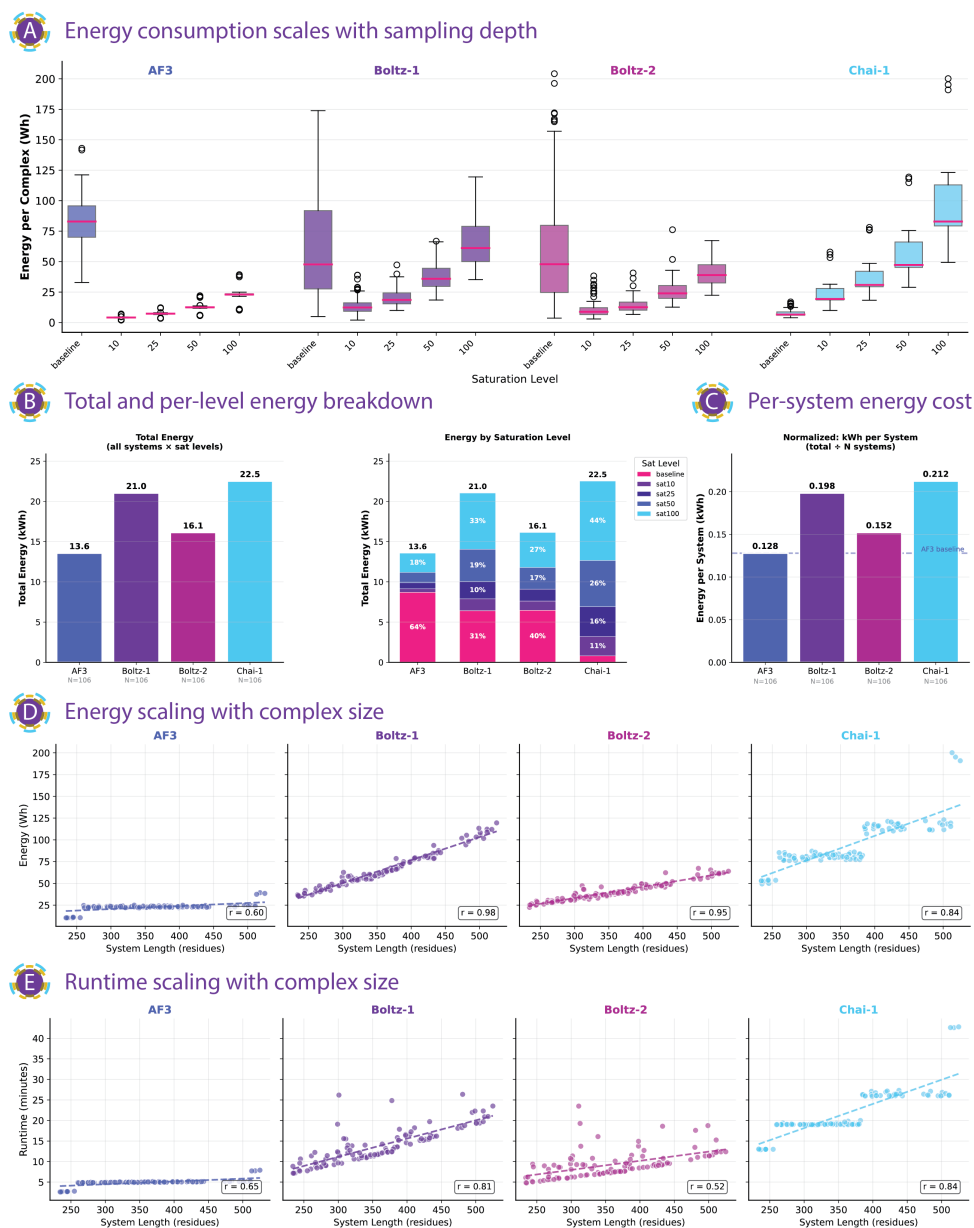
Supplementary Figure 1: **Saturation sampling improves docking quality but does not realign confidence scores across models.** (A) DockQ improvement across models. Distribution of DockQ improvement following saturation sampling, shown as  $\Delta\text{DockQ} (\text{Best} - \text{Sample}_0)$  for AF3 ( $N = 106$ ), Boltz-2 ( $N = 106$ ), and Chai-1 ( $N = 103$ ). All three models exhibit right-skewed improvement distributions, indicating that additional sampling frequently refines docking geometry. (B) Quality category shifts after saturation sampling. DockQ quality category distributions (Incorrect, Acceptable, Medium, High) for the initial single-sample prediction ( $\text{Sample}_0$ ) and the best structure obtained after saturation sampling (Best). Saturation sampling substantially increases the fraction of acceptable or higher-quality models for AF3 and Boltz-2, with more limited gains for Chai-1. (C) Paired DockQ improvement per system. Paired comparison of DockQ scores between  $\text{Sample}_0$  and Best predictions for each system. Lines connect paired predictions; boxplots summarize score distributions. All three models show significant improvements in DockQ after saturation sampling (Wilcoxon signed-rank test,  $p < 0.001$ ). The dashed line indicates the DockQ = 0.23 acceptability threshold. (D) ipTM remains stable despite quality gains. Paired comparison of ipTM confidence scores between  $\text{Sample}_0$  and Best predictions. In contrast to DockQ, ipTM shows no consistent improvement with saturation sampling for AF3 or Chai-1, and only a variable decrease for Boltz-2 (p-values shown), indicating weak coupling between confidence and structural refinement. (E) Cross-model confidence agreement. Scatter plots show pairwise ipTM values between model pairs (AF3 vs. Boltz-2, AF3 vs. Chai-1, Boltz-2 vs. Chai-1), with point shapes indicating the model achieving higher DockQ for each system and dashed lines indicating equal confidence. Pearson correlation coefficients and winner counts are reported for each comparison, revealing modest cross-model agreement and frequent mismatches between higher confidence and better structural quality.

