

EVALUATING REPRESENTATION LEARNING ON THE PROTEIN STRUCTURE UNIVERSE

Arian R. Jamasb^{*,1,†}, Alex Morehead^{*,2}, Chaitanya K. Joshi^{*,1}, Zuobai Zhang^{*,3}, Kieran Didi¹, Simon Mathis¹, Charles Harris¹, Jian Tang³, Jianlin Cheng², Pietro Liò¹, Tom L. Blundell¹

¹University of Cambridge, ²University of Missouri, ³Mila - Québec AI Institute

ABSTRACT

We introduce *ProteinWorkshop*, a comprehensive benchmark suite for representation learning on protein structures with Geometric Graph Neural Networks. We consider large-scale pre-training and downstream tasks on both experimental and predicted structures to enable the systematic evaluation of the quality of the learned structural representation and their usefulness in capturing functional relationships for downstream tasks. We find that: (1) large-scale pretraining on AlphaFold structures and auxiliary tasks consistently improve the performance of both rotation-invariant and equivariant GNNs, and (2) more expressive equivariant GNNs benefit from pretraining to a greater extent compared to invariant models.

We aim to establish a common ground for the machine learning and computational biology communities to rigorously compare and advance protein structure representation learning. Our open-source codebase reduces the barrier to entry for working with large protein structure datasets by providing: (1) storage-efficient dataloaders for large-scale structural databases including AlphaFoldDB and ESM Atlas, as well as (2) utilities for constructing new tasks from the entire PDB. *ProteinWorkshop* is available at: github.com/a-r-j/ProteinWorkshop.

1 INTRODUCTION

Modern protein structure prediction methods have led to an explosion in the availability of structural data (Jumper et al., 2021; Baek et al., 2021). While many sequence-based functional annotations can be directly mapped to structures, this has resulted in a significantly-increasing gap between structures and meaningful *structural* annotations (Varadi et al., 2021). Recent work has focused on developing methods to draw biological insights from large volumes of structural data by either determining representative structures that can be used to provide such annotations (Holm, 2022) or representing structures in a simplified and compact manner such as sequence alphabets (van Kempen et al., 2023) or graph embeddings (Greener & Jamali, 2022). These works have significantly reduced the computational resources required to process and analyse such structural representatives at scale. Nonetheless, it remains to be shown how such results can help us better understand the relationship between protein sequence, structure, and function through the use of deep learning algorithms.

Several deep learning methods have been developed for protein structures. In particular, Geometric Graph Neural Networks (GNNs) (Duval et al., 2023) have emerged as the architecture of choice for learning structural representations of biomolecules (Schütt et al., 2018; Gasteiger et al., 2020; Jing et al., 2020; Schütt et al., 2021; Morehead et al., 2022; Zhang et al., 2023b). Methods can be categorised according to (1) the featurisation schemes and level of granularity of input structures ($C\alpha$, backbones, all-atom); as well as (2) the enforcement of physical symmetries and inductive biases (invariant or equivariant representations) (Joshi et al., 2023a). However, there remains a need for a robust, standardised benchmark to track the progress of new and established methods with greater granularity and relevance to downstream applications.

*Equal contribution

†Current affiliation: Prescient Design, Genentech

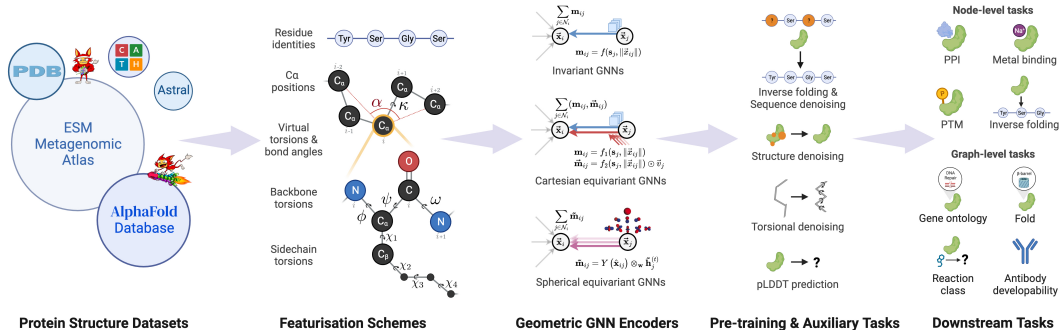


Figure 1: Overview of *ProteinWorkshop*, a comprehensive benchmark suite for evaluating pre-training and representation learning of Geometric GNNs on large-scale protein structure data.

In this work, we develop a unified and rigorous framework for evaluating protein structural encoders, providing pretraining corpora that span known foldspace and tasks that assess the ability of models to learn informative representations at different levels of structural granularity. Previous works in protein structure representation learning have focused on learning effective *global* (i.e. graph-level) representations of protein structure, typically evaluating the methods on function or fold classification tasks (Gligorijević et al., 2021; Zhang et al., 2023b). However, there has been comparatively little investigation into the ability of different methods to learn informative local (*node-level*) representations. Good node-level representations are important for a variety of annotation tasks, such as binding or interaction site prediction (Gainza et al., 2020), as well as providing conditioning signals in structure-conditioned molecule design methods (Schneuing et al., 2022; Corso et al., 2023). Understanding the structure-function relationship at this granular level can drive progress in protein design by revealing structural motifs that underlie desirable properties, enabling them to be incorporated into designs.

Our contributions are as follows:

- We curate numerous *structure-based* pretraining and fine-tuning datasets from the literature with a focus on tasks that can enable structural annotation of predicted structures. We compile a highly-modular benchmark, enabling the community to rapidly evaluate protein representation learning methods across tasks, models, and pretraining setups.
- We benchmark Geometric GNNs for representation learning of proteins at different levels of structural granularity ($C\alpha$, backbones, sidechain) and across several classes of models, ranging from general-purpose (Schütt et al., 2018; Satorras et al., 2021) to protein-specific architectures (Morehead & Cheng, 2024; Zhang et al., 2023b). We are the first to evaluate higher order equivariant GNNs (Thomas et al., 2018; Batatia et al., 2022) for proteins.
- We pretrain and evaluate models on, to our knowledge, the *largest non-redundant* protein structure corpus containing 2.27 million structures from AlphaFoldDB.
- Our benchmarks show that sequence and structure denoising-based auxiliary tasks and structure denoising-based pretraining consistently improve Geometric GNNs. Moreover, we surprisingly observe that sequence-based pretrained ESM-2-650M (Lin et al., 2022) augmented with our structural featurisation matches state-of-the-art GNNs on (super)family fold and gene ontology prediction.

2 PROTEINWORKSHOP

The overarching goal of *ProteinWorkshop* is to effectively cover the design space of protein structure representation learning methods. To achieve this, the benchmark is highly modular by design, enabling evaluation of different combinations of structural encoders, protein featurisation schemes, and auxiliary tasks over a wide range of both supervised and unsupervised tasks. A user manual is available in Appendix D, containing detailed listings and descriptions of all components.

2.1 FEATURISATION SCHEMES

Protein structures are typically represented as geometric graphs, with researchers opting to use a coarse-grained C_α atoms graph as full atom representations can quickly become computationally intractable due to a large number of nodes. However, this is a lossy representation, with much of the structural detail, such as orientation of the backbone and sidechain structure, being only implicitly encoded. Due to the computational burden incurred by operating on full-atom node representations, we focus primarily on C_α -based graph representations, investigating featurisation strategies to incorporate higher-level structural information. Note that we do provide utilities to enable users to work with backbone and full-atom graphs in the benchmark. Details about different featurisation schemes are provided in Appendix D.4 and Table 5.

2.2 PRE-TRAINING TASKS

The benchmark contains a comprehensive suite of pretraining tasks. Broadly, these can be categorised into: masked-attribute prediction, denoising-based and contrastive learning-based tasks. These can be used as both a pretraining objective or as auxiliary tasks in a downstream supervised task.

Sequence Denoising. The benchmark contains two variations based on two sequence corruption processes $C(\hat{S}|\mathcal{S}, \nu)$ that receive an amino acid sequence $\mathcal{S} \in [0, 1]^{|\nu| \times 23}$ and return a sequence $\hat{S} \in [0, 1]^{|\nu| \times 23}$ with fraction ν of its positions corrupted. The first scheme is based on mutating a fraction of the residues to a random amino acid and tasking the model with recovering the uncorrupted sequence. The second is a masked residue prediction task, where a fraction of the residues are altered to a mask value and the model is tasked to recover the uncorrupted sequence.

Structure Denoising. We provide two structure-based denoising tasks: coordinate denoising and torsional denoising. In the coordinate denoising task, noise is sampled from a normal or uniform distribution and scaled by noise factor, $\nu \in \mathbb{R}$, and applied to each of the atom coordinates in the structure to ensure structural features, such as backbone or sidechain torsion angles, are also corrupted. The model is then tasked with predicting either the per-node noise or the original uncorrupted coordinates. For torsional denoising, the noise is applied to the backbone torsion angles and Cartesian coordinates are recomputed using pNeRF (AlQuraishi, 2019) and the uncorrupted bond lengths and angles prior to feature computation. Similarly to the coordinate denoising task, the model is then tasked with predicting either the per-residue angular noise or the original dihedral angles.

Sequence-Structure Co-Denoising. This is a multitask formulation of the previously described structure and sequence denoising tasks, with separate output heads for denoising each modality.

Masked Attribute Prediction. We use inverse folding (Section 2.3.1) as a pretraining task. The benchmark additionally incorporates the distance, angle and dihedral angle masked-attribute prediction (Zhang et al., 2023b) as well as a backbone dihedral angle prediction task.

pLDDT Prediction. Structure prediction models typically provide per-residue pLDDT (predicted Local Distance Difference Test) scores as local confidence measures which have been shown to correlate well with disordered regions (Wilson et al., 2022). We formulate a node-level regression task on predicted structures, somewhat analogous to structure quality assessment, where the model is tasked with predicting the scaled per-residue pLDDT $y \in [0, 1]$ values.

2.3 DOWNSTREAM TASKS

We curate several structure-based and sequence-based datasets from the literature and existing benchmarks[†], summarised in Table 1. The tasks are selected to evaluate not only the *global* structure representational power of each method, but also to evaluate the ability of each method to learn informative *local* representations for residue-level prediction and annotation tasks.

The raw structures are, where possible and accounting for obsolescence, retrieved directly from the PDB (or another structural source) as several processed datasets used by the community discard full atomic coordinates in favour of retaining only C_α positions, making them unsuitable for in-depth

[†]To retain focus on *protein* representation learning, we deliberately exclude commonly-used tasks based on protein-small molecule interactions as it is hard to disentangle the effect of the small molecule representation and the potential for bias (Boyles et al., 2019).

Table 1: Overview of supervised tasks and datasets.

	Task	Dataset Origin	Structures	# Train	# Validation	# Test	Metric
Node-level	Inverse Folding	Ingraham et al. (2019)	Experimental	3.9 M	105 K	180 K	Perplexity
	PPI Site Prediction	Gainza et al. (2020)	Experimental	478 K	53 K	117 K	AUPRC
	Metal Bind Site Prediction		Experimental	1.1 M	13.7 K	29.8 K	Accuracy
	PTM Site Prediction	Yan et al. (2023)	Predicted	44 K	2.4 K	2.5 K	ROC-AUC
Graph-level	Fold Prediction	Hou et al. (2017)	Experimental	12.3 K	0.7 K	1.3/0.7/1.3 K	Accuracy
	Gene Ontology Prediction	Glgorijević et al. (2021)	Experimental	27.5 K	3.1 K	3.0 K	F_{\max}
	Reaction Class Prediction	Hermosilla et al. (2020)	Experimental	29.2 K	2.6 K	5.6 K	Accuracy
	Antibody Dev. Prediction	Huang et al. (2021)	Experimental	1.7 K	0.24 K	0.48 K	AUPRC

experimentation. This provides an entry point for users to apply a custom sequence of pre-processing steps, such as deprotonation or fixing missing regions which are common in experimental data.

2.3.1 NODE-LEVEL TASKS

Inverse Folding. Many generative methods for protein design produce backbone structures that require the design of an associated sequence. As a result, inverse folding is an important part of *de novo* design pipelines for proteins (Dauparas et al., 2022). Formally, this is a node-level classification task where the model learns a mapping for each residue r_i to an amino acid type $y \in \{1, \dots, n\}$, where n is the vocabulary size ($n = 20$ for the canonical set of amino acids). Inverse folding is a generic task that can be applied to any dataset in the benchmark. In the literature, it is commonly evaluated on the CATH dataset (Section 2.4) compiled by Ingraham et al. (2019).

PPI Site Prediction. Identifying protein-protein interaction sites has important applications in developing refined protein-protein interaction networks and docking tools, providing biological context to guide protein engineering and target identification in drug discovery campaigns (Jamasb et al., 2021). This task is a node-level binary classification task where the goal is to predict whether or not a residue is involved in a protein-protein interaction interface. We use the dataset of experimental structures curated from the PDB by Gainza et al. (2020) and retain the original splits, though we modify the labelling scheme to be based on inter-atomic proximity (3.5 Å), which can be user-defined, rather than solvent exclusion. The dataset is curated from the PDB by preprocessing such as the presence of one of the seven specified ligands (e.g., ADP or FAD), clustering based on 30% sequence identity and random subsampling. It contains 1,459 structures, which are randomly assigned to training (72%), validation (8%) and test set (20%). 12 (Å) radius patches were extracted from the generated structures and a patch labelled as part of a binding pocket if its centre point was < 3 (Å) away from an atom of the corresponding ligand.

