

000 001 UNAAGI: ATOM-LEVEL DIFFUSION FOR GENERAT- 002 ING NON-CANONICAL AMINO ACID SUBSTITUTIONS 003 004

005 **Anonymous authors**

006 Paper under double-blind review

007 008 ABSTRACT 009

010 Proposing beneficial amino acid substitutions, whether for mutational effect pre-
011 diction or protein engineering, remains a central challenge in structural biology.
012 Recent inverse folding models, trained to reconstruct sequences from structure,
013 have had considerable impact in identifying functional mutations. However, cur-
014 rent approaches are constrained to designing sequences composed exclusively
015 of natural amino acids (NAAs). The larger set of non-canonical amino acids
016 (NCAAs), which offer greater chemical diversity, and are frequently used in in-
017 vivo protein engineering, remain largely inaccessible for current variant effect pre-
018 diction methods.

019 To address this gap, we introduce **UNAAGI**, a diffusion-based generative model
020 that reconstructs residue identities from atomic-level structure using an $E(3)$ -
021 equivariant framework. By modeling side chains in full atomic detail rather than
022 as discrete tokens, UNAAGI enables the exploration of both canonical and non-
023 canonical amino acid substitutions within a unified generative paradigm. We
024 evaluate our method on experimentally benchmarked mutation effect datasets and
025 demonstrate that it achieves substantially improved performance on NCAA substi-
026 tutions compared to the current state-of-the-art. Furthermore, our results suggest
027 a shared methodological foundation between protein engineering and structure-
028 based drug design, opening the door for a unified training framework across these
029 domains.

030 031 1 INTRODUCTION

032 Proteins are linear chains of amino acids that fold into specific three-dimensional structures to per-
033 form a wide range of biological functions. Chemically, an amino acid consists of a central α -carbon
034 bonded to a hydrogen atom, an amino group (NH_2), a carboxyl group ($COOH$), and a variable side
035 chain denoted as R . While the R group can take on diverse chemical forms, only 20 distinct side-
036 chain structures are commonly found in naturally occurring proteins. These are referred to as the
037 genetically encoded, canonical, or Natural Amino Acids (NAAs) (Branden & Tooze, 1991). Several
038 studies suggest that this set of 20 NAAs has been evolutionarily optimized to achieve functional
039 completeness and broad chemical diversity (Ilardo et al., 2015; Doig, 2017; Ilardo et al., 2019).
040 Nevertheless, Non-Canonical Amino Acids (NCAAs) offer side-chain chemistries and func-
041 tionalities absent from the standard set, such as enhanced metal coordination, tailored electrostatics, or
042 novel catalytic properties. Incorporating NCAAs enables protein engineers to expand the functional
043 repertoire of proteins beyond evolutionary constraints, unlocking new possibilities in synthetic biol-
044 ogy and therapeutics (Link et al., 2003; Chin, 2017; Rogers et al., 2018).

045 In recent years, machine learning – particularly deep learning – has achieved remarkable success in
046 protein research, advancing tasks such as structure prediction, design, and mutational effect estima-
047 tion (Jumper et al., 2021; Dauparas et al., 2022; Abramson et al., 2024). Variant effect prediction
048 models are often trained in an unsupervised fashion, where amino acid propensities in a protein are
049 predicted conditioned on a sequence (Riesselman et al., 2018) or structural (Dauparas et al., 2022)
050 context. Such procedure generally model amino acids using a discrete distribution, typically limited
051 to the 20 naturally occurring amino acids, based on the propensities observed in our vast databases
052 on evolutionary data. While such models have shown impressive zero-shot predictive performance
053 for various protein properties, the discrete modeling prohibits generalization beyond the naturally
occurring amino acids, representing a fundamental limitation in the current modeling paradigm.

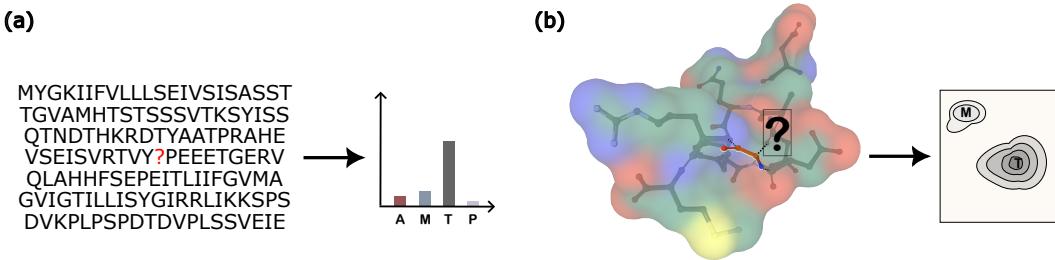


Figure 1: An illustration of the motivation behind UNAAGI. **(a)** Previous models predicted amino acid identities from a fixed vocabulary, limiting sampling to natural amino acids. **(b)** UNAAGI directly samples side chains at the atom level, enabling generation of more diverse structures, including non-canonical amino acids (NCAs).

To address this gap, we propose a novel generative framework for residue-level sequence design based on atom-level side-chain generation via equivariant molecular diffusion. We term our approach **Uncanonical Novel Amino Acid Generative Inference (UNAAGI)**. By modeling the atomic coordinates of side chains directly, UNAAGI implicitly infers residue identity and is capable of proposing plausible non-canonical substitutions; a visual illustration of our motivation is shown in Figure 1. We evaluate our model by comparing its generative likelihoods to experimentally measured mutational effect (Rogers et al., 2018; Notin et al., 2023), and find that the performance on non-canonical amino acids is improved considerably over the current state-of-the-art. Importantly, our model maintains consistent performance across both canonical and non-canonical residues.

Specifically, our contributions are as follows:

- We introduce a novel $E(3)$ -equivariant diffusion framework for atom-wise side-chain generation, offering a new approach to residue identity inference.
- We investigate the zero-shot performance of UNAAGI on variant effect prediction. While our model does not fully match the state-of-the-art for naturally occurring amino acids, we are the first method to demonstrate meaningful levels of correlation for non-canonical substitutions.
- We discuss the broader implications of our approach, suggesting that it may bridge structure-based drug design (SBDD) and protein engineering, as both fields can benefit from shared modeling tools such as energy-based design.

In the following sections, we formalize the generation task at the atom level, describe our $E(3)$ -equivariant diffusion framework UNAAGI, evaluate it on mutational benchmarks, and discuss its implications for protein design and synthetic biology.

2 RELATED WORK

Our approach has connections to prior work in molecular diffusion, variant effect prediction and peptide design.

2.1 MOLECULAR DIFFUSION

Hoogeboom et al. (2022) introduced the first framework for $E(3)$ -equivariant diffusion, generating molecular conformations in 3D space via a diffusion process. However, their approach suffered from low generative quality and often produced molecules with disconnected atoms. Subsequent works, including Peng et al. (2023) and Vignac et al. (2023), proposed refinements that incorporated inductive biases from molecular bonding structures, leading to improved sampling quality.

Le et al. (2024) provided a comparative study of several molecular diffusion frameworks, analyzing the impact of different modeling choices on sample quality. Their work identified combinations of design choices that optimize performance for molecular generative tasks.

Beyond unconditional molecular generation, a critical application is the conditional generation of compounds that bind to specific protein targets—known as Structure-Based Drug Design (SBDD). Schneuing et al. (2024) extended the framework of Hoogeboom et al. (2022) to conditional ligand generation in 3D space, while Guan et al. (2023) proposed similar models for target-aware ligand sampling. More recently, Schneuing et al. (2025) advanced this line of research by introducing models capable of sampling ligands with variable atom counts and incorporating protein pocket side-chain conformations during generation.

2.2 MUTATIONAL EFFECT PREDICTION

Mutational effect prediction is the task of predicting the impact of amino acid substitutions in a protein on a particular property of interest (e.g. stability, affinity or function). Multiple state-of-the-art protein models—whether sequence-based, structure-based, or hybrids of the two—have demonstrated sensitivity to the mutational landscape of proteins.

Statistical models of multiple sequence alignments remain a strong baseline when many sequences are available for a given protein family. A notable example is the DeepSequence VAE model (Rieselman et al., 2018). Language-based models such as Meier et al. (2021) and Lin et al. (2023) obviate the need for alignments, leveraging conservation and coevolutionary signals learned from massive protein sequence datasets. While sequence-based model can predict mutational effect in many contexts, their performance declines when evolutionary signals are sparse, such as for shallow MSAs, antibodies, or de novo designed peptides.

Structure-based approaches replace the sequence context of a mutation with a structural context. Early models considered the structural environment around one amino acid at a time (Boomsma & Frellsen, 2017; Torng & Altman, 2017), which was later generalized to entire sequences in inverse-folding models (Ingraham et al., 2019; Dauparas et al., 2022; Hsu et al., 2022). Although primarily designed for structure reconstruction, these models also show utility in mutational effect prediction Frellsen et al. (2025). Our approach mirrors the early structure-based models by considering only single amino acids at a time, with the crucial difference that UNAAGI extends mutational effect prediction to include non-canonical amino acids.

2.3 MODELING NON-CANONICAL AMINO ACID SUBSTITUTIONS

PepINVENT (Geylan et al., 2025) is a transformer-based language model trained on atomic-level peptide representations using CHUCKLES (Siani et al., 1994), which converts amino acid residue tokens into SMILES strings. This enables the model to generate amino acid SMILES for masked positions. The strategy allows for generalization to non-canonical amino acids, but the original study does not validate the sample weights with experimental data. We include this method as a baseline.

NCFflow (Lee & Kim, 2025) and RareFold (Li et al., 2025) are based on AlphaFold3-style architectures. NCFflow adopts the Flow Matching framework and evaluates mutational effect prediction on the same benchmark as ours, allowing for direct comparisons (see Section 4). RareFold introduces new tokens for NCAs and focuses on predicting protein structures containing NCAA substitutions or designing peptides with NCAs. Their study reports on successful experimentally validated designed peptides with NCAs, but the model itself cannot predict affinity scores or mutational effect in silico, making direct comparison on NCAA mutational effect prediction infeasible.

3 UNAAGI

We introduce UNAAGI, a diffusion-based model that bridges molecular generation and protein mutational effect prediction by constructing amino acid side chains atom by atom, thereby extending the mutational landscape beyond the natural amino acid space. By learning from the chemistry of canonical residues, UNAAGI can generalize and interpolate to non-canonical amino acids. In what follows, we outline the general idea of equivariant graph diffusion, describe the model architecture, training setup, and evaluation protocol, and finally present the baselines used for comparison.

