

TaxDiff: Taxonomic-guided diffusion model for protein sequence generation

Zongying Lin^{1,3}, Hao Li^{1,2,3}, Liuzhenghao Lv^{1,3}, Yu Wang^{1,3}, Bin Lin¹, Junwu Zhang¹,
Zijun Chen^{1,3}, Calvin Yu-Chian Chen^{1,2,3,4}, Li Yuan^{1,2,3*} & Yonghong Tian^{1,2,3*}

¹School of Electronic and Computer Engineering, Peking University, 518055, China

²Peng Cheng Laboratory, Shenzhen, 518000, China

³AI for Science (AI4S)-Preferred Program, Peking University Shenzhen Graduate School, Shenzhen, 518055, China

⁴School of Chemical Biology and Biotechnology, Peking University, 518055, China

Appendix A Introduction

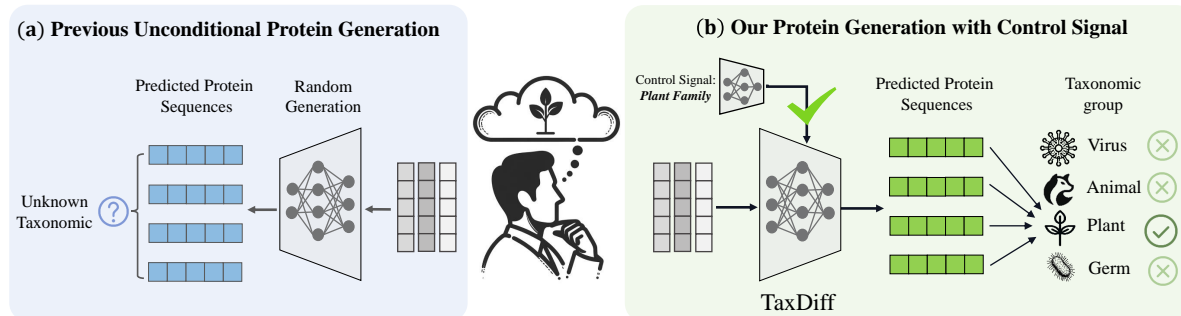


Figure A1 (a). Traditional protein sequence generation models operate without control signals, thus researchers can only randomly generate sequences and subsequently filter them until they fulfill the desired criteria. (b). Our TaxDiff takes species features as guidance for controllable protein sequence generation, meeting the need of biology downstream tasks.

Protein design [1–4] aims to generate protein variants with targeted biological functions, which is significant in multiple biological areas, including enzyme reaction catalysis [5,6], vaccine design [7–9], and fluorescence intensity [10,11].

Protein design contains two paradigms: sequence generation [5, 12, 13] and structure generation [14, 15]. Recently, EvoDiff [12] proposed a universal designing paradigm, combining structure and sequence generation using the diffusion framework [16], which improves the protein design efficiency. Despite the success of EvoDiff [12] and other sequences generative models [5,13,17] that are widely used for designing biologically plausible protein sequences, these protein design models are limited to unconditional generation. As shown in Figure A1.(a), in practical scenes, biological researchers need to filter the randomly generated proteins to fulfil the desired criteria [18], which is time-consuming and labor-intensive. Thus, unconditional protein generation, which can not control protein properties, is still some way from practical application.

To address the uncontrollable challenge, we propose a taxonomic-guided diffusion model, TaxDiff, to design target proteins with the biological-species control signals. Specifically, TaxDiff inserts the taxonomic control features into each Denoise Transformer block of the diffusion model to achieve controllable generation. Unlike the protein family in ProGen [19], which focuses on sequence and structural similarities, TaxDiff utilizes a broader taxonomic classification to reflect evolutionary lineage and shared characteristics among all living organisms. For fine-grained protein sequence generation, we also propose the patchify attention mechanism in the denoise transformer block to capture the protein feature on global and local scales. Furthermore, we reclassify protein sequences at the family and species levels to consolidate the overly detailed classification units within UniProt [20]. Our TaxDiff follows the protein design paradigm of EvoDiff [12]. Thus, TaxDiff is capable of generating both protein sequences and structures in a shared space.