Metal Binding Site Prediction. Many proteins coordinate transition metal ions to carry out their functions. Predicting the binding sites of metal ions can elucidate the role of metal binding on protein function. This is a binary node classification task where each residue is mapped to a label $y \in \{0, 1\}$ indicating whether the residue (or its constituent atoms) is within 3.5 (Å) of a user-defined metal ion or ligand heteroatom, respectively. We provide tooling to curate a dataset of experimental structures from the PDB for this task, where binding site assignments for each residue are computed on-the-fly. While the benchmark supports this task on arbitrary subsets of the PDB and ligands, we provide the Zinc-binding dataset from Dürr et al. (2023) specifically for this task. The dataset is constructed by sequence-based clustering of the PDB at 30% sequence identity to remove sequence and structural redundancy. Clusters with a member shorter than 3000 residues, containing at least one zinc atom with resolution better than 2.5 (Å) determined by x-ray crystallography and not containing nucleic acids are used to compose the dataset. If multiple structures fulfil these criteria, the highest resolution structure is used. The train (2,085) / validation (26) / test (59) splits are constructed such that proteins in the validation and test sets have no partial overlap with any protein in the training data.

Post-Translational Modification Site Prediction. Identifying the precise sites where post-translational modifications (PTMs) occur is essential for understanding protein behaviour and designing targeted therapeutic interventions. We frame prediction of PTM sites as a multilabel classification task where each residue is mapped to a label $y \in \{1, \dots, 13\}$ distinguishing between modifications on different amino acids (e.g. phosphorylation on S/T/Y and N-linked glycosylation on N). We use a dataset of 48,811 AlphaFold2-predicted structures curated by Yan et al. (2023),

where each structure contains the PTM metadata necessary to construct residue-wise site prediction labels. The dataset is split into training (43,907, validation (2,393) and test (2,511) sets based on 50% sequence identity and 80% coverage.

2.3.2 GRAPH-LEVEL TASKS

Fold Prediction. The utility of this task is that it serves as a litmus test for the ability of a model to distinguish different structural folds. It stands to reason that models that perform poorly on distinguishing fold classes likely learn limited or low-quality structural representations. This is a multiclass graph classification task where each protein, \mathcal{G} , is mapped to a label $y \in \{1, \dots, 1195\}$ denoting the fold class. We adopt the fold classification dataset originally curated from SCOP 1.75 by (Hou et al., 2017). This provides three different test sets stratified based on topological similarity: Fold, in which proteins originating from the same superfamily are absent during training; Superfamily, in which proteins originating from the same family are absent during training; and Family, in which proteins from the same family are present during training.

Gene Ontology Prediction. Predicting protein function in the form of functional annotations such as GO terms has important applications in protein analysis and engineering, providing researchers with the ability to cluster functionally-related structures or to guide protein generation methods to design new proteins with desired functional properties. This is a multilabel classification task, assigning functional Gene Ontology (GO) annotation to structures. GO annotations are assigned within three ontologies: biological process (BP), cellular component (CC) and molecular function (MF). We use the dataset of experimental structures curated from the PDB by Gligorijević et al. (2021) and retain the original multi-cutoff based splits, with cutoff at 30% sequence similarity.

Reaction Class Prediction. As proteins’ reaction classifications are based on their enzyme-catalyzed reaction according to all four levels of the standard Enzyme Commission (EC) number, methods that predict such classifications may help elucidate the function of newly-designed proteins as they are developed. This is a multiclass graph classification task where each protein, \mathcal{G} , is mapped to a label $y \in \{1, \dots, 384\}$ denoting which class of reactions a given protein catalyzes; all four levels of the EC assignment are employed to define the reaction class label. We adopt the reaction class prediction dataset originally curated from the PDB by Hermosilla et al. (2020), split on the basis of sequence similarity using a 50% threshold.

Antibody Developability Prediction. Therapeutic antibodies must be optimised for favourable physicochemical properties in addition to target binding affinity and specificity to be viable development candidates. Consequently, we frame prediction of antibody developability as a binary graph classification task indicating whether a given antibody is developable. We adopt the antibody developability dataset originally curated from SabDab (Dunbar et al., 2014) by Chen et al. (2020). This dataset contains 2,426 antibodies that have both sequences and PDB structures available, where each example contains both a heavy chain and a light chain with resolution < 3 (Å). The label is based on thresholding the developability index (DI) (Lauer et al., 2012) as computed by BIOVIA’s platform (Systèmes, 2016), which relies on an antibody’s hydrophobic and electrostatic interactions. This task is interesting from a benchmarking perspective as it enables targeted performance assessment of models on a specific (immunoglobulin) fold, providing insight into whether general-purpose structure-based encoders can be applicable to fold-specific tasks.

2.4 PRE-TRAINING DATASETS

The benchmark contains several large corpora of both experimental and predicted structures that can be used for pretraining or inference. We provide utilities for configuring supervised tasks and splits directly from the PDB. Additionally, we build storage-efficient dataloaders for large pretraining corpora of predicted structures (AlphaFoldDB, ESM Atlas). We believe our codebase will considerably reduce the barrier to entry for working with large structure-based datasets.

2.4.1 EXPERIMENTAL STRUCTURES

PDB. We provide utilities for curating datasets directly from the Protein Data Bank (Berman, 2000). In addition to using the collection in its entirety, users can define filters to subset and split the data using a combination of structural similarity, sequence similarity or temporal strategies. Structures can

be filtered by length, number of chains, resolution, deposition date, presence/absence of particular ligands and structure determination method.

CATH. We provide the dataset derived from CATH 4.2 40% (Knudsen & Wiuf, 2010) non-redundant chains developed by Ingraham et al. (2019) as an additional, smaller, pretraining dataset.

ASTRAL. ASTRAL (Brenner, 2000) provides protein *domain* structures which are regions of proteins that can maintain their structure and function independently of the rest of the protein. Domains typically exhibit highly-specific functions and can be considered structural building blocks.

2.4.2 PREDICTED STRUCTURES

AlphaFoldDB Representative Structures. This dataset contains 2.27 million representative structures, identified through large-scale structural-similarity-based clustering of the 214 million structures contained in the AlphaFold Database (Varadi et al., 2021) using FoldSeek (van Kempen et al., 2023). We additionally provide a subset of this collection — the so-called dark proteome — corresponding to the 31% of the representative structures that lack annotations.

ESM Atlas, ESM High Quality. These datasets are compressed collections of predicted structures produced by ESMFold. ESM Atlas is the full collection of all 772m predicted structures for the MGnify 2023 release (Richardson et al., 2022). ESM High Quality is a curated subset of high confidence (mean pLDDT) structures from the collection.

3 METHODS AND EXPERIMENTAL SETUP

Overview. To demonstrate the utility of our benchmark, we investigate how combinations of protein structure representation, architecture choice and pretraining/auxiliary tasks affect predictive performance across a range of tasks. The tasks are selected to focus on important real-world structural annotation tasks and such that we can evaluate these combinations in terms of both the local and global representational power. To this end, we select state-of-the-art protein structure encoders and generic geometric GNN architectures that span the design space of geometric GNN models with regard to both message passing body order and tensor order (Joshi et al., 2023a). We evaluate several structural representations that, to varying degrees, capture the full detail of the protein structure.

Architectures. We provide a unified implementation of several rotation invariant and equivariant architectures. We benchmark 4 general purpose models: SchNet (Schütt et al., 2018), EGNN (Satorras et al., 2021), TFN (Thomas et al., 2018), MACE (Batatia et al., 2022); and 2 protein-specific architectures: GCPNet (Morehead & Cheng, 2024), GearNet (Zhang et al., 2023b). We also compare geometric GNNs to the pretrained sequence-based language model ESM (Lin et al., 2022) augmented with structural featurisation. We chose the 650M pretrained ESM-2 because this is the scale at which significant structure-related abilities were observed for ESM.

Featurisation Schemes. We consider five featurisation schemes, progressively increasing the amount of structural information provided to the model by incorporating sequence positional information, virtual dihedral and bond angles over the $C\alpha$ trace, backbone torsion angles, and sidechain torsion angles. Featurisation schemes are detailed in Table 5 in the Appendix.

Pretraining Dataset. For all pretraining experiments we use AlphaFoldDB (Barrio-Hernandez et al., 2023). This dataset provides a rich diversity of 2.27 million non-redundant protein structures and, to our knowledge, is substantially more diverse than any other previously used structure-based pretraining corpus, whilst remaining of a size that is amenable to experimentation. Models pretrained on AlphaFoldDB should, in principle, exhibit strong generalisation to the currently known (and predicted) natural protein structure universe as it would have ‘seen’ the same protein fold during pretraining. To facilitate working with large-scale AlphaFoldDB and ESM Atlas, we developed storage-efficient dataloaders based on FoldComp (Kim et al., 2023), described in Appendix D.6.

Pretraining and Auxiliary Tasks. In our evaluation, we focus predominantly on denoising-based pretraining and auxiliary tasks as these are comparatively less explored than contrastive or masked-attribute prediction tasks (Zhang et al., 2023b). We consider five pretraining tasks: (1) structure-based denoising, (2) sequence denoising, (3) torsional denoising, (4) inverse folding and (5) pLDDT prediction. Structure and sequence denoising are also used as auxiliary tasks in our experiments. We

also investigate an inverse folding pre-training task which we subsequently finetune on the CATH dataset for benchmarking inverse folding as a downstream task (see below).

Noising Schemes. For structure-based denoising we draw i.i.d. noise samples from a Gaussian distribution and scale by $\sigma = 0.1$ to corrupt the input coordinates or dihedral angles. Geometric scalar and vector-valued features are computed from the noised structure, *i.e.* $\tilde{\mathcal{G}} = (\mathcal{V}, \tilde{\mathcal{E}}, \tilde{\mathbf{X}}, \tilde{\mathbf{S}}, \tilde{\mathbf{V}})$, where $\tilde{\mathbf{x}}_i = \mathbf{x}_i + \sigma \boldsymbol{\epsilon}_i$ and $\boldsymbol{\epsilon}_i \sim \mathcal{N}(0, I_3)$. For sequence-based denoising, we use the mutation strategy and corrupt 25% of the residues in each protein. When sequence denoising is used as an auxiliary task, we weight the loss with a coefficient $\lambda = 0.1$, similar to NoisyNodes (Godwin et al., 2021).

Training. As we are interested in benchmarking large-scale datasets and models, we try to consistently use six layers for all models, each with 512 hidden channels. For equivariant GNNs, we reduced the number of layers and hidden channels to fit 80GB of GPU memory on one NVIDIA A100 GPU. For downstream tasks, we set the maximum number of epochs to 150 and use the Adam optimizer with a batch size of 32 and ReduceLROnPlateau learning rate scheduler, monitoring the validation metric with patience of 5 epochs and reduction of 0.6. See Appendix D.10 for details on hyperparameter tuning for optimal learning rates and dropouts for each architecture. We train models to convergence, monitoring the validation metric and performing early stopping with a patience of 10 epochs. Pretraining is performed for 10 epochs using a linear warm-up with cosine schedule. We report standard deviations over three runs across three random seeds.

4 RESULTS & DISCUSSIONS

4.1 AUXILIARY DENOISING CONSISTENTLY IMPROVES BASELINE PERFORMANCE

In Table 2, we first set out to determine the following questions about architectural choices in conjunction with denoising auxiliary tasks and *without* pretraining:

- Whether invariant or equivariant models perform better?** Across 10 tasks, equivariant models such as EGNN and GCPNet attain the best performance on 5. Notably, sequence-based ESM-2-650M augmented with our structural featurisation matches state-of-the-art protein-specific GNNs (Fan et al., 2023) on (super)family and gene ontology prediction.
- Which input representation is the best for each respective task?** Featurising models with $C\alpha$ atoms, virtual angles, and backbone torsions provides the best performance overall on 22 out of 60 combinations of models and tasks. This suggests that letting models implicitly learn about side chain orientation and flexibility by using backbone-only featurisation may prevent overfitting on crystallisation artifacts (Dauparas et al., 2022).
- Whether auxiliary denoising tasks improve model performance?** Both sequence and structure denoising are particularly useful auxiliary tasks for training protein structure encoders, until sufficient structural detail makes the tasks trivial, improving results over not using auxiliary tasks for 50 out of 60 combinations of models and primary tasks. Notably, structure denoising helped stabilise the training of MACE models on the GO and Reaction tasks, where other runs did not converge.

4.2 INCORPORATING MORE STRUCTURAL DETAIL IMPROVES PRE-TRAINING PERFORMANCE

We then investigated protein structure pre-training in Table 3 to determine:

- Which input representation is best for pre-training?** Incorporating greater structural detail with dihedral angles generally improves validation metrics on pre-training tasks, more so than architecture.
- Which GNNs benefit from which pre-training task?** Inverse folding, structure denoising, sequence denoising, and torsional denoising benefit equivariant models the most in the context of pre-training, whereas pLDDT prediction benefits invariant models the most, suggesting that certain pre-training tasks benefit certain classes of models more than other tasks. Unfortunately, we were currently unable to pre-train spherical equivariant GNNs (TFN, MACE) due to the high computational requirements of these models.

Table 3: **Validation performance for pretraining tasks on AlphaFoldDB.** Metrics: Inverse Folding: perplexity; pLDDT, Structure Denoising, Torsional Denoising: RMSE; Seq. Denoising: Accuracy. Best (second-best) results are **bolded** (underlined). **Key takeaway:** Incorporating **backbone structural features** (i.e., adding torsion angles ϕ, ψ, ω), in general, improves pretraining performance compared to using only **virtual angles** along the sequence.

Method	Task					
	Inverse Folding (\downarrow)	pLDDT Pred. (\downarrow)	Structure Denoising (\downarrow)	Seq. Denoising (\uparrow)	Torsional Denoising (\downarrow)	
$C\alpha + \kappa, \alpha$	SchNet	7.791	0.2397	0.0704	36.81	0.0586
	GearNet-Edge	6.596	0.2326	0.0672	43.76	0.0615
	EGNN	6.016	0.2406	0.0700	40.51	0.0586
	GCPNet	6.243	0.2395	0.0679	44.81	0.0562
$C\alpha + \phi, \psi, \omega$	SchNet	5.562	<u>0.2388</u>	0.0603	45.61	0.0489
	GearNet-Edge	<u>5.324</u>	0.2402	<u>0.0562</u>	50.15	0.0538
	EGNN	5.962	0.2403	0.0593	<u>53.80</u>	<u>0.0487</u>
	GCPNet	3.839	0.2399	0.0561	59.54	0.0443

4.3 PRE-TRAINING AND GREATER STRUCTURAL DETAIL BENEFIT DOWNSTREAM TASKS

Following the observation that more fine-grained input representations improve pretraining performance, Table 4 explores finetuning on downstream tasks:

- **Whether these lessons from pretraining translate to downstream tasks?** Equivariant GNNs outperform invariant GNNs in the majority of cases and generally benefit the most from pretraining on structure-based tasks, particularly when provided with greater structural detail in input features.
- **Which combination of parameters performs best on downstream tasks?** Overall, providing a greater amount of structural detail compared to a strict $C\alpha$ atom representation benefits downstream performance for *both* invariant and equivariant models. Notably, structure denoising generally improves downstream performance for *both* types of models.