162 3.1 MULTI-MODAL DIFFUSION FOR MOLECULAR GRAPHS
163

164 Diffusion models have emerged as a powerful class of generative models (Sohl-Dickstein et al.,
165 2015; Ho et al., 2020; Song & Ermon, 2020; Song et al., 2021a;b). In the discrete-time setting, these
166 models learn to reverse a Markovian forward process that gradually corrupts data with noise until the
167 samples resemble a tractable prior distribution (typically Gaussian). The generative model is trained
168 to approximate the reverse process, which recovers data via a parameterized denoising chain.
169

170 Originally developed for image generation, diffusion models have since been extended to molecular
171 and protein structure generation (Hoogeboom et al., 2022). Let $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ denote a sample from
172 the data distribution. The forward process is defined as
173

174
$$q(\mathbf{x}_{1:T} \mid \mathbf{x}_0) = \prod_{t=1}^T q(\mathbf{x}_t \mid \mathbf{x}_{t-1}), \quad (1)$$

175

176 which incrementally adds noise over T steps until \mathbf{x}_T becomes indistinguishable from Gaussian
177 noise. The reverse process is then modeled as
178

179
$$p_\theta(\mathbf{x}_{0:T}) = p(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1} \mid \mathbf{x}_t), \quad (2)$$

180

181 which gradually denoises \mathbf{x}_T back to a clean sample.
182

183 Molecular data is inherently *multi-modal*: atomic coordinates are continuous, while atom types,
184 bond types, and other chemical features are discrete. To handle this, diffusion over discrete variables
185 has been explored through categorical transition kernels—discrete analogues of Gaussian noise in
186 the forward process (Hoogeboom et al., 2021; Austin et al., 2023). Recent molecular diffusion
187 models combine continuous and categorical processes to jointly model the heterogeneous modalities
188 of molecules (Peng et al., 2023; Vignac et al., 2023; Guan et al., 2023; Le et al., 2024).
189

190 The forward noise process for continuous and discrete modalities can be expressed as
191

192
$$\begin{aligned} q(\mathbf{x}_t \mid \mathbf{x}_0) &= \mathcal{N}(\mathbf{x}_t \mid \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I}), \\ q(\mathbf{c}_t \mid \mathbf{c}_0) &= \mathcal{C}(\mathbf{c}_t \mid \bar{\alpha}_t \mathbf{c}_0 + (1 - \bar{\alpha}_t) \tilde{\mathbf{c}}), \end{aligned} \quad (3)$$

193 where $\bar{\alpha}_t = \prod_{k=1}^t (1 - \beta_k) \in (0, 1)$ is the cumulative noise schedule. For discrete features, $\tilde{\mathbf{c}}$ denotes
194 the categorical prior (e.g., uniform or empirical distribution). Following D3PM (Austin et al., 2023),
195 we adopt a *mask diffusion* scheme, which gradually converts labels into an absorbing state during
196 the forward process.
197

198 In our setting, we perturb both atomic coordinates \mathbf{X} and atom-wise categorical features \mathbf{H} . The fea-
199 tures in \mathbf{H} include atom types, formal charges, hybridization states, aromaticity, ring membership,
200 and atomic degree. Each modality is diffused independently using either Gaussian or categorical
201 noise. This chemically-aware perturbation strategy is inspired by Le et al. (2024), who showed that
202 injecting chemically meaningful noise improves sample quality and stability.
203

204 The training objective follows the standard ELBO on the data log-likelihood:
205

206
$$\log p(\mathbf{x}_0) \geq \mathcal{L}_0 + \mathcal{L}_{\text{prior}} + \sum_{t=1}^{T-1} \mathcal{L}_t, \quad (4)$$

207 where $\mathcal{L}_0 = \log p(\mathbf{x}_0 \mid \mathbf{x}_1)$ is the reconstruction term and $\mathcal{L}_{\text{prior}} = -\text{KL}(q(\mathbf{x}_T) \parallel p(\mathbf{x}_T))$ matches the
208 prior. In practice, these terms are often omitted. The main learning signal comes from minimizing
209 the per-timestep KL divergence:
210

211
$$\mathcal{L}_t = \text{KL}[q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1} \mid \mathbf{x}_t)], \quad (5)$$

212 which has closed-form solutions under both Gaussian and categorical noise. This is typically op-
213 timized by predicting either the clean data $\hat{\mathbf{x}}_0$, the noise ϵ , or categorical logits, depending on the
214 parameterization (Ho et al., 2020; Austin et al., 2023). We adopt the $\hat{\mathbf{x}}_0$ -parameterization to di-
215 rectly predict the clean data. The loss function is mean squared error for atomic coordinates and
216 cross-entropy for other categorical features.
217

216 3.2 E(3)-EQUIVARIANT GRAPH NEURAL NETWORK
217218 Molecular structures are inherently three-dimensional and symmetric under Euclidean transfor-
219 mations such as rotation and translation. To respect these symmetries, we incorporate *equivariance* as
220 an inductive bias.221 Formally, a function $f : \mathcal{X} \rightarrow \mathcal{Y}$ is equivariant to a group G if
222

223
$$f(g \cdot x) = g \cdot f(x), \quad \forall g \in G.$$

224 For molecular data, the relevant group is $E(3)$, which includes all 3D rotations and translations.
225 Specifically, a denoising function on atomic coordinates $\mathbf{X} \in \mathbb{R}^{N \times 3}$ should satisfy
226

227
$$f(\mathbf{X}Q + \mathbf{t}) = f(\mathbf{X})Q + \mathbf{t},$$

228 for any orthogonal transformation $Q \in O(3)$ and translation $\mathbf{t} \in \mathbb{R}^3$.
229230 We adopt an $E(3)$ -equivariant Graph Neural Network, closely following the EQGAT-diff model of
231 Le et al. (2024), to learn the score function over atom-level graphs. Each node represents an atom
232 and carries scalar features h_i (e.g., atom type, charge) and vector features $\mathbf{v}_i \in \mathbb{R}^3$, while edges
233 encode bond types and interatomic distances.234 Message passing is performed on a fully connected graph with attention-based aggregation. The
235 directional information between atoms i and j is encoded as a unit vector
236

237
$$\mathbf{x}_{ji,n} = \frac{\mathbf{x}_j - \mathbf{x}_i}{\|\mathbf{x}_j - \mathbf{x}_i\|},$$

238

239 which is used to construct equivariant vector features. Messages are computed via MLPs over
240 concatenated scalar and edge features, then split into components to update scalar, vector, and po-
241 sitional embeddings in an equivariant manner. This design ensures that all updates preserve $E(3)$ -
242 equivariance and enables joint modeling of geometry, connectivity, and chemistry during denoising.
243244 3.3 GRAPH TOPOLOGY AND VIRTUAL NODE PADDING
245246 Unlike traditional Structure-Based Drug Design (SBDD), where generative models sample ligands
247 conditioned on a protein pocket, our task focuses on reconstructing amino acid side chains. A
248 key distinction is that each generated side chain is covalently bonded to a fixed protein backbone,
249 imposing local structural constraints around the masked region. To account for this, we reformulate
250 the graph construction. Specifically, the positions and types of backbone atoms for the residue under
251 reconstruction are fixed, anchoring generation to the known structural context.252 A further challenge is unified sampling across natural amino acids (NAA) and non-canonical amino
253 acids (NCAA), which vary in atom counts. In standard SBDD models, the atom count must be fixed
254 before diffusion, since the graph topology must be predetermined. We address this using a virtual
255 node strategy inspired by DrugFlow (Schneuing et al., 2025). Unlike empirical sampling schemes
256 (Hoogeboom et al., 2022; Schneuing et al., 2024; Guan et al., 2023) or size-estimation networks
257 (Igashov et al., 2024), this approach is fully end-to-end. This allows sampling from a smoother
258 distribution over amino acid identities across side chains of varying sizes.259 Virtual nodes are assigned a special atom type NOATOM and are disconnected from all other atoms,
260 with edges labeled NOBOND. During training, a random number of virtual nodes $n_{\text{virt}} \sim U(0, N_{\text{max}})$
261 are added to the side chain. Following Schneuing et al. (2025), these nodes are initialized at the side-
262 chain center of mass. The total number of nodes is capped by a predefined upper bound N_{max} , while
263 the effective side-chain size becomes variable and learnable. At sampling time, the model denoises a
264 graph with N_{max} nodes. Nodes denoised to the virtual type are removed post hoc, allowing variable-
265 sized side chains to be generated within a unified diffusion framework.266 4 EXPERIMENTS
267268 We now describe the experimental procedure for UNAAGI, including training details, evaluation
269 methods, and comparisons against baselines.

270 4.1 DATASET COMPOSITION AND PREPROCESSING
271

272 We use a dataset of 1,000 protein structures submitted to the Protein Data Bank (PDB) for training.
 273 For each structure, we extract every residue along with its local environment, defined as surrounding
 274 residues whose center of mass lies within a 10 Å radius of the target residue. To ensure independence
 275 between training and evaluation, we verified that none of these PDB structures overlap with the
 276 benchmark assays.

277 In addition to canonical amino acids, we include non-canonical amino acids (NCAs) to enhance
 278 the model’s capacity to generalize to novel chemical environments and to prevent overfitting. For
 279 NCAs, we use datasets from Ilardo et al. (2019) and SwissSidechain (Gfeller et al., 2012), al-
 280 though these consist of NCAs without associated proteomic context. To capture diverse physical
 281 interactions, we also augment the training set with protein–ligand complexes from PDBBind (Wang
 282 et al., 2005). Graph construction differs slightly between sources: - For independent NCAs, the
 283 four backbone atoms (C, C, O, N) are fixed during training, and only side-chain atoms are diffused.
 284 - For PDBBind, we fix pocket atoms and apply noise to all ligand atoms, following the procedure of
 285 SBDD models in Le et al. (2024).

286 4.2 EVALUATION ON DEEP MUTATIONAL SCANNING (DMS)
287

288 During sampling, UNAAGI generates diverse side-chain conformations given a fixed environment.
 289 We interpret the sampling frequencies of side-chain identities and conformations as a proxy for the
 290 probability density learned by the diffusion model.

291 To evaluate predictive power, we compare this learned distribution against experimental Deep Mu-
 292 tational Scanning (DMS) data. Specifically, for each site in the DMS benchmark, we perform 100
 293 sampling iterations and record frequencies of each sampled amino acid. While larger sample sizes
 294 would provide smoother frequency estimates, we adopt 100 iterations as a trade-off between statis-
 295 tical stability and computational feasibility. Sampled side chains are matched to known natural or
 296 non-canonical amino acids using graph isomorphism (Weisfeiler & Lehman, 1968).

297 From the sampling frequencies, we estimate the likelihood of each mutant or wild-type residue.
 298 We then compute the negative log-likelihood (NLL) for each sampled amino acid, calculating the
 299 differential log-likelihood as:

$$301 \quad \Delta \log \mathcal{L} = -\log P(\text{mutant}) + \log P(\text{wild-type}),$$

302 where $P(\text{mutant})$ and $P(\text{wild-type})$ denote the estimated probabilities of mutant and wild-type
 303 residues, respectively. This log-likelihood difference is correlated with experimental $\Delta\Delta G$ values
 304 from the DMS benchmark.

305 As a sanity check on the ability to recover the natural amino acids, we selected assays with fewer
 306 than 100 residues from Notin et al. (2023)¹. The full list of assays and results is provided in
 307 Supplementary A.1.

308 To evaluate performance on NCAs, we use peptide–protein complexes from Rogers et al. (2018),
 309 specifically PDB ID 5LY1 (JMJD2A/KDM4A bound to a macrocyclic peptide CP2) and PDB ID
 310 2ROC (Mcl-1 bound to PUMA). This dataset includes DMS measurements of $\Delta\Delta G$ for substitu-
 311 tions across 20 canonical and 20 non-canonical amino acids. For fairness, we restrict evaluation to
 312 substitutions where both the mutant and wild-type residues appear in the sampled data.

314 4.3 RESULTS
315316 4.3.1 PROTEINGYM
317

318 We evaluate predictive performance by computing the Spearman correlation between likelihood
 319 differences and ground-truth experimental values, following the reporting standard of ProteinGym,
 320 with more detailed per-assay results in Supplementary A.1.

321
 322 ¹Note that the computational cost of UNAAGI scales linearly with protein size (since only local structural
 323 environments are considered). Scaling to larger proteins is therefore not fundamentally problematic, but given
 a limited computational budget, we prioritized protein diversity over protein length in this experiment.

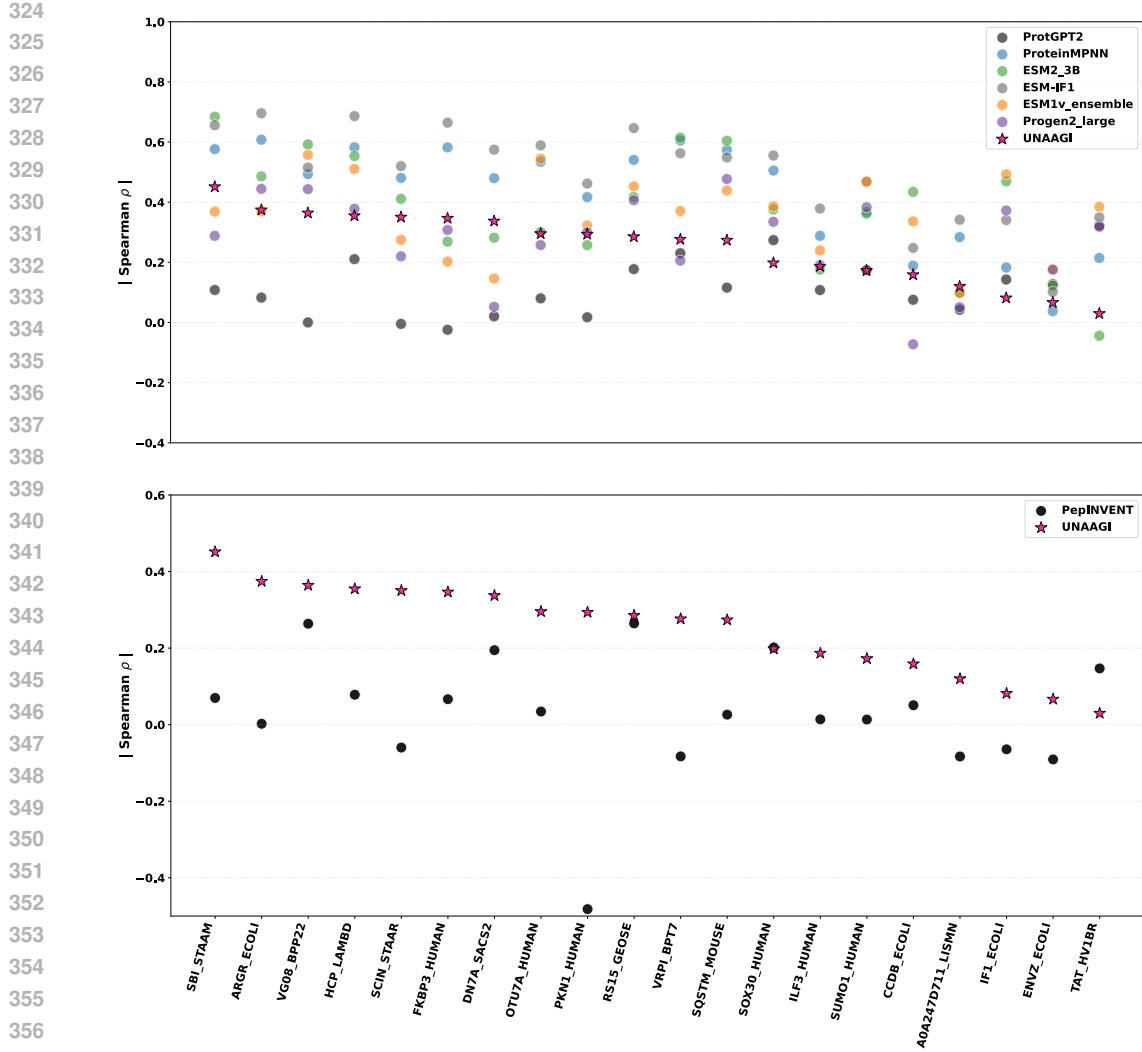


Figure 2: **Top:** Results of UNAAGI on 20 ProteinGym assays, evaluated using Spearman correlation against experimental DMS measurements. **Down:** Results of UNAAGI on ProteinGym, comparing to PepINVENT.

For reference, we also evaluate PepINVENT as a baseline in this setting, although it was not originally tested on DMS benchmarks. Following the repository guidelines, we sample by masking one residue at a time and generating candidate substitutions. This is repeated for all positions in each assay. For peptide–protein complexes, we concatenate peptide and protein sequences as input to the language model. Sampled canonical SMILES are mapped back to amino acids if they are in the benchmark to estimate frequencies.

While UNAAGI does not achieve state-of-the-art performance on any single assay, it still exhibits substantial correlation with experimental mutational effect across most assays (Figure 2). Importantly, UNAAGI is distinct from all other baselines in Figure 2, as it generates residues in an atom-wise manner rather than sampling from a fixed vocabulary. The experiment is somewhat adversarial, since we for this dataset know *a priori* that we are restricted to the 20 amino acids. It is encouraging that the models produces reasonable performance despite the more expressive output distribution.

For a fair comparison, we evaluate UNAAGI against its closest methodological baseline, PepINVENT, which represents sequences and samples residues atom by atom. The comparison, in terms of Spearman correlation on ProteinGym, is also shown in Figure 2. UNAAGI consistently outper-

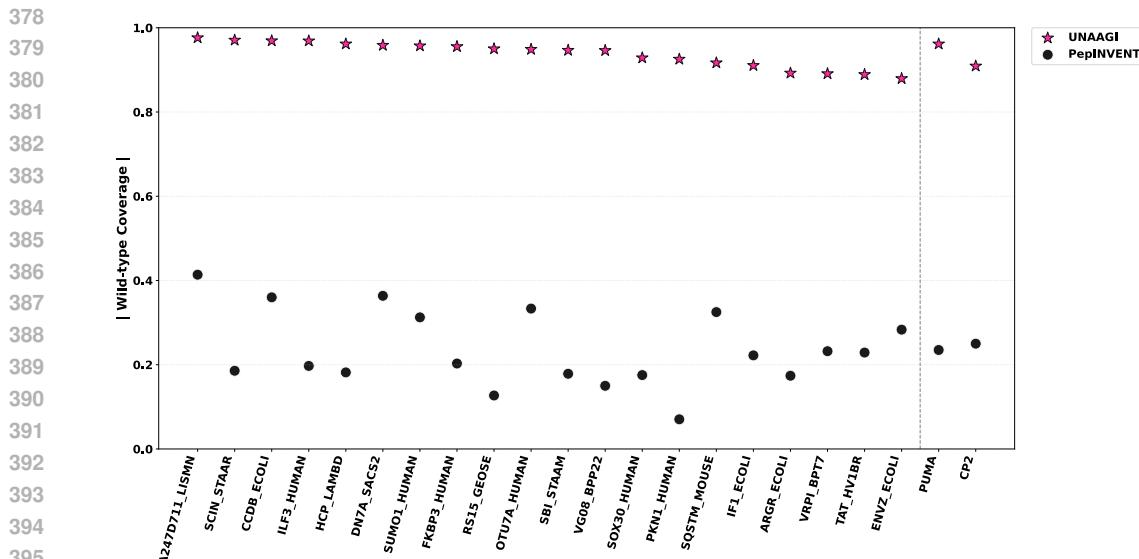


Figure 3: Wild-type coverage rate of UNAAGI and PepINVENT across benchmark assays.

forms PepINVENT on nearly all assays, with the exception of TAT_HV1BR. Notably, PepINVENT shows little to no correlation on most assays, whereas UNAAGI achieves substantial predictive performance.

Since both models operate in an atom-wise sampling regime, they must demonstrate sufficient capacity to reconstruct wild-type residues. To quantify this, we define the *wild-type coverage rate*—the frequency with which the wild-type residue is successfully sampled across all positions in an assay. Results are shown in Figure 3, where UNAAGI achieves consistently high coverage across sites, while PepINVENT rarely recovers the wild-type residue. This limitation restricts PepINVENT to evaluation on only a small subset of benchmark positions, and even under this reduced setting, its performance remains below that of UNAAGI, further emphasizing UNAAGI’s advantage.

4.4 DMS FOR NCAA SUBSTITUTIONS

We further evaluate UNAAGI on a DMS benchmark containing substitutions to non-canonical amino acids (NCAs). For comparison, we include both PepINVENT and NCFlow, as NCFlow was evaluated on the same benchmark in its original work.

As shown in Figure 4, UNAAGI achieves consistent performance across both canonical and non-canonical substitution benchmarks. The correlations observed on NCAA substitutions are comparable to those obtained on canonical amino acids, demonstrating that UNAAGI can generalize beyond the natural amino acid space. In contrast, PepINVENT shows limited signal on NCAA substitutions, and NCFlow performs poorly across the NCAA benchmark, with no or negative correlations regardless of the affinity prediction module used. These results indicate that UNAAGI is the first diffusion-based approach to provide reliable predictive power on NCAA mutational effect benchmarks.

5 CONCLUSION

In this paper, we introduced UNAAGI, a molecular diffusion model that generates amino acid side chains in an atom-wise manner and can sample across both canonical and non-canonical amino acids. We evaluated UNAAGI on Deep Mutational Scanning (DMS) benchmarks and found that it achieves meaningful correlations in variant effect prediction. Crucially, this correlation extends to benchmarks containing NCAA substitutions, making UNAAGI the first machine learning method to provide measurable signal on this problem.

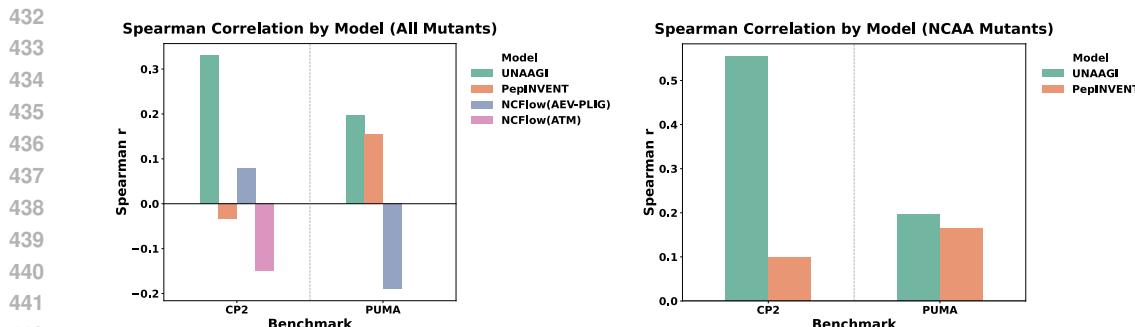


Figure 4: Results of UNAAGI on NCAA DMS benchmarks. NCFlow (ATM) was reported in the original work to fail entirely on the PUMA assay. **Left:** Spearman correlations compared to PepINVENT and NCFlow variants with different affinity modules. **Right:** Performance on NCAA substitutions specifically.