756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809



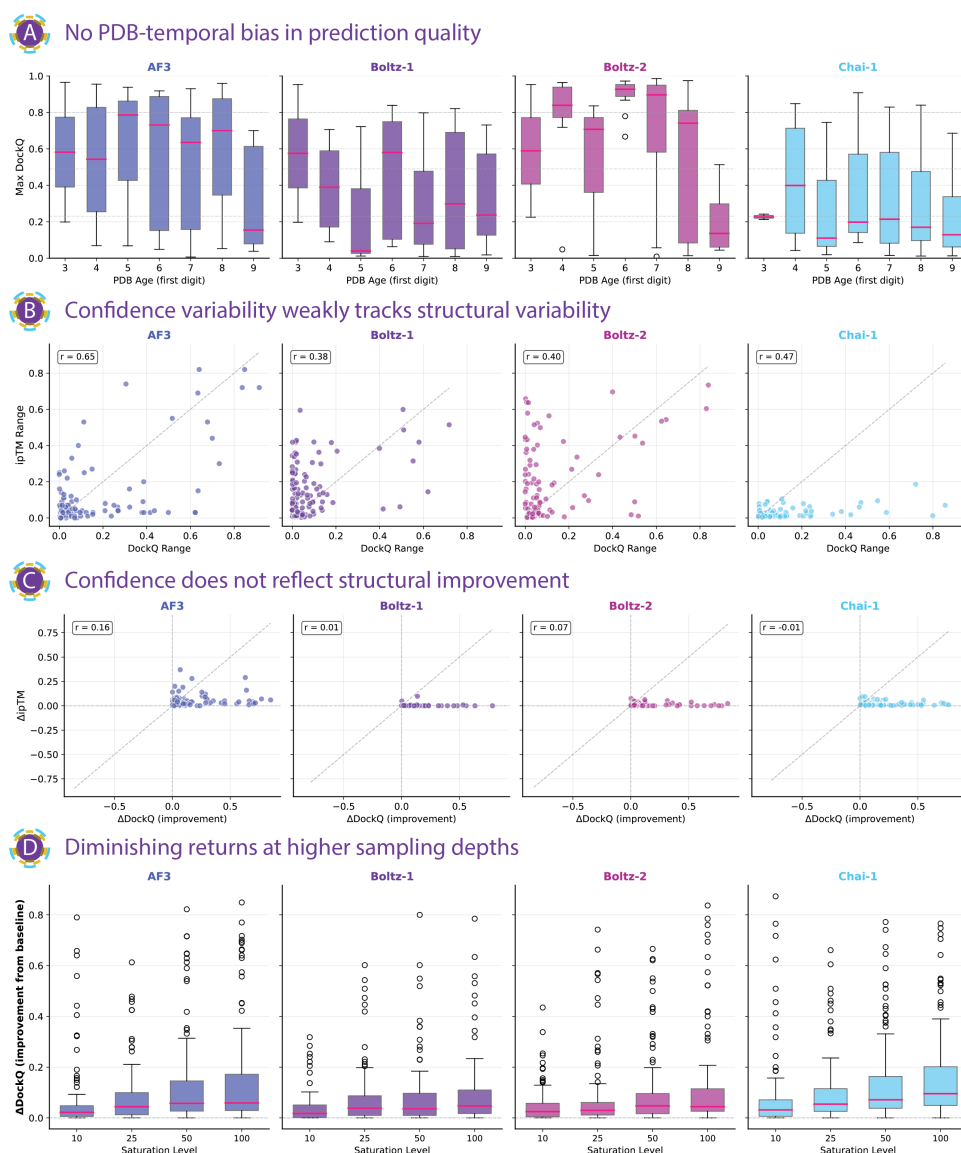
Supplementary Figure 2: **Robustness of docking predictions across VHH variants.** (A) Per-system DockQ across VHH variants. Best DockQ scores obtained after saturation sampling for 15 antigen systems, each evaluated with two distinct VHH binders (VHH1 and VHH2). For each antigen (PDB ID shown above each panel), bars show the best DockQ achieved by AF3, Boltz-2, and Chai-1 for each VHH replicate. The dashed horizontal line indicates the DockQ = 0.23 threshold for acceptable docking quality. Letter annotations denote relative ranking of predictions within each antigen system. Consistent high DockQ across both VHHs indicates robust recovery of the antigen epitope despite paratope variation, whereas discordant outcomes highlight sensitivity to VHH identity. (B) Model consistency across VHH variants. Bars show the number of antigen systems (out of 15) for which both VHH replicates achieved acceptable docking quality (DockQ  $\geq$  0.23). AF3 succeeds in 10/15 systems, compared to 6/15 for Boltz-2 and 2/15 for Chai-1, indicating superior consistency of AF3 in recovering correct antibody-antigen interactions across VHH variants.

810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863



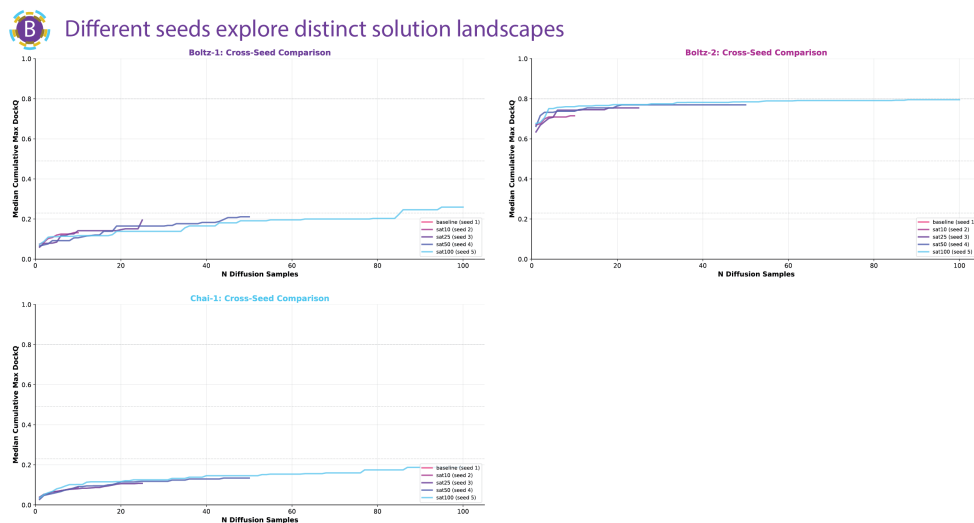
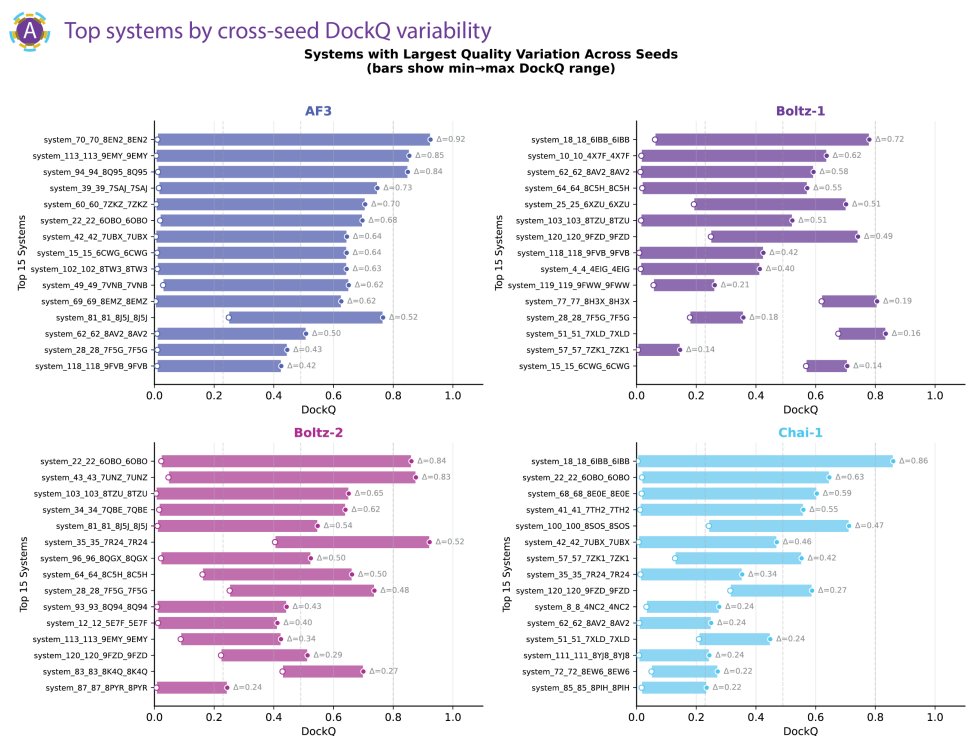
Supplementary Figure 3: **Computational cost of diffusion-based antibody-antigen docking across models.**

(A) Energy consumption per complex as a function of sampling depth. Boxplots show GPU energy usage (Wh) per system for AF3, Boltz-1, Boltz-2, and Chai-1 across increasing numbers of diffusion samples (baseline, 10, 25, 50, and 100). Each point corresponds to a single system; boxes summarize medians and interquartile ranges. (B) Aggregate energy consumption across all systems and sampling levels. Left: total GPU energy consumed (kWh) across all 106 systems and all sampling depths for each model. Middle: contribution of each sampling level to total energy usage, shown as stacked bars with percentages indicating relative contribution. (C) Normalized energy usage per system. Bars show total energy expenditure per system (kWh) aggregated across all sampling levels, enabling direct comparison of average computational cost between models. (D) Scaling of energy consumption with system size. Scatter plots show GPU energy usage (Wh) versus total system length (number of residues) for each model, with dashed lines indicating linear fits. Pearson correlation coefficients quantify scaling strength. (E) Scaling of runtime with system size. Scatter plots show wall-clock runtime (minutes) versus system length for each model, with linear fits and Pearson correlation coefficients shown. Runtime scaling mirrors energy scaling, with model-specific differences in slope and variance.



Supplementary Figure 4: **Temporal effects, confidence variability, and sampling-dependent improvements across models.** (A) Prediction quality as a function of target age. Boxplots show maximum DockQ distributions for AF3, Boltz-1, Boltz-2, and Chai-1 grouped by the first digit of the PDB deposition year. No systematic degradation in performance is observed for older structures. (B) Relationship between structural variability and confidence variability. Scatter plots show the range of ipTM values versus the range of DockQ scores across stochastic replicates for each system. Pearson correlation coefficients indicate weak-to-moderate coupling depending on the model. (C) Coupling between structural improvement and confidence change. Scatter plots show changes in ipTM ( $\Delta\text{ipTM}$ ) versus improvements in DockQ ( $\Delta\text{DockQ}$ ) from baseline to best sampled structure. Near-zero correlations indicate confidence scores do not track refinement. (D) Distribution of DockQ improvement as a function of sampling depth. Boxplots show  $\Delta\text{DockQ}$  relative to baseline across increasing sample counts (10, 25, 50, 100). Deeper sampling expands the upper tail of improvement, while median gains show diminishing returns.

918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971



Supplementary Figure 5: **Seed-dependent variability and cross-seed convergence of docking quality.** (A) Systems with the largest seed-dependent variability. Horizontal bar plots show the top 15 systems per model (AF3, Boltz-1, Boltz-2, Chai-1) ranked by the DockQ range across five independent seeds ( $\Delta = \max - \min$  DockQ). Some systems span nearly the full DockQ scale. (B) Cross-seed convergence of sampling trajectories. Line plots show median cumulative maximum DockQ as a function of diffusion sample count for each seed. Curves illustrate how seeds explore distinct solution landscapes and converge to different quality plateaus, revealing partial convergence and persistent seed dependence.

## C METHODS

### C.1 VHH-ANTIGEN DATASET CURATION

For benchmarking antibody-antigen complex prediction, VHH-antigen structures were curated from two complementary sources: SAbDab-nano (Schneider et al., 2022) (downloaded March 2025; 90%

972 sequence redundancy cutoff; bound complexes only) and the Antigen-Antibody Complex Database  
973 (AACDB) (Zhou et al., 2025). To minimize overlap with the training data of AlphaFold3 (AF3),  
974 Boltz-2, and Chai-1, only post-October 2021 depositions were retained from SAbDab-nano, cor-  
975 responding to the earliest training cutoff among the evaluated tools (Boltz-1). To explicitly probe  
976 potential memorization effects, pre-cutoff structures from AACDB were retained and later used to  
977 define train-leakage subsets.

978 All structures were filtered to include protein antigens only, with crystallographic resolution  $\leq$   
979 3.0 Å, VHH length between 110-150 amino acids, and antigen length between 100-400 amino acids.  
980 The two datasets were merged and exact PDB-chain duplicates were removed. For epitope-paratope  
981 variation analysis, we defined a set of distinct VHs binding the same antigen within a single PDB  
982 entry (Supplementary Table 1). This yielded 15 PDB entries containing such replicates (30 systems).

983 To maintain computational feasibility, the total number of unique systems was capped below 125.  
984 After manual inspection, 17 systems were excluded due to structural artifacts, VHs contacting  
985 multiple antigen chains, or ambiguous chain annotation. The final benchmark comprised 106 VHH-  
986 antigen systems from 91 unique PDB entries.

987 Antigen secondary structure was assigned using DSSP, amino-acid composition was computed di-  
988 rectly from sequences, and PDB novelty was inferred from the first character of the PDB accession  
989 code.  
990

## 991 C.2 REAL VS. SHUFFLED COMPLEX GENERATION

992  
993 To assess interaction specificity, we constructed a full combinatorial pairing matrix between the  
994 curated VHs and antigens. For each VHH, the experimentally observed pairing with its cognate  
995 antigen was designated as the *real* complex, while all other VHH-antigen pairings were designated  
996 as *shuffled* complexes. This design preserves realistic molecular interfaces while systematically  
997 breaking biological specificity.  
998

999 Unless otherwise stated, all analyses were performed on the complete  $106 \times 106$  VHH-antigen  
1000 matrix, enabling direct comparison between real and shuffled predictions under identical modeling  
1001 conditions.  
1002

## 1003 C.3 STRUCTURE PREDICTION

1004  
1005 Structure predictions for the full combinatorial dataset (106 VHs  $\times$  106 antigens; 50 samples per  
1006 complex) were performed using three independent structure prediction frameworks: AlphaFold3  
1007 (AF3), Boltz-2, and Chai-1. Predictions were executed across multiple high-performance computing  
1008 environments. Exact hardware specifications, runtime stacks, and software versions were logged at  
1009 job start for reproducibility.  
1010

### 1011 C.3.1 CHAI-1

1012 Chai-1 predictions were performed using version 0.6.1. All complexes were run on  
1013 the Immunohub cluster using one GPU per job, starting on August 14, 2025 at  
1014 12:00 PM. Each system was predicted with 50 models using the following param-  
1015 eters: `num_trunk_samples=5`, `num_diffn_samples=10`, `num_trunk_recycles=3`,  
1016 `num_diffn_timesteps=200`, `seed=42`. No MSAs were computed; instead, ESM-based em-  
1017 beddings were used under Chai-1’s default inference mode.  
1018

### 1019 C.3.2 BOLTZ-2

1020  
1021 Predictions were performed using the Boltz CLI v2.2.0. Boltz-1 and Boltz-2 were  
1022 selected via `-model boltz1` or `-model boltz2`, respectively. Unless otherwise  
1023 noted, all runs used `-use_msa_server`, `recycling_steps=3`, `sampling_steps=200`,  
1024 `diffusion_samples=50`, and a fixed random seed (42). Outputs were written in PDB or mm-  
1025 CIF format, with full per-residue and interaction-level scores exported using `-write_full_pde`.  
Boltz-2 jobs ran on Immunohub and Boltz-1 on VEGA.

### 1026 C.3.3 ALPHAFOLD3

1027  
1028 AlphaFold3 predictions were performed using `alphafold3-3.0.1` via Singularity containers.  
1029 Each system was predicted with 50 diffusion samples using `num_diffusion_samples=50`,  
1030 `run_data_pipeline=true`, and `run_inference=true`. Each job was run with a fixed ran-  
1031 dom seed (`seed=1`). MSAs were computed using the built-in data pipeline.

### 1032 C.4 TRAIN, TEST, AND MIXED LABELING

1033  
1034 Each VHH-antigen system was assigned to *train*, *test*, or *both* categories based on the training cutoffs  
1035 of the corresponding model: (i) AF3: ~30 September 2021, (ii) Chai-1: ~12 January 2021, and (iii)  
1036 Boltz-2: ~1 June 2023. A system was labeled *train* if both the VHH and antigen were present in the  
1037 model’s training set, *test* if neither was present, and *both* if one component was present while the  
1038 other was not.

### 1039 C.5 STRUCTURAL SIMILARITY BETWEEN REPLICATES

1040  
1041 To quantify structural reproducibility, geometric and contact-based similarity metrics were com-  
1042 puted across independently generated predictions. All analyses were performed on  $C\alpha$  atomic  
1043 coordinates extracted from PDB files using custom Python scripts (NumPy, SciPy, pandas). For each  
1044 complex, all pairwise comparisons among replicate structures were computed and aggregated to  
1045 yield per-system metrics.  
1046

#### 1047 C.5.1 ROOT MEAN SQUARE DEVIATION (RMSD)

1048  
1049 Global structural similarity was quantified using RMSD after superposition on the antigen via the  
1050 Kabsch algorithm:

$$1051 \text{RMSD} = \sqrt{\frac{1}{N} \sum_{i=1}^N \|\mathbf{X}_i - (\mathbf{R}\mathbf{Y}_i + \mathbf{t})\|^2},$$

1052  
1053 where  $\mathbf{X}$  and  $\mathbf{Y}$  are corresponding  $C\alpha$  coordinates, and  $\mathbf{R}$  and  $\mathbf{t}$  are the optimal rotation and trans-  
1054 lation. Lower RMSD values indicate higher structural consistency. RMSD values were averaged  
1055 across all replicate pairs to obtain mean  $\pm$  SD estimates per system.  
1056

### 1057 C.6 ENERGY MONITORING

1058  
1059 GPU telemetry (power draw, utilization, memory usage) was sampled every 5 seconds using  
1060 `nvidia-smi` and written to per-run CSV files. Wall-clock timing and system metrics were cap-  
1061 tured via `/usr/bin/time -v`. Total energy consumption (Wh) was calculated as:  
1062

$$1063 E_{\text{run}} = \sum_i P_i \Delta t / 3600,$$

1064  
1065 with  $\Delta t = 5$  s. Runs were aggregated by model and diffusion sample count and reported as mean  
1066  $\pm$  SD.  
1067

### 1068 C.7 COMPUTATIONAL COST AND SAMPLING EFFICIENCY

1069  
1070 Structural quality was assessed using DockQ (Mirabello & Wallner, 2024), which integrates in-  
1071 terface RMSD, ligand RMSD, and fraction of native contacts into a single score ranging from 0  
1072 (incorrect) to 1 (perfect). Quality categories followed CAPRI conventions.  
1073

1074 For saturation analysis, each system was evaluated using five independent random seeds, with each  
1075 seed assigned to a single saturation level ( $N = 1, 10, 25, 50, 100$  diffusion samples). Marginal  
1076 quality gains were computed as per-system differences in maximum DockQ between consecutive  
1077 saturation levels and summarized as medians across systems. Efficiency frontiers were generated by  
1078 plotting cumulative median  $\Delta$ DockQ against cumulative median energy expenditure.

1079 Correlations between confidence-score changes ( $\Delta$ ipTM) and quality changes ( $\Delta$ DockQ) were as-  
sessed using Pearson correlation coefficients.

## 1080 D DATA AVAILABILITY

1081  
1082 Supplementary tables will be released upon publication.  
1083

## 1084 E DISCUSSION

1085  
1086 This work benchmarks whether modern AI complex-prediction tools can discriminate cognate  
1087 antibody-antigen binding and recover correct paratope-epitope interfaces. Our central findings are  
1088 twofold. First, across tools, high-level confidence scores frequently fail to separate real from mis-  
1089 matched ("shuffled") complexes. Second, while increased sampling can improve best-case geom-  
1090 etry, it does not resolve the upstream problem of selecting the correct binding mode. Importantly,  
1091 even when structural quality improves substantially under aggressive sampling (as measured, for  
1092 example, by DockQ for real complexes), model confidence scores often remain unchanged, reveal-  
1093 ing a disconnect between structural refinement and confidence calibration. Below, we place these  
1094 results in the broader AI biologics landscape and outline practical implications for drug discovery  
1095 and model development.  
1096

### 1097 E.1 WHERE THE FIELD STANDS: DISCONNECT BETWEEN *de novo* GENERATION AND *de novo* 1098 PREDICTION

1099 The AI biologics ecosystem is now shaped by strong competition among major players (e.g., La-  
1100 tent, Chai, Nabla, Isomorphic, Boltz), with rapid iteration cycles. Despite headline claims, although  
1101 *de novo* generation success rates are increasing into the double-digit range (approximately 10-15%  
1102 on difficult interface tasks (Team et al., 2025; Stark et al., 2025)), our results indicate a disconnect  
1103 between *de novo* design ("generate an antibody for a given target without exploiting other informa-  
1104 tion") and *de novo* prediction ("identify all antigens that can be bound by a given antibody, and vice  
1105 versa"). While such success rates may appear encouraging in isolation, they remain low in absolute  
1106 terms for practical screening scenarios, where even modest false-positive rates translate into sub-  
1107 stantial experimental burden. Our real-versus-shuffled discrimination benchmark demonstrates that  
1108 models can generate interfaces that appear plausible across many pairings, producing an abundance  
1109 of false positives that are challenging to triage.

1110 An emerging industry screening workflow is to generate thousands to millions of candidates and  
1111 then filter using internal confidence scores (e.g., pLDDT, ipTM, PAE). Our results suggest that this  
1112 strategy is fragile for antibody-antigen binding, because confidence does not equate to accurate biol-  
1113 ogy. In a many-VHH-versus-many-antigen screening setting, shuffled complexes frequently achieve  
1114 confidence comparable to real complexes, while some high-quality structures remain underconfident  
1115 depending on the tool. This breaks the assumption that "high-confidence docking implies correct  
1116 binding" and helps explain why confidence-guided selection can yield many experimental failures  
1117 even when predicted structures appear polished. Consistent with recent proposals for alternative  
1118 evaluation metrics (Xu et al., 2025; Almeida et al., 2025), our results underscore the need for scores  
1119 that more directly reflect biological meaning rather than internal structural self-consistency.

1120 Crucially, drug viability depends on properties that these scores do not directly encode, including  
1121 induced-fit compatibility, entropic penalties, off-target propensity, developability constraints, and  
1122 tolerance to antigen dynamics or conformational selection. Confidence metrics were not designed  
1123 as surrogates for molecular functionality, and our data reinforce that treating them as such inflates  
1124 false positives, especially in interface-dependent problems.

### 1125 E.2 WHY EVOLUTIONARY SIGNALS AND SAMPLING HELP LESS THAN HOPED FOR 1126 ANTIBODIES

1127  
1128 A common intuition is that better MSAs and richer evolutionary information should improve AI  
1129 docking and design. However, antibody paratopes-especially CDR loops-are among the most vari-  
1130 able regions in biology, and much of binding specificity arises from flexibility, conformational di-  
1131 versity, and context-dependent loop rearrangements rather than sequence conservation (Smorodina  
1132 et al., 2025b; Spondlin et al., 2025; Fernández-Quintero et al., 2021). This supports the view  
1133 that the utility of MSAs in antibody-focused tasks is more nuanced than in general protein struc-  
ture prediction. While MSAs are critical for accurate monomer modeling (removal of MSAs in

1134 AF3 leads to substantial degradation in structural accuracy (Ali et al., 2025)), their contribution to  
1135 antibody-antigen docking is constrained by the extreme diversity of CDRs, which limits informative  
1136 evolutionary signal. As a result, MSA-derived confidence estimates may be inherently less reliable  
1137 for antibody-antigen interfaces than for other protein complexes. Alternative MSA strategies may  
1138 therefore provide better guidance, particularly for evolutionary complex or rapidly mutating systems  
1139 such as antibodies.

1140 On the other hand, sampling helps but exposes a deeper bottleneck: path selection versus refine-  
1141 ment. Saturation or diffusion sampling improves best-case DockQ in many systems, confirming that  
1142 sampling is valuable as a refinement mechanism. However, we observe essentially zero correla-  
1143 tion between changes in DockQ ( $\Delta\text{DockQ}$ ) and changes in ipTM ( $\Delta\text{ipTM}$ ) across all three models  
1144 (Pearson  $r = -0.03, -0.04,$  and  $-0.02$  for AF3, Boltz-2, and Chai-1, respectively), indicating that  
1145 confidence scores do not track structural improvement. The persistence of "non-improver" systems  
1146 supports a two-stage view: (i) path or seed selection determines the qualitative docking mode (often  
1147 incorrect), and (ii) saturation refines within that chosen mode.

1148 In this sense, current models appear to have a "fixed mindset": once a confidence level is assigned  
1149 to a given seed or trajectory, it remains largely invariant even when the resulting structure improves  
1150 substantially. This likely reflects an architectural constraint: confidence scores such as ipTM are  
1151 derived from MSA-based pairwise representations computed in the trunk network, which remains  
1152 fixed during diffusion sampling. If the model commits early to an incorrect binding mode, additional  
1153 diffusion samples often cannot rescue it. This has immediate practical implications: simply sampling  
1154 more may be an expensive path to diminishing returns, and large seed sweeps, while sometimes  
1155 recommended, can be prohibitive at screening scale without guaranteeing specificity or providing  
1156 confidence-based validation that sampling helped. Developing confidence measures sensitive to  
1157 binding-mode quality, or orthogonal strategies such as consensus across independent seeds, remains  
1158 an important challenge.

### 1159 1160 1161 E.3 STRUCTURAL BIOLOGY PERSPECTIVE: FLEXIBILITY IS NOT OPTIONAL 1162

1163  
1164 Antibody recognition is not merely geometric matching; it is often a dynamic process in which  
1165 CDR loops adapt to the antigen surface and exploit transient conformations. This is particularly  
1166 relevant for targets with disordered regions or multiple accessible states. A major unresolved chal-  
1167 lenge is therefore distinguishing hallucinated (geometrically plausible but physically or functionally  
1168 implausible) complexes from viable ones that can maintain binding under realistic dynamics.

1169 This points to a missing ingredient in many AI pipelines: standardized representations of flexibility  
1170 and experimentally grounded dynamics. A key practical bottleneck in current AI biologics work-  
1171 flows is selection. Generating candidates is increasingly easy; selecting the biologically functional  
1172 subset is not. We argue that molecular-dynamics-informed filtering is a promising bridge between  
1173 AI-generated hypotheses and functional plausibility (Bashour et al., 2024; Park & Izadi, 2024).  
1174 However, today's AI-generated "MD-like" outputs are often not physically faithful enough to serve  
1175 as ground truth.

1176 Standardized, cross-method training data for dynamics therefore matter. Efforts such as the DINO  
1177 proposal (Smorodina et al., 2025c; Anonymous, 2025), aimed at becoming a "PDB-equivalent for  
1178 dynamics," directly address the lack of unified MD data across force fields and simulation proto-  
1179 cols. By enabling consistent learning from flexibility and improving computational-experimental  
1180 alignment, such resources could underpin a new generation of screening criteria that go beyond  
1181 static confidence scores, including interface stability under perturbation, persistence of key contacts,  
1182 accessibility of binding-competent states, and robustness to antigen motion.

1183 As many AI startups increasingly partner with pharmaceutical companies, the core industrial ques-  
1184 tion is not "which model has the highest benchmark score?" but "which pipeline yields validated  
1185 candidates efficiently?" Our benchmark suggests that confidence-score-based filtering alone is in-  
1186 sufficient for antibody-antigen interface prediction. The field may need to shift from (i) generating  
1187 more candidates and ranking by confidence to (ii) ranking by biophysical plausibility and functional  
robustness, potentially using MD-informed filters and interface-specific mechanistic constraints.

1188 E.4 LIMITATIONS AND FUTURE DIRECTIONS  
1189

1190 Our study focuses on VHH-antigen complexes under a controlled real-versus-shuffled pairing design  
1191 that explicitly tests specificity rather than docking plausibility alone. As such, it quantifies the  
1192 feasibility of generalized binding prediction. This design may overestimate performance limitations  
1193 in narrower settings where the antigen epitope is known, constraints are available, or the target is  
1194 rigid. Additionally, similar benchmarks should be extended to more complex antibody formats to  
1195 assess generalizability beyond VHHs.

1196 Large-scale resources such as AbSet (Almeida et al., 2025), which contains hundreds of thousands of  
1197 antibody structures, highlight an additional challenge: volume does not imply reliability. Systematic  
1198 assessment of how many computational antibody structures are structurally sound, interface-correct,  
1199 or functionally meaningful remains limited. Our results support the view that antibody-antigen  
1200 complex prediction is not yet a solved problem despite massive *in silico* datasets.

1201 Looking forward, progress will likely require (1) interface-specific training objectives incorporating  
1202 hard, mutation-diverse negatives across affinity ranges (Ursu et al., 2024); (2) confidence estimates  
1203 calibrated to specificity and functional plausibility rather than structural self-consistency; and (3)  
1204 standardized dynamic datasets enabling models and filters to account for flexibility and induced fit.  
1205 The most impactful near-term advance may not be a new confidence score, but a robust, scalable  
1206 post-prediction filtering layer grounded in biophysics and dynamics.

1207 In summary, this study addresses a central challenge in immunotherapy: accurate modeling of  
1208 antibody-antigen interactions. Despite major advances in structure prediction, current tools fall  
1209 short in reliably identifying correct paratope-epitope interfaces. Through systematic benchmarking,  
1210 we expose key failure modes in interface identification and binding-mode selection. These findings  
1211 suggest that the next leap in AI biologics may come not from increasingly sophisticated architectures  
1212 alone, but from moving beyond static confidence proxies toward biochemical, mechanistic,  
1213 and flexibility-aware evaluation of binding-and from building the data infrastructure needed to make  
1214 that possible.

1215  
1216 REFERENCES

- 1217 Website. <https://arxiv.org/html/2512.11892v1>, 2025.
- 1218
- 1219 J. Abramson et al. Accurate structure prediction of biomolecular interactions with alphafold 3.  
1220 *Nature*, 630:493–500, 2024.
- 1221
- 1222 R. Akbar et al. Progress and challenges for the machine learning-based design of fit-for-purpose  
1223 monoclonal antibodies. *MAbs*, 14:2008790, 2022.
- 1224
- 1225 M. Ali et al. Improving nanobody structure prediction with self-distillation. *bioRxiv*, 2025. doi:  
1226 10.64898/2025.12.01.691162.
- 1227
- 1228 D. S. Almeida et al. Abset: A standardized data set of antibody structures for machine learning  
1229 applications. *Journal of Chemical Information and Modeling*, 65:4767–4774, 2025.
- 1230
- 1231 F. Ambrosetti, B. Jiménez-García, J. Roel-Touris, and A. M. J. J. Bonvin. Modeling antibody-  
1232 antigen complexes by information-driven docking. *Structure*, 28:119–129.e2, 2020.
- 1233
- 1234 Anonymous. Probabilistic modeling of antibody structural dynamics. In *2nd edition of Frontiers  
1235 in Probabilistic Inference: Learning meets Sampling*, 2025. URL [https://openreview.  
1236 net/forum?id=26C7kYxxlp](https://openreview.net/forum?id=26C7kYxxlp).
- 1237
- 1238 H. Bashour et al. Biophysical cartography of the native and human-engineered antibody landscapes  
1239 quantifies the plasticity of antibody developability. *Communications Biology*, 7:922, 2024.
- 1240
- 1241 S. Basu and B. Wallner. Dockq: A quality measure for protein-protein docking models. *PLoS One*,  
11:e0161879, 2016.
- W. Bielska et al. Applying computational protein design to therapeutic antibody discovery – current  
state and perspectives. *arXiv*, 2025. arXiv:q-bio.BM.

- 1242 J. Boitreau et al. Chai-1: Decoding the molecular interactions of life. *Synthetic Biology*, 2024.
- 1243 P. Bradley. Structure-based prediction of t cell receptor:peptide-mhc interactions. *eLife*, 12, 2023.
- 1244 X. Dang et al. Epitope mapping of monoclonal antibodies: a comprehensive comparison of different  
1245 technologies. *MAbs*, 15:2285285, 2023.
- 1246 C. Discovery et al. Chai-1: Decoding the molecular interactions of life. *bioRxiv*, 2024. doi:  
1247 10.1101/2024.10.10.615955.
- 1248 R. Evans et al. Protein complex prediction with alphafold-multimer. *bioRxiv*, 2021. doi: 10.1101/  
1249 2021.10.04.463034.
- 1250 M. L. Fernández-Quintero, G. Georges, J. M. Varga, and K. R. Liedl. Ensembles in solution as a  
1251 new paradigm for antibody structure prediction and design. *MAbs*, 13:1923122, 2021.
- 1252 M. L. Fernández-Quintero et al. Challenges in antibody structure prediction. *MAbs*, 15:2175319,  
1253 2023.
- 1254 S. Fromm, M. Ludaic, and A. Elofsson. Evaluating deep learning based structure prediction methods  
1255 on antibody-antigen complexes. *bioRxiv*, 2025. doi: 10.1101/2025.07.11.662141.
- 1256 L. R. Genz, S. Nair, N. Nagar, and M. Topf. Assessing scoring metrics for alphafold2 and alphafold3  
1257 protein complex predictions. *Protein Science*, 34:e70327, 2025.
- 1258 V. Greiff, G. Yaari, and L. G. Cowell. Mining adaptive immune receptor repertoires for biological  
1259 and clinical information using machine learning. *Current Opinion in Systems Biology*, 24:109–  
1260 119, 2020.
- 1261 M. Grieswelle et al. A new benchmark for deep learning based affinity prediction: Solving the  
1262 inter-protein scoring noise problem. *ChemRxiv*, 2025. doi: 10.26434/chemrxiv-2025-sf3cs.
- 1263 C. Guo, G. Pleiss, Y. Sun, and K. Q. Weinberger. On calibration of modern neural networks. *arXiv*,  
1264 2017. doi: 10.48550/ARXIV.1706.04599. arXiv:1706.04599 [cs.LG].
- 1265 A. Harmalkar, S. Lyskov, and J. J. Gray. Reliable protein-protein docking with alphafold, rosetta,  
1266 and replica-exchange. *bioRxiv*, 2023. doi: 10.1101/2023.07.28.551063.
- 1267 F. N. Hitawala and J. J. Gray. What does alphafold3 learn about antigen and nanobody docking, and  
1268 what remains unsolved? *Bioengineering*, 2024.
- 1269 F. N. Hitawala and J. J. Gray. What does alphafold3 learn about antibody and nanobody docking,  
1270 and what remains unsolved? *MAbs*, 17:2545601, 2025.
- 1271 C. M. Holt et al. Contrastive learning enables epitope overlap predictions for targeted antibody  
1272 discovery. *bioRxiv*, 2025. doi: 10.1101/2025.02.25.640114.
- 1273 I. Johansson-Åkhe and B. Wallner. Improving peptide-protein docking with alphafold-multimer  
1274 using forced sampling. *Frontiers in Bioinformatics*, 2:959160, 2022.
- 1275 J. Jumper et al. Highly accurate protein structure prediction with alphafold. *Nature*, 596:583–589,  
1276 2021.
- 1277 M. Kull et al. Beyond temperature scaling: Obtaining well-calibrated multiclass probabilities with  
1278 dirichlet calibration. *arXiv*, 2019. doi: 10.48550/ARXIV.1910.12656. arXiv:1910.12656 [cs.LG].
- 1279 D. Kuroda and J. J. Gray. Shape complementarity and hydrogen bond preferences in protein-protein  
1280 interfaces: implications for antibody modeling and protein-protein docking. *Bioinformatics*, 32:  
1281 2451–2456, 2016.
- 1282 L. Lannelongue and M. Inouye. Environmental impacts of machine learning applications in protein  
1283 science. *Cold Spring Harbor Perspectives in Biology*, 15, 2023.
- 1284 A. H. Laustsen, V. Greiff, A. Karatt-Vellatt, S. Muyldermans, and T. P. Jenkins. Animal immu-  
1285 nization, in vitro display technologies, and machine learning for antibody discovery. *Trends in*  
1286 *Biotechnology*, 39:1263–1273, 2021.

- 1296 S. Levine et al. Origin-1: a generative ai platform for de novo antibody design against novel epitopes.  
1297 *bioRxiv*, 2026. doi: 10.64898/2026.01.14.699389.
- 1298
- 1299 F. Meng et al. A comprehensive overview of recent advances in generative models for antibodies.  
1300 *Computational and Structural Biotechnology Journal*, 23:2648–2660, 2024.
- 1301
- 1302 C. Mirabello and B. Wallner. Dockq v2: improved automatic quality measure for protein multimers,  
1303 nucleic acids, and small molecules. *Bioinformatics*, 40, 2024.
- 1304 C. Morand, A.-L. Ligozat, and A. Névéol. How green can ai be? a study of trends in machine learn-  
1305 ing environmental impacts. *arXiv*, 2024. doi: 10.48550/ARXIV.2412.17376. arXiv:2412.17376  
1306 [cs.LG].
- 1307 M. D. Overath et al. Predicting experimental success in de novo binder design: A meta-analysis of  
1308 3,766 experimentally characterised binders. *bioRxiv*, 2025. doi: 10.1101/2025.08.14.670059.
- 1309
- 1310 M. D. Overath et al. Accelerating multi-objective vhh discovery via integrated high-throughput  
1311 selection and alphafold3-guided structure prediction. *bioRxiv*, 2026. doi: 10.64898/2026.01.19.  
1312 700436.
- 1313 E. Park and S. Izadi. Molecular surface descriptors to predict antibody developability: sensitivity to  
1314 parameters, structure models, and conformational sampling. *MABs*, 16:2362788, 2024.
- 1315
- 1316 S. Passaro et al. Boltz-2: Towards accurate and efficient binding affinity prediction. *bioRxiv*, 2025.  
1317 doi: 10.1101/2025.06.14.659707.
- 1318 J. A. Ruffolo, L.-S. Chu, S. P. Mahajan, and J. J. Gray. Fast, accurate antibody structure prediction  
1319 from deep learning on massive set of natural antibodies. *Nature Communications*, 14:2389, 2023.
- 1320
- 1321 C. Schneider, M. I. J. Raybould, and C. M. Deane. Sabdab in the age of biotherapeutics: updates in-  
1322 cluding sabdab-nano, the nanobody structure tracker. *Nucleic Acids Research*, 50:D1368–D1372,  
1323 2022.
- 1324 E. Smorodina et al. Structural informatic study of determined and alphafold2 predicted molecular  
1325 structures of 13 human solute carrier transporters and their water-soluble qty variants. *Scientific*  
1326 *Reports*, 12:20103, 2022.
- 1327
- 1328 E. Smorodina et al. Structural modeling of antibody variant epitope specificity with complemen-  
1329 tary experimental and computational techniques. In *ICLR 2025 Workshop on Generative and*  
1330 *Experimental Perspectives for Biomolecular Design*, 2025a.
- 1331
- 1332 E. Smorodina et al. Structural modeling of antibody variant epitope specificity with complemen-  
1333 tary experimental and computational techniques. In *ICLR 2025 Workshop on Generative and*  
1334 *Experimental Perspectives for Biomolecular Design*, 2025b.
- 1335 Eva Smorodina, Victor Greiff, and Rahmad Akbar. DINO: dynamics-informed dataset to overcome  
1336 the limitations of static molecular data in AI-driven drug discovery. In *NeurIPS 2025 AI for*  
1337 *Science Workshop*, 2025c. URL <https://openreview.net/forum?id=ENVitzhEhh>.
- 1338
- 1339 F. C. Spoenlin et al. Predicting the conformational flexibility of antibody and t cell receptor  
1340 complementarity-determining regions. *Nature Machine Intelligence*, 7:1755–1767, 2025.
- 1341
- 1342 H. Stark et al. Boltzgen: Toward universal binder design. *bioRxiv*, 2025. doi: 10.1101/2025.11.20.  
1343 689494.
- 1344
- 1345 Chai Discovery Team et al. Zero-shot antibody design in a 24-well plate. *bioRxiv*, 2025. doi:  
1346 10.1101/2025.07.05.663018.
- 1347
- 1348 S. Unsal, B. Holland, I. Sardag, and E. Timucin. Confidence scoring for ai-predicted antibody-  
1349 antigen complexes: Anticonf as a precision-driven metric. *bioRxiv*, 2025. doi: 10.1101/2025.07.  
25.666870.
- 1349
- E. Ursu et al. Training data composition determines machine learning generalization and biological  
rule discovery. *bioRxiv*, 2024. doi: 10.1101/2024.06.17.599333.

1350 S. Weeratunga et al. Interrogation and validation of the interactome of neuronal munc18-interacting  
1351 mint proteins with alphafold2. *Journal of Biological Chemistry*, 300:105541, 2024.  
1352

1353 B. D. Weitzner et al. Modeling and docking of antibody structures with rosetta. *Nature Protocols*,  
1354 12:401–416, 2017.

1355 J. Wohlwend et al. Boltz-1: Democratizing biomolecular interaction modeling. *Biophysics*, 2024.  
1356

1357 X. Xu, I. Coratella, V. Reys, and A. M. Bonvin. Deeprank-ab: a scoring function for antibody-  
1358 antigen complexes based on geometric deep learning. *bioRxiv*, 2025. doi: 10.64898/2025.12.03.  
1359 691974.

1360 R. Yin and B. G. Pierce. Evaluation of alphafold antibody-antigen modeling with implications for  
1361 improving predictive accuracy. *Protein Science*, 33:e4865, 2024.  
1362

1363 Y. Zhou et al. A comprehensive antigen-antibody complex database unlocking insights into interac-  
1364 tion interface. *eLife*, 14, 2025.  
1365  
1366  
1367  
1368  
1369  
1370  
1371  
1372  
1373  
1374  
1375  
1376  
1377  
1378  
1379  
1380  
1381  
1382  
1383  
1384  
1385  
1386  
1387  
1388  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1400  
1401  
1402  
1403