Table 4: **Pretrained model benchmark results.** Results for each model and featurisation pair are given as: **no pretraining** / **sequence denoising** / **structure denoising**, except for inverse folding on CATH, which is pretrained with **inverse folding** on AFDB. **Key takeaway:** The equivariant GCPNet model benefits most from pretraining and maximum structural detail.

Method	Features	GO-BP (\uparrow)	GO-MF (\uparrow)	GO-CC (\uparrow)	Fold (\uparrow)			Reaction (\uparrow)	Inverse Folding (\downarrow)
					Fold	Family	Superfamily		
GearNet	$C\alpha + \kappa, \alpha$	0.393 / 0.342 / 0.376	0.476 / 0.481 / 0.490	0.436 / 0.446 / 0.457	32.79 / 29.33 / 34.00	95.35 / 91.72 / 96.37	47.56 / 43.40 / 49.63	77.45 / 78.64 / 79.37	12.35 / 7.84
	$C\alpha + \kappa, \alpha, \phi, \psi, \omega$	0.397 / 0.378 / 0.392	0.480 / 0.479 / 0.497	0.441 / 0.445 / 0.457	33.75 / 31.20 / 36.47	94.35 / 94.58 / 94.02	46.60 / 45.83 / 48.23	76.61 / 77.22 / 80.31	11.61 / 7.29
GCPNet	$C\alpha + \kappa, \alpha$	0.364 / 0.358 / 0.336	0.465 / 0.484 / 0.451	0.427 / 0.414 / 0.404	36.97 / 37.57 / 41.81	95.65 / 96.82 / 96.51	47.35 / 53.45 / 54.95	76.46 / 78.89 / 78.84	8.80 / 7.37
	$C\alpha + \kappa, \alpha, \phi, \psi, \omega$	0.362 / 0.348 / 0.334	0.466 / 0.501 / 0.502	0.424 / 0.409 / 0.404	38.34 / 41.14 / 42.81	95.94 / 96.09 / 96.61	49.81 / 52.60 / 57.17	75.49 / 77.97 / 79.18	7.56 / 6.55

5 CONCLUSIONS

This work focuses on building a comprehensive and multi-task benchmark for protein structure representation learning. *ProteinWorkshop* provides a unified implementation of large pretraining corpora, featurisation schemes, Geometric GNN models and benchmarking tasks to evaluate the effectiveness of protein structure encoding methods. Key findings include that structural pretraining, as well as auxiliary sequence and structure denoising tasks, improve performance on a wide range of downstream tasks and that incorporating more structural detail in featurisation improves performance. Our benchmark is flexible for including new tasks and datasets and is open to the wider research community.

Availability. The *ProteinWorkshop* codebase is available under a permissive MIT License at github.com/a-r-j/ProteinWorkshop and accompanying documentation, preprocessed datasets and pretrained model weights are hosted publicly at proteins.sh. Preprocessed datasets and pretrained model weights are deposited on Zenodo at the following URLs, respectively: zenodo.org/record/8282470 and zenodo.org/record/8287754.

ACKNOWLEDGEMENTS

We acknowledge that this work was supported by a variety of institutions. ARJ was funded by a Biotechnology and Biological Sciences Research Council (BBSRC) DTP studentship (BB/M011194/1). AM and JC were supported by a U.S. NSF grant (DBI2308699) and two U.S. NIH grants (R01GM093123 and R01GM146340). CKJ was supported by the A*STAR Singapore National Science Scholarship (PhD). This work was performed using resources provided by the Cambridge Service for Data Driven Discovery (CSD3) operated by the University of Cambridge Research Computing Service (www.csd3.cam.ac.uk), provided by Dell EMC and Intel using Tier-2 funding from the Engineering and Physical Sciences Research Council (capital grant EP/T022159/1), and DiRAC funding from the Science and Technology Facilities Council (www.dirac.ac.uk). Additionally, this work was performed using high performance computing infrastructure provided by Research Support Services at the University of Missouri-Columbia (DOI: 10.32469/10355/97710). We thank Martin Steinegger and Do-Yoon Kim for allowing us to use the illustrations in Figure 1.

BROADER IMPACTS

Our benchmark unifies protein representation learning tasks, large-scale pre-training datasets, featurisation schemes, and models. The wide range of tasks studied in our benchmark can enable us to develop insight into effective pre-training strategies, and whether pre-trained protein structural representations can have material impact in real-world computational biology and drug discovery. It is not lost on us that these models can play a role in developing harmful chemical matter in the hands of a bad actor. Additionally, training very large models can contribute to climate change. However, we hope that developing highly effective structural representations will have broad, positive implications across biology and medicine that significantly outweigh the potential for misuse.

REFERENCES

- Mohammed AlQuraishi. Parallelized natural extension reference frame: Parallelized conversion from internal to cartesian coordinates. *Journal of Computational Chemistry*, 2019. (Cited on page 3, 24)
- Minkyung Baek, Frank DiMaio, Ivan Anishchenko, Justas Dauparas, et al. Accurate prediction of protein structures and interactions using a three-track neural network. *Science*, 2021. (Cited on page 1, 23)
- Inigo Barrio-Hernandez, Jingsi Yeo, Jürgen Jänes, Tanita Wein, Mihaly Varadi, Sameer Velankar, Pedro Beltrao, and Martin Steinegger. Clustering predicted structures at the scale of the known protein universe. *Bioinformatics*, 2023. doi: 10.1101/2023.03.09.531927. (Cited on page 6)
- Ilyes Batatia, Dávid Péter Kovács, Gregor NC Simm, Christoph Ortner, and Gábor Csányi. Mace: Higher order equivariant message passing neural networks for fast and accurate force fields. In *NeurIPS*, 2022. (Cited on page 2, 6, 23)
- H. M. Berman. The protein data bank. *Nucleic Acids Research*, 2000. doi: 10.1093/nar/28.1.235. (Cited on page 5, 23)
- Fergus Boyles, Charlotte M Deane, and Garrett M Morris. Learning from the ligand: using ligand-based features to improve binding affinity prediction. *Bioinformatics*, 2019. (Cited on page 3, 25)
- Anthony R. Bradley, Alexander S. Rose, Antonín Pavelka, Yana Valasatava, Jose M. Duarte, Andreas Prlić, and Peter W. Rose. MMTF—an efficient file format for the transmission, visualization, and analysis of macromolecular structures. *PLOS Computational Biology*, 2017. (Cited on page 23)
- S. E. Brenner. The ASTRAL compendium for protein structure and sequence analysis. *Nucleic Acids Research*, 28(1):254–256, January 2000. (Cited on page 6, 24)
- Martin Buttenschoen, Garrett M Morris, and Charlotte M Deane. Posebusters: Ai-based docking methods fail to generate physically valid poses or generalise to novel sequences. *arXiv preprint arXiv:2308.05777*, 2023. (Cited on page 20)

- Henriette Capel, Robin Weiler, Maurits Dijkstra, Reinier Vleugels, Peter Bloem, and K. Anton Feenstra. ProteinGLUE multi-task benchmark suite for self-supervised protein modeling. *Scientific Reports*, 12(1), September 2022. (Cited on page 19)
- Chen Chen, Xiao Chen, Alex Morehead, Tianqi Wu, and Jianlin Cheng. 3d-equivariant graph neural networks for protein model quality assessment. *Bioinformatics*, 39(1):btad030, 2023. (Cited on page 19)
- Xingyao Chen, Thomas Dougherty, Chan Hong, Rachel Schibler, Yi Cong Zhao, Reza Sadeghi, Naim Matasci, Yi-Chieh Wu, and Ian Kerman. Predicting antibody developability from sequence using machine learning. *bioRxiv*, June 2020. (Cited on page 5, 30)
- Gabriele Corso. Modeling molecular structures with intrinsic diffusion models. *arXiv preprint arXiv:2302.12255*, 2023. (Cited on page 19)
- Gabriele Corso, Hannes Stärk, Bowen Jing, Regina Barzilay, and Tommi Jaakkola. Diffdock: Diffusion steps, twists, and turns for molecular docking. In *ICLR*, 2023. (Cited on page 2, 23)
- Bowen Dai and Chris Bailey-Kellogg. Protein interaction interface region prediction by geometric deep learning. *Bioinformatics*, 37(17):2580–2588, 2021. (Cited on page 19)
- Christian Dallago, Jody Mou, Jody Mou, Kadina Johnston, Bruce Wittmann, Nicholas Bhattacharya, Samuel Goldman, Ali Madani, and Kevin Yang. Flip: Benchmark tasks in fitness landscape inference for proteins. In *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks*, volume 1. Curran, 2021. (Cited on page 19)
- J. Dauparas, I. Anishchenko, N. Bennett, H. Bai, R. J. Ragotte, L. F. Milles, B. I. M. Wicky, A. Courbet, R. J. de Haas, N. Bethel, P. J. Y. Leung, T. F. Huddy, S. Pellock, D. Tischer, F. Chan, B. Koepnick, H. Nguyen, A. Kang, B. Sankaran, A. K. Bera, N. P. King, and D. Baker. Robust deep learning-based protein sequence design using ProteinMPNN. *Science*, 2022. (Cited on page 4, 7, 19, 20)
- DeepMind-Isomorphic. Performance and structural coverage of the latest, in-development alphafold model. *DeepMind*, 2023. (Cited on page 20)
- James Dunbar, Konrad Krawczyk, Jinwoo Leem, Terry Baker, Angelika Fuchs, Guy Georges, Jiye Shi, and Charlotte M Deane. Sabdab: the structural antibody database. *Nucleic acids research*, 42(D1):D1140–D1146, 2014. (Cited on page 5, 30)
- Simon L. Dürr, Andrea Levy, and Ursula Rothlisberger. Metal3d: a general deep learning framework for accurate metal ion location prediction in proteins. *Nature Communications*, 14(1), May 2023. (Cited on page 4, 28)
- Alexandre Duval, Simon V Mathis, Chaitanya K Joshi, Victor Schmidt, Santiago Miret, Fragkiskos D Malliaros, Taco Cohen, Pietro Liò, Yoshua Bengio, and Michael Bronstein. A hitchhiker’s guide to geometric gnns for 3d atomic systems. *arXiv preprint*, 2023. (Cited on page 1, 21)
- Vijay Prakash Dwivedi, Chaitanya K Joshi, Anh Tuan Luu, Thomas Laurent, Yoshua Bengio, and Xavier Bresson. Benchmarking graph neural networks. *Journal of Machine Learning Research*, 2023. (Cited on page 19)
- Stephan Eismann, Patricia Suriana, Bowen Jing, Raphael JL Townshend, and Ron O Dror. Protein model quality assessment using rotation-equivariant, hierarchical neural networks. *arXiv preprint arXiv:2011.13557*, 2020. (Cited on page 19)
- Sara El-Gebali, Jaina Mistry, Alex Bateman, Sean R Eddy, Aurélien Luciani, Simon C Potter, Matloob Qureshi, Lorna J Richardson, Gustavo A Salazar, Alfredo Smart, Erik L L Sonnhammer, Layla Hirsh, Lisanna Paladin, Damiano Piovesan, Silvio C E Tosatto, and Robert D Finn. The pfam protein families database in 2019. *Nucleic Acids Research*, 47(D1):D427–D432, October 2018. (Cited on page 19)
- William A Falcon. Pytorch lightning. *GitHub*, 2019. (Cited on page 20)

- Hehe Fan, Zhangyang Wang, Yi Yang, and Mohan Kankanhalli. Continuous-discrete convolution for geometry-sequence modeling in proteins. In *The Eleventh International Conference on Learning Representations*, 2023. (Cited on page 7)
- Matthias Fey and Jan Eric Lenssen. Fast graph representation learning with pytorch geometric. *arXiv preprint*, 2019. (Cited on page 20)
- Pablo Gainza, Freyr Sverrisson, Frederico Monti, Emanuele Rodola, D Boscaini, MM Bronstein, and BE Correia. Deciphering interaction fingerprints from protein molecular surfaces using geometric deep learning. *Nature Methods*, 17(2):184–192, 2020. (Cited on page 2, 4, 19, 26, 27)
- Johannes Gasteiger, Janek Groß, and Stephan Günnemann. Directional message passing for molecular graphs. In *ICLR*, 2020. (Cited on page 1, 21)
- Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach, Hal Daumé Iii, and Kate Crawford. Datasheets for datasets. *Communications of the ACM*, 2021. (Cited on page 25)
- Mario Geiger and Tess Smidt. e3nn: Euclidean neural networks. *arXiv preprint arXiv:2207.09453*, 2022. (Cited on page 20)
- Vladimir Gligorijević, P. Douglas Renfrew, Tomasz Kosciolk, Julia Koehler Leman, Daniel Berenberg, Tommi Vatanen, Chris Chandler, Bryn C. Taylor, Ian M. Fisk, Hera Vlamakis, Ramnik J. Xavier, Rob Knight, Kyunghyun Cho, and Richard Bonneau. Structure-based protein function prediction using graph convolutional networks. *Nature Communications*, 2021. (Cited on page 2, 4, 5, 19, 26, 29)
- Jonathan Godwin, Michael Schaarschmidt, Alexander Gaunt, Alvaro Sanchez-Gonzalez, Yulia Rubanova, Petar Veličković, James Kirkpatrick, and Peter Battaglia. Simple gnn regularisation for 3d molecular property prediction & beyond. *arXiv preprint*, 2021. (Cited on page 7, 19)
- Joe G Greener and Kiarash Jamali. Fast protein structure searching using structure graph embeddings. *bioRxiv*, 2022. (Cited on page 1)
- Charles Harris, Kieran Didi, Arian Jamasb, Chaitanya Joshi, Simon Mathis, Pietro Lio, and Tom Blundell. Posecheck: Generative models for 3d structure-based drug design produce unrealistic poses. In *NeurIPS 2023 Generative AI and Biology (GenBio) Workshop*, 2023. (Cited on page 20)
- Pedro Hermosilla, Marco Schäfer, Matěj Lang, Gloria Fackelmann, Pere Pau Vázquez, Barbora Kozlíková, Michael Krone, Tobias Ritschel, and Timo Ropinski. Intrinsic-extrinsic convolution and pooling for learning on 3d protein structures. *arXiv preprint arXiv:2007.06252*, 2020. (Cited on page 4, 5, 19, 26, 30)
- Liisa Holm. Dali server: structural unification of protein families. *Nucleic acids research*, 2022. (Cited on page 1)
- Jie Hou, Badri Adhikari, and Jianlin Cheng. DeepSF: deep convolutional neural network for mapping protein sequences to folds. *Bioinformatics*, 2017. (Cited on page 4, 5, 26, 29)
- Kexin Huang, Tianfan Fu, Wenhao Gao, Yue Zhao, Yusuf Roohani, Jure Leskovec, Connor Coley, Cao Xiao, Jimeng Sun, and Marinka Zitnik. Therapeutics data commons: Machine learning datasets and tasks for drug discovery and development. In *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks*, volume 1. Curran, 2021. (Cited on page 4, 19, 26)
- Yufei Huang, Lirong Wu, Haitao Lin, Jiangbin Zheng, Ge Wang, and Stan Z Li. Data-efficient protein 3d geometric pretraining via refinement of diffused protein structure decoy. *arXiv preprint arXiv:2302.10888*, 2023. (Cited on page 19)
- John Ingraham, Vikas Garg, Regina Barzilay, and Tommi Jaakkola. Generative models for graph-based protein design. In *NeurIPS*, 2019. (Cited on page 4, 6, 23, 26, 27)

- Arian R. Jamasb, Ben Day, Cătălina Cangea, Pietro Liò, and Tom L. Blundell. Deep learning for protein–protein interaction site prediction. In *Methods in Molecular Biology*, pp. 263–288. Springer US, 2021. (Cited on page [4](#), [27](#))
- Arian Rokkum Jamasb, Ramon Viñas Torné, Eric J Ma, Yuanqi Du, Charles Harris, Kexin Huang, Dominic Hall, Pietro Lio, and Tom Leon Blundell. Graphein - a python library for geometric deep learning and network analysis on biomolecular structures and interaction networks. In *NeurIPS*, 2022. (Cited on page [20](#))
- Wengong Jin, Jeremy Wohlwend, Regina Barzilay, and Tommi Jaakkola. Iterative refinement graph neural network for antibody sequence-structure co-design. *arXiv preprint arXiv:2110.04624*, 2021. (Cited on page [19](#))
- Bowen Jing, Stephan Eismann, Patricia Suriana, Raphael John Lamarre Townshend, and Ron Dror. Learning from protein structure with geometric vector perceptrons. In *ICLR*, 2020. (Cited on page [1](#), [19](#))
- Chaitanya K. Joshi, Cristian Bodnar, Simon V. Mathis, Taco Cohen, and Pietro Liò. On the expressive power of geometric graph neural networks. In *International Conference on Machine Learning*, 2023a. (Cited on page [1](#), [6](#), [21](#))
- Chaitanya K Joshi, Arian R Jamasb, Ramon Viñas, Charles Harris, Simon V Mathis, Alex Morehead, and Pietro Liò. grnade: Geometric deep learning for 3d rna inverse design. *arXiv preprint*, 2023b. (Cited on page [20](#))
- John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, et al. Highly accurate protein structure prediction with alphafold. *Nature*, 2021. (Cited on page [1](#), [24](#))
- Hyunbin Kim, Milot Mirdita, and Martin Steinegger. Foldcomp: a library and format for compressing and indexing large protein structure sets. *Bioinformatics*, 39(4), March 2023. (Cited on page [6](#), [24](#))
- Michael Knudsen and Carsten Wiuf. The CATH database. *Human Genomics*, 4(3):207, 2010. (Cited on page [6](#), [23](#))
- Rohith Krishna, Jue Wang, Woody Ahern, Pascal Sturmfels, Preetham Venkatesh, Indrek Kalvet, Gyu Rie Lee, Felix S Morey-Burrows, Ivan Anishchenko, Ian R Humphreys, et al. Generalized biomolecular modeling and design with rosettafold all-atom. *bioRxiv*, 2023. (Cited on page [20](#))
- Timothy M. Lauer, Neeraj J. Agrawal, Naresh Chennamsetty, Kamal Egodage, Bernhard Helk, and Bernhardt L. Trout. Developability index: A rapid in silico tool for the screening of antibody aggregation propensity. *Journal of Pharmaceutical Sciences*, 101(1):102–115, January 2012. (Cited on page [5](#), [30](#))
- Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Nikita Smetanin, Allan dos Santos Costa, Maryam Fazel-Zarandi, Tom Sercu, Sal Candido, et al. Language models of protein sequences at the scale of evolution enable accurate structure prediction. *bioRxiv*, 2022. (Cited on page [2](#), [6](#), [24](#))
- Shengchao Liu, Weitao Du, Yanjing Li, Zhuoxinran Li, Zhiling Zheng, Chenru Duan, Zhiming Ma, Omar Yaghi, Anima Anandkumar, Christian Borgs, et al. Symmetry-informed geometric representation for molecules, proteins, and crystalline materials. *arXiv preprint arXiv:2306.09375*, 2023a. (Cited on page [20](#))
- Shengchao Liu, Hongyu Guo, and Jian Tang. Molecular geometry pretraining with SE(3)-invariant denoising distance matching. In *The Eleventh International Conference on Learning Representations*, 2023b. (Cited on page [20](#))
- Sajid Mahmud, Alex Morehead, and Jianlin Cheng. Accurate prediction of protein tertiary structural changes induced by single-site mutations with equivariant graph neural networks. *bioRxiv*, pp. 2023–10, 2023. (Cited on page [19](#))
- Alex Morehead and Jianlin Cheng. Geometry-complete diffusion for 3d molecule generation. In *ICLR Machine Learning for Drug Discovery 2023*, 2023. (Cited on page [22](#))

- Alex Morehead and Jianlin Cheng. Geometry-complete perceptron networks for 3d molecular graphs. *Bioinformatics*, 2024. (Cited on page 2, 6, 19, 22)
- Alex Morehead, Chen Chen, and Jianlin Cheng. Geometric transformers for protein interface contact prediction. In *International Conference on Learning Representations*, 2022. (Cited on page 1)
- Alex Morehead, Chen Chen, Ada Sedova, and Jianlin Cheng. Dips-plus: The enhanced database of interacting protein structures for interface prediction. *Scientific Data (Nature)*, 2023a. (Cited on page 19)
- Alex Morehead, Jeffrey A Ruffolo, Aadyot Bhatnagar, and Ali Madani. Towards joint sequence-structure generation of nucleic acid and protein complexes with $se(3)$ -discrete diffusion. In *NeurIPS 2023 Workshop on Machine Learning in Structural Biology*, 2023b. (Cited on page 20)
- Alex Morehead, Jian Liu, and Jianlin Cheng. Protein structure accuracy estimation using geometry-complete perceptron networks. *Protein Science*, 2024. (Cited on page 19)
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems 32*, pp. 8024–8035. Curran Associates, Inc., 2019. (Cited on page 20)
- Roshan Rao, Nicholas Bhattacharya, Neil Thomas, Yan Duan, Peter Chen, John Canny, Pieter Abbeel, and Yun Song. Evaluating protein transfer learning with tape. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. (Cited on page 19)
- Lorna Richardson, Ben Allen, Germana Baldi, Martin Beracochea, Maxwell L Bileschi, Tony Burdett, Josephine Burgin, Juan Caballero-Pérez, Guy Cochrane, Lucy J Colwell, Tom Curtis, Alejandra Escobar-Zepeda, Tatiana A Gurbich, Varsha Kale, Anton Korobeynikov, Shriya Raj, Alexander B Rogers, Ekaterina Sakharova, Santiago Sanchez, Darren J Wilkinson, and Robert D Finn. MGnify: the microbiome sequence data analysis resource in 2023. *Nucleic Acids Research*, 2022. (Cited on page 6, 24)
- Victor Garcia Satorras, Emiel Hoogeboom, and Max Welling. E (n) equivariant graph neural networks. In *ICML*, 2021. (Cited on page 2, 6, 22)
- Arne Schneuing, Yuanqi Du, Charles Harris, Arian Jamasb, Ilia Igashov, Weitao Du, Tom Blundell, Pietro Lió, Carla Gomes, Max Welling, et al. Structure-based drug design with equivariant diffusion models. *arXiv preprint*, 2022. (Cited on page 2)
- Kristof Schütt, Oliver Unke, and Michael Gastegger. Equivariant message passing for the prediction of tensorial properties and molecular spectra. In *ICML*, 2021. (Cited on page 1)
- Kristof T Schütt, Huziel E Sauceda, P-J Kindermans, Alexandre Tkatchenko, and K-R Müller. SchNet—a deep learning architecture for molecules and materials. *The Journal of Chemical Physics*, 2018. (Cited on page 1, 2, 6, 21)
- Vignesh Ram Somnath, Charlotte Bunne, and Andreas Krause. Multi-scale representation learning on proteins. *Advances in Neural Information Processing Systems*, 34, 2021. (Cited on page 19)
- Freyr Sverrisson, Jean Feydy, Bruno E Correia, and Michael M Bronstein. Fast end-to-end learning on protein surfaces. In *CVPR*, 2021. (Cited on page 19)
- Dassault Systèmes. Biovia discovery studio, 2016. (Cited on page 5, 30)
- Nathaniel Thomas, Tess Smidt, Steven Kearnes, Lusann Yang, Li Li, Kai Kohlhoff, and Patrick Riley. Tensor field networks: Rotation-and translation-equivariant neural networks for 3d point clouds. *arXiv preprint*, 2018. (Cited on page 2, 6, 23)

- Raphael Townshend, Martin Vögele, Patricia Suriana, Alex Derry, Alexander Powers, Yianni Laloudakis, Sidhika Balachandar, Bowen Jing, Brandon Anderson, Stephan Eismann, Risi Kondor, Russ Altman, and Ron Dror. Atom3d: Tasks on molecules in three dimensions. In *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks*, volume 1. Curran, 2021. (Cited on page 19)
- Michel van Kempen, Stephanie S. Kim, Charlotte Tumescheit, Milot Mirdita, Jeongjae Lee, Cameron L. M. Gilchrist, Johannes Söding, and Martin Steinegger. Fast and accurate protein structure search with foldseek. *Nature Biotechnology*, 2023. (Cited on page 1, 6, 24)
- Mihaly Varadi, Stephen Anyango, Mandar Deshpande, Sreenath Nair, Cindy Natassia, Galabina Yordanova, David Yuan, Oana Stroe, Gemma Wood, Agata Laydon, Augustin Židek, Tim Green, Kathryn Tunyasuvunakool, Stig Petersen, John Jumper, Ellen Clancy, Richard Green, Ankur Vora, Mira Lutfi, Michael Figurnov, Andrew Cowie, Nicole Hobbs, Pushmeet Kohli, Gerard Kleywegt, Ewan Birney, Demis Hassabis, and Sameer Velankar. AlphaFold protein structure database: massively expanding the structural coverage of protein-sequence space with high-accuracy models. *Nucleic Acids Research*, 2021. (Cited on page 1, 6, 24)
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. (Cited on page 21)
- Limei Wang, Haoran Liu, Yi Liu, Jerry Kurtin, and Shuiwang Ji. Learning hierarchical protein representations via complete 3d graph networks. In *ICLR*, 2023. (Cited on page 19)
- Carter J Wilson, Wing-Yiu Choy, and Mikko Karttunen. Alphafold2: a role for disordered protein/region prediction? *International Journal of Molecular Sciences*, 23(9):4591, 2022. (Cited on page 3, 25)
- Minghao Xu, Zuobai Zhang, Jiarui Lu, Zhaocheng Zhu, Yangtian Zhang, Ma Chang, Runcheng Liu, and Jian Tang. Peer: A comprehensive and multi-task benchmark for protein sequence understanding. In *NeurIPS*, 2022. (Cited on page 19)
- Omry Yadan. Hydra - a framework for elegantly configuring complex applications. Github, 2019. (Cited on page 20)
- Yu Yan, Jyun-Yu Jiang, Mingzhou Fu, Ding Wang, Alexander R. Pelletier, Dibakar Sigdel, Dominic C.M. Ng, Wei Wang, and Peipei Ping. MIND-s is a deep-learning prediction model for elucidating protein post-translational modifications in human diseases. *Cell Reports Methods*, 2023. (Cited on page 4, 26, 28)
- Jiaxuan You, Zhitao Ying, and Jure Leskovec. Design space for graph neural networks. *Advances in Neural Information Processing Systems*, 33, 2020. (Cited on page 17)
- Sheheryar Zaidi, Michael Schaarschmidt, James Martens, Hyunjik Kim, Yee Whye Teh, Alvaro Sanchez-Gonzalez, Peter Battaglia, Razvan Pascanu, and Jonathan Godwin. Pre-training via denoising for molecular property prediction. In *The Eleventh International Conference on Learning Representations*, 2023. (Cited on page 19)
- Zuobai Zhang, Minghao Xu, Vijil Chenthamarakshan, Aurélie Lozano, Payel Das, and Jian Tang. Enhancing protein language models with structure-based encoder and pre-training. *arXiv preprint arXiv:2303.06275*, 2023a. (Cited on page 20)
- Zuobai Zhang, Minghao Xu, Arian Rokkum Jamasb, Vijil Chenthamarakshan, Aurelie Lozano, Payel Das, and Jian Tang. Protein representation learning by geometric structure pretraining. In *The Eleventh International Conference on Learning Representations*, 2023b. (Cited on page 1, 2, 3, 6, 8, 19, 21, 25)
- Zuobai Zhang, Minghao Xu, Aurélie Lozano, Vijil Chenthamarakshan, Payel Das, and Jian Tang. Physics-inspired protein encoder pre-training via siamese sequence-structure diffusion trajectory prediction, 2023c. (Cited on page 25)
- Zhaocheng Zhu, Chence Shi, Zuobai Zhang, Shengchao Liu, Minghao Xu, Xinyu Yuan, Yangtian Zhang, Junkun Chen, Huiyu Cai, Jiarui Lu, et al. Torchdrug: A powerful and flexible machine learning platform for drug discovery. *arXiv preprint*, 2022. (Cited on page 19, 20)

APPENDICES

A	Illustrated Results	17
B	Related Work	19
C	Discussion and Future Work	19
D	<i>ProteinWorkshop</i> User Manual	20
	D.1 Dependencies	20
	D.2 Usage	20
	D.3 Computational Resources	20
	D.4 Featurisation Schemes	20
	D.5 Protein Structure Encoder Architectures	21
	D.6 Pretraining Datasets	23
	D.7 Pretraining Tasks	24
	D.8 Downstream Tasks	25
	D.9 SE(3) Equivariant Noise Predictor	25
	D.10 Hyperparameter Selection	26
E	Documentation for Datasets	26
	E.1 CATH - Inverse Folding	26
	E.2 MaSIF-Site - PPI Site Prediction	27
	E.3 ccPDB - Metal Binding Site Prediction	28
	E.4 PTM - Post-Translational Modification Site Prediction	28
	E.5 FOLD - Fold Prediction	29
	E.6 GO - Gene Ontology Prediction	29
	E.7 EC Reaction - Reaction Class Prediction	30
	E.8 TDC - Antibody Developability Prediction	30

A ILLUSTRATED RESULTS

Following [You et al. \(2020\)](#), in this section, we provide an alternative means of interpreting the main results in Table 2 of the main text. Across Figures 2-10, we illustrate the test metric stability of each encoder, featurization scheme, and auxiliary task across each dataset included in Table 2 via a ranking analysis. For each configuration, we rank design choices by their performance, deeming performance to be tied when the difference $\epsilon < 0.02$ ($\epsilon < 0.001$ for PPI site prediction). We collect rankings over all configurations and report the mean ranking and their smoothed distribution for each task with bar and violin plots, respectively (lower is better). An average ranking score of 1 indicates the design choice invariably results in the highest performance over other possibilities in its category. The smoothed distributions provide further insight into the broader behaviour of a design choice.

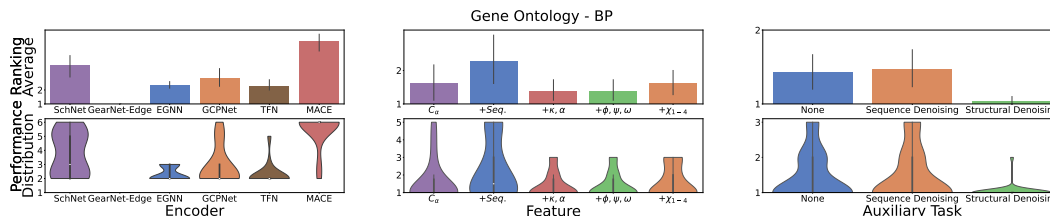


Figure 2: Ranking analysis of Gene Ontology-Biological Process (GO-BP) test performance across different encoders, feature sets, and auxiliary tasks.

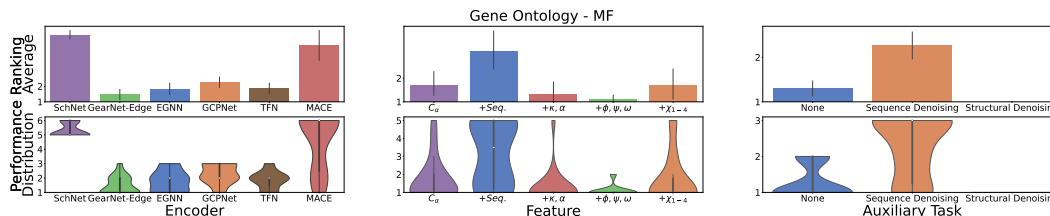


Figure 3: Ranking analysis of Gene Ontology-Molecular Function (GO-MF) test performance across different encoders, feature sets, and auxiliary tasks.

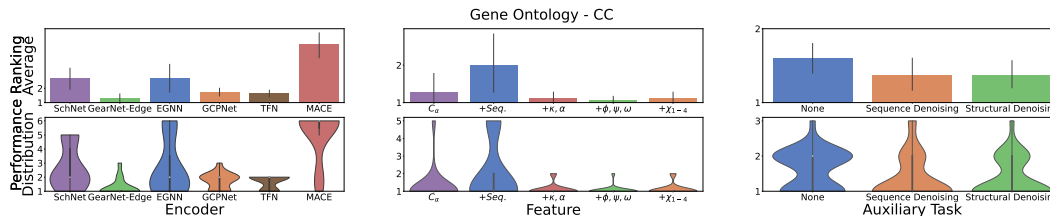


Figure 4: Ranking analysis of Gene Ontology-Cellular Component (GO-CC) test performance across different encoders, feature sets, and auxiliary tasks.

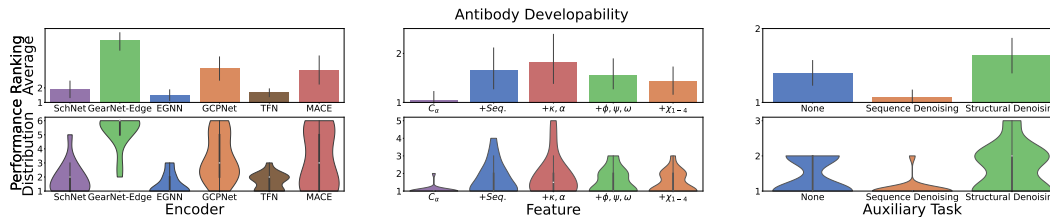


Figure 5: Ranking analysis of Antibody Developability test performance across different encoders, feature sets, and auxiliary tasks.

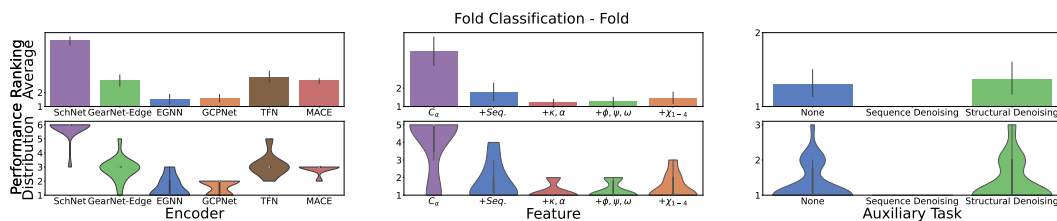


Figure 6: Ranking analysis of Fold Classification-Fold test performance across different encoders, feature sets, and auxiliary tasks.

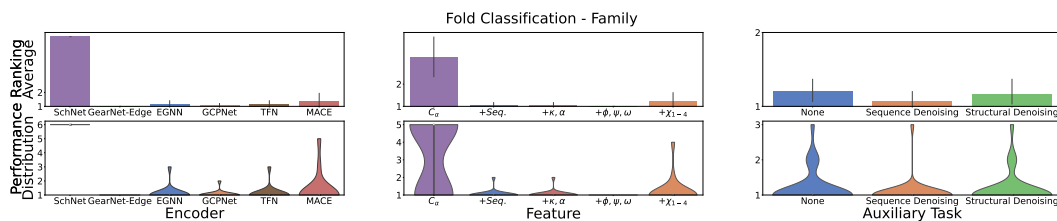


Figure 7: Ranking analysis of Fold Classification-Family test performance across different encoders, feature sets, and auxiliary tasks.

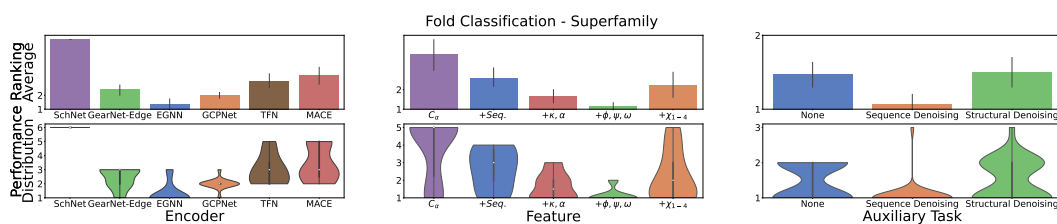


Figure 8: Ranking analysis of Fold Classification-Superfamily test performance across different encoders, feature sets, and auxiliary tasks.

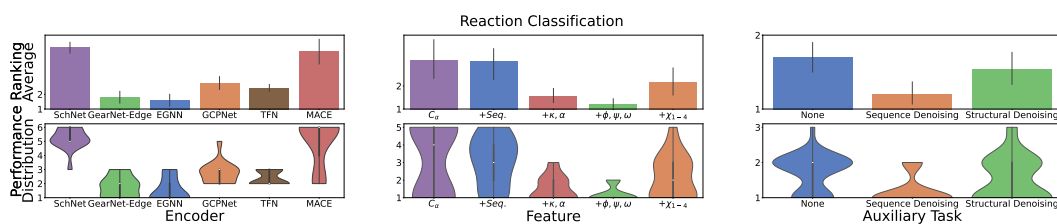


Figure 9: Ranking analysis of Reaction Classification test performance across different encoders, feature sets, and auxiliary tasks.

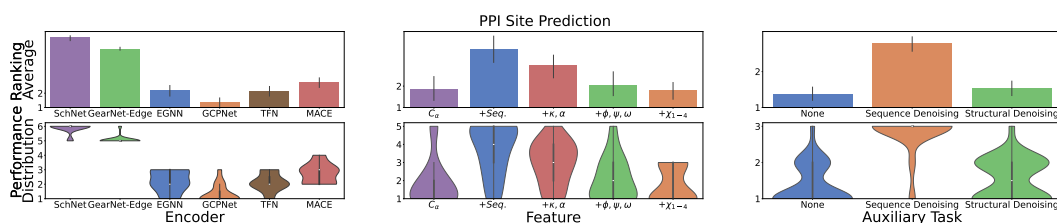


Figure 10: Ranking analysis of Protein-Protein Interaction (PPI) Site Prediction test performance across different encoders, feature sets, and auxiliary tasks.

B RELATED WORK

Protein Structure Representation Learning. Several structure-based encoders for proteins have been designed to extract information from different levels of granularity, such as residue-level, atom-level and surfaces. Previous works have aimed to encode protein structural priors directly within architectures to model proteins hierarchically (Somnath et al., 2021; Hermosilla et al., 2020), as computationally efficient point clouds (Gainza et al., 2020; Sverrisson et al., 2021), or as geometric graphs (Jing et al., 2020; Jin et al., 2021; Morehead & Cheng, 2024; Wang et al., 2023; Zhang et al., 2023b; Mahmud et al., 2023) for tasks such as protein function prediction (Gligorijević et al., 2021), protein model quality assessment (Eismann et al., 2020; Chen et al., 2023; Morehead et al., 2024), and protein interaction region prediction (Dai & Bailey-Kellogg, 2021; Morehead et al., 2023a).

Protein Benchmarks. Several benchmarks have been proposed for evaluating the efficacy of learnt protein *sequence* representations. However, *structure-based* benchmarks are comparatively unaddressed. Rao et al. (2019) developed TAPE (Tasks Assessing Protein Embeddings), providing a large pretraining corpus of protein sequences curated from Pfam (El-Gebali et al., 2018), as well as a collection of five supervised benchmark tasks assessing the ability of protein language models to predict structural qualities (contact prediction and secondary structure prediction), and functional properties (fluorescence and stability prediction). Xu et al. (2022) proposed PEER (Protein Sequence Understanding), focussing on multitask evaluation of protein sequence models. Therapeutic Data Commons (Huang et al., 2021) provide several datasets relevant to therapeutic development, however the few protein structure-derived datasets it contains are cast as sequence-based tasks. Dallago et al. (2021) developed FLIP, a sequence-based benchmark of protein fitness landscapes. ProteinGLUE (Capel et al., 2022) is another sequenced-based benchmark focussing on per-residue tasks.

To our best knowledge, the only protein structure-benchmark to date is ATOM3D (Townshend et al., 2021), which proposes a collection of tasks largely assessing geometric methods at predicting graph-level properties of protein structures. TorchProtein (Zhu et al., 2022) also provides a small collection of global-structural datasets. Most existing benchmarks do not exhaustively evaluate both the local and global representation learning power of proposed methods. As the field develops, we identify a need for a consistent benchmarking framework of diverse tasks to ensure improving results reported in the literature map on to progress in the downstream problems we hope to address. Similar benchmarking efforts for general purpose GNNs have provided experimental rigour to architectural research (Dwivedi et al., 2023).

Denoising-Based Pre-training and Regularisation. Several methods have been developed for pre-training GNNs, predominantly focussing on cases where 3D coordinate information is only implicitly encoded in the graph structures. In this work, we build on work by Godwin et al. (2021) and Zaidi et al. (2023) to investigate whether denoising-based auxiliary and pre-training tasks are effective methods for pre-training geometric GNNs operating on protein structures, similar to concurrent works bridging the gap between denoising objectives for geometric neural networks and diffusion generative modeling for biomolecules (Huang et al., 2023; Corso, 2023).

C DISCUSSION AND FUTURE WORK

Completeness of input featurisation. The extent to which providing complete information about all atoms and side chain orientations of each residue in a protein is debatable, as the exact coordinates from PDB files are known to contain artifacts from structure determination via crystallography. This was most recently noted by Dauparas et al. (2022) in the context of developing and experimentally validating an inverse folding model. Our benchmarking results provides similar insights – at present, letting models implicitly learn about side chain orientation by using backbone-only featurisation performs better or equally well as explicitly providing complete side chain atomic information across both global and node level tasks.

On the choice of pre-training tasks. We currently focussed on pre-training tasks that roughly fall under the category of denoising, i.e. corrupting information in the input (sequence identity, coordinates) and tasking the model with producing the uncorrupted input. We were particularly interested in self-supervised objectives that were (1) extremely scalable, so as to pre-train on the large-scale AlphaFoldDB of 2.4M structures; and (2) train protein representations at the fine-grained node level, so as to be general-purpose across the downstream tasks considered. We did not benchmark other

tasks from the literature, such as contrastive learning and generative modelling-inspired objectives (Liu et al., 2023b;a; Zhang et al., 2023a). Such tasks are generally (1) computationally heavier and more cumbersome to set up than corruption-type objectives, making them harder to scale up, and (2) only train protein representations for the global/graph level and do not operate at the node level. Naturally, we would like to continue exploring more pre-training strategies as we continue to expand *ProteinWorkshop*.

Beyond protein-only representation learning. Recently, Krishna et al. (2023) and DeepMind-Isomorphic (2023) generalised protein structure prediction models to full biological assemblies including proteins, small molecules, nucleic acids, and other ligands. Geometric graphs are a natural choice for representation learning across biomolecular systems: Advances in geometric GNN methodology should, in principle, be adaptable to modelling other biomolecules (Joshi et al., 2023b), their complexes (Morehead et al., 2023b), and quaternary structures.

We currently focus on protein representation learning because (1) large scale datasets for self-supervised learning, as well as well-defined downstream tasks, are readily available and accepted by the community; and (2) we see protein representation learning as a fundamental task, improving upon which should also advance the modelling of proteins in complex with other biomolecules. Comparatively, the scale of data available for biomolecular complexes is smaller and there is less consensus among the community on evaluation (Harris et al., 2023; Buttenschoen et al., 2023).

D *ProteinWorkshop* USER MANUAL

D.1 DEPENDENCIES

The benchmark is developed using PyTorch (Paszke et al., 2019), PyTorch Geometric (Fey & Lenssen, 2019), PyTorch Lightning (Falcon, 2019), and Graphein (Jamash et al., 2022). Experiment configuration is performed using Hydra (Yadan, 2019). Certain architectures introduce additional dependencies, such as TorchDrug (Zhu et al., 2022) and e3nn (Geiger & Smidt, 2022).

D.2 USAGE

The modular design of our benchmark means it can be readily adapted into different workflows easily. Firstly, the benchmark is pip-installable from PyPI and contains several importable modules, such as dataloaders, featurisers and models, that can be imported into new projects. This will aid in standardising the datasets and workflows used in protein representation learning. Secondly, the benchmark serves as an easily extendable template, which users can fork and work directly in, thereby reducing implementation burden. Lastly, we provide a CLI that can be used to quickly run single experiments and hyperparameter sweeps with minimal development time overhead.

D.3 COMPUTATIONAL RESOURCES

All models are trained on 80GB NVIDIA A100 GPUs. All baseline and finetuning results are performed using a single GPU while pre-training is performed using four GPUs.

D.4 FEATURISATION SCHEMES

Protein structures are typically represented as graphs, with researchers typically opting to use a coarse-grained $C\alpha$ atoms graph as full atom representations can quickly become computationally intractable due to a large number of nodes. The extent to which coarse-graining schemes are ‘complete’ representation of the geometry and structure of the protein residue is variable. For instance, backbone-only features ignore the orientations of the side chain atoms in the residue, so models must account for this information implicitly. However, providing complete information about all atoms and side chain orientations is debatable as exact coordinates from PDB files are known to contain crystallography artefacts (Dauparas et al., 2022). Due to the computational burden incurred by operating on full-atom node representations, we focus primarily on $C\alpha$ -based graph representations, investigating featurisation strategies to incorporate higher-level structural information. Note that we do provide utilities to enable users to work with backbone and full-atom graphs in the benchmark.

We represent protein structures as geometric graphs, $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \vec{\mathbf{X}}, \mathbf{S}, \vec{\mathbf{V}})$, where \mathcal{V} is a set of nodes, \mathcal{E} is a set of edges, $\vec{\mathbf{X}} \in \mathbb{R}^{|\mathcal{V}| \times 3}$ is a matrix of Cartesian node coordinates, $\mathbf{S} \in \mathbb{R}^{|\mathcal{V}| \times d}$ is a matrix of d -dimension scalar node features, and $\vec{\mathbf{V}} \in \mathbb{R}^{|\mathcal{V}| \times d \times 3}$ is a tensor of vector-valued features.

The five scalar featurisation schemes considered in the baselines are provided in Table 5. Pretraining experiments are performed on the featurisation schemes in rows three and four.

Table 5: **Structural featurisation schemes.** Residue type is a one-hot encoding of the amino acid type for each node; positional encoding is a 16-dimensional transformer-like positional encoding (Vaswani et al., 2017); and $\phi, \psi, \omega \in \mathbb{R}^6$ and $\chi_{1-4} \in \mathbb{R}^8$ are backbone dihedral angles and sidechain torsion angles, respectively, embedded on the unit circle. Similarly, $\kappa, \alpha \in \mathbb{R}^4$ are *virtual* torsion and bond angles defined over $C\alpha$ atoms. In our experimental evaluation, we consistently use k -NN edge construction with $k = 16$.

Granularity	$C\alpha$ Features	Backbone	Sidechain
$C\alpha$	Residue Type		
$C\alpha$	Residue Type, Positional Encoding		
$C\alpha$	Residue Type, Positional Encoding, κ, α		
$C\alpha$	Residue Type, Positional Encoding, κ, α	ϕ, ψ, ω	
$C\alpha$	Residue Type, Positional Encoding, κ, α	ϕ, ψ, ω	$\chi_1, \chi_2, \chi_3, \chi_4$

Additionally, vector message passing methods such as GCPNet receive node orientation vectors (i.e., $\mathbf{v}_{o^{i-1}}$ and $\mathbf{v}_{o^{i+1}}$) and edge directional vectors (i.e., $\mathbf{v}_{d^{ij}}$) as vector input features for nodes and edges, respectively. These vector features, prior to being normalised as unit vectors, are defined as:

$$\vec{\mathbf{v}}_{o^{i-1}} = \vec{\mathbf{x}}_{i-1} - \vec{\mathbf{x}}_i, \vec{\mathbf{v}}_{o^{i+1}} = \vec{\mathbf{x}}_{i+1} - \vec{\mathbf{x}}_i, \text{ and } \vec{\mathbf{v}}_{d^{ij}} = \vec{\mathbf{x}}_i - \vec{\mathbf{x}}_j. \quad (1)$$

D.5 PROTEIN STRUCTURE ENCODER ARCHITECTURES

We provide a unified implementation of several rotation invariant and equivariant geometric GNNs, spanning the range of message passing body order and tensor order (Joshi et al., 2023a), including general-purpose as well as protein-specific models. See Duval et al. (2023) for a self-contained introduction to geometric GNNs.

D.5.1 INVARIANT GNNs

SchNet (Schütt et al., 2018) SchNet is one of the most popular and simplest instantiation of E(3) invariant message passing GNNs. SchNet constructs messages through element-wise multiplication of scalar features modulated by a radial filter conditioned on the pairwise distance $\|\vec{\mathbf{x}}_{ij}\|$ between two neighbours. Scalar features are update from iteration t to $t + 1$ via:

$$\mathbf{s}_i^{(t+1)} := \mathbf{s}_i^{(t)} + \sum_{j \in \mathcal{N}_i} f_1 \left(\mathbf{s}_j^{(t)}, \|\vec{\mathbf{x}}_{ij}\| \right) \quad (2)$$

Hyperparameters: number of layers = 6, hidden channels = 512.

DimeNet++ (Gasteiger et al., 2020) DimeNet is an E(3) invariant GNN which uses both distances $\|\vec{\mathbf{x}}_{ij}\|$ and angles $\vec{\mathbf{x}}_{ij} \cdot \vec{\mathbf{x}}_{ik}$ to perform message passing among triplets, as follows:

$$\mathbf{s}_i^{(t+1)} := \sum_{j \in \mathcal{N}_i} f_1 \left(\mathbf{s}_i^{(t)}, \mathbf{s}_j^{(t)}, \sum_{k \in \mathcal{N}_i \setminus \{j\}} f_2 \left(\mathbf{s}_j^{(t)}, \mathbf{s}_k^{(t)}, \|\vec{\mathbf{x}}_{ij}\|, \vec{\mathbf{x}}_{ij} \cdot \vec{\mathbf{x}}_{ik} \right) \right) \quad (3)$$

Hyperparameters: number of layers = 6, hidden channels = 512.

GearNet-Edge (Zhang et al., 2023b) GearNet-Edge is an SE(3) invariant architecture leveraging relational graph convolutional layers and edge message passing. The original GearNet-Edge formulation presented in Zhang et al. (2023b) operates on multirelational protein structure graphs making use of several edge construction schemes (k -NN, euclidean distance and sequence distance based). Our

benchmark contains full capabilities for working with multirelational graphs but use a single edge type (i.e. $|\mathcal{R}| = 1$) in our experiments to enable more direct architectural comparisons.

The relational graph convolutional layer is defined for relation type r as:

$$\mathbf{s}_i^{(t+1)} := \mathbf{s}_i^{(t)} + \sigma \left(\text{BN} \left(\sum_{r \in \mathcal{R}} \mathbf{W}_r \sum_{j \in \mathcal{N}_r(i)} \mathbf{s}_j^{(t)} \right) \right) \quad (4)$$

The edge message passing layer is defined for relation type r as:

$$\mathbf{m}_{(i,j,r_1)}^{(t+1)} := \sigma \left(\text{BN} \left(\sum_{r \in |\mathcal{R}|'} \mathbf{W}'_r \sum_{(w,k,r_2) \in \mathcal{N}'_r((i,j,r_1))} \mathbf{m}_{(w,k,r_2)}^{(t)} \right) \right) \quad (5)$$

$$\mathbf{s}_i^{(t+1)} := \sigma \left(\text{BN} \left(\sum_{r \in |\mathcal{R}|} \mathbf{W}_r \sum_{j \in \mathcal{N}_r(i)} \left(\mathbf{s}_j^{(t)} + \text{FC}(\mathbf{m}_{(i,j,r)}^{(t+1)}) \right) \right) \right), \quad (6)$$

where $\text{FC}(\cdot)$ denotes a linear transformation upon the message function.

Hyperparameters: number of layers = 6, hidden channels = 512.

D.5.2 EQUIVARIANT GNNs IN CARTESIAN COORDINATES

EGNN (Satorras et al., 2021) We consider E(3) equivariant GNN layers proposed by Satorras et al. (2021) which updates both scalar features \mathbf{s}_i as well as node coordinates $\vec{\mathbf{x}}_i$, as follows:

$$\mathbf{s}_i^{(t+1)} := f_2 \left(\mathbf{s}_i^{(t)}, \sum_{j \in \mathcal{N}_i} f_1 \left(\mathbf{s}_i^{(t)}, \mathbf{s}_j^{(t)}, \|\vec{\mathbf{x}}_{ij}^{(t)}\| \right) \right) \quad (7)$$

$$\vec{\mathbf{x}}_i^{(t+1)} := \vec{\mathbf{x}}_i^{(t)} + \sum_{j \in \mathcal{N}_i} \vec{\mathbf{x}}_{ij}^{(t)} \odot f_3 \left(\mathbf{s}_i^{(t)}, \mathbf{s}_j^{(t)}, \|\vec{\mathbf{x}}_{ij}^{(t)}\| \right) \quad (8)$$

Hyperparameters: number of layers = 6, hidden channels = 512.

GCPNet (Morehead & Cheng, 2024) GCPNet is an SE(3) equivariant architecture that jointly learns scalar and vector-valued features from geometric protein structure inputs and, through the use of geometry-complete frame embeddings, sensitises its predictions to account for potential changes induced by the effects of molecular chirality on protein structure. In contrast to the original GCPNet formulation presented in Morehead & Cheng (2024), the implementation we provide in the benchmark incorporates the architectural enhancements proposed in Morehead & Cheng (2023) which include the addition of a scalar message attention gate (i.e., $f_a(\cdot)$) and a simplified structure for the model’s geometric graph convolution layers (i.e., $f_n(\cdot)$). With geometry-complete graph convolution in mind, for node i and layer t , scalar edge features $\mathbf{s}_{e^{ij}}^{(t)}$ and vector edge features $\mathbf{v}_{e^{ij}}^{(t)}$ are used along with scalar node features $\mathbf{s}_{n^i}^{(t)}$ and vector node features $\mathbf{v}_{n^i}^{(t)}$ to update each node feature type as:

$$(\mathbf{s}_{m^{ij}}^{(t+1)}, \mathbf{v}_{m^{ij}}^{(t+1)}) := f_e^{(t+1)} \left((\mathbf{s}_{n^i}^{(t)}, \mathbf{v}_{n^i}^{(t)}), (\mathbf{s}_{n^j}^{(t)}, \mathbf{v}_{n^j}^{(t)}), (f_a^{(t+1)}(\mathbf{s}_{e^{ij}}^{(t)}), \mathbf{v}_{e^{ij}}^{(t)}), \mathcal{F}_{ij} \right) \quad (9)$$

$$(\mathbf{s}_{n^i}^{(t+1)}, \mathbf{v}_{n^i}^{(t+1)}) := f_n^{(t+1)} \left((\mathbf{s}_{n^i}^{(t)}, \mathbf{v}_{n^i}^{(t)}), \sum_{j \in \mathcal{N}(i)} (\mathbf{s}_{m^{ij}}^{(t+1)}, \mathbf{v}_{m^{ij}}^{(t+1)}) \right), \quad (10)$$

where the geometry-complete and chirality-sensitive local frames for node i (i.e., its edges) are defined as $\mathcal{F}_{ij} = (\mathbf{a}_{ij}, \mathbf{b}_{ij}, \mathbf{c}_{ij})$, with $\mathbf{a}_{ij} = \frac{\mathbf{x}_i - \mathbf{x}_j}{\|\mathbf{x}_i - \mathbf{x}_j\|}$, $\mathbf{b}_{ij} = \frac{\mathbf{x}_i \times \mathbf{x}_j}{\|\mathbf{x}_i \times \mathbf{x}_j\|}$, and $\mathbf{c}_{ij} = \mathbf{a}_{ij} \times \mathbf{b}_{ij}$, respectively.