We carry out extensive experiments to evaluate TaxDiff across multiple benchmarks, encompassing both unconditional and taxonomic-guided controllable protein sequence generation. In unconditional protein sequence generation, the sequence-based TaxDiff demonstrated comparable structural modeling capabilities to structure-based protein generation models, even significantly outperforming them in common metrics such as TM-score, RMSD, and Fident, with improvements of 11.93%, 5.4552, and 7.13% respectively. In taxonomic-guided controllable protein sequence generation, the pLDDT scores from

* Corresponding author (email: yuanli-ece@pku.edu.cn, yhtian@pku.edu.cn)

protein structure prediction model OmegaFold [21] far surpassed other sequence generation models, nearing the levels of natural protein sequences. Empirical studies also indicated that due to the patchify attention mechanism, the efficiency of TaxDiff was markedly enhanced, requiring as little as 24 minutes to generate 1,000 protein sequences, which is only 1/4 to 2/3 of the time required by other models. All experimental results demonstrate that TaxDiff possesses superior capabilities in exploring protein sequence space and producing structurally coherent proteins. The main contributions of our study are outlined as follows:

- To the best of our knowledge, our TaxDiff is the first controllable protein generation model utilizing guidance from taxonomies.
- Our TaxDiff proposes a taxonomic-guided framework that fits all diffusion-based protein design models. We propose a patchify attention mechanism that enhances protein design by reducing training and inference time while improving model performance.
- Experiments demonstrate that our TaxDiff achieves state-of-the-art results in both taxonomic-guided controllable and unconditional protein sequence generation, excelling in structural modeling scores and sequence consistency.

Appendix B Preliminary

In this section, we first introduce the problem setting of controllable protein sequence generation in Section Appendix B.1, then describe the Diffusion Models (Section Appendix B.2), which is utilized as our main generation framework.

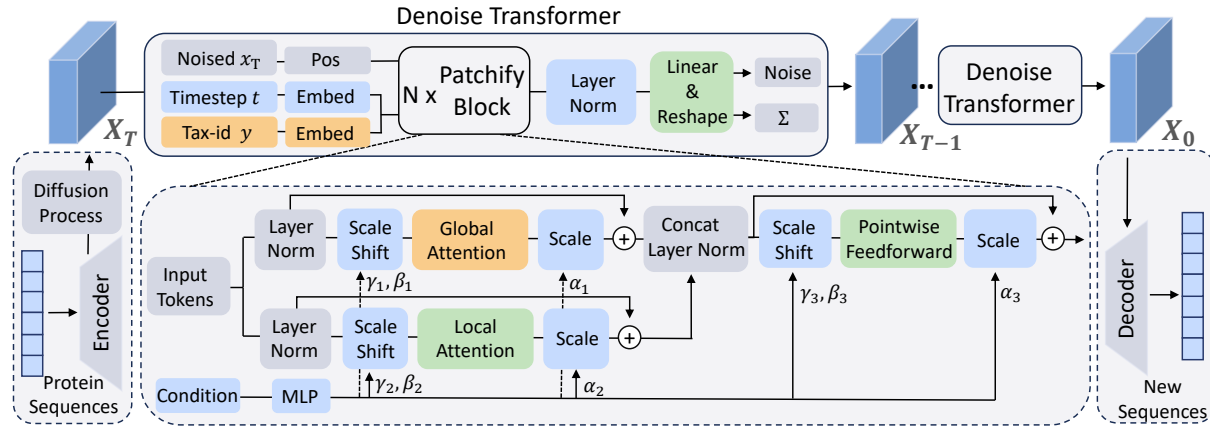


Figure B1 Overall architecture of the proposed TaxDiff. This framework delineates how we fuse the Denoise Transformer into the denoising process of the Diffusion model. For a taxonomic-guided controllable generation, we additionally accept a Tax-id y and embed it with Timestep t into the Patchify Blocks. The bottom middle of this framework elaborates on the details in Patchify Block. Σ is the predicted diagonal covariance.

Appendix B.1 Protein sequence generation

In this paper, we consider generating protein sequences under the guidance of taxonomies. Protein sequence space can be represented as $\mathbf{S} = \{\mathbf{x}, \mathbf{y}\}$ where $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_N) \in \mathbf{R}^{N \times L}$ are the protein sequences of length L and $\mathbf{y} = (\mathbf{y}_1, \dots, \mathbf{y}_M) \in \mathbf{R}^{M \times 1}$ represents the biological taxonomic category to which the protein belongs, such as Bacteria, Eukaryota, Archaea and so on. N and M , respectively represent the total number of protein sequences and categories.

We consider the following two generative tasks:

(I) Controllable generation. With a collection of protein sequences x with label y , we build a conditional generative model $p_\theta(x|y)$ that is capable of controllable protein sequences generation given desired biological taxonomic category y .

(II) Unconditional generation. Using the set of proteins x , unconditional generation train parameterized generative models $p_\theta(x)$ which can randomly generate diverse and realistic protein sequences without other additional labels. In a sense, unconditional generation also belongs to controllable generation, it just requires us to set the control signal to null: $p_\theta(x|\emptyset)$.

Appendix B.2 Diffusion model

We introduce the Diffusion Models (DMs) [22, 23], which is the generation framework of our TaxDiff. DMs are latent variable models [24] that model the pure data x_0 as Markov chain $x_0 \dots x_T$. In TaxDiff, the x_0 represents the initial and pure protein sequence, while x_T is the noisy protein sequence formed after adding Gaussian noise. DMs can be described with two Markovian processes: a forward diffusion process $q(x_{1:T}|x_0) = \prod_{t=1}^T q(x_t|x_{t-1})$ and a reverse denoising process $p_\theta(x_{0:T}) = p(x_T) \prod_{t=1}^T p_\theta(x_{t-1}|x_t)$. The variables x_t and x_{t-1} represent the protein sequence data with noise at timestep t and timestep $t-1$, respectively. The forward process gradually adds Gaussian noise to data x_t :

$$q(x_t|x_{t-1}) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)\mathbf{I}) \quad (\text{B1})$$

The hyperparameter $\bar{\alpha}_t$ in Eq. B1 controls the amount of noise added to the pure protein sequence x_0 at each timestep t . The $\bar{\alpha}_t$ are chosen so that the samples x_t can approximately converge to standard Gaussians $\mathcal{N}(0, \mathbf{I})$. Typically, this forward process q is a pure noise adder without trainable parameters.

The generation process of DMs is defined as learning a parameterized reverse denoising process p_θ , which aims to incrementally denoise the noisy variables $x_{T,1}$ to approximate initial data x_0 in the target data distribution. The denoising process p_θ in TaxDiff is represented by the Denoise Transformer, as illustrated in Figure B1, and can be formally expressed as:

$$p_\theta(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \Sigma_\theta(x_t, t)) \quad (\text{B2})$$

where the noise distribution $p(x_t)$ is defined as standard Gaussians $\mathcal{N}(0, \mathbf{I})$. The means μ_θ and variances Σ_θ typically are neural networks such as U-Nets [25] for images or Transformers for text [26]. However, TaxDiff uses a Denoise Transformer based on Patchify attention (described in Appendix D.4) to predict the means and variances.

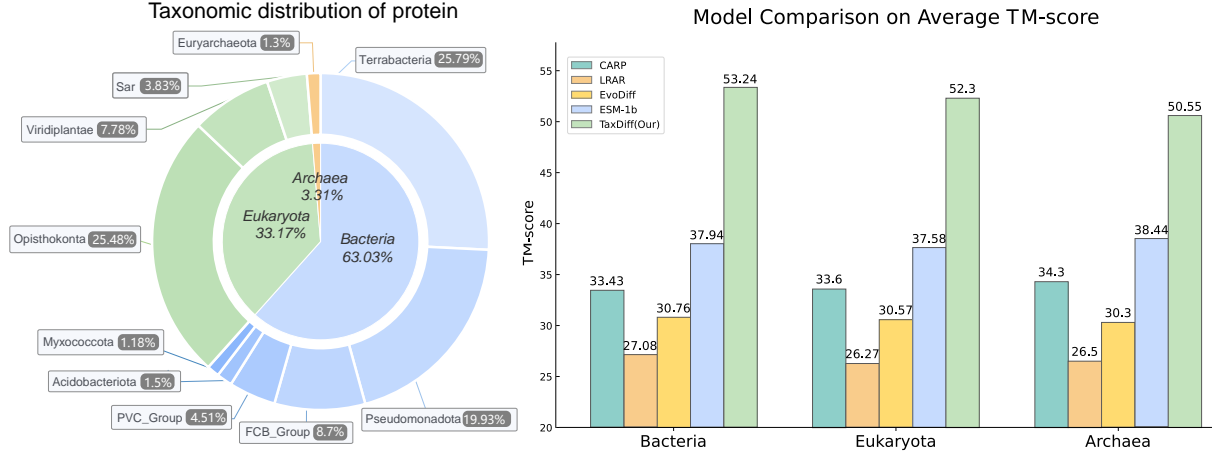


Figure B2 Taxonomic distribution and comparison of taxonomic-guided models. *Left:* We display the distribution of species classification from the second and third taxonomic levels, showcasing the top 10 categories. *Right:* We compare the performance of different models in terms of TM-score under the condition of the second taxonomic level.

Appendix C Method

We begin by introducing the framework and data flow of TaxDiff in Section Appendix C.1. Next, we elucidate the taxonomic-guided controllable generation process in Section Appendix C.2, followed by a discussion of the optimization of the denoising capability of the diffusion model using the patchify attention mechanism in Section Appendix C.3. Finally, we describe the training procedure, with a specific focus on the design of the loss function, in Section Appendix C.4. The overall architecture of the proposed TaxDiff is illustrated in Figure B1.

Appendix C.1 TaxDiff framework

In recognition of the diversity of amino acids, we introduce an additional dimension D to enrich features at the amino acid level. Through Encoder, feature-augmented x can thus be represented as $x \in \mathbb{R}^{L \times D}$. In the Denoise Transformer block, three different types of inputs are processed: the data x_T formed by the forward process in DMs that gradually adds Gaussian noise, the timestep t , and the protein taxonomic identifier y (tax-id). x_T undergoes standard Transformer-based frequency position embedding (sine-cosine version) [27], while the timestep t and tax-id y are individually embedded, resulting in two distinct conditional tokens that are concatenated with the x_T . Conditional tokens are designed for seamless integration, rendering them indistinguishable from protein sequence tokens. After passing through the terminal Patchify block, these conditional tokens are removed from the sequence. This approach enables the use of standard Transformer blocks without modification.

After the patchify block, the sequence tokens must be decoded into a predicted noise and diagonal covariance (Σ). Both of them retain the shape equivalent to that of the original input. To facilitate this, adaptive layer normalization (adaLN) [28] is applied, and each token is linearly decoded into a tensor of dimensions $L \times 2D$. In the decoder, the denoised final result $x_0 \in \mathbb{R}^{L \times D}$ is subjected to an argmax layer: $\text{argmax}(x_0) \in \mathbb{R}^L$. The output is then parsed and segmented at each padding or stopping sign, thereby generating the protein sequences.

Appendix C.2 Taxonomic-guided generation

Traditional conditional diffusion models take the class labels or text as extra information [29,30], but our TaxDiff model encodes tax-id y as a condition for controllable generation. In this case, the reverse process is formalized as $p_\theta(x_{t-1}|x_t, y)$, with both means μ_θ and variances Σ_θ being conditional by y . In this way, the taxonomic-guided encourages the sampling procedure towards maximizing the conditional log-likelihood $\log p(y|x)$ [31]. Invoking Bayes' Rule, we have $\log p(y|x) \propto \log p(x|y) - \log p(x)$ and hence $\nabla_x \log p(y|x) \propto \nabla_x \log p(x|y) - \nabla_x \log p(x)$. DMs based on score functions guides model

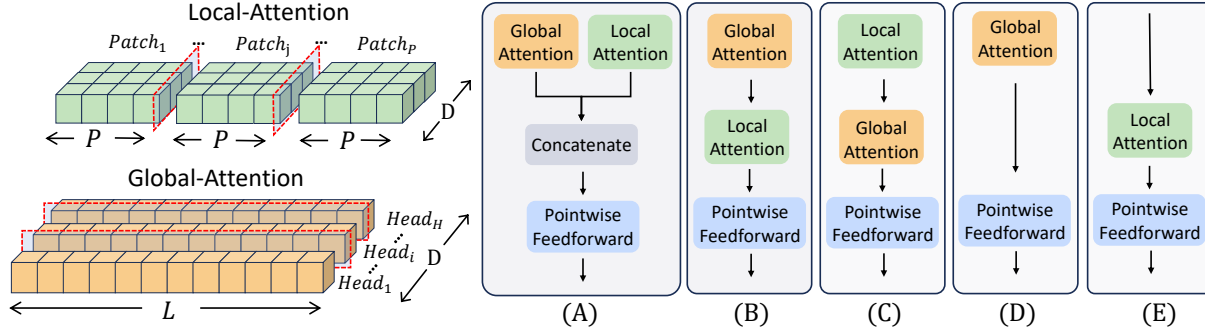


Figure C1 The patchify attention mechanism. *Left:* The division between Local-Attention and Global-Attention. *Right:* we present five different approaches to combining them. The division of the feature dimension in Global-Attention is determined by the number of heads H , while the number of tokens generated by Patchify is determined by the patchify size P .

sample x with an increased probability $p(x|y)$ by adjusting noise prediction:

$$\begin{aligned} \hat{\mu}_\theta(x_t, y) &= \mu_\theta(x_t, \emptyset) + s \cdot \nabla_x \log p(x|y) \\ &\propto \mu_\theta(x_t, \emptyset) + s \cdot (\mu_\theta(x_t, y) - \mu_\theta(x_t, \emptyset)) \end{aligned} \quad (\text{C1})$$

where the scalar $s > 1$ determines the scale of guidance. Notably, setting $s = 1$ reverts the process to standard sampling.

For diffusion models where $y = \emptyset$, this is achieved by randomly dropping out y with a certain probability and replacing it with a 'null' embedding \emptyset . This technique of controllable guidance significantly improves samples over generic sampling techniques [31, 32]. This observed improvement is also demonstrated in Figure B2 and discussed in Experiments Appendix D.4.

The classification of taxonomic identifiers (tax-id) y within UniProt is highly detailed, often resulting in exceedingly fine-grained categories that encompass only a limited number of sequences. Such fine-grained categories hinder effective feature extraction across the broader Taxonomic domain. To address this issue, we reclassified the ninth layer of original taxonomic lineages, which corresponds to the family and species levels [33]. Specifically, we align the protein sequences from UniProt with the taxonomic data from NCBI [34] by recursive tracing from the terminal child nodes up to the root. Utilizing the nine layers of classification subordinate to the root node, we assigned a novel tax-id y to each sequence. This reorganization effectively condenses the categories to a total of 23,427. Within this refined classification, cellular organisms constitute a dominant 98% of the protein sequences, with the remaining 2% attributed to viruses and other entities. Predominant domains within cellular organisms include Bacteria (63%), Eukaryota (33%), and Archaea (3%).

To further examine the impact of taxonomic guidance, we calculated the average performance of different models under the guidance of these three taxonomic domains. The performance of the different models in the TM-score is shown in Figure B2 Right.

Appendix C.3 Patchify Attention

To enhance the model's capacity to extract global features from protein sequences while simultaneously preserve the salient of amino-acid local features, we have implemented patchify attention mechanism upon the input protein sequences. This procedure is delineated in Figure C1.

At the protein sequence level, Global-Attention is employed to capture the intricate relationships between different amino acids. For each head i within the Global-Attention, with a total of H heads, we calculate the $Q_i = x \times W_i^Q$, $K_i = x \times W_i^K$ and $V_i = x \times W_i^V$ as linear transformations of the input matrix $x \in \mathbb{R}^{L \times D}$, where W_i^Q , W_i^K and W_i^V are the weight matrices unique to each head. The attention weights are then computed using the scaled dot-product attention mechanism:

$$\text{Attention}_{\text{head}}(Q_i, K_i, V_i) = \text{softmax} \left(\frac{Q_i K_i^T}{\sqrt{d_k}} \right) V_i \quad (\text{C2})$$

Here, d_k is the dimension of the K_i . The final output of the Global-Attention obtained by concatenating the individual heads' outputs and applying a subsequent linear transformation:

$$\text{Global-Attention}(x) = \text{concat}(\text{head}_1, \dots, \text{head}_H) W^O \quad (\text{C3})$$

where each head_i represents an attention block, and W^O serves as an output weight matrix to synthesize all Heads.

At the amino-acids level, Local-Attention divides the sequences of length L based on the patchify-size P . For each patch j , the queries (Q_j), keys (K_j) and values (V_j) are deduced by:

$$Q_j = \text{Patch}_j(x) \times W^Q \quad (\text{C4})$$

$$K_j = \text{Patch}_j(x) \times W^K \quad (\text{C5})$$

$$V_j = \text{Patch}_j(x) \times W^V \quad (\text{C6})$$

where $Patch_j(x)$ denotes the partitioning of x into its corresponding patch, and W^Q , W^K and W^V are the shared weight matrices. Attention weights within each patch are similarly computed utilizing the scaled dot-product attention:

$$Attention_{\text{patch}}(Q_j, K_j, V_j) = \text{softmax} \left(\frac{Q_j K_j^T}{\sqrt{d_k}} \right) V_j \quad (\text{C7})$$

where d_k is the key’s dimension within the patch. The Local-Attention’s ultimate output is derived by concatenating the outputs from all patches:

$$Local\text{-}Attention(x) = \text{concat}(Patch_1, \dots, Patch_P) \quad (\text{C8})$$

where $Patch_j$ equating to $Attention_{\text{patch}}(Q_j, K_j, V_j)$.

Additionally, we derive the scaling parameter α from residual connections within the Patchify block, while the scaling parameter γ and bias parameter β are regressed from the conditioning embedding vectors of t and y . The adaptive layer normalization (adaLN) uniformly applies the same function across all tokens.

To effectively integrate Global and Local-Attention, we explored five methodologies and subjected them to extensive experiments. The combinatorial strategies are illustrated on the left side of Figure C1: (A) Global-Attention and Local-Attention are used in parallel, followed by concatenation and fusion of the extracted features via a Pointwise Feedforward network; (B) and (C) Global-Attention and Local-Attention are deployed in a sequential format, capped with Pointwise Feedforward to predict noise and variance; (D) and (E) isolate the utilization to either Global-Attention or Local-Attention exclusively. The comparison of these five methods, as well as the impact of varying patchify size on local feature extraction capabilities, are examined in Section Appendix D.5.

Appendix C.4 Training Procedure

The training of diffusion models is aimed to learn the reverse process, which is expressed as $p_\theta(x_{t-1}|x_t) = \mathcal{N}(\mu_\theta(x_t), \Sigma_\theta(x_t))$, while the Denoise Transformer is used to estimate p_θ . The model training within the variational lower bound of the log-likelihood of x_0 , with the exclusion of an additional term irrelevant for training, the loss function can be represented to:

$$L(\theta) = -p(x_0|x_1) + \sum_t D_{KL}(q^*(x_{t-1}|x_t, x_0)||p_\theta(x_{t-1}|x_t)) \quad (\text{C9})$$

Given that both q^* and p_θ are Gaussian distributions, the Kullback–Leibler divergence (D_{KL}) can be evaluated by the means and covariances of these distributions.

For simple model training, μ_θ is reparameterized as a noise prediction network ϵ_θ , then the model can be trained using a simple mean-squared error function (MSE loss) between the predicted noise $\epsilon_\theta(x_t)$ and the actual sampled Gaussian noise ϵ_t :

$$L_{MSE}(\theta) = ||\epsilon_\theta(x_t) - \epsilon_t||^2 \quad (\text{C10})$$

However, for training diffusion models with a learned reverse process covariance Σ_θ , it becomes imperative to optimize the entire D_{KL} term. Followed DiT [26], we train ϵ_θ with L_{MSE} and Σ_θ with the full loss function $L(\theta)$. Upon the successful training of p_θ , new protein sequences can be generated by initializing $x_t \sim \mathcal{N}(0, I)$ and subsequently sampling $x_{t-1} \sim p_\theta(x_{t-1}|x_t)$ via the reparameterization trick.

Appendix D Experiments

In this section, we demonstrate the advantages of TaxDiff through a series of comprehensive experiments. We begin by presenting our experimental setup and Evaluation in Section Appendix D.1 and Section Appendix D.2. Following this, we present and analyze the results of unconditional and controllable generation in Section Appendix D.3 and Section Appendix D.4, respectively. Then, We explore various methods of combining Global and Local-Attention, and examine the impact of different patchify-size in Section Appendix D.5. Finally, We show the scaling ability and generative of TaxDiff in Section Appendix D.6 and Section Appendix D.7.

Appendix D.1 Experiment Setup

Datasets: For our model training, we utilized the Uniref50(2023-04 release) dataset from Uniprot [20], a protein database formed through clustering. Within TaxDiff, we retained protein sequences in Uniref50 that were less than 256 amino acids in length. Sequences falling short of 256 amino acids undergo zero-padding to equalize their length, thereby standardizing the representation of proteins as sequences with a uniform length $L = 256$. The focus on sequences shorter than 256 amino acids, covering 62% of sequences in Uniref50 (approximately 37.88 million out of 60 million sequences), was deliberate. This range not only represents a dominant subset of sequences but also encompasses small molecules [35] and peptides [36] that are pivotal in life sciences and pharmaceutical research, providing targeted relevance to our study. It is noteworthy that the evaluation datasets AFDB and PDB, which are derived from the Foldseek [37]. We used directly these two original databases without any selection criteria or filtering steps.

Setting: Following the DiT [26], we replace the standard layer norm in the Transformer block with adaLN and initialize batch normalization to zero within each Patchify block (adaLN-zero) [38] and employed AdamW [39] for training all models. Furthermore, we employ a linear beta scheduler from DDPM [23], adhering to the original DDPM sampling strategy across 1,000 steps. A constant learning rate of 1×10^{-4} without weight decay is used for training, while the batch size is set to 512 on eight 4090 GPUs. Results in D1 and D2 have the same setup with patchify-size $P = 16$, the layer of patchify block

$N = 12$. Following common practices in generative modeling, we maintained an exponential moving average (EMA) of the TaxDiff weights during training, with a decay factor of 0.9999. The same training hyperparameters were applied across all TaxDiff models and patch sizes.

Baselines: We compare TaxDiff to several competitive baseline models. The left-to-right autoregressive (LRAR) model and convolutional autoencoding representations of proteins (CARP) [40] were both trained using the same dilated convolutional neural network architectures on the UniRef50 dataset. For LRAR, a causal mask is applied to the convolutional modules to prevent information leakage. FoldingDiff [14] and RFdiffusion [15] are recent progress on diffusion models for protein structure generation. Notably, the RFdiffusion and FoldingDiff directly produce protein structures; hence, we first unconditionally select structures generated by these two and then use ESM-IF to design their corresponding sequences. EvoDiff [12] is a diffusion model that leverages evolutionary-scale data based on programmable, sequence-first design. The results reported in our result are based on the config640M configuration. ESM-1b [41] and state-of-the-art ESM-2 [42] are the protein masked language models, which were trained on different releases of UniRef50. Specifically, ESM-2 version is ESM2-t33-650M-UR50D and ESM1b is ESM1b-t33-650M-UR50S. ProtGPT2 [43] and ProGen2 [19] are autoregressive large protein language models based on GPT2 [44], which have undergone pre-training in UniRef50 and UniRef90 respectively. ProGen2 refers specifically to ProGen2-large.

It is worth noting that using ESM models to predict sequences is not common. However, there are already some applications in EvoDiff [12] and human antibodies [45] due to the strong protein sequence understanding capability of the ESM method. Recognizing the importance of clarity in explaining non-standard methodologies, we have included the relevant content in our work to elucidate our approach. To be specific, the process begins with a series of `<mask>_seqlen` tokens as input to simulate a completely unknown sequence of the desired length. This method allows the model to iteratively predict the missing amino acids at each position. To ensure the biological relevance of the generated sequences, we implement penalties on specific tokens (e.g., X, pad, eos, and cls) to prevent their generation, addressing the unique considerations necessary for protein sequence design. The sequence generation is facilitated by performing random sampling based on the predicted probability distribution of amino acids, thus ensuring variability and fidelity to plausible protein sequences. Finally, the tokens are converted back into amino acid sequences, culminating in the designed protein sequences.

Appendix D.2 Evaluation

We measure model performance through sequence consistency and structural foldability. For the feasibility of the sequences, We employed OmegaFold [21] to predict structures and calculate the average Predicted Local Distance Difference Test (pLDDT) across the entire structure, which reflects OmegaFold’s confidence in its structure prediction for each residue on sequences level. The pLDDT is calculated by a function based on the difference between the inter-residue distances predicted by the model and the true distance. It has values ranging from 0 to 100 and indicates the confidence level of each residue in the prediction which was first used in AlphaFold2 [46]. OmegaFold performs structure prediction without the need for homologous sequences or evolutionary information, relying solely on a single sequence for prediction. However, due to inherent noise and errors in OmegaFold’s structure predictions, which only consider the foldability of individual sequences, we further measure our results using Foldseek [37]. Foldseek facilitates the pairing of the queried protein p^{query} with structurally similar proteins from an existing protein database, yielding pairs represented as (p^{query}, p^{target}) . Here, p^{target} denotes the protein in the database with a significant structural similarity to p^{query} . The magnitude of the average Template Modeling score (TM-score [47]) value and Root-Mean-Square Deviation (RMSD) reflects the degree of structural similarity. TM-score takes into account the overall topological structure of proteins, focusing more on the protein’s overall structure. RMSD calculates the square root of the average position deviation of corresponding atoms between two protein structures, being highly sensitive to the size of the protein structure and local variations. Additionally, Foldseek also calculates the sequence identity (Fident) between p^{query} and p^{target} , reflecting their sequence-level similarity. In Foldseek, we choose natural protein structures from the Protein Data Bank (PDB Dataset) [48] to verify the natural-like degree of our sequences and high-confidence protein structures predicted by AlphaFold (AFDB Dataset) [1] to expand the comparison range and verify the broad validity of our model.

TM-score is a quantitative measure used to assess the similarity between two protein structures. It is defined as:

$$\text{TM-score} = \max \left(\frac{1}{L} \sum_{i=1}^L \frac{1}{1 + \left(\frac{d_i}{d_0(L)} \right)^2} \right) \quad (\text{D1})$$

where L represents the length of the target protein, d_i denotes the distance between the i th pair of aligned residues, and $d_0(L)$ is a length-dependent scale for normalization. The TM-score ranges from 0 to 1, with values closer to 1 indicating a higher degree of structural similarity. Typically, a TM-score greater than 0.5 suggests that two protein structures share the same fold.

RMSD is a standard measure used to quantify the average distance between atoms (usually backbone atoms) of superimposed proteins. The RMSD is calculated as:

$$\text{RMSD} = \sqrt{\frac{1}{N} \sum_{i=1}^N (r_i - r'_i)^2} \quad (\text{D2})$$

where N is the number of aligned atoms, r_i represents the position vector of the i th atom in the reference structure, and r'_i is the corresponding position in the model structure. RMSD provides a measure of the structural deviation between two protein conformations, with lower values indicating higher similarity.

Fident quantifies the proportion of residues that are identical between the sequences, providing insight into the degree of homology. The Fident value is calculated as follows:

$$\text{Fident} = \frac{\text{Number of identical residues}}{\text{Total number of aligned residues}} \times 100\% \quad (\text{D3})$$

where the numerator represents the number of residues that are identical between the aligned sequences, and the denominator is the total number of residues that have been aligned.

Appendix D.3 Unconditional Sequence Generation

In the unconditional generation, TaxDiff sets the condition y to be \emptyset and generates 1,000 protein sequences within the length range of 10 to 256 residues. For sequences lengths shorter than 10, we consider them to be invalid protein sequences and, therefore remove them, which is a practice also applicable to sequences generated by all other models.

Table D1 Unconditional generation result comparison on AFDB and PDB datasets. Metrics are calculated with 1,000 samples generated from each model.

Architecture	Method	pLDDT \uparrow	AFDB Dataset			PDB Dataset		
			TM-score($\%$) \uparrow	RMSD \downarrow	Fident($\%$) \uparrow	TM-score($\%$) \uparrow	RMSD \downarrow	Fident($\%$) \uparrow
CNN	CARP	34.66	25.31	18.6365	14.29	31.02	12.6905	15.23
	LRAR	48.88	26.33	18.6102	14.91	30.50	13.6436	15.45
Encoder	ESM-1b	59.25	32.69	19.9506	16.82	36.67	16.1619	17.58
	ESM-2	51.01	24.56	23.7637	17.18	28.92	20.3257	17.56
Diffusion	EvoDiff	44.29	24.22	20.0326	15.01	29.58	13.9564	15.64
	FoldingDiff	67.44	34.96	12.2538	19.04	37.32	9.8115	19.91
Decoder	ProