
BLOSUM Is All You Learn — Generative Antibody Models Reflect Evolutionary Priors

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Generative models have emerged as powerful tools for antibody sequence design,
2 with recent studies demonstrating that log-likelihood scores from these models
3 can correlate with binding affinity and potentially serve as effective ranking met-
4 rics. In this work, we investigate the biochemical basis of these model-derived
5 log-likelihoods by comparing them with classical evolutionary similarity met-
6 rics. We find that BLOSUM similarity scores between designed and parental
7 antibody sequences correlate strongly with measured binding affinity—on par with
8 the predictive performance of a state-of-the-art diffusion-based generative model.
9 Moreover, these BLOSUM scores also align closely with log-likelihoods from
10 multiple generative models, suggesting that such models may be implicitly learning
11 evolutionary priors encoded in substitution matrices. In contrast, similarity scores
12 based on position weight matrices (PWMs) and position-specific scoring matrices
13 (PSSMs) that do not require the knowledge of the parental sequence show weaker
14 and less consistent alignment with binding affinity, with performance depending
15 on the source of the background sequence data. Additionally, using consensus
16 sequences in place of parental sequences to compute BLOSUM scores largely
17 eliminates the observed correlation with affinity, underscoring the context-specific
18 nature of the correlations. These findings highlight the potential of interpretable,
19 evolution-inspired metrics to complement generative modeling in antibody design,
20 offering insights into both model behavior and biological relevance.

21 1 Introduction

22 Antibody design continues to be a critical area of research for both therapeutic and diagnostic
23 applications, where the goal is to target antibodies toward a predetermined epitope of interest,
24 optimize binding affinity and specificity, and improve overall developability of antibody candidates.
25 Computational methods, particularly generative models trained on large-scale antibody sequence
26 datasets, have shown significant promise in accelerating this process. These models can propose
27 candidate sequences that adhere to learned patterns of natural antibody repertoires, potentially
28 capturing features relevant to biological function and molecular recognition. A major challenge,
29 however, lies in the development of reliable *in silico* metrics to rank designed sequences by their
30 likelihood of exhibiting high binding affinity.

31 Among available metrics, log-likelihood scores derived from generative models trained solely on
32 structural and/or sequence data have recently been shown to correlate with experimentally measured
33 binding affinity [Ucar et al., 2024], indicating their potential as a practical tool for prioritization. Yet
34 the biochemical basis of this correlation remains poorly understood. It is not clear whether these
35 scores reflect meaningful biophysical or evolutionary constraints, or whether they serve merely as
36 model-specific heuristics.

In this work, we investigate the biochemical grounding of model-derived scores by examining their relationship with classical sequence similarity metrics. Specifically, we assess whether these scores align with metrics grounded in evolutionary substitution dynamics, including BLOcks SUBstitution Matrix (BLOSUM) similarity [Henikoff and Henikoff, 1992], position weight matrices (PWM) [Stormo et al., 1982], and position-specific scoring matrices (PSSM) [Gribskov et al., 1987]. These metrics are computed using different strategies: BLOSUM similarity is evaluated between designed sequences and various types of reference sequences, while PWM and PSSM scores are derived from amino acid frequency profiles obtained from large antibody datasets.

By comparing these scores with both measured binding affinities and log-likelihoods from multiple generative models—including a diffusion-based model DiffAbXL-A [Ucar et al., 2024], an inverse folding model AntiFold [Høie et al., 2023], and a language model IgLM [Shuai et al., 2023], we aim to clarify whether such models capture biologically meaningful substitution patterns. Our results highlight the importance of context-specific similarity metrics and suggest that interpretable, evolution-inspired scores can complement generative models by improving the transparency and reliability of model-guided antibody design pipelines.

2 Related Work

Generative models have become increasingly important in protein and antibody design, leveraging advancements in deep learning to generate novel sequences with desirable properties. These models typically fall into three broad categories: large language models (LLMs) trained on sequence data [Malherbe and Uçar, 2024], graph-based models that capture spatial and topological properties [Kong et al., 2023], and diffusion-based models that simulate denoising processes over sequence and/or structure spaces [Ucar et al., 2024]. Each approach differs in how it represents and generates protein information—ranging from purely sequence-based generation to more complex structure-aware co-design frameworks.

While much of the recent focus has been on improving generative quality and integrating structure prediction tools (e.g., AlphaFold [Jumper et al., 2021], RoseTTAFold [Baek et al., 2021]), a parallel challenge lies in evaluating and ranking generated designs. Commonly used *in silico* metrics include structural scores (e.g., RMSD, ipTM, pAE) [Abramson et al., 2024] and sequence-based scores such as amino acid recovery (AAR) or log-likelihood under pretrained models [Ucar et al., 2024, Luo et al., 2022]. However, these metrics are not explicitly optimized to reflect functional outcomes like binding affinity, which limits their utility in candidate prioritization.

Recent studies have begun to explore whether log-likelihood scores from generative models can serve as more functionally meaningful evaluation metrics. For instance, Shanehsazzadeh et al. [2023b] observed that higher-likelihood sequences generated by IgMPNN yielded higher proportions of binders, though the correlation was indirect and assessed via enrichment metrics. Other work in general protein fitness prediction has used zero-shot likelihood ranking across mutational scans [Truong Jr and Bepler, 2023], but the results have varied across assay types. Specifically for antibodies, Chungyoun et al. [2024] found that log-likelihood does not always align with functional readouts such as binding or expression, suggesting that model scores alone may be insufficient.

More recently, Ucar et al. [2024] systematically evaluated log-likelihood scores from a diffusion-based antibody generative model (DiffAbXL-A), finding that these scores consistently and significantly correlate with measured binding affinity across diverse datasets. This supports the view that generative models can internalize biophysical and evolutionary constraints relevant to antigen binding, even without explicit supervision on affinity labels.

Moreover, BLOSUM matrices have long been used to quantify functional similarity between sequences and are known to reflect evolutionary constraints. Prior efforts have used such metrics for sequence alignment and clustering but not explicitly for affinity ranking in design contexts [Altschul et al., 1990]. This work builds on prior studies by directly investigating the biochemical basis of log-likelihood as a scoring function. We examine whether classical evolutionary similarity metrics—such as those derived from BLOSUM, PWM, and PSSM—are predictive of binding affinity and whether they correlate with the scores produced by generative models.

88 3 Method

89 We evaluate the relationship between log-likelihood scores from generative models and classical evolutionary metrics for predicting antibody binding affinity. Specifically, we compare four approaches: 90 (i) log-likelihoods from three representative generative models — a diffusion-based model, an inverse 91 folding model, and a language model [Ucar et al., 2024, Høie et al., 2023, Shuai et al., 2023], (ii) 92 BLOSUM-based scoring [Henikoff and Henikoff, 1992], (iii) PWM-derived scores [Stormo et al., 93 1982], and (iv) PSSM-derived scores [Gribskov et al., 1987]. All methods are evaluated across 94 multiple datasets with experimentally measured affinities. 95

96 **Scoring Region Masking.** To ensure consistent comparison across scoring methods, we define 97 a dataset-specific mutation mask \mathcal{M} that identifies positions relevant for scoring. For BLOSUM, 98 PWM, and PSSM-based evaluations, we first align all sequences using the AHO numbering scheme 99 [Honegger and PluÈckthun, 2001], and include a position in \mathcal{M} if it differs between the parental 100 sequence and any of its variants. In contrast, for log-likelihood-based scoring using the DiffAbXL-A 101 and IgLM models, no alignment is performed; the mask is computed directly from raw sequence 102 positions that differ from the parental sequence. For AntiFold, which uses IMGT numbering internally 103 to define complementarity-determining regions (CDRs), we provide the model with the specific list of 104 CDRs designed in each library to define \mathcal{M} . This approach maintains compatibility with how these 105 generative models process input and ensures that each method is evaluated in its appropriate context.

106 **Log-likelihood scoring.** We evaluate log-likelihoods using three different generative models: 107 DiffAbXL-A, AntiFold, and IgLM.

108 **DiffAbXL-A.** DiffAbXL-A is a scaled variant of the diffusion-based model DiffAb [Luo et al., 109 2022], trained to generate all six complementarity-determining regions (CDRs) of an antibody given 110 structural context. The model is trained on an expanded synthetic dataset with longer input lengths, 111 improving its generalization ability [Ucar et al., 2024]. Log-likelihoods are computed in **De Novo** 112 (**DN**) mode, where both sequence and structure are masked over the scoring region. We refer to this 113 mode as DiffAbXL-A-DN. For a designed sequence s , and positions j in the mutation mask \mathcal{M} , we 114 define the score:

$$LL = \sum_{j \in \mathcal{M}} \log (P_j(s_j | \mathcal{U}) + \varepsilon), \quad (1)$$

115 where $P_j(s_j | \mathcal{U})$ is the probability of residue s_j predicted by the model given unmasked context \mathcal{U} , 116 and ε is a small constant (e.g., 10^{-9}) for numerical stability.

117 **AntiFold.** AntiFold is an inverse folding model based on ESM-IF1 [Høie et al., 2023], which 118 autoregressively predicts sequence from structure. To compute log-likelihoods, we first pass either the 119 parental or mutant sequence along with its backbone structure into the model to obtain per-position 120 log-probabilities. We then gather the log-probabilities corresponding to the mutant amino acids at the 121 mutated positions \mathcal{M} . In the main analysis, we use the parental sequence for context (**AntiFold**_{PA}). 122 The log-likelihood is defined as:

$$LL_{PA} = \sum_{j \in \mathcal{M}} \log (P_j(s_j^{\text{mut}} | s_{<j}^{\text{pa}}, \mathcal{S}) + \varepsilon), \quad (2)$$

123 where \mathcal{S} is the structure, s_j^{mut} is the mutant residue at position j , and $s_{<j}^{\text{pa}}$ is the prefix of the parental 124 sequence.

125 **IgLM.** IgLM is a decoder-only language model trained with a masked span infilling objective on 126 antibody sequences [Shuai et al., 2023]. We evaluate log-likelihood using two scoring modes: 127 preceding context only ([pre]) and bidirectional context ([bi]). In both cases, we use the parental 128 sequence as context in the main experiments. Mutation regions may consist of one or more contiguous 129 spans (e.g., individual or multiple CDRs), each of which is scored independently in bidirectional 130 context case.

131 *Preceding context only ([pre]):*

$$\text{IgLM[pre]}_{PA}(s_{\text{mut}}; s_{\text{pa}}) = \sum_{t \in \mathcal{M}} \log p(s_t^{\text{mut}} | s_{<t}^{\text{pa}}), \quad (3)$$

132 where the parental sequence s_{pa} is passed to the model to obtain logits, and the mutant sequence s_{mut} 133 is used to select log-probabilities at mutation sites $t \in \mathcal{M}$.

134 *Bidirectional context ([bi]):*

$$\text{IgLM[bi]}_{\text{PA}}(\mathbf{s}^{\text{mut}}; \mathbf{s}_{\text{pa}}) = \sum_k \sum_{t \in S_k} \log p(s_t^{\text{mut}} | s_{<t}^{\text{pa}}, s_{>t}^{\text{pa}}), \quad (4)$$

135 where each mutated span S_k is a contiguous region (e.g., a CDR) that is masked in the parental se-
 136 quence, and the model predicts the mutant residues using bidirectional context. Multiple spans can be
 137 masked and evaluated independently if the mutation mask covers disjoint regions. Furthermore, since
 138 IgLM is trained to model masked spans using bidirectional context, scoring based on bidirectional
 139 information is better aligned with its training objective and has been empirically shown to produce
 140 lower perplexity than scoring based solely on preceding context [Shuai et al., 2023].

141 Results for scoring using the mutant sequence as context (IgLM[pre]_{MUT}, IgLM[bi]_{MUT}, and
 142 AntiFold_{MUT}), along with additional details on how the scores are computed, are provided in Sections
 143 D - G of the Appendix.

144 **BLOSUM similarity scoring.** BLOSUM similarity is computed using substitution matrices (e.g.,
 145 BLOSUM45, 62, 80, and 90). For each designed sequence, similarity is calculated with respect to
 146 one of the following reference sequences: (1) the parental sequence from the same dataset (denoted
 147 as BLOSUM_{PA}); (2) a global consensus sequence derived from two antibody datasets—either
 148 human antibodies from the Observed Antibody Space (OAS) or antibody-antigen complexes from
 149 SAbDab—denoted as BLOSUM_{GOAS} and BLOSUM_{GSAbDab}, respectively; or (3) a consensus sequence
 150 derived from the antibody variants within each dataset, referred to as dataset-specific (DS) and
 151 denoted as BLOSUM_{DS}. For (2), separate consensus sequences are constructed for heavy, kappa,
 152 and lambda chains using AHO numbering [Honegger and PluÈckthun, 2001]. The similarity score is
 153 defined as:

$$\text{ScoreBLOSUM} = \frac{1}{|\mathcal{M}|} \sum j \in \mathcal{M} B(s_j^{\text{ref}}, s_j), \quad (5)$$

154 where s_j^{ref} is the residue at position j in the reference sequence, s_j is the corresponding residue in the
 155 designed sequence, and $B(a, b)$ is the substitution score from the chosen BLOSUM matrix.

156 **PWM-based similarity scoring.** Position weight matrices (PWMs) are constructed from aligned
 157 human antibody sequences in the OAS and SAbDab databases. Sequences are split by chain type
 158 (heavy, kappa, and lambda) and aligned using AHO numbering. From these alignments, we compute
 159 amino acid frequency distributions at each position, producing normalized matrices in which each
 160 column sums to 1.

161 For each designed sequence, we compute the PWM score as the sum of amino acid frequencies at the
 162 mutated positions, matched by chain type:

$$\text{ScorePWM} = \sum j \in \mathcal{M}_H f_H(j, s_j) + \sum j \in \mathcal{M}_L f_L(j, s_j), \quad (6)$$

163 where $f_H(j, s_j)$ and $f_L(j, s_j)$ are the amino acid frequencies at position j in the heavy and light
 164 chain PWMs, respectively, and $\mathcal{M}_H, \mathcal{M}_L$ are the subsets of \mathcal{M} corresponding to the heavy and light
 165 chains.

166 **PSSM-based similarity scoring.** Position-specific scoring matrices (PSSMs) are computed from
 167 aligned antibody sequences, where each entry reflects the log-odds score of observing amino acid a
 168 at position j relative to background. We derive PSSMs from three sources: (1) human antibodies
 169 from the OAS repertoire, (2) the SAbDab database, and (3) the sequences within each experimental
 170 dataset, where a separate PSSM is constructed for each dataset using only its constituent sequences.
 171 We refer to this third approach as dataset-specific PSSMs (**PSSM_{DS}**). As with PWMs, sequences are
 172 aligned using AHO numbering and split by chain type.

173 For each designed sequence, we compute the PSSM score as:

$$\text{ScorePSSM} = \sum j \in \mathcal{M}_H S_H(j, s_j) + \sum j \in \mathcal{M}_L S_L(j, s_j), \quad (7)$$

174 where $S_H(j, s_j)$ and $S_L(j, s_j)$ denote the log-odds substitution scores from the PSSMs for the heavy
 175 and light chains, respectively. Higher scores indicate higher evolutionary preference for the observed
 176 amino acids at those positions. Additional details on the construction of PSSMs and PWMs can be
 177 found in Section C of the Appendix.

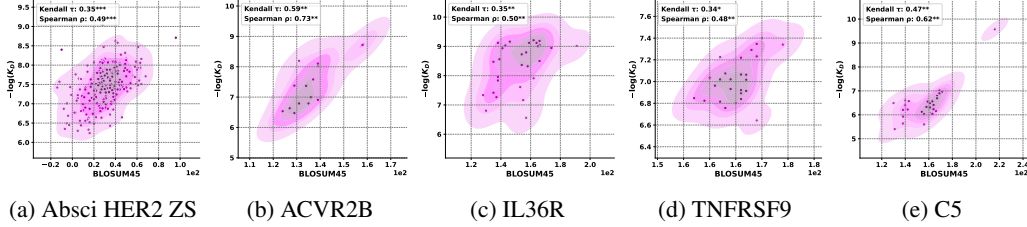


Figure 1: **Correlation between BLOSUM45 and $-\log K_D$:** a) HER2 Zero-Shot (ZS), b) ACVR2B, c) IL36R, d) TNFRSF9, e) C5. *, **, *** indicate p-values under 0.05, 0.01 and 1e-4 respectively.

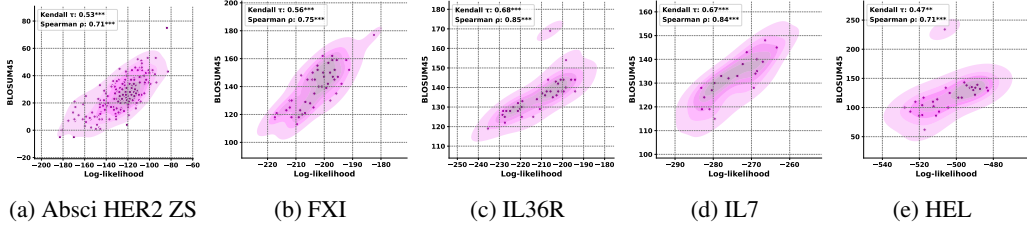


Figure 2: **Correlation between BLOSUM45 and the log-likelihood scores of DiffAbXL-A-DN:** a) HER2 Zero-Shot (ZS), b) FXI, c) IL36R, d) IL7, e) HEL. *, **, *** indicate p-values under 0.05, 0.01 and 1e-4 respectively.

178 **Hydrophobicity and rigidity scoring.** To evaluate coarse-grained biophysical trends, we compute
 179 the average hydrophobicity and rigidity of mutated residues using the Kyte-Doolittle [Kyte and
 180 Doolittle, 1982] and Karplus-Schulz [Karplus and Schulz, 1985] scales, respectively:

$$\text{ScoreHydro} = \frac{1}{|\mathcal{M}|} \sum_{j \in \mathcal{M}} H(s_j), \quad \text{ScoreRigid} = \frac{1}{|\mathcal{M}|} \sum_{j \in \mathcal{M}} \frac{1}{F(s_j)}, \quad (8)$$

181 where $H(s_j)$ and $F(s_j)$ denote the hydrophobicity and flexibility value of amino acid s_j .

182 **Affinity correlation analysis.** After computing all scores, we evaluate their correlation with exper-
 183 imental binding affinities, expressed as $-\log(K_D)$ or equivalent (e.g., $-\log(IC_{50})$). Spearman's
 184 rank correlation coefficient (ρ) and Kendall's tau (τ) are computed to assess the predictive power and
 185 biological relevance of each scoring method. We also examine correlations between statistical simi-
 186 larity scores and log-likelihoods to probe whether generative models implicitly capture evolutionary
 187 constraints.

188 4 Empirical Evaluation

189 4.1 Datasets

190 In this study, we use fourteen datasets drawn from four sources: Absci HER2 [Shanehsazzadeh et al.,
 191 2023b], IgDesign [Shanehsazzadeh et al., 2023a], Nature [Porebski et al., 2024], and proprietary
 192 datasets from AstraZeneca (AZ).

193 **Absci HER2.** These datasets involve HCDR re-designs of Trastuzumab, an antibody that targets
 194 HER2. Sequence generation was performed using a two-step pipeline: first, machine learning models
 195 were used to predict HCDR loop structures conditioned on the HER2 backbone (PDB:1N8Z, Chain
 196 C), the Trastuzumab framework, and the known epitope; second, sequences were generated via
 197 inverse folding on the predicted structures. While HCDR3 lengths ranged from 9 to 17 residues, we
 198 focus on sequences with HCDR3 length 13, consistent with the native antibody. HCDR1 and HCDR2
 199 were fixed at 8 residues. Binding affinities (K_D) were measured using a FACS-based ACE assay. We
 200 analyze two datasets: (1) a "zero-shot binders" set, and (2) an SPR-validated "control" set containing
 201 both binders and non-binders.

202 **IgDesign.** This set includes seven antigen targets—FXI, IL36R, C5, TSLP, IL17A, ACVR2B, and
 203 TNFRSF9. Antibodies were designed by mutating either the HCDR3 alone or all three CDRs on the
 204 heavy chain. Each design library was synthesized and experimentally tested using SPR.

Table 1: Spearman correlations of log-likelihood scores of DiffAbXL-A-DN, biophysical features, PWM, PSSM, and BLOSUM scores with binding affinity. *, **, *** indicate p-values under 0.05, 0.01, and 1e-4, respectively. The measurements are qAC_{50} for AZ Target-1, IC_{50} for IL7, and K_D for the rest.

Approach	Model	Absci HER2		Nature			AZ		Absci IgDesign						
		Zero Shot	Control	HEL	IL7	HER2	Target-1	Target-2	IL17A	ACVR2B	FXI	TNFRSF9	IL36R	C5	TSLP
Diffusion	DiffAbXL-A-DN	0.43***	0.22***	0.62**	-0.79***	0.37*	-0.11	0.41**	0.62**	0.54*	0.18	0.18	0.14	-0.32	-0.02
Biophysical	Hydrophobicity	-0.17*	0.04	-0.13	0.57*	-0.38	0.30	0.14	0.49	0.49	-0.39*	0.41*	0.70***	0.32	-0.05
	Rigidity	0.09	0.45***	-0.18	0.33	-0.37	0.18	0.06	0.49	0.49	-0.22	-0.16	0.15	0.19	-0.12
	PWM _{OAS}	0.30***	0.12*	0.29	0.34	0.17	-0.23	0.17	-0.24	0.13	-0.14	0.02	-0.16	-0.03	-0.15
	PWM _{Sabab}	0.30***	0.10*	0.35*	-0.49*	0.43*	0.24	0.16	-0.10	0.42	-0.23	0.29	-0.32*	-0.31	0.03
	PSSM _{OAS}	0.36***	0.22**	0.28	-0.07	0.20	-0.02	0.03	-0.13	0.11	-0.20	0.06	-0.29	-0.10	-0.103
	PSSM _{Sabab}	0.40***	0.29***	0.19	-0.40	0.31	0.38	0.10	-0.10	0.48*	-0.18	0.28	-0.39*	-0.32	-0.05
Statistical	PSSM _{DS}	0.50***	0.47***	0.52**	-0.18	-0.16	0.29	0.24*	-0.01	0.24	0.53**	0.05	0.26	0.00	-0.11
	BLOSUM45 _{GOAS}	0.36***	0.03	0.26	-0.59**	0.26	-0.52**	-0.07	0.10	0.03	0.29*	0.03	-0.27	-0.25	-0.36***
	BLOSUM45 _{GSabab}	0.39***	0.095	0.22	-0.80***	-0.34	-0.21	0.16	-0.39	-0.006	-0.42**	0.26	-0.27	-0.32	-0.34**
	BLOSUM45 _{DS}	0.46***	0.31***	0.52**	0.69***	-0.32	0.43*	0.30**	0.01	0.47*	0.52**	0.17	0.38*	0.10	-0.29**
	BLOSUM90 _{PA}	0.48***	0.25***	0.57**	-0.86***	-0.74***	-0.09	0.26*	0.77	0.63**	0.32	0.41*	0.57**	0.57**	0.19*
	BLOSUM180 _{PA}	0.48***	0.26***	0.58**	-0.85***	-0.73***	-0.10	0.24*	0.77	0.73**	0.32	0.42*	0.58**	0.64**	0.20*
	BLOSUM62 _{PA}	0.48***	0.26***	0.57**	-0.85***	-0.72***	-0.06	0.26*	0.77	0.71**	0.32	0.38*	0.58**	0.60**	0.21*
	BLOSUM45 _{PA}	0.50***	0.29***	0.59**	-0.87***	-0.71***	-0.08	0.30**	0.77	0.73**	0.33	0.48**	0.50**	0.62**	0.22*

Nature. We also include datasets reported by Porebski et al. [2024], covering HER2, IL7, and HEL. Mutations in anti-HER2 are limited to HCDR3, while anti-IL7 involves LCDR1 and LCDR3. The HEL dataset consists of nanobodies with mutations across all three CDRs. Dataset sizes range from 19 to 38 sequences. We use K_D values for HER2 and HEL, and IC_{50} for IL7. For structure-based methods, parental structures were predicted using ImmuneBuilder2 (HER2), IgFold (IL7), and NanoBodyBuilder2 (HEL) [Abanades et al., 2023, Ruffolo et al., 2023] as described in [Ucar et al., 2024].

AZ These proprietary datasets consist of two antibody libraries targeting separate antigens. The first target includes 24 variants, generated via rational design across four regions (HCDR1–3 and LCDR3). The second comprises 85 sequences drawn from three design strategies: two rationally designed libraries (one mutating heavy chain CDRs, the other light chain CDRs) and a third created using a machine learning model introducing changes across all six CDRs. Binding measurements are reported as qAC_{50} for Target-1 and K_D for Target-2. For models requiring structure, we use the corresponding crystal structures for both targets.

4.2 Results

Benchmarking predictive power across scoring methods. We assess the correlation between several scoring methods—including DiffAbXL-A log-likelihoods, BLOSUM similarity, PWM similarity, and PSSM similarity—and experimentally measured binding affinities across fourteen benchmark datasets. DiffAbXL-A was selected based on prior findings that its log-likelihood scores exhibit the strongest correlation with experimental binding affinity in [Ucar et al., 2024]. Spearman’s rank correlation coefficients (ρ) are summarized in Table 1, and correlation statistics between log-likelihoods and BLOSUM or statistical similarity scores are shown in Table 2. Log-likelihood scores derived from the DiffAbXL-A model in De Novo (DN) mode show consistent and often promising correlations with binding affinity across diverse design tasks. For example, we observe $\rho = 0.43$ on Absci HER2 Zero Shot, $\rho = 0.62$ on Nature HEL, and $\rho = 0.62$ on IgDesign IL17A. These results suggest that the model is capturing sequence features that are predictive of functional binding, even though it was not trained on the specific antibody libraries present in these datasets or on binding affinity prediction tasks. On the Nature IL7 dataset, where inhibition rather than binding affinity is measured (via IC_{50}), a strong negative correlation is observed ($\rho = -0.79$). However, IC_{50} reflects the concentration needed to achieve 50% inhibition of a biological response, and is influenced by multiple factors beyond binding—such as receptor expression levels, signaling kinetics, and assay-specific artifacts. Unlike K_D , which directly measures molecular interaction strength, IC_{50} integrates downstream effects and may vary substantially even when two molecules have similar affinities. As a result, correlations involving IC_{50} should be interpreted cautiously.

BLOSUM similarity scores align strongly with binding affinity. Across the board, BLOSUM-based similarity—particularly using BLOSUM45—shows good correlation with binding affinity. This trend holds across nearly all datasets, including Absci HER2, Nature HEL, and several IgDesign targets. For instance, BLOSUM45 achieves $\rho = 0.50$ on Absci HER2 Zero Shot, $\rho = 0.59$ on Nature HEL, and $\rho = 0.73$ on IgDesign ACVR2B. This finding supports the idea that evolutionary closeness to the parental (reference) sequence is a strong indicator of retained binding functionality. The consistent performance of BLOSUM matrices, especially those calibrated for more distant

Table 2: Spearman correlations between DiffAbXL-A-DN log-likelihoods and sequence-similarity scores (BLOSUM45, PWM, and PSSM). *, **, *** indicate p-values under 0.05, 0.01, and 1e-4, respectively.

Method	Absci HER2		Nature			AZ		Absci IgDesign						
	ZS	Ctrl	HEL	IL7	HER2	T-1	T-2	IL17A	ACVR2B	FXI	TNFRSF9	IL36R	C5	TSLP
BLOSUM45 _{PA}	0.71***	0.81***	0.71***	0.85***	-0.24	0.48*	0.56***	0.88***	0.64**	0.75***	0.41*	0.85***	-0.61**	0.46***
PSSM _{DS}	0.80***	0.73***	0.79***	0.07	-0.12	-0.47*	0.23*	-0.25	0.49*	0.38**	-0.01	0.63***	-0.73***	0.27**
PSSM _{SAbDab}	0.77***	0.55***	0.45**	0.46*	0.38	-0.05	0.70***	0.09	0.24	-0.61***	-0.18	-0.69***	0.36*	0.27**
PSSM _{OAS}	0.76***	0.51***	0.57**	0.12	0.40*	-0.16	0.47***	-0.11	0.13	-0.63***	0.02	-0.69***	-0.41*	0.38***
PWM _{SAbDab}	0.80***	0.52***	0.55**	0.58**	0.25	0.02	0.56***	0.19	0.68**	-0.17	-0.12	-0.71***	0.48**	0.55***
PWM _{OAS}	0.77***	0.48***	0.56**	-0.18	0.12	0.48*	0.31**	0.20	0.41	-0.04	0.10	-0.72***	-0.48**	0.33**

homologs (e.g., BLOSUM45), suggests they capture robust patterns relevant to antigen recognition and molecular stability. We note that across all datasets examined, the parental antibody—used as the reference for computing BLOSUM scores and thus assigned the highest similarity score by design—is consistently among the strongest binders. We expected that if this were not true, the correlation between BLOSUM similarity and binding affinity would be substantially weaker.

Consensus-based BLOSUM scores lose predictive power. In contrast, when BLOSUM similarity is computed between designed sequences and global consensus sequences—either from OAS or SAbDab—the correlation with binding affinity largely disappears. For example, BLOSUM45_{OAS} shows weak or negative correlations on many datasets, with significant drops observed on most datasets. This suggests that global evolutionary priors do not substitute well for dataset-specific reference sequences in affinity prediction.

PWM similarity shows modest utility, with SAbDab outperforming OAS. PWM scores based on the OAS repertoire show weak and inconsistent correlation with binding affinity, with meaningful results limited to the Absci HER2 datasets (e.g., $\rho = 0.30$ on Zero Shot). However, when PWMs are computed from the SAbDab dataset, the correlation improves slightly. In addition to Absci HER2, we observe positive correlations on Nature HEL ($\rho = 0.35$) and Nature HER2 ($\rho = 0.43$), suggesting that structure-based databases such as SAbDab may better reflect the selective pressures acting on antibody binding regions since they contain antibody-antigen complexes. Nonetheless, performance remains below that of BLOSUM and DiffAbXL-A. This may be due to the local and repertoire-specific nature of the PWM used, which reflects background amino acid usage rather than target-specific substitution effects. Because PWMs are constructed from observed frequencies rather than substitution dynamics, they may fail to capture functionally relevant mutations, especially when applied outside their source distribution.

PSSM similarity improves with data-specific priors. PSSM scores show variable performance depending on how the matrices are derived. Global PSSMs constructed from the OAS or SAbDab datasets yield modest correlations with binding affinity, with improvements observed for SAbDab-based PSSMs in some datasets (e.g., Absci HER2 and IgDesign ACVR2B). However, when PSSMs are constructed specifically from the dataset under evaluation (PSSM_{DS}), correlation improves substantially. For example, PSSM_{DS} yields $\rho = 0.50$ on Absci HER2 Zero Shot and $\rho = 0.52$ on Nature HEL, both surpassing the performance of global PSSMs. These results highlight the importance of local sequence context in capturing meaningful constraints for affinity prediction and suggest that dataset-specific PSSMs can serve as useful tools when sufficient in-distribution sequence data are available.

Log-likelihood scores correlate with BLOSUM, PSSM, and PWM similarity metrics. We observe strong correlations between DiffAbXL-A log-likelihood scores and BLOSUM45 similarity across most datasets (Table 2), with 12 out of 14 cases exceeding a Spearman ρ of 0.4, and several surpassing 0.8 (e.g., IL17A, IL36R, IL7). This suggests that the generative model is implicitly learning amino acid substitution patterns that align closely with established evolutionary priors. Similar, though generally weaker, correlations are observed with PWM-based scores. Notably, PWMs derived from SAbDab show stronger alignment with log-likelihoods than those from OAS, especially in datasets such as IL7 ($\rho = 0.58$) and ACVR2B ($\rho = 0.68$). PSSM similarity also correlates with model scores, particularly when computed from dataset-specific (PSSM_{DS}) or SAbDab-based matrices, supporting the view that the model internalizes substitution preferences.

Biophysical scores show target-specific utility. Two coarse-grained biophysical metrics—Kyte-Doolittle hydrophobicity and Karplus-Schulz rigidity—were also evaluated. Hydrophobicity content correlates positively with binding affinity on Nature IL7 ($\rho = 0.57$) and IgDesign IL36R ($\rho = 0.70$), suggesting a potential link between hydrophobic residues and binding affinity in these datasets.

Table 3: Spearman correlations between the model log-likelihoods (LLs) and K_D values as well as between LLs and BLOSUM45_{PA} scores for DiffAbXL-A-DN, AntiFold and IgLM. *, **, *** indicate p-values under 0.05, 0.01, and 1e-4, respectively.

Model	Metric	Absci HER2		Nature		Absci IgDesign						
		ZS	Ctrl	HEL	HER2	IL17A	ACVR2B	FXI	TNFRSF9	IL36R	C5	TSLP
DiffAbXL-A-DN	LL vs K_D	0.43***	0.22***	0.62**	0.37*	0.62**	0.54*	0.18	0.18	0.14	-0.32	-0.02
	LL vs BLOSUM45 _{PA}	0.71***	0.81***	0.71***	-0.24	0.88***	0.64**	0.75***	0.41*	0.85***	-0.61**	0.46***
AntiFold _{PA}	LL vs K_D	0.45***	0.15**	0.64***	-0.38	-0.79**	0.29	0.02	0.04	0.18	-0.01	0.24**
	LL vs BLOSUM45 _{PA}	0.61***	0.61***	0.76***	0.60**	-0.39	-0.04	0.14	0.65***	-0.22	0.01	0.80***
IgLM <pre>PA</pre>	LL vs K_D	-0.09	0.22***	0.35*	0.02	0.12	0.16	0.31*	0.12	-0.63***	-0.26	0.23*
	LL vs BLOSUM45 _{PA}	-0.22**	-0.07	0.57**	-0.33	0.16	-0.36	0.60***	0.11	-0.43**	-0.82***	0.51***
IgLM[bi] _{PA}	LL vs K_D	0.26**	-0.05	0.59**	-0.12	0.14	-0.01	-0.11	-0.09	0.24	-0.21	-0.08
	LL vs BLOSUM45 _{PA}	0.35***	0.14**	0.73***	0.43*	0.56*	0.48*	-0.58***	0.38*	-0.04	-0.59**	0.36***

Rigidity shows a moderate correlation only on Absci HER2 Control ($\rho = 0.45$), with little signal elsewhere. These results indicate that such physicochemical scores may capture some target-specific trends but are not general-purpose predictors of affinity.

Consistent evolutionary signatures in log-likelihoods from diverse generative models. As with DiffAbXL-A, log-likelihoods from both AntiFold and IgLM show strong correlations with BLOSUM45 similarity scores (Table 3). For example, AntiFold_{PA} log-likelihoods correlate with BLOSUM similarity at $\rho = 0.76$ on Nature HEL and $\rho = 0.80$ on IgDesign TSLP. Similarly, IgLM[bi]_{PA} achieves $\rho = 0.73$ on Nature HEL and shows moderate-to-strong alignment on multiple other datasets. These findings further support the hypothesis that generative models implicitly learn substitution preferences that align with classical evolutionary priors. We also note that when log-likelihoods (LLs) correlate with binding affinity (K_D), they almost always exhibit an even stronger correlation with BLOSUM similarity. However, the converse does not necessarily hold: a strong LL-BLOSUM correlation does not imply a meaningful LL- K_D correlation.

5 Conclusion

We examined the connection between generative model scores and classical evolutionary similarity metrics in the context of antibody design. Across fourteen datasets of experimentally characterized antibody variants, BLOSUM similarity to the wild type (WT)—particularly using matrices such as BLOSUM45—showed strong and consistent correlation with binding affinity, often rivaling or exceeding the performance of model-based log-likelihood scores. Among generative models, the diffusion-based model DiffAbXL-A showed the highest correlation with affinity values, consistent with prior findings in [Ucar et al., 2024]. Moreover, evaluations of AntiFold and IgLM revealed that their log-likelihood scores also align with BLOSUM similarities, even when their correlation with binding affinity is weaker. This reinforces the idea that generative models may implicitly learn substitution patterns shaped by evolutionary pressure, even without explicit supervision.

In contrast, PWM-based scores exhibited limited and variable performance, particularly when derived from general antibody repertoires. Slight improvements were observed with PWMs constructed from antibodies in complex with antigens, though these still underperformed relative to BLOSUM similarities over the parental antibody sequence in each dataset. PSSM-based scores showed somewhat stronger and more stable correlations than PWMs, especially when constructed from dataset-specific alignments, occasionally matching the predictive power of BLOSUM and log-likelihood scores. However, their performance was less consistent when based on broad repertoires such as OAS or structure databases such as SAbDab. Simple biophysical descriptors such as hydrophobicity and rigidity captured some signal in a few cases but lacked generalizability.

These findings highlight the utility of interpretable, evolution-derived metrics such as BLOSUM and underscore the importance of contextual information—such as reference sequence choice—in score interpretation. They also suggest that generative models encode evolutionary signals that can be leveraged for scoring and prioritization tasks. As generative models continue to improve, understanding the extent to which their internal representations align with biological priors will be essential for advancing robust and interpretable generative models for designing therapeutic antibodies.

References

- Brennan Abanades, Wing Ki Wong, Fergus Boyles, Guy Georges, Alexander Bujotzek, and Charlotte M Deane. Immunebuilder: Deep-learning models for predicting the structures of immune proteins. *Communications Biology*, 6(1):575, 2023.
- Josh Abramson, Jonas Adler, Jack Dunger, Richard Evans, Tim Green, Alexander Pritzel, Olaf Ronneberger, Lindsay Willmore, Andrew J Ballard, Joshua Bambrick, et al. Accurate structure prediction of biomolecular interactions with alphafold 3. *Nature*, pages 1–3, 2024.
- Stephen F Altschul, Warren Gish, Webb Miller, Eugene W Myers, and David J Lipman. Basic local alignment search tool. *Journal of molecular biology*, 215(3):403–410, 1990.
- Minkyung Baek, Frank DiMaio, Ivan Anishchenko, Justas Dauparas, Sergey Ovchinnikov, Gyu Rie Lee, Jue Wang, Qian Cong, Lisa N Kinch, R Dustin Schaeffer, et al. Accurate prediction of protein structures and interactions using a three-track neural network. *Science*, 373(6557):871–876, 2021.
- Michael Chungyoun, Jeffrey Ruffolo, and Jeffrey Gray. Flab: Benchmarking deep learning methods for antibody fitness prediction. *BioRxiv*, pages 2024–01, 2024.
- Michael Gribskov, Andrew D McLachlan, and David Eisenberg. Profile analysis: detection of distantly related proteins. *Proceedings of the National Academy of Sciences*, 84(13):4355–4358, 1987.
- Steven Henikoff and Jorja G Henikoff. Amino acid substitution matrices from protein blocks. *Proceedings of the National Academy of Sciences*, 89(22):10915–10919, 1992.
- Magnus Høie, Alissa Hummer, Tobias Olsen, Morten Nielsen, and Charlotte Deane. Antifold: Improved antibody structure design using inverse folding. In *NeurIPS 2023 Generative AI and Biology (GenBio) Workshop*, 2023.
- Annemarie Honegger and Andreas Plückthun. Yet another numbering scheme for immunoglobulin variable domains: an automatic modeling and analysis tool. *Journal of molecular biology*, 309(3): 657–670, 2001.
- John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, et al. Highly accurate protein structure prediction with AlphaFold. *nature*, 596(7873):583–589, 2021.
- Paul Andrew Karplus and GE Schulz. Prediction of chain flexibility in proteins: a tool for the selection of peptide antigens. *Naturwissenschaften*, 72(4):212–213, 1985.
- Xiangzhe Kong, Wenbing Huang, and Yang Liu. End-to-end full-atom antibody design. *arXiv preprint arXiv:2302.00203*, 2023.
- Jack Kyte and Russell F Doolittle. A simple method for displaying the hydropathic character of a protein. *Journal of molecular biology*, 157(1):105–132, 1982.
- Shitong Luo, Yufeng Su, Xingang Peng, Sheng Wang, Jian Peng, and Jianzhu Ma. Antigen-specific antibody design and optimization with diffusion-based generative models for protein structures. *Advances in Neural Information Processing Systems*, 35:9754–9767, 2022.
- Cedric Malherbe and Talip Uçar. Igbblend: Unifying 3d structures and sequences in antibody language models. *bioRxiv*, pages 2024–10, 2024.
- Benjamin T Porebski, Matthew Balmforth, Gareth Browne, Aidan Riley, Kiarash Jamali, Maximilian JLJ Fürst, Mirko Velic, Andrew Buchanan, Ralph Minter, Tristan Vaughan, et al. Rapid discovery of high-affinity antibodies via massively parallel sequencing, ribosome display and affinity screening. *Nature biomedical engineering*, 8(3):214–232, 2024.
- Jeffrey A Ruffolo, Lee-Shin Chu, Sai Pooja Mahajan, and Jeffrey J Gray. Fast, accurate antibody structure prediction from deep learning on massive set of natural antibodies. *Nature communications*, 14(1):2389, 2023.

379 Amir Shanehsazzadeh, Julian Alverio, George Kasun, Simon Levine, Ido Calman, Jibran A Khan,
380 Chelsea Chung, Nicolas Diaz, Breanna K Luton, Ysis Tarter, Cailen McCloskey, Katherine B
381 Bateman, Hayley Carter, Dalton Chapman, Rebecca Consbruck, Alec Jaeger, Christa Kohnert,
382 Gaelin Kopec-Belliveau, John M Sutton, Zheyuan Guo, Gustavo Canales, Kai Ejan, Emily Marsh,
383 Alyssa Ruelos, Rylee Ripley, Brooke Stoddard, Rodante Caguiat, Kyra Chapman, Matthew
384 Saunders, Jared Sharp, Douglas Ganini da Silva, Audree Feltner, Jake Ripley, Megan E Bryant,
385 Danni Castillo, Joshua Meier, Christian M Stegmann, Katherine Moran, Christine Lemke, Shaheed
386 Abdulhaqq, Lillian R Klug, and Sharrol Bachas. Igdesign: In vitro validated antibody design
387 against multiple therapeutic antigens using inverse folding. December 2023a. doi: 10.1101/
388 2023.12.08.570889. URL <http://dx.doi.org/10.1101/2023.12.08.570889>.

389 Amir Shanehsazzadeh, Sharrol Bachas, Matt McPartlon, George Kasun, John M Sutton, Andrea K
390 Steiger, Richard Shuai, Christa Kohnert, Goran Rakocovic, Jahir M Gutierrez, et al. Unlocking de
391 novo antibody design with generative artificial intelligence. *bioRxiv*, pages 2023–01, 2023b.

392 Richard W Shuai, Jeffrey A Ruffolo, and Jeffrey J Gray. Iglm: Infilling language modeling for
393 antibody sequence design. *Cell Systems*, 14(11):979–989, 2023.

394 Gary D Stormo, Thomas D Schneider, Larry Gold, and Andrzej Ehrenfeucht. Use of the ‘percep-
395 tron’ algorithm to distinguish translational initiation sites in e. coli. *Nucleic acids research*, 10(9):
396 2997–3011, 1982.

397 Timothy Truong Jr and Tristan Bepler. Poet: A generative model of protein families as sequences-of-
398 sequences. *Advances in Neural Information Processing Systems*, 36:77379–77415, 2023.

399 Talip Ucar, Cedric Malherbe, and Ferran Gonzalez. Exploring log-likelihood scores for ranking
400 antibody sequence designs. In *NeurIPS 2024 Workshop on AI for New Drug Modalities*, 2024.

401 Appendix

402 A Impact Statement

403 This research advances the use of machine learning in antibody engineering by investigating the
404 predictive power of statistical models such as BLOSUM and PWM for ranking antibody designs
405 based on binding affinity as well as studying their relationship with the log-likelihood scores of
406 diffusion-based generative models. By demonstrating an association between derived scores and
407 real-world experimental data, this work provides a pathway to accelerate therapeutic antibody
408 discovery while minimizing costly trial-and-error experimentation. Ethical and societal impacts,
409 such as improved healthcare outcomes and broader accessibility of life-saving treatments, mirror the
410 established considerations in the broader field of machine learning-driven drug discovery.

411 B License Information

412 IgDesign datasets [Shanehsazzadeh et al., 2023a] are released under MIT license. Absci Her2 datasets
413 [Shanehsazzadeh et al., 2023b] are released under BSD License. SAbDab and OAS datasets are
414 available under a CC-BY 4.0 license. We will release our code upon the acceptance of our paper with
415 Apache 2.0 license.

416 C Construction of PWM and PSSM

417 To analyze amino acid preferences at each structurally equivalent position, we first aligned all
418 sequences using the AHO numbering scheme, which provides a consistent positional framework
419 across antibody variable domains. Each sequence was mapped into an AHO-labeled position-residue
420 dictionary, and alignment matrices were constructed accordingly.

421 We computed the Position Weight Matrix (PWM) by counting the occurrences of each amino
422 acid (including gaps) at each AHO-defined position across all aligned sequences. To prevent zero
423 probabilities and ensure numerical stability, a Laplace pseudocount of 1 was added to each residue
424 count. The resulting frequency $f_{i,a}$ of amino acid a at position i was computed as:

$$f_{i,a} = \frac{n_{i,a} + 1}{n_i + k}$$

425 where $n_{i,a}$ is the count of amino acid a at position i , n_i is the total number of observations at that
426 position (including gaps), and $k = 21$ is the number of possible residue types (20 amino acids plus
427 the gap character). The PWM thus represents a normalized probability distribution over residues at
428 each position.

429 To construct the Position-Specific Scoring Matrix (PSSM), we first estimated background frequencies
430 q_a for each amino acid a across the entire aligned dataset, again applying Laplace pseudocounts:

$$q_a = \frac{n_a + 1}{N + k}$$

431 where n_a is the total count of amino acid a in the full alignment, N is the total number of observed
432 residues across all positions and sequences (including gaps), and $k = 21$ as before.

433 The log-odds score $S_{i,a}$ for each residue a at position i was then calculated as:

$$S_{i,a} = \log_2 \left(\frac{f_{i,a}}{q_a} \right)$$

434 This score reflects how much more (or less) likely a residue is to appear at a specific position compared
435 to its global background expectation. The final PSSM was stored as a position-by-residue matrix of
436 log-odds scores. To summarize the most likely residue at each position, a consensus sequence was
437 derived by selecting the amino acid with the highest frequency in the PWM.

438 D Additional Results

Table 4: Spearman correlations between the model log-likelihoods (LLs) and K_D values as well as between LLs and BLOSUM45_{PA} scores for DiffAbXL-A-DN, AntiFold and IgLM. *, **, *** indicate p-values under 0.05, 0.01, and 1e-4, respectively.

Model	Metric	Absci HER2		Nature		Absci IgDesign						
		ZS	Ctrl	HEL	HER2	IL17A	ACVR2B	FXI	TNFRSF9	IL36R	C5	TSLP
DiffAbXL-A-DN	LL vs K_D	0.43***	0.22***	0.62**	0.37*	0.62**	0.54*	0.18	0.18	0.14	-0.32	-0.02
	LL vs BLOSUM45 _{PA}	0.71***	0.81***	0.71***	-0.24	0.88***	0.64**	0.75***	0.41*	0.85***	-0.61**	0.46***
AntiFold _{PA}	LL vs K_D	0.45***	0.15**	0.64***	-0.38	-0.79**	0.29	0.02	0.04	0.18	-0.01	0.24**
	LL vs BLOSUM45 _{PA}	0.61***	0.61***	0.76***	0.60**	-0.39	-0.04	0.14	0.65***	-0.22	0.01	0.80***
AntiFold _{MUT}	LL vs K_D	-0.31***	-0.31***	0.33	0.05	0.17	0.34	-0.13	-0.34*	-0.05	0.39*	0.29**
	LL vs BLOSUM45 _{PA}	-0.29***	-0.27***	0.42*	0.18	0.12	0.25	0.28	0.19	-0.02	0.39*	0.00
IgLM[pre] _{PA}	LL vs K_D	-0.09	0.22***	0.35*	0.02	0.12	0.16	0.31*	0.12	-0.63***	-0.26	0.23*
	LL vs BLOSUM45 _{PA}	-0.22**	-0.07	0.57**	-0.33	0.16	-0.36	0.60***	0.11	-0.43**	-0.82***	0.51***
IgLM[pre] _{MUT}	LL vs K_D	-0.33***	-0.43***	0.05	-0.20	-0.46	-0.20	-0.00	-0.30	0.47**	0.03	-0.06
	LL vs BLOSUM45 _{PA}	-0.17*	-0.17**	0.22	0.36	-0.64*	-0.06	-0.27	-0.36*	0.74***	0.56**	-0.73***
IgLM[bi] _{PA}	LL vs K_D	0.26**	-0.05	0.59**	-0.12	0.14	-0.01	-0.11	-0.09	0.24	-0.21	-0.08
	LL vs BLOSUM45 _{PA}	0.35***	0.14**	0.73***	0.43*	0.56*	0.48*	-0.58***	0.38*	-0.04	-0.59**	0.36***
IgLM[bi] _{MUT}	LL vs K_D	0.25**	-0.08	-0.53***	-0.12	-0.13	-0.10	-0.15	-0.24	-0.23	-0.29	0.19*
	LL vs BLOSUM45 _{PA}	0.39***	0.25***	-0.66***	0.55**	0.48	0.56**	-0.54***	0.57**	-0.37*	-0.54**	0.53***

439 E Log-Likelihood Scoring

440 E.1 Setup and Notation

441 Let an antibody sequence be denoted by

$$\mathbf{s} = (s_1, s_2, \dots, s_T),$$

442 where each s_t represents one of the 20 canonical amino acids. A left-to-right autoregressive language
443 model (e.g., GPT-2) defines the sequence probability via the factorization

$$p(\mathbf{s}) = \prod_{t=1}^T p(s_t | s_{<t}), \quad s_{<t} = (s_1, \dots, s_{t-1}).$$

444 Let $\ell_t(\cdot) \in \mathbb{R}^{20}$ denote the unnormalized logits output by the model at position t . The corresponding
445 conditional log-probability is computed as

$$\log p(s_t | s_{<t}) = \log[\text{softmax}(\ell_{t-1})]_{s_t}, \quad t \geq 2.$$

446 For $t = 1$, either a special start-of-sequence token or a uniform prior may be used.

447 E.2 Full Sequence Log-Likelihood

448 The total log-likelihood of a sequence \mathbf{s} is given by

$$\log p(\mathbf{s}) = \sum_{t=1}^T \log p(s_t | s_{<t}).$$

449 This quantity can be computed from a single forward pass through the model. The following
450 pseudocode illustrates the computation using transformers-style APIs:

```

451 # Tokenize and run one forward pass
452 input_ids = tokenizer.encode(seq)           # shape [1, T]
453 logits    = model(input_ids).logits         # shape [1, T, 20]
454 log_probs = softmax(logits, dim=-1).log()
455
456 # Shift so that log_probs[:, t-1, *] = log p(s_t | s_{<t})
457 shifted = log_probs[:, :-1, :]
458 labels  = input_ids[:, 1:]
459
460 # Gather and sum
461 token_ll = shifted.gather(-1, labels.unsqueeze(-1)).squeeze(-1)
462 total_ll = token_ll.sum()
463

```

465 E.3 Single Contiguous Masked Region

466 To compute the log-likelihood of a contiguous span (s_a, \dots, s_b) conditioned on its left context, we
 467 can use

$$\log p(s_a, \dots, s_b \mid s_{<a}) = \sum_{t=a}^b \log p(s_t \mid s_{<t}), \quad 1 \leq a \leq b \leq T.$$

468 Given the vector of log-probabilities computed above, the relevant entries are summed as follows:

```
469 # token_ll[i] = log p(s_{i+1} | s_{<=i})
470 span_ll = token_ll[(a - 1) : b].sum()
472
```

473 E.4 Multiple Disjoint Masked Regions

474 Consider a collection of K disjoint spans $\{[a_k, b_k]\}_{k=1}^K$. The total log-likelihood over these regions is

$$\sum_{k=1}^K \sum_{t=a_k}^{b_k} \log p(s_t \mid s_{<t}).$$

475 This is equivalent to summing over the union of all token positions in the selected spans:

$$\mathcal{M} = \bigcup_{k=1}^K \{a_k, \dots, b_k\}.$$

```
476 positions = []
477 for (a, b) in spans:
478     positions.extend(range(a - 1, b))
479 multi_ll = token_ll[positions].sum()
480
```

482 E.5 Context Choice for Library Scoring: LL_{MUT} vs. LL_{PA}

483 For a designed antibody library derived from a common parental sequence, the log-likelihood of each
 484 designed sequence can be computed using one of two distinct approaches:

485 **Mutation-context likelihood (LL_{MUT}).** In this setting, the model is conditioned directly on each
 486 designed (mutant) sequence to compute its own log-likelihood:

$$LL_{\text{MUT}}(\mathbf{s}_{\text{mut}}) = \sum_{t \in P} \log p(s_t^{\text{mut}} \mid s_{<t}^{\text{mut}}),$$

487 where P denotes the set of mutated positions. This approach reflects the model's confidence in the
 488 mutant sequence given the full autoregressive context of the design.

489 **Parent-context likelihood (LL_{PA}).** Alternatively, model logits may be obtained from the original
 490 parental sequence \mathbf{s}_{pa} , and the log-likelihood is then evaluated using the designed sequence \mathbf{s}_{mut} by
 491 gathering log-probabilities only at mutated positions:

$$LL_{\text{PA}}(\mathbf{s}_{\text{mut}}; \mathbf{s}_{\text{pa}}) = \sum_{t \in P} \log p(s_t^{\text{mut}} \mid s_{<t}^{\text{pa}}).$$

492 In practice, the parental sequence is passed to the model to obtain the sequence of conditional
 493 distributions, and the designed sequence is used only to select which token probabilities to score at
 494 positions $t \in P$.

495 **Implementation notes.** Let `input_pa` denote the tokenized parental sequence, and `input_mut`
 496 denote the mutant sequence. Then:

```

497 # For LL_MUT
498 logits_mut = model(input_mut).logits
499 log_probs_mut = softmax(logits_mut, dim=-1).log()
500 # Use positions P on token_ll_mut to compute LL_MUT
501
502 # For LL_PA
503 logits_pa = model(input_pa).logits
504 log_probs_pa = softmax(logits_pa, dim=-1).log()
505 # Gather log_probs_pa at positions t in P,
506 # using tokens from input_mut for indexing
507

```

509 **Use cases.** LL_{MUT} reflects how likely a full designed sequence is under the model’s distribution, incorporating all mutated residues and their autoregressive influence. In contrast, LL_{PA} measures how well the mutated residues are supported by the local sequence context inherited from the parent, isolating the evaluation to mutation sites without considering their downstream impact.

513 Summary.

- 514 • LL_{MUT} : evaluates mutant sequence under its own autoregressive context.
- 515 • LL_{PA} : evaluates mutant residues in the fixed parental context.

516 F IgLM Log-Likelihood Scoring

517 F.1 Preceding vs. Bidirectional Context Scoring

518 IgLM is a decoder-only Transformer model trained with an infilling objective designed for antibody sequence modeling [Shuai et al., 2023]. Instead of standard left-to-right language modeling, IgLM is trained to reconstruct masked spans within a sequence using both left and right flanking context. During training, a contiguous span of amino acids is removed and replaced with a special [MASK] token. The remaining prefix and suffix of the sequence are concatenated, separated by a [SEP] token, and the removed span is appended after this context, followed by an [ANS] token to indicate the end of the span. This reordered sequence is used as input to the model, which is then trained to autoregressively predict the span tokens (and the [ANS] terminator), conditioned on the entire flanking context.

527 Formally, for a sequence $s = (s_1, \dots, s_T)$ with a masked span $S = (s_s, \dots, s_e)$, IgLM constructs an input of the form:

[CHAIN] [SPECIES] s_1, \dots, s_{s-1} , [MASK], s_{e+1}, \dots, s_T , [SEP], s_s, \dots, s_e , [ANS]

529 where [CHAIN] and [SPECIES] are fixed metadata tokens indicating chain type and species.

530 During evaluation, IgLM supports two log-likelihood computation strategies:

- 531 • Preceding context (autoregressive) scoring ([pre]): left-to-right log-likelihoods are computed over the full sequence.
- 533 • Bidirectional context (infilling-based) scoring ([bi]): log-likelihoods are computed using bidirectional context, consistent with the training setup.

535 Both strategies can be used with either the mutant or parental sequence as input context. For example, IgLM[pre]_{MUT} uses the mutant sequence for autoregressive scoring, while IgLM[bi]_{MUT} uses the mutant sequence for bidirectional infilling.

538 **Preceding-Context Scoring.** In the [pre] setting, the full sequence (mutant or parental) is passed to the model, and token-level log-likelihoods are computed left-to-right. For a given context sequence s^{CTX} , and mutant target s^{mut} , the log-likelihood is:

$$LL_{pre}(s^{mut}; s^{CTX}) = \sum_{t \in \mathcal{M}} \log p(s_t^{mut} | s_{<t}^{CTX}),$$

541 where \mathcal{M} is the set of mutated positions, and $CTX \in \{MUT, PA\}$ indicates the context source.

Bidirectional (Infilling-Based) Scoring. In the [bi] setting, IgLM uses its span-infilling mechanism to evaluate masked spans with bidirectional context. For each mutated span $S_k = (s_{a_k}, \dots, s_{b_k})$, the span is masked in the context sequence, and the corresponding mutant residues are passed as the target:

$$LL_{bi}(s^{\text{mut}}; s^{\text{CTX}}) = \sum_k \sum_{t=a_k}^{b_k} \log p(s_t^{\text{mut}} | s_{<t}^{\text{CTX}}, s_{>t}^{\text{CTX}}).$$

As before, CTX indicates whether the context is the parental or mutant sequence.

Implementation Considerations. The IgLM implementation supports both evaluation protocols. For autoregressive scoring, the model processes the entire context sequence once and gathers log-probabilities at mutation sites. For infilling-based scoring, the model must be called separately for each mutated span: the corresponding region is masked in the context sequence, and the mutant tokens are appended as the infill segment. The model then computes log-probabilities for these tokens after the [SEP] marker, consistent with its training objective.

As IgLM was trained to model masked spans using bidirectional context, the infilling-based scoring ([bi]) is more aligned with the model’s architecture and is empirically reported to yield lower perplexity than autoregressive scoring. However, both modes are supported and yield useful comparisons when paired with either mutant or parental input.

Summary.

- LL_{pre} : computes log-likelihood using the full sequence as input, relying only on autoregressive left context.
- LL_{bi} : computes log-likelihood of the designed regions using IgLM’s infilling mechanism, conditioned on bidirectional context from the sequence.

G AntiFold Log-Likelihood Scoring:

AntiFold is an antibody-specific inverse folding model based on the ESM-IF1 architecture, trained to predict sequences given fixed backbone structures [Høie et al., 2023]. Given a structure input, AntiFold generates amino acid sequences autoregressively from N- to C-terminus using a decoder-only Transformer architecture with *causal attention*. This means each position attends only to its preceding sequence positions and not to future residues. However, AntiFold conditions globally on the full backbone structure, which is processed separately and fed as a contextual embedding at each decoding step.

During evaluation, AntiFold outputs the log-probability assigned to each of the 20 amino acids at each sequence position. These per-position log-probabilities can be used to compute the log-likelihood of any specified subset of residues, including disjoint masked regions.

Mutation-context likelihood (LL_{MUT}). In the first evaluation strategy, the model is run using the full mutant sequence as input, along with the associated structure¹. Per-position log-probabilities are extracted from the model output, and the log-likelihood of the mutant residues is computed over the specified masked positions:

$$LL_{\text{MUT}}(s_{\text{mut}}) = \sum_{t \in P} \log p(s_t^{\text{mut}} | s_{<t}^{\text{mut}}, \text{structure}),$$

where P denotes the indices of mutated residues. Because the input sequence is the mutant itself, the decoder conditions on the correct mutated left context for each position in P .

Parent-context likelihood (LL_{PA}). Alternatively, log-probabilities can be computed using the parental sequence as input. In this case, the model is conditioned on the original (pre-mutation) left context, and the mutant amino acids are scored using the model’s output:

$$LL_{\text{PA}}(s_{\text{mut}}; s_{\text{pa}}) = \sum_{t \in P} \log p(s_t^{\text{mut}} | s_{<t}^{\text{pa}}, \text{structure}).$$

¹Backbone structure of mutant sequence is assumed to be same as the parental sequence

582 This formulation assesses how well the mutant residues fit into the structural and sequential context
583 defined by the parent. However, since the decoder is causal, any mutations occurring at early positions
584 can affect the correctness of conditioning for later positions if the full mutant context is not used.

585 **Implementation details.** AntiFold outputs a matrix of per-position log-probabilities in CSV format.
586 Given a mutant sequence and set of mutated positions P , the log-likelihood is computed by gathering
587 the model's log-probability for each mutant residue at the corresponding position:

```
588 # For LL_MUT  
589 log_probs_mut = antifold(model_input=mutant.pdb)  
590 ll_mut = sum([log_probs_mut[t][s_mut[t]] for t in P])  
591  
592 # For LL_PA  
593 log_probs_pa = antifold(model_input=parent.pdb)  
594 ll_pa = sum([log_probs_pa[t][s_mut[t]] for t in P])  
595
```

597 Here, $s_mut[t]$ refers to the amino acid at position t in the mutant sequence, and
598 $log_probs[t][aa]$ gives the log-probability assigned to amino acid aa at position t .

599 Summary.

- 600 • LL_{MUT} : scores the mutant using its own autoregressive context.
- 601 • LL_{PA} : scores the mutant residues using log-probabilities from the parent context.
- 602 • In both cases, log-likelihood is computed by summing over specified mutated regions.