Hyperparameters: number of layers = 6, hidden scalar channels = 128, hidden vector channels = 16.

D.5.3 EQUIVARIANT GNNs IN SPHERICAL COORDINATES

Tensor Field Network (Thomas et al., 2018) Tensor Field Networks are E(3) or SE(3) equivariant GNNs that have been successfully used in protein structure prediction (Baek et al., 2021) and protein-ligand docking (Corso et al., 2023). These models use higher order spherical tensors $\tilde{\mathbf{h}}_{i,l} \in \mathbb{R}^{2l+1 \times f}$ as node features, starting from order $l = 0$ up to arbitrary $l = L$. The first two orders correspond to scalar features s_i and vector features $\tilde{\mathbf{v}}_i$, respectively. The higher order tensors $\tilde{\mathbf{h}}_i$ are updated via tensor products \otimes of neighbourhood features $\tilde{\mathbf{h}}_j$ for all $j \in \mathcal{N}_i$ with the higher order spherical harmonic representations Y of the relative displacement $\frac{\tilde{\mathbf{x}}_{ij}}{\|\tilde{\mathbf{x}}_{ij}\|} = \hat{\mathbf{x}}_{ij}$:

$$\tilde{\mathbf{h}}_i^{(t+1)} := \tilde{\mathbf{h}}_i^{(t)} + \sum_{j \in \mathcal{N}_i} Y(\hat{\mathbf{x}}_{ij}) \otimes_{\mathbf{w}} \tilde{\mathbf{h}}_j^{(t)}, \quad (11)$$

where the weights \mathbf{w} of the tensor product are computed via a learnt radial basis function of the relative distance, *i.e.* $\mathbf{w} = f(\|\tilde{\mathbf{x}}_{ij}\|)$.

Hyperparameters: choice of symmetry = SO(3), number of layers = 4, spherical harmonic tensor order = 2, hidden irreps per tensor type = 64x0e + 64x0o + 8x1e + 8x1o + 4x2e + 4x2o. We were particularly interested in benchmarking the impact of higher order tensors and SO(3) equivariance.

MACE MACE (Batatia et al., 2022) is a higher order E(3) or SE(3) equivariant GNN originally developed for molecular dynamics simulations. MACE provides an efficient approach to computing high body order equivariant features in the Tensor Field Network framework via Atomic Cluster Expansion: They first aggregate neighbourhood features analogous to equation 11 (the A functions in Batatia et al. (2022) (eq. 9)) and then take $k - 1$ repeated self-tensor products of these neighbourhood features. In our formalism, this corresponds to:

$$\tilde{\mathbf{h}}_i^{(t+1)} := \underbrace{\tilde{\mathbf{h}}_i^{(t+1)} \otimes_{\mathbf{w}} \dots \otimes_{\mathbf{w}} \tilde{\mathbf{h}}_i^{(t+1)}}_{k-1 \text{ times}}, \quad (12)$$

Hyperparameters: choice of symmetry = O(3), number of layers = 2, spherical harmonic tensor order = 2, hidden irreps per tensor type = 32x0e + 32x1o + 32x2e, body order = 4. Note that the number of channels for all tensor types must be the same for MACE, which is restrictive for scaling the depth and number of parameters.

D.6 PRETRAINING DATASETS

The benchmark contains several large corpora of both experimental and predicted structural data that can be used for pretraining or inference. We provide utilities for configuring supervised tasks and splits directly from the PDB (Berman, 2000). Additionally, we build storage-efficient dataloaders for large pretraining corpora of predicted structures including the AlphaFoldDB and ESM Atlas. We believe our codebase will considerably reduce the barrier to entry for pretraining and working with large structure-based datasets.

D.6.1 EXPERIMENTAL STRUCTURES

PDB. We provide utilities for curating datasets directly from the Protein Data Bank (Berman, 2000). In addition to using the collection in its entirety, users can define filters to subset and split the data using a combination of structural similarity, sequence similarity or temporal strategies. Structures can be filtered by length, number of chains, resolution, deposition date, presence/absence of particular ligands and structure determination method. The benchmark supports working with PDB structures in both `.pdb` and `.mmCIF` format (Bradley et al., 2017), which significantly reduces the requirements for data storage.

CATH. We provide the dataset derived from CATH 4.2 40% (Knudsen & Wiuf, 2010) non-redundant chains developed by Ingraham et al. (2019) as an additional, smaller, pretraining dataset. These data are split based on random assignment of the CATH topology classifications based on an 80/10/10 split.

ASTRAL. ASTRAL (Brenner, 2000) provides compendia of protein *domain* structures, regions of proteins that can maintain their structure and function independently of the rest of the protein. Domains typically exhibit highly-specific functions and can be considered structural building blocks of proteins.

D.6.2 PREDICTED STRUCTURES

We provide ready-to-go dataloaders for several large-scale collections of predicted structures derived from both AlphaFold2 (Jumper et al., 2021) and ESMFold (Lin et al., 2022). This is facilitated by FoldComp (Kim et al., 2023), a (minimally) lossy compression scheme for predicted protein structures. FoldComp stores protein structures as a collection of discretised dihedral and bond angles which can be used to reconstruct the whole structure using fixed bond lengths and canonical amino acid geometry. FoldComp achieves a disk-space reduction of almost an order of magnitude, describing a residue with only 13 bytes – down from 97 bytes per-residue in a traditional uncompressed format. Whilst lossy, this procedure results in 0.08 Å and 0.14 Å RMSD for backbone and all-atom reconstruction, making it highly suitable for pretraining tasks which use input representations complete up to the backbone. Furthermore, this lightweight format enables the dataloaders in the benchmark to read structures *directly from disk* with no pre-processing or caching required.

AlphaFoldDB Representative Structures. This dataset contains 2.27 million representative structures, identified through large-scale structural-similarity-based clustering of the 214 million structures contained in the AlphaFold Database (Varadi et al., 2021) using FoldSeek (van Kempen et al., 2023). We additionally provide a subset of this collection — the so-called dark proteome — corresponding to the 31% of the representative structures that lack annotations.

ESM Atlas, ESM High Quality. These datasets are compressed collections of predicted structures produced by ESMFold. ESM Atlas is the full collection of all 772m predicted structures for the MGnify 2023 release (Richardson et al., 2022). ESM High Quality is a curated subset of high confidence (mean pLDDT) structures from the collection.

D.7 PRETRAINING TASKS

The benchmark contains a comprehensive suite of pretraining tasks. Broadly, these tasks can be categorised into: masked-attribute prediction, denoising-based and contrastive learning-based tasks. In most cases, these tasks can be used as both a pretraining objective or as auxiliary tasks in a downstream supervised task.

Sequence Denoising. The benchmark contains two variations based on two sequence corruption processes $C(\hat{S}|\mathcal{S}, \nu)$ that receive an amino acid sequence $\mathcal{S} \in [0, 1]^{|\nu| \times 23}$ and return a sequence $\hat{S} \in [0, 1]^{|\nu| \times 23}$ with fraction ν of its positions corrupted. The first scheme is based on mutating a fraction of the residues to a random amino acid and tasking the model with recovering the uncorrupted sequence. The second is a masked residue prediction task, where a fraction of the residues are altered to a mask value and the model is tasked to recover the uncorrupted sequence.

Structure Denoising. We provide two structure-based denoising tasks: coordinate denoising and torsional denoising. In the coordinate denoising task, noise is sampled from a normal or uniform distribution and scaled by noise factor, $\nu \in \mathbb{R}$, and applied to each of the atom coordinates in the structure to ensure structural features, such as backbone or sidechain torsion angles, are also corrupted. The model is then tasked with predicting either the per-node noise or the original uncorrupted coordinates. For the torsional denoising variant, the noise is applied to the backbone torsion angles and Cartesian coordinates are recomputed using pNeRF (AlQuraishi, 2019) and the uncorrupted bond lengths and angles prior to feature computation. Similarly to the coordinate denoising task, the model is then tasked with predicting either the per-residue angular noise or the original dihedral angles.

Sequence-Structure Co-Denoising. This is a multitask formulation of the previously described structure and sequence denoising tasks, with separate output heads for denoising each modality.

Masked Attribute Prediction Tasks We use inverse folding (Section 2.3.1) as a pretraining task. The benchmark additionally incorporates the distance, angle and dihedral angle masked-attribute

prediction tasks proposed in Zhang et al. (2023b) as well as a backbone dihedral angle prediction task.

pLDDT Prediction. Protein structure prediction models typically provide per-residue pLDDT (predicted Local Distance Difference Test) scores as local confidence measures in the quality of the prediction shown to correlate well with disordered regions (Wilson et al., 2022). We formulate a self-supervised node-level regression task on predicted structures, somewhat analogous to structure quality assessment (QA), where the model is tasked with predicting the scaled per-residue pLDDT $y \in [0, 1]$ values.

D.8 DOWNSTREAM TASKS

We curate several structure-based and sequence-based datasets from the literature and existing benchmarks[†]. The tasks are selected to evaluate not only the *global* structure representational power of each method, but also to evaluate the ability of each method to learn informative *local* representations for residue-level prediction and annotation tasks.

The raw structures are, where possible and accounting for obsolescence, retrieved directly from the PDB (or another structural source) as several processed datasets used by the community discard full atomic coordinates in favour of retaining only C_α positions making them unsuitable for in-depth experimentation. This provides an entry point for users to apply a custom sequence of pre-processing steps, such as deprotonation or fixing missing regions which are common in experimental data.

The following downstream tasks are available in our benchmark. Detailed documentation including composition, splitting details, metrics, and data sheets (Gebru et al., 2021) are available in Appendix E

Node-level Tasks

- CATH Inverse folding.
- PPI Site Prediction.
- Metal Binding Site Prediction.
- Post-Translational Modification Site Prediction.

Graph-level Tasks

- Fold Prediction.
- Gene Ontology Prediction.
- Reaction Class Prediction.
- Antibody Developability Prediction.

D.9 SE(3) EQUIVARIANT NOISE PREDICTOR

Similar to Zhang et al. (2023c), for structure-based denoising tasks we use an SE(3) equivariant noise predictor network to predict per-residue perturbations from SE(3) invariant scalar embeddings and corrupted atomic coordinates $\tilde{\mathbf{X}}$. Each edge e_{ij} is featurised by concatenating the two adjoining scalar node representations \tilde{s}_i, \tilde{s}_j and the euclidean distance between them $\|\tilde{\mathbf{x}}_i - \tilde{\mathbf{x}}_j\|_2$. We use a two-layer MLP to produce a score \mathbf{m}_{ij} :

$$\mathbf{m}_{ij} = \text{MLP}(\tilde{s}_i, \tilde{s}_j, \text{MLP}(\|\tilde{\mathbf{x}}_i - \tilde{\mathbf{x}}_j\|_2)). \quad (13)$$

Subsequently, Equation 13 is used to aggregate normalised directional edge vectors over the neighbourhood \mathcal{N}_i of each node:

$$\epsilon_\theta(\tilde{\mathcal{G}}) = \sum_{j \in \mathcal{N}_i} \mathbf{m}_{ij} \cdot \frac{\tilde{\mathbf{x}}_i - \tilde{\mathbf{x}}_j}{\|\tilde{\mathbf{x}}_i - \tilde{\mathbf{x}}_j\|_2} \quad (14)$$

[†]To retain focus on *protein* representation learning, we deliberately exclude commonly-used tasks based on protein-small molecule interactions as it is hard to disentangle the effect of the small molecule representation and the potential for bias (Boyles et al., 2019)

Table 6: Overview of supervised tasks and datasets.

	Task	Dataset Origin	Structures	# Train	# Validation	# Test
Node-level	Inverse Folding	Ingraham et al. (Ingraham et al., 2019)	Experimental	3.9 M	105 K	180 K
	PPI Site Prediction	Gainza et al. (2020)	Experimental	478 K	53 K	117 K
	Metal Binding Site Prediction		Experimental	1.1 M	13.7 K	29.8 K
	Post-Trans. Modification Site Prediction	Yan et al. (2023)	Predicted	44 K	2.4 K	2.5 K
Graph-level	Fold Prediction	Hou et al. (Hou et al., 2017)	Experimental	12.3 K	0.7 K	1.3/0.7/1.3 K
	Gene Ontology Prediction	Gligorijević et al. (2021)	Experimental	27.5 K	3.1 K	3.0 K
	Reaction Class Prediction	Hermosilla et al. (2020)	Experimental	29.2 K	2.6 K	5.6 K
	Antibody Developability Prediction	Huang et al. (Huang et al., 2021)	Experimental	1.7 K	0.24 K	0.48 K

D.10 HYPERPARAMETER SELECTION

Given the large number of models and featurisation schemes, we try our best to do a consistent and fair hyperparameter search. We fix a consistently high number of layers and large hidden dimension across models, as we wanted to focus on scaling model size and dataset size via pre-training. We use the Fold Classification task to select the best learning rate and dropout per model and featurisation scheme for downstream tasks. While our best performing models sometimes do not outperform the best reported results in the literature, we have obtained these results in a consistent experimental setup. Our goal was to demonstrate the utility of our benchmarking framework and uncover the impact of architectural considerations such as featurisation schemes, geometric GNN models, and pre-training/auxiliary tasks under fair and rigorous experimental settings.

Protein Structure Encoders

1. We use a consistent batch size of 32.
2. For all models, we try to consistently use six layers, each with 512 hidden channels. For tensor-based equivariant GNNs (TFN, MACE), we reduced the number of layers and hidden channels to fit 80GB of GPU memory on one NVIDIA A100 GPU.
3. For each encoder and featurisation baseline, we search over learning rates: 0.00001, 0.0001, 0.0003, 0.001 and select the best based on the validation performance on the fold classification task.
4. For each finetuning configuration, we conduct a small sweep over learning rates: 0.0001, 0.0003, 0.00001 and perform model selection based on the validation performance.
5. We train all models to convergence or to a maximum of 24 hours on a single A100 GPU.

Output Heads

1. All primary output heads use a three-layer MLP with 512 as the hidden dimension.
2. For all auxiliary tasks we use a two-layer MLP with 128 as the hidden dimension.
3. For all structure denoising tasks we use a two-layer SE(3) equivariant noise predictor network (Section D.9). The message and distance MLPs each consist of two layers of 128 hidden units.
4. For each encoder and featurisation, we search over decoder dropout: 0.0, 0.1, 0.3, 0.5, and select the best based on the validation performance on the Fold Classification task.

E DOCUMENTATION FOR DATASETS

Below, we provide detailed documentation for each dataset included in our benchmark, summarised in Table 6. Each dataset will be made available for download in processed and raw forms from Zenodo upon publication. Note that, for all datasets, we authors bear all responsibility in case of any violation of rights regarding the usage of such datasets, whether they were compiled from existing sources or curated from scratch.

E.1 CATH - INVERSE FOLDING

This is a common protein engineering task where the goal is to recover an amino acid sequence given a structure up to backbone completeness. Formally, this is a node-level classification task where the

model learns a mapping for each residue r_i to an amino acid type $y \in \{1, \dots, n\}$, where n is the vocabulary size ($n = 20$ for the canonical set of amino acids).

Note that inverse folding is a generic task that can be applied to any dataset. In the literature, it is commonly evaluated on the CATH dataset compiled by Ingraham et al. (2019). We additionally use inverse folding on AlphaFoldDB as a pretraining task.

- **Motivation** Several generative methods for protein design produce backbone structures that require the design of an associated sequence. As a result, inverse folding is an important part of *de novo* design pipelines for proteins.
- **Collection** For this dataset, we adopt the commonly-used CATH dataset originally compiled by Ingraham et al. (2019).
- **Composition** The dataset consists of protein structures randomly split into training, validation and test sets such that proteins in different sets do not share the same CATH topology classification (i.e., CAT code).
- **Hosting** A preprocessed version of the dataset can be downloaded from the benchmark’s Zenodo data record.
- **Licensing** We have released a preprocessed version of the dataset under a Creative Commons Attribution 4.0 International license. The original dataset is available under a Creative Commons Attribution 4.0 International license at <http://cathdb.info>.
- **Maintenance** We will announce any errata discovered in or changes made to the dataset using the benchmark’s GitHub repository.
- **Uses** This dataset can be used for multilabel node classification tasks where a model learns a mapping for each residue r_i to an amino acid type $y \in \{1, \dots, n\}$, where n is the vocabulary size (e.g., $n = 20$ for the canonical set of amino acids).
- **Metric** Perplexity.

E.2 MASIF-SITE - PPI SITE PREDICTION

This task is a node-level binary classification task where the goal is to predict whether or not a residue is involved in a protein-protein interaction interface.

- **Motivation** Identifying protein-protein interaction (PPI) sites has important applications in developing improved protein-protein interaction networks and docking tools, providing biological context to guide protein engineering and target identification in drug discovery campaigns (Jamasb et al., 2021).
- **Collection** We adopt the dataset of experimental structures curated from the PDB by Gainza et al. (2020) and retain the original splits, though we modify the labelling scheme to be based on inter-atomic proximity (3.5 Å), which can be user-defined, rather than solvent exclusion.
- **Composition** The dataset is curated from the PDB by preprocessing such as the presence of one of the seven specified ligands (e.g., ADP or FAD), clustering based on 30% sequence identity and random subsampling. It contains 1,459 structures, which are randomly assigned to training (72%), validation (8%) and test set (20%). 12 (Å) radius patches were extracted from the generated structures and a patch labelled as part of a binding pocket if its centre point was < 3 (Å) away from an atom of the corresponding ligand.
- **Hosting** The original dataset is made available by the authors on Zenodo.
- **Licensing** We have released a preprocessed version of the dataset under a Creative Commons Attribution 4.0 International license. The original dataset is available under an Apache 2.0 license at <https://github.com/LPDI-EPFL/masif/blob/master/LICENSE>.
- **Maintenance** We will announce any errata discovered in or changes made to the dataset using the benchmark’s GitHub repository.
- **Uses** This dataset can be used for binary node classification tasks where the goal is to predict whether or not a residue is involved in a protein-protein interaction interface.
- **Metric** AUPRC.

E.3 CCPDB - METAL BINDING SITE PREDICTION

This is a binary node classification task where each residue is mapped to a label $y \in \{0, 1\}$ indicating whether the residue (or its constituent atoms) is within 3.5 (Å) of a user-defined metal ion or ligand heteroatom, respectively.

- **Motivation** Several proteins coordinate transition metal ions to carry out their functions. As such, predicting the binding sites of metal ions can elucidate the role of metal binding on protein function.
- **Collection** The dataset is constructed from experimental structures curated from the PDB, where binding site assignments for each residue are computed on-the-fly. While the benchmark supports this task on arbitrary subsets of the PDB and ligands, we provide the Zinc-binding dataset from Dürr et al. (2023) specifically for this task.
- **Composition** The dataset is constructed by sequence-based clustering of the PDB at 30% sequence identity to remove sequence and structural redundancy. Clusters with a member shorter than 3000 residues, containing at least one zinc atom with resolution better than 2.5 (Å) determined by x-ray crystallography and not containing nucleic acids are used to compose the dataset. If multiple structures fulfill these criteria, the highest resolution structure is used. The train (2,085) / validation (26) / test (59) splits are constructed such that proteins in the validation and test sets have no partial overlap with any protein in the training data.
- **Hosting** A preprocessed version of the dataset can be downloaded from the benchmark’s Zenodo data record.
- **Licensing** We have released a preprocessed version of the dataset under a Creative Commons Attribution 4.0 International license. The original dataset is freely available without a license at <https://academic.oup.com/database/article/doi/10.1093/database/bay142/5298333#130010908>.
- **Maintenance** We will announce any errata discovered in or changes made to the dataset using the benchmark’s GitHub repository.
- **Uses** This dataset can be used for binary node classification tasks where each residue is mapped to a label $y \in \{0, 1\}$ indicating whether the residue (or its constituent atoms) is within 3.5 (Å) of a user-defined metal ion or ligand heteroatom, respectively.
- **Metric** Accuracy.

E.4 PTM - POST-TRANSLATIONAL MODIFICATION SITE PREDICTION

We frame prediction of post-translational modification (PTM) sites as a multilabel classification task where each residue is mapped to a label $y \in \{1, \dots, 13\}$ distinguishing between modifications on different amino acids (e.g. phosphorylation on S/T/Y and N-linked glycosylation on N).

- **Motivation** Identifying the exact sites where post-translational modifications (PTMs) occur is essential for understanding protein behaviour and designing targeted therapeutic interventions.
- **Collection** We adopt a dataset of 48,811 AlphaFold2-predicted structures curated by Yan et al. (2023), where each structure contains the PTM metadata necessary to construct residue-wise site prediction labels.
- **Composition** The dataset is split into training (43,907), validation (2,393) and test (2,511) sets based on 50% sequence identity and 80% coverage. In total, there are 240,090 PTMs present in the dataset compared to 3,391,208 residues where PTMs could be possible but are not present. The most common PTMs are phosphorylations on serine (93,734) and N-linked glycosylation at asparagine (59,143) which together account for around 70% of the PTMs.
- **Hosting** A preprocessed version of the dataset can be downloaded from the benchmark’s Zenodo data record.
- **Licensing** We have released a preprocessed version of the dataset under a Creative Commons Attribution 4.0 International license. The original dataset is available under a Creative Commons Attribution 4.0 International license at <https://zenodo.org/record/7655709>.

- **Maintenance** We will announce any errata discovered in or changes made to the dataset using the benchmark’s GitHub repository.
- **Uses** This dataset can be used for multilabel node classification tasks where each residue is mapped to a label $y \in \{1, \dots, 13\}$ distinguishing between modifications on different amino acids (e.g., phosphorylation on S/T/Y and N-linked glycosylation on N).
- **Metric** ROC-AUC.

E.5 FOLD - FOLD PREDICTION

This is a multiclass graph classification task where each protein, \mathcal{G} , is mapped to a label $y \in \{1, \dots, 1195\}$ denoting the fold class.

- **Motivation** The utility of fold prediction is that it serves as a litmus test for the ability of a model to distinguish different structural folds. It stands to reason that models that perform poorly on distinguishing fold classes likely learn limited or low-quality structural representations.
- **Collection** We adopt the fold classification dataset originally curated from SCOP 1.75 by (Hou et al., 2017).
- **Composition** This dataset provides three different test sets stratified based on topological similarity: Fold, in which proteins originating from the same superfamily are absent during training; Superfamily, in which proteins originating from the same family are absent during training; and Family, in which proteins from the same family are present during training.
- **Hosting** A preprocessed version of the dataset can be downloaded from the benchmark’s Zenodo data record.
- **Licensing** We have released a preprocessed version of the dataset under a Creative Commons Attribution 4.0 International license. The original dataset is available under a Creative Commons Attribution 4.0 International license at <https://academic.oup.com/bioinformatics/article/34/8/1295/4708302>.
- **Maintenance** We will announce any errata discovered in or changes made to the dataset using the benchmark’s GitHub repository.
- **Uses** This dataset can be used for multilabel graph classification tasks where each protein, \mathcal{G} , is mapped to a label $y \in \{1, \dots, 1195\}$ denoting the fold class.
- **Metric** Micro-averaged accuracy.

E.6 GO - GENE ONTOLOGY PREDICTION

This is a multilabel classification task, assigning functional Gene Ontology (GO) annotation to structures. GO annotations are assigned within three ontologies: biological process (BP), cellular component (CC) and molecular function (MF).

- **Motivation** Predicting protein function in the form of functional annotations such as gene ontology (GO) terms has important applications in protein analysis and engineering, providing researchers with the ability to cluster functionally-related structures or to guide protein generation methods to design new proteins with desired functional properties.
- **Collection** We adopt the dataset of experimental structures originally curated from the PDB by (Gligorijević et al. (2021)).
- **Composition** We retain the original multi-cutoff based dataset splits proposed by (Gligorijević et al., 2021), with cutoff at 30% sequence similarity.
- **Hosting** A preprocessed version of the dataset can be downloaded from the benchmark’s Zenodo data record.
- **Licensing** We have released a preprocessed version of the dataset under a Creative Commons Attribution 4.0 International license. The original dataset is available under a Creative Commons Attribution 4.0 International license at <https://www.nature.com/articles/s41467-021-23303-9>.

- **Maintenance** We will announce any errata discovered in or changes made to the dataset using the benchmark’s GitHub repository.
- **Uses** This dataset can be used for multilabel graph classification tasks, assigning a functional Gene Ontology (GO) annotation to protein structures.
- **Metric** F_{\max} score.

E.7 EC REACTION - REACTION CLASS PREDICTION

This is a multiclass graph classification task where each protein, \mathcal{G} , is mapped to a label $y \in \{1, \dots, 384\}$ denoting which class of reactions a given protein catalyzes; all four levels of the EC assignment are employed to define the reaction class label.

- **Motivation** As proteins’ reaction classifications are based on their enzyme-catalyzed reaction according to all four levels of the standard Enzyme Commission (EC) number, methods that predict such classifications can help elucidate the function of newly-designed proteins as they are developed.
- **Collection** We adopt the reaction class prediction dataset originally curated from the PDB by [Hermosilla et al. \(2020\)](#).
- **Composition** The dataset is split on the basis of sequence similarity using a 50% threshold.
- **Hosting** A preprocessed version of the dataset can be downloaded from the benchmark’s Zenodo data record.
- **Licensing** We have released a preprocessed version of the dataset under a Creative Commons Attribution 4.0 International license. The original dataset is available under an MIT license at https://github.com/phermosilla/IEConv_proteins/blob/master/LICENCE.
- **Maintenance** We will announce any errata discovered in or changes made to the dataset using the benchmark’s GitHub repository.
- **Uses** This dataset can be used for multilabel graph classification tasks where each protein, \mathcal{G} , is mapped to a label $y \in \{1, \dots, 384\}$ denoting which class of reactions a given protein catalyzes.
- **Metric** Accuracy.

E.8 TDC - ANTIBODY DEVELOPABILITY PREDICTION

Therapeutic antibodies must be optimised for favourable physicochemical properties in addition to target binding affinity and specificity to be viable development candidates. Consequently, we frame prediction of antibody developability as a binary graph classification task indicating whether a given antibody is developable.

- **Motivation** From a benchmarking perspective, predicting the developability of a given antibody is important as it enables targeted performance assessment of models on a specific (immunoglobulin) fold, providing insight into whether general-purpose structure-based encoders can be applicable to fold-specific tasks.
- **Collection** We adopt the antibody developability dataset originally curated from SabDab ([Dunbar et al., 2014](#)) by [Chen et al. \(2020\)](#).
- **Composition** This dataset contains 2,426 antibodies that have both sequences and PDB structures available, where each example contains both a heavy chain and a light chain with resolution < 3 (Å). Labels are based on thresholding the developability index (DI) ([Lauer et al., 2012](#)) as computed by BIOVIA’s platform ([Systèmes, 2016](#)), which relies on an antibody’s hydrophobic and electrostatic interactions.
- **Hosting** A preprocessed version of the dataset can be downloaded from the benchmark’s Zenodo data record.
- **Licensing** We have released a preprocessed version of the dataset under a Creative Commons Attribution 4.0 International license. The original dataset is available under a Creative Commons Attribution 3.0 Unported license at https://tdcommons.ai/single_pred_tasks/develop/#sabdad-chen-et-al.

- **Maintenance** We will announce any errata discovered in or changes made to the dataset using the benchmark's GitHub repository.
- **Uses** This dataset can be used for binary graph classification tasks indicating whether a given antibody is developable.
- **Metric** AUPRC.