Methodologically, UNAAGI follows the molecular diffusion paradigm developed for Structure-Based Drug Discovery (SBDD), learning the density of feasible chemical combinations in protein structures. Unlike protein language models (Lin et al., 2023), UNAAGI does not rely on coevolutionary information from protein sequences, yet still shows predictive ability on mutational effect tasks. This suggests a promising connection between SBDD and variant effect prediction, as both are governed by the same non-covalent interaction principles. Building on this insight, future work could integrate additional SBDD techniques—such as pharmacophore-based conditioning (Peng et al., 2025) or flexible side-chain modeling (Schneuing et al., 2025)—to further improve variant effect prediction.

At the same time, UNAAGI has important limitations. It remains far from solving the problem of NCAA variant effect prediction and cannot yet propose high-fitness NCAA substitutions efficiently. As shown in Supplementary A.1, although UNAAGI achieves substantial correlations on the NCAA mutants it successfully samples, coverage remains limited to a small subset of the 20 NCAs in the benchmark. This limitation is largely driven by the scarcity of relevant NCAA data. Moreover, UNAAGI tends to interpolate between canonical-like structures, while sampling chemically distinct NCAs remains an out-of-distribution challenge under the current setup.

We also observe substantial variability in predictive performance across different assays, a phenomenon shared by all methods. This highlights the intrinsic complexity of the variant effect prediction problem: no single model performs universally well, and outcomes depend strongly on model design, training strategy, and data availability.

Looking forward, UNAAGI opens several avenues for exploration. A natural next step is scaling the model, both in parameter count (beyond the current 3.6M) and in training data (from the 1,000 PDB structures used here toward the full database). Incorporating protein structures with NCAA substitutions—although sparse in the PDB—could further expand coverage of the NCAA chemical space. Finally, applying guidance or auxiliary correctors may enhance the generation of chemically diverse NCAs and enable biasing toward desirable properties, such as synthetic accessibility.

REFERENCES

- Josh Abramson, Jonas Adler, Jack Dunger, Richard Evans, Tim Green, Alexander Pritzel, Olaf Ronneberger, Lindsay Willmore, Andrew J. Ballard, Joshua Bambrick, Sebastian W. Bodenstein, David A. Evans, Chia-Chun Hung, Michael O’Neill, David Reiman, Kathryn Tunyasuvanakool, Zachary Wu, Akvilė Žemgulytė, Eirini Arvaniti, Charles Beattie, Ottavia Bertolli, Alex Bridgland, Alexey Cherepanov, Miles Congreve, Alexander I. Cowen-Rivers, Andrew Cowie, Michael Figurnov, Fabian B. Fuchs, Hannah Gladman, Rishabh Jain, Yousuf A. Khan, Caroline M. R. Low, Kuba Perlin, Anna Potapenko, Pascal Savy, Sukhdeep Singh, Adrian Stecula, Ashok Thillaisundaram, Catherine Tong, Sergei Yakneen, Ellen D. Zhong, Michal Zielinski, Augustin Žídek, Victor Bapst, Pushmeet Kohli, Max Jaderberg, Demis Hassabis, and John M. Jumper.

- 486 Accurate structure prediction of biomolecular interactions with alphafold 3. *Nature*, 630(8016):
 487 493–500, 2024. doi: 10.1038/s41586-024-07487-w. URL <https://doi.org/10.1038/s41586-024-07487-w>.
- 489
 490 Jacob Austin, Daniel D. Johnson, Jonathan Ho, Daniel Tarlow, and Rianne van den Berg. Structured
 491 denoising diffusion models in discrete state-spaces, 2023. URL <https://arxiv.org/abs/2107.03006>.
- 493
 494 Wouter Boomsma and Jes Frellsen. Spherical convolutions and their application in molecular mod-
 495 ellng. *Advances in neural information processing systems*, 30, 2017.
- 496 Carl Branden and John Tooze. *Introduction to Protein Structure*. Garland Publishing, New York,
 497 1991.
- 499 Jason W. Chin. Expanding and reprogramming the genetic code. *Nature*, 550:53–60, 2017. doi:
 500 10.1038/nature24031. URL <https://doi.org/10.1038/nature24031>.
- 501 J. Dauparas, I. Anishchenko, N. Bennett, H. Bai, R. J. Ragotte, L. F. Milles, B. I. M. Wicky,
 502 A. Courbet, R. J. de Haas, N. Bethel, P. J. Y. Leung, T. F. Huddy, S. Pellock, D. Tischer, F. Chan,
 503 B. Koepnick, H. Nguyen, A. Kang, B. Sankaran, A. K. Bera, N. P. King, and D. Baker. Robust
 504 deep learning-based protein sequence design using proteinmpnn. *Science*, 378(6615):49–56,
 505 2022. doi: 10.1126/science.add2187. URL <https://www.science.org/doi/abs/10.1126/science.add2187>.
- 507
 508 Alan J. Doig. Frozen, but no accident – why the 20 standard amino acids were selected. *FEBS
 509 Journal*, 284(9):1296–1305, 2017. doi: 10.1111/febs.13982. URL <https://doi.org/10.1111/febs.13982>.
- 511
 512 Jes Frellsen, Maher M Kassem, Tone Bengtsen, Lars Olsen, Kresten Lindorff-Larsen, Jesper
 513 Ferkinghoff-Borg, and Wouter Boomsma. Zero-shot protein stability prediction by inverse folding
 514 models: a free energy interpretation. *arXiv preprint arXiv:2506.05596*, 2025.
- 515
 516 G. Geylan, J. P. Janet, A. Tibo, J. He, A. Patronov, M. Kabeshov, W. Czechtizky, F. David, O. En-
 517 gkvist, and L. De Maria. Pepinvent: Generative peptide design beyond natural amino acids.
Chemical Science, 16(20):8682–8696, Apr 2025. doi: 10.1039/d4sc07642g.
- 518
 519 David Gfeller, Olivier Michelin, and Vincent Zoete. Expanding molecular modeling and design
 520 tools to non-natural sidechains. *Journal of Computational Chemistry*, 33(19):1525–1535, 2012.
 521 doi: 10.1002/jcc.22982.
- 522
 523 Jiaqi Guan, Wesley Wei Qian, Xingang Peng, Yufeng Su, Jian Peng, and Jianzhu Ma. 3d equivariant
 524 diffusion for target-aware molecule generation and affinity prediction, 2023. URL <https://arxiv.org/abs/2303.03543>.
- 525
 526 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models, 2020. URL
 527 <https://arxiv.org/abs/2006.11239>.
- 528
 529 Emiel Hoogeboom, Didrik Nielsen, Priyank Jaini, Patrick Forré, and Max Welling. Argmax flows
 530 and multinomial diffusion: Learning categorical distributions, 2021. URL <https://arxiv.org/abs/2102.05379>.
- 532
 533 Emiel Hoogeboom, Victor Garcia Satorras, Clément Vignac, and Max Welling. Equivariant diffu-
 534 sion for molecule generation in 3d, 2022. URL <https://arxiv.org/abs/2203.17003>.
- 535
 536 Chloe Hsu, Robert Verkuil, Jason Liu, Zeming Lin, Brian Hie, Tom Sercu, Adam Lerer, and
 537 Alexander Rives. Learning inverse folding from millions of predicted structures. In Kamalika
 538 Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato
 539 (eds.), *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of
Proceedings of Machine Learning Research, pp. 8946–8970. PMLR, 17–23 Jul 2022. URL
<https://proceedings.mlr.press/v162/hsu22a.html>.

- 540 Ilia Igashov, Hannes Stärk, Clément Vignac, Arne Schneuing, Victor Garcia Satorras, Pascal
 541 Frossard, Max Welling, Michael Bronstein, and Bruno Correia. Equivariant 3d-conditional dif-
 542 fusion model for molecular linker design. *Nature Machine Intelligence*, 6(4):417–427, 2024.
 543 ISSN 2522-5839. doi: 10.1038/s42256-024-00815-9. URL <https://doi.org/10.1038/s42256-024-00815-9>.
- 544
 545 Melissa Ilardo, Markus Meringer, Stephen Freeland, Bakhtiyor Rasulev, and H. James Cleaves.
 546 Extraordinarily adaptive properties of the genetically encoded amino acids. *Scientific Reports*, 5:
 547 9414, 2015. doi: 10.1038/srep09414. URL <https://doi.org/10.1038/srep09414>.
- 548
 549 Melissa Ilardo, Raghav Bose, Markus Meringer, Bakhtiyor Rasulev, Oleksandr Zhulyn, Ian Carrick,
 550 and H. James Cleaves. Adaptive properties of the genetically encoded amino acid alphabet are
 551 inherited from its subsets. *Scientific Reports*, 9:12468, 2019. doi: 10.1038/s41598-019-47574-x.
 552 URL <https://doi.org/10.1038/s41598-019-47574-x>.
- 553 John Ingraham, Vikas Garg, Regina Barzilay, and Tommi Jaakkola. Generative models for graph-
 554 based protein design. *Advances in neural information processing systems*, 32, 2019.
- 555
 556 John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger,
 557 Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, Alex Bridgland,
 558 Clemens Meyer, Simon A. A. Kohl, Andrew J. Ballard, Andrew Cowie, Bernardino Romera-
 559 Paredes, Stanislav Nikolov, Rishabh Jain, Jonas Adler, Trevor Back, Stig Petersen, David Reiman,
 560 Ellen Clancy, Michal Zielinski, Martin Steinegger, Michalina Pacholska, Tamas Berghammer,
 561 Sebastian Bodenstein, David Silver, Oriol Vinyals, Andrew W. Senior, Koray Kavukcuoglu,
 562 Pushmeet Kohli, and Demis Hassabis. Highly accurate protein structure prediction with al-
 563 phafold. *Nature*, 596(7873):583–589, 2021. doi: 10.1038/s41586-021-03819-2. URL <https://doi.org/10.1038/s41586-021-03819-2>.
- 564
 565 Tuan Le, Julian Cremer, Frank Noe, Djork-Arné Clevert, and Kristof T Schütt. Navigating the
 566 design space of equivariant diffusion-based generative models for de novo 3d molecule gen-
 567 eration. In *The Twelfth International Conference on Learning Representations*, 2024. URL
 568 <https://openreview.net/forum?id=kzGuIRXZrQ>.
- 569
 570 Jin Sub Lee and Philip M. Kim. Design of peptides with non-canonical amino acids using flow
 571 matching. *bioRxiv*, 2025. doi: 10.1101/2025.07.31.667780. URL <https://www.biorxiv.org/content/early/2025/07/31/2025.07.31.667780>.
- 572
 573 Quizhen Li, Diandra Daumiller, and Patrick Bryant. Rarefold: Structure prediction and de-
 574 sign of proteins with noncanonical amino acids. *bioRxiv*, 2025. doi: 10.1101/2025.05.19.
 575 654846. URL <https://www.biorxiv.org/content/early/2025/05/23/2025.05.19.654846>.
- 576
 577 Zeming Lin, Haydar Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wei Lu, Tom Sercu, Sergio
 578 Candido, and Alexander Rives. Evolutionary-scale prediction of atomic-level protein structure
 579 with a language model. *Science*, 379:1123–1130, 2023. doi: 10.1126/science.ade2574.
- 580
 581 A. J. Link, M. L. Mock, and D. A. Tirrell. Non-canonical amino acids in protein engineering.
 582 *Current Opinion in Biotechnology*, 14(6):603–609, 2003. doi: 10.1016/j.copbio.2003.10.011.
 583 URL <https://doi.org/10.1016/j.copbio.2003.10.011>.
- 584
 585 Joshua Meier, Roshan Rao, Robert Verkuil, Jason Liu, Tom Sercu, and Alexander Rives. Language
 586 models enable zero-shot prediction of the effects of mutations on protein function. *bioRxiv*,
 587 2021. doi: 10.1101/2021.07.09.450648. URL <https://doi.org/10.1101/2021.07.09.450648>.
- 588
 589 Pascal Notin, Aaron Kollasch, Daniel Ritter, Lood van Niekerk, Steffanie Paul, Han Spin-
 590 ner, Nathan Rollins, Ada Shaw, Rose Orenbuch, Ruben Weitzman, Jonathan Frazer,
 591 Mafalda Dias, Dinko Franceschi, Yarin Gal, and Debora Marks. Proteingym: Large-
 592 scale benchmarks for protein fitness prediction and design. In A. Oh, T. Naumann,
 593 A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural Informa-
 594 tion Processing Systems*, volume 36, pp. 64331–64379. Curran Associates, Inc., 2023.
 595 URL https://proceedings.neurips.cc/paper_files/paper/2023/file/

- 594 cac723e5ff29f65e3fcbb0739ae91bee-Paper-Datasets_and_Benchmarks.
 595 pdf.
 596
- 597 Jian Peng, Jun-Lin Yu, Zeng-Bao Yang, Yi-Ting Chen, Si-Qi Wei, Fan-Bo Meng, Yao-Geng
 598 Wang, Xiao-Tian Huang, and Guo-Bo Li. Pharmacophore-oriented 3d molecular generation to-
 599 ward efficient feature-customized drug discovery. *Nature Computational Science*, Aug 2025.
 600 ISSN 2662-8457. doi: 10.1038/s43588-025-00850-5. URL <https://doi.org/10.1038/s43588-025-00850-5>.
- 601
- 602 Xingang Peng, Jiaqi Guan, Qiang Liu, and Jianzhu Ma. Moldiff: Addressing the atom-bond in-
 603 consistency problem in 3d molecule diffusion generation, 2023. URL <https://arxiv.org/abs/2305.07508>.
- 604
- 605 Adam J Riesselman, John B Ingraham, and Debora S Marks. Deep generative models of genetic
 606 variation capture the effects of mutations. *Nature methods*, 15(10):816–822, 2018.
- 607
- 608 J. M. Rogers, T. Passioura, and H. Suga. Nonproteinogenic deep mutational scanning of linear and
 609 cyclic peptides. *Proceedings of the National Academy of Sciences of the United States of America*,
 610 115(43):10959–10964, 2018. doi: 10.1073/pnas.1809901115. URL <https://doi.org/10.1073/pnas.1809901115>.
- 611
- 612 Arne Schneuing, Charles Harris, Yuanqi Du, Kieran Didi, Arian Jamasb, Ilia Igashov, Weitao Du,
 613 Carla Gomes, Tom L. Blundell, Pietro Liò, Max Welling, Michael Bronstein, and Bruno Correia.
 614 Structure-based drug design with equivariant diffusion models. *Nature Computational Science*,
 615 4(12):899–909, 2024. ISSN 2662-8457. doi: 10.1038/s43588-024-00737-x. URL <https://doi.org/10.1038/s43588-024-00737-x>.
- 616
- 617 Arne Schneuing, Ilia Igashov, Adrian W. Dobbeltstein, Thomas Castiglione, Michael M. Bronstein,
 618 and Bruno Correia. Multi-domain distribution learning for de novo drug design. In *The Thirteenth
 619 International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=g3VCIM94ke>.
- 620
- 621
- 622 M. A. Siani, D. Weininger, and J. M. Blaney. CHUCKLES: a method for representing and search-
 623 ing peptide and peptoid sequences on both monomer and atomic levels. *Journal of Chemi-
 624 cal Information and Computer Sciences*, 34(3):588–593, May 1994. ISSN 0095-2338. doi:
 625 10.1021/ci00019a017.
- 626
- 627 Jascha Sohl-Dickstein, Eric A. Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsuper-
 628 vised learning using nonequilibrium thermodynamics, 2015. URL <https://arxiv.org/abs/1503.03585>.
- 629
- 630 Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution,
 631 2020. URL <https://arxiv.org/abs/1907.05600>.
- 632
- 633 Yang Song, Conor Durkan, Iain Murray, and Stefano Ermon. Maximum likelihood training of score-
 634 based diffusion models. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan (eds.),
 635 *Advances in Neural Information Processing Systems*, 2021a. URL <https://openreview.net/forum?id=AklttWFnxS9>.
- 636
- 637 Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and Ben
 638 Poole. Score-based generative modeling through stochastic differential equations, 2021b. URL
 639 <https://arxiv.org/abs/2011.13456>.
- 640
- 641 Wen Torn and Russ B Altman. 3D deep convolutional neural networks for amino acid environment
 642 similarity analysis. *BMC Bioinformatics*, 18(1):302, 2017.
- 643
- 644 Clement Vignac, Nagham Osman, Laura Toni, and Pascal Frossard. Midi: Mixed graph and 3d
 645 denoising diffusion for molecule generation, 2023. URL <https://arxiv.org/abs/2302.09048>.
- 646
- 647 Renxiao Wang, Xueguang Fang, Ying Lu, Chuangye Yang, and Shaomeng Wang. The pdbsbind
 648 database: methodologies and updates. *Journal of Medicinal Chemistry*, 48(12):4111–4119, 2005.
 649 doi: 10.1021/jm048957q.

648 Boris Weisfeiler and A. A. Lehman. A Reduction of a Graph to a Canonical Form and an Algebra
 649 Arising During This Reduction. *Nauchno-Technicheskaya Informatsia*, Ser. 2(N9):12–16, 1968.
 650

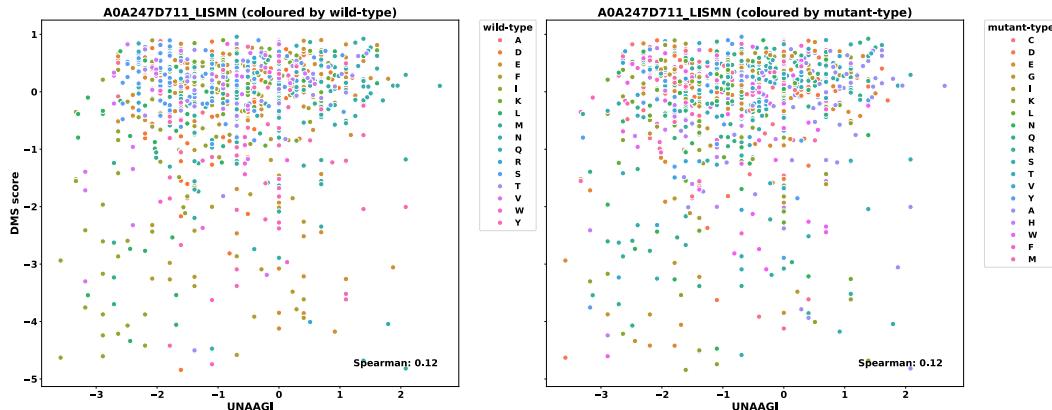
651 A APPENDIX

652 A.1 DETAILED RESULTS ON DMS BENCHMARKS

653 In this section, we present the detailed performance of UNAAGI on the DMS benchmarks. Section A.1.1 reports results on ProteinGym, while Section A.1.2 reports results on DMS assays containing NCAAs.

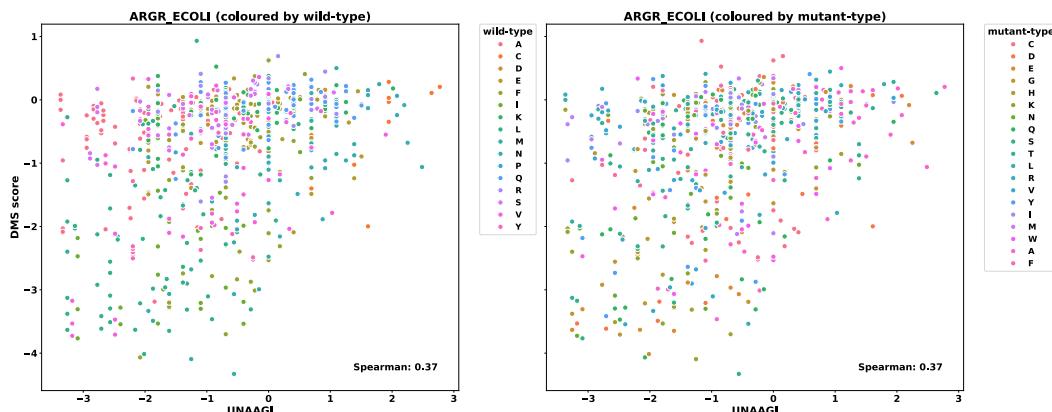
654 A.1.1 DETAILED RESULTS ON PROTEINGYM

655 **A0A247D711_LISMN**



656 Figure 5: Scatterplot of UNAAGI vs. DMS score

657 **ARGR_ECOLI**



658 Figure 6: Scatterplot of UNAAGI vs. DMS score

659 **CCDB_ECOLI**

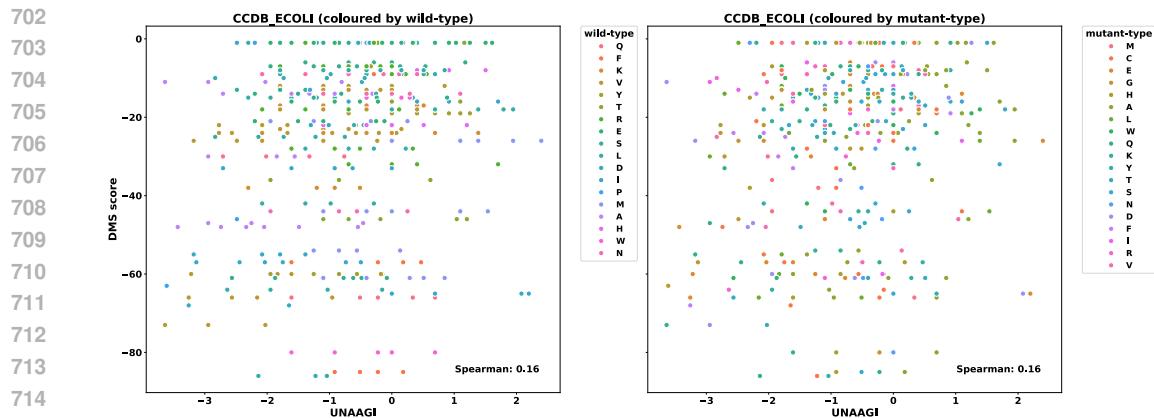


Figure 7: Scatterplot of UNAAGI vs. DMS score

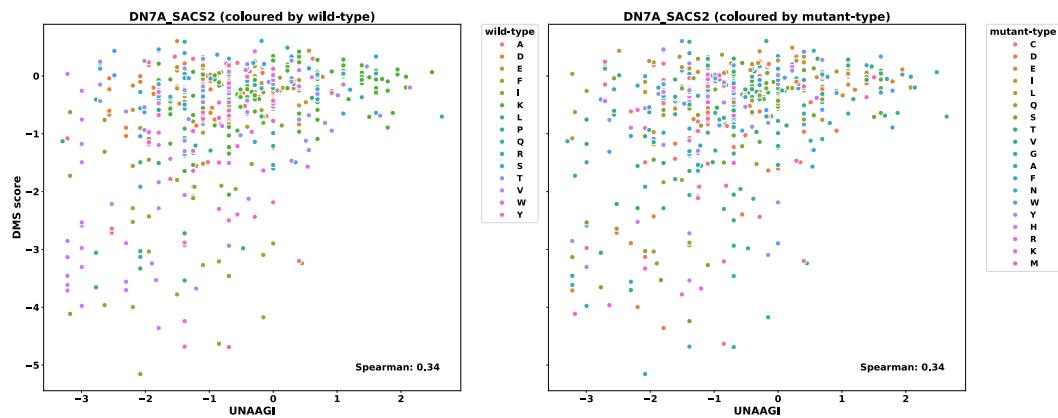
DN7A_SACS2

Figure 8: Scatterplot of UNAAGI vs. DMS score

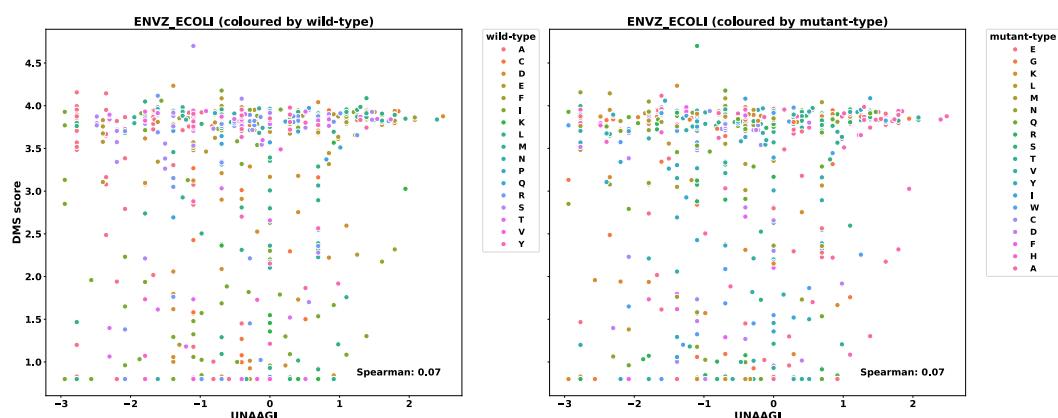
ENVZ_ECOLI

Figure 9: Scatterplot of UNAAGI vs. DMS score

FKBP3_HUMAN

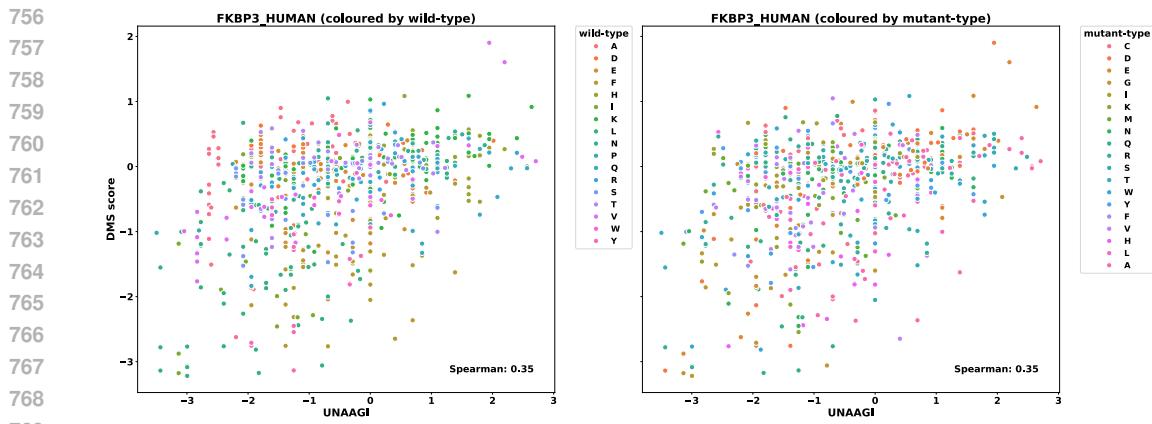


Figure 10: Scatterplot of UNAAGI vs. DMS score

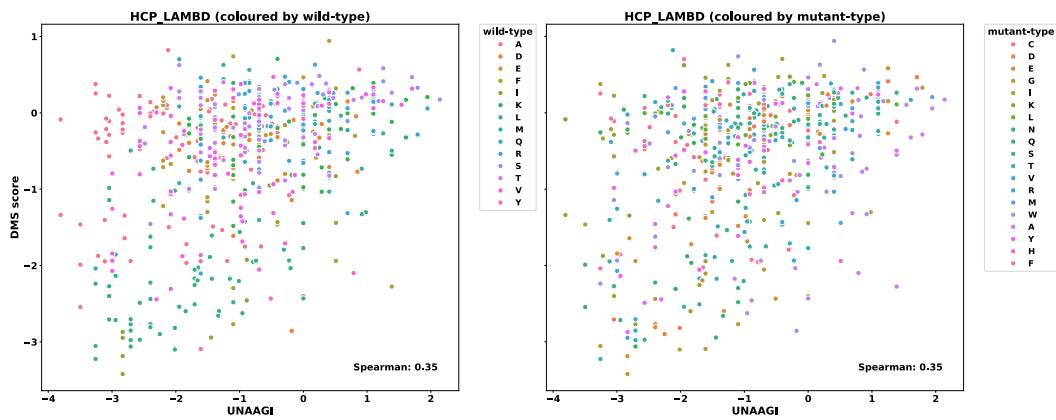
HCP_LAMBD

Figure 11: Scatterplot of UNAAGI vs. DMS score

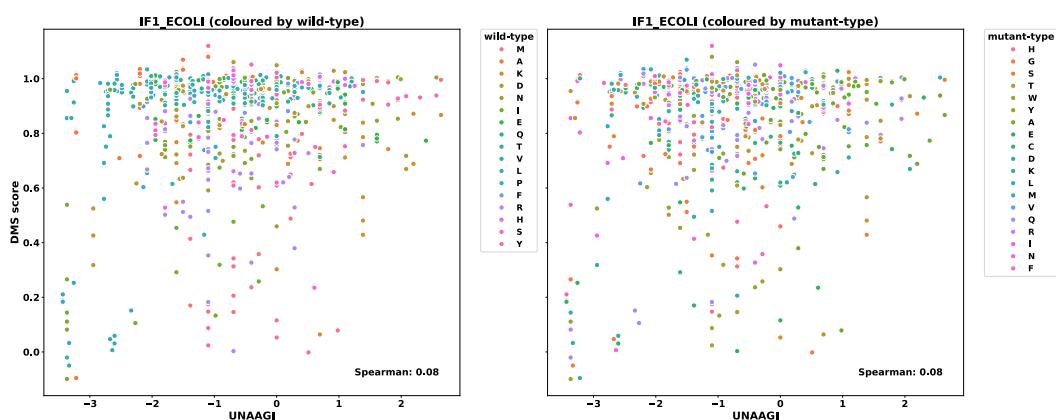
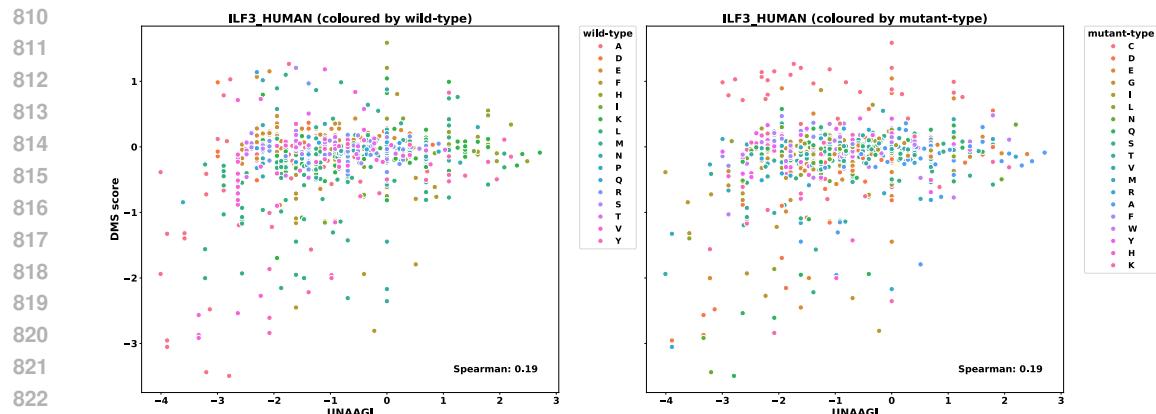
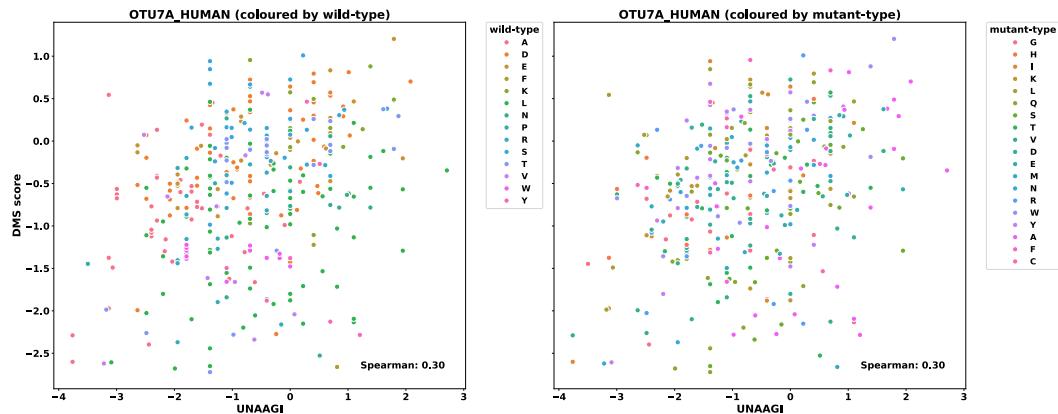
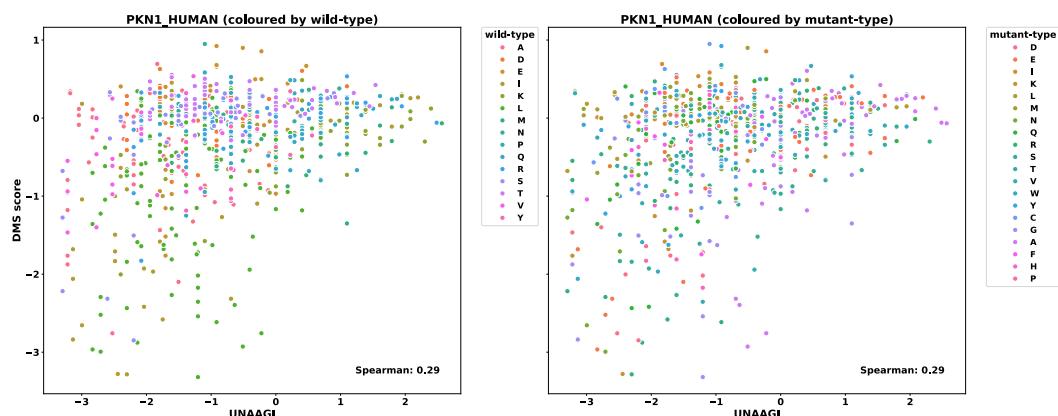
IF1_ECOLI

Figure 12: Scatterplot of UNAAGI vs. DMS score

ILF3_HUMAN

**OTU7A_HUMAN****PKN1_HUMAN****RS15_GEOSE**

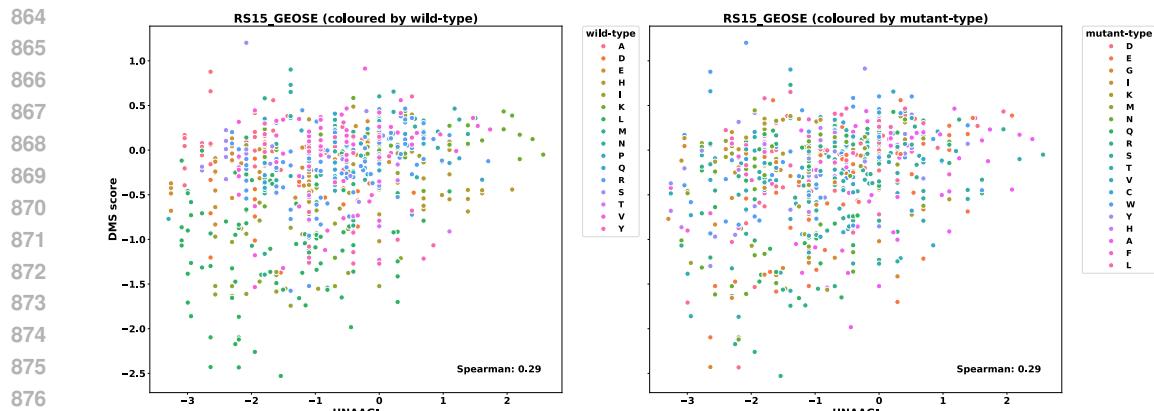


Figure 16: Scatterplot of UNAAGI vs. DMS score

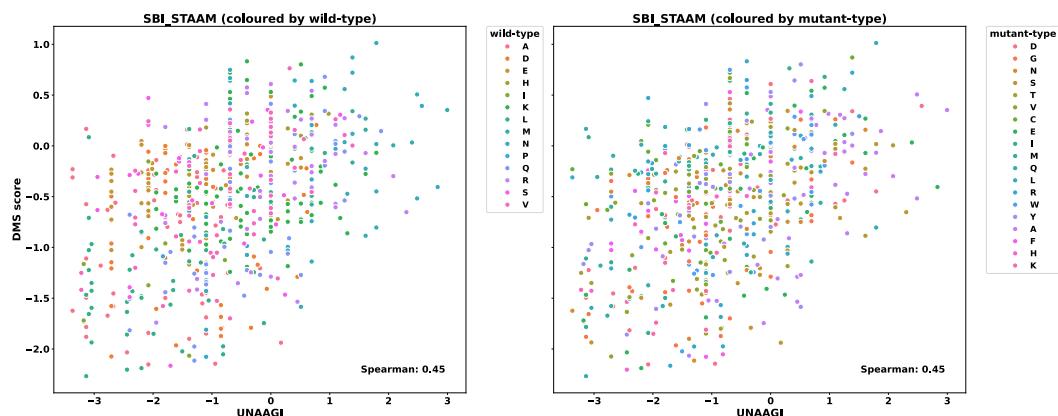
SBI_STAAM

Figure 17: Scatterplot of UNAAGI vs. DMS score

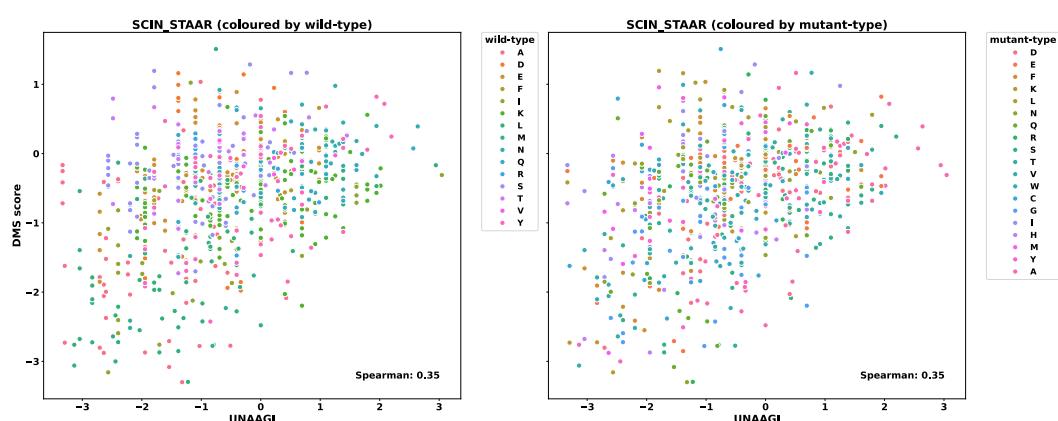
SCIN_STAAR

Figure 18: Scatterplot of UNAAGI vs. DMS score

SOX30_HUMAN

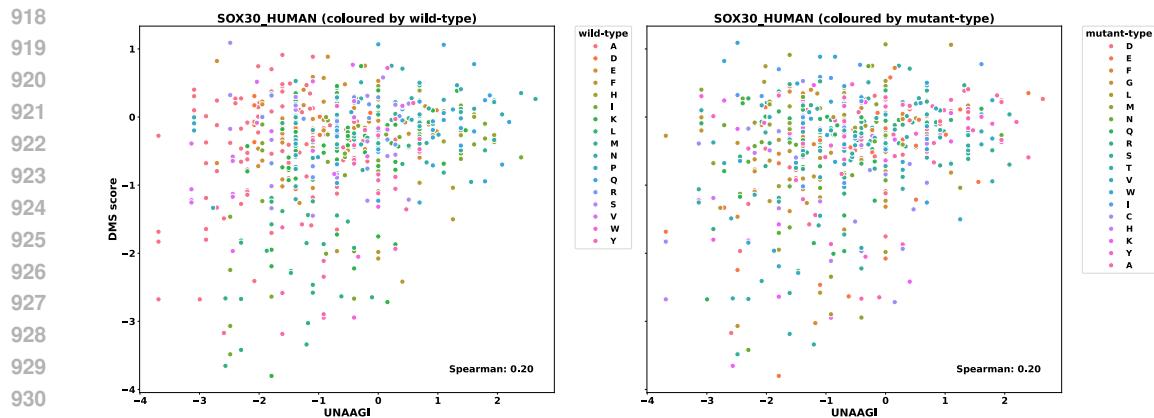


Figure 19: Scatterplot of UNAAGI vs. DMS score

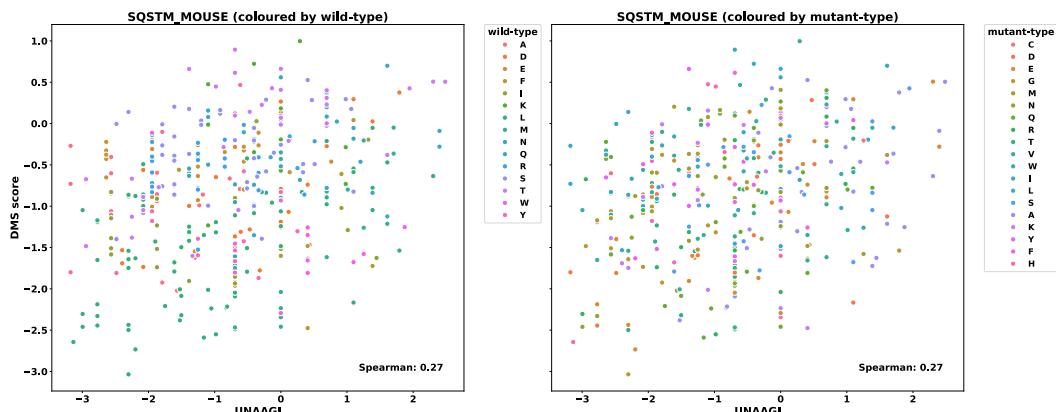
SQSTM_MOUSE

Figure 20: Scatterplot of UNAAGI vs. DMS score

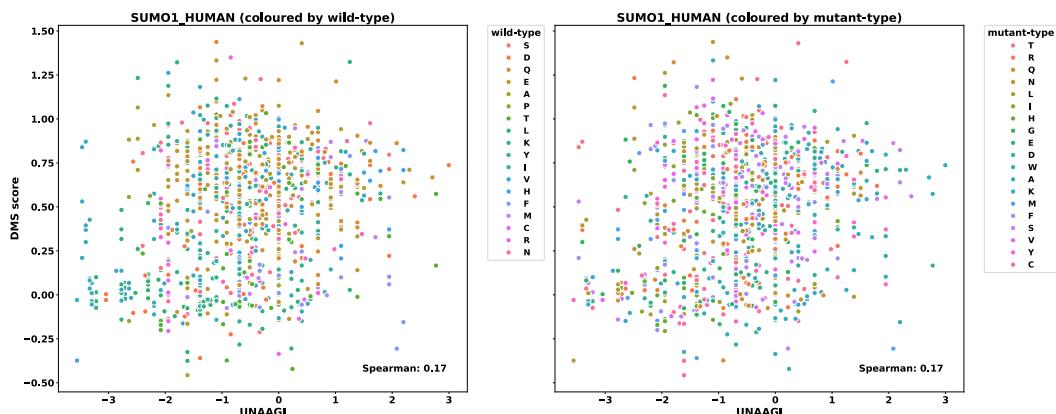
SUMO1_HUMAN

Figure 21: Scatterplot of UNAAGI vs. DMS score

TAT_HV1BR

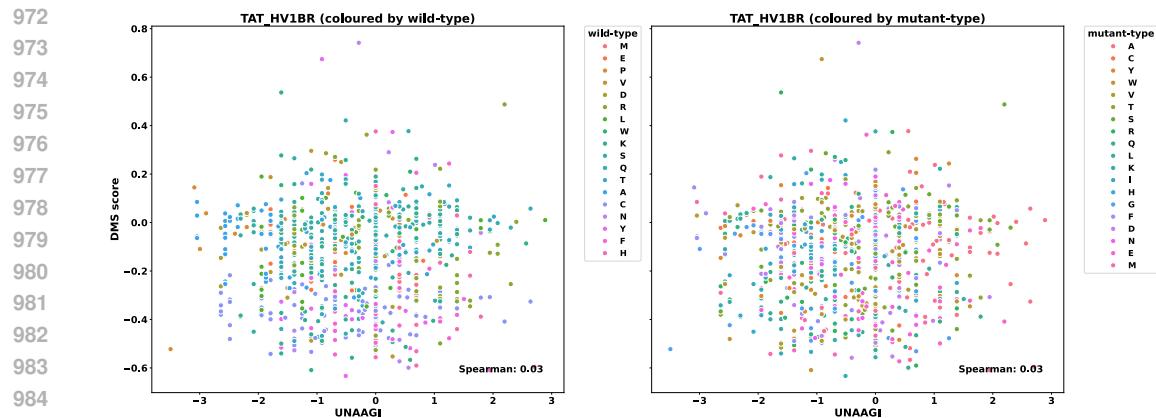


Figure 22: Scatterplot of UNAAGI vs. DMS score

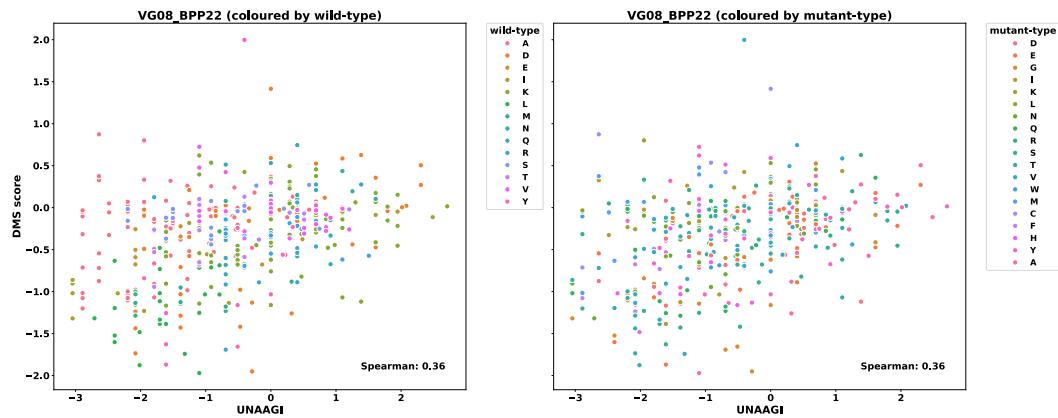
VG08_BPP22

Figure 23: Scatterplot of UNAAGI vs. DMS score

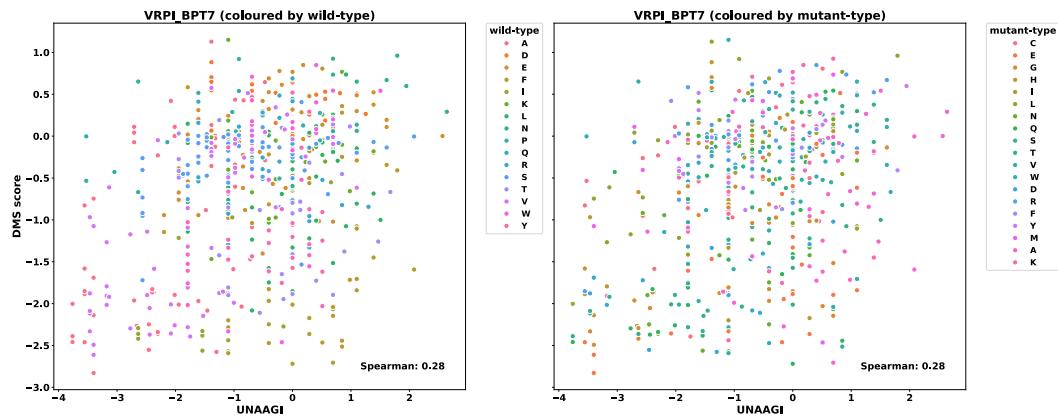
VRPI_BPT7

Figure 24: Scatterplot of UNAAGI vs. DMS score

1026 A.1.2 DETAILED RESULTS ON DMS OF NCAAs
1027

1028 For results on DMS of NCAAs, we report the spearman correlation over all the mutants, and for the
1029 NCAA mutants specifically, as depicted in Figure 4.

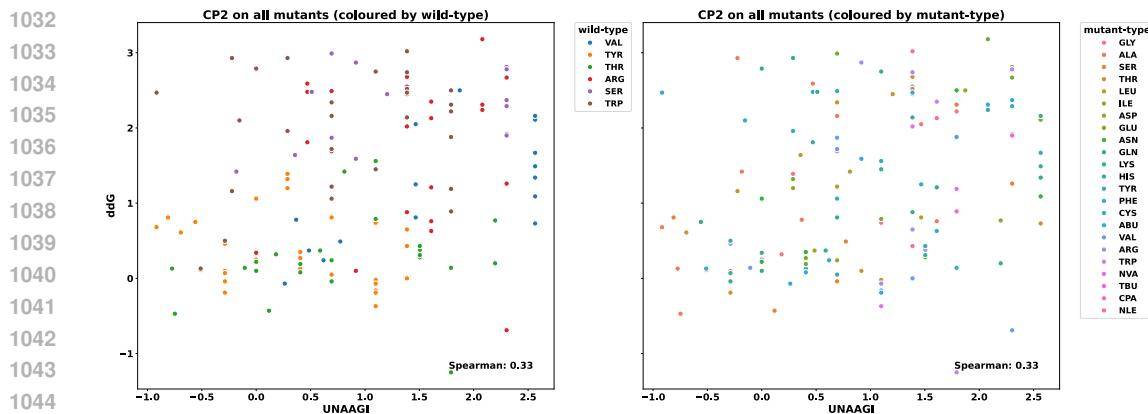
1030 **CP2**
1031

Figure 25: Scatterplot of UNAAGI vs. ddG, for all mutants

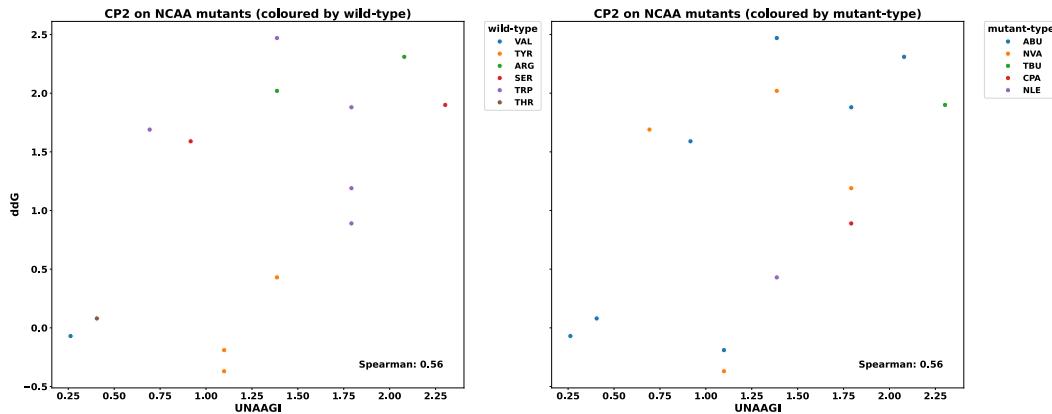


Figure 26: Scatterplot of UNAAGI vs. ddG, for NCAA mutants

1065 **PUMA**

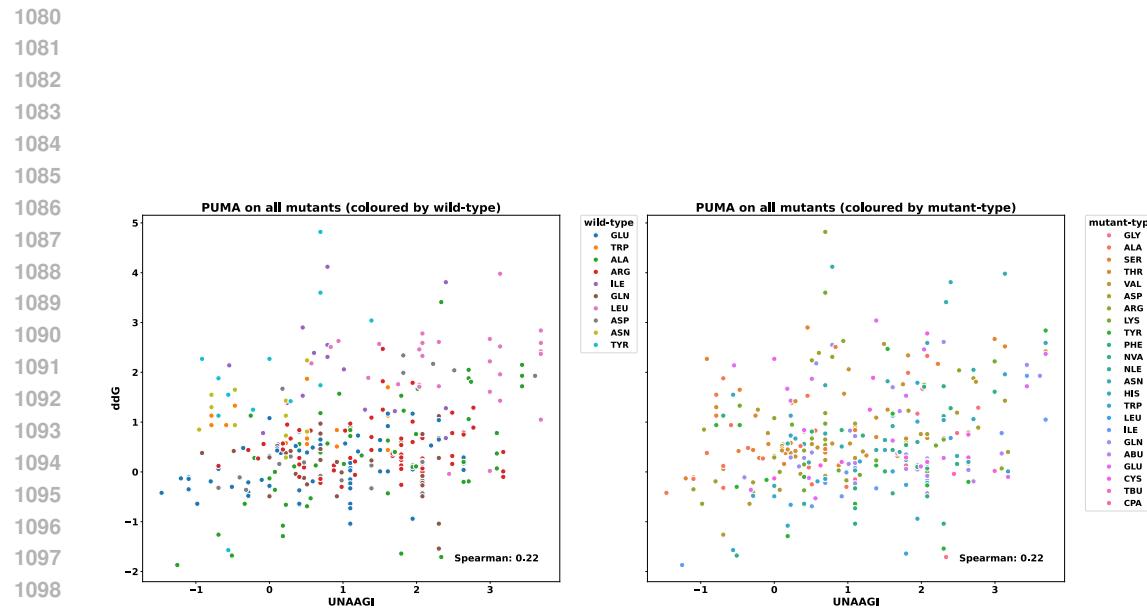


Figure 27: Scatterplot of UNAAGI vs. ddG, for all mutants

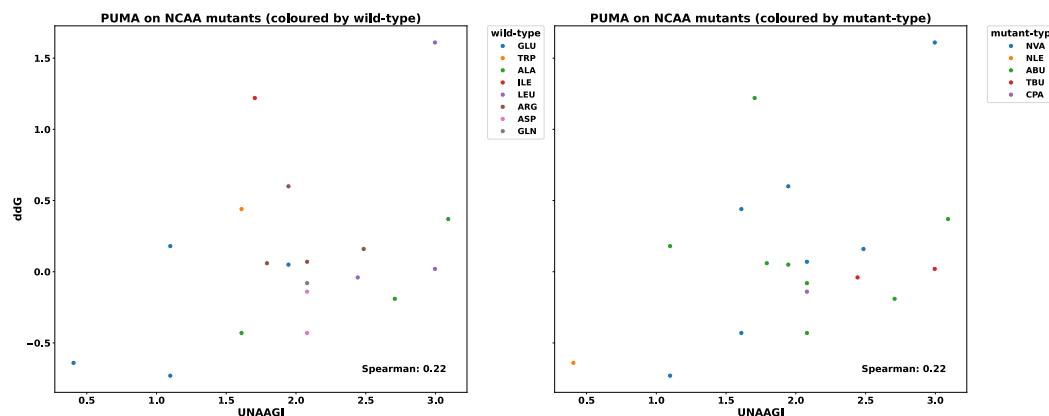


Figure 28: Scatterplot of UNAAGI vs. ddG, for NCAA mutants

1134
1135

A.2 THE USE OF LARGE LANGUAGE MODELS (LLMs)

1136
1137
1138
1139

In accordance with ICLR requirements, we disclose the role of LLMs in preparing this work. Large language models were not used for any experimental design, data analysis, or implementation. Their use was limited to assisting with minor language editing, such as correcting grammar and polishing phrasing.

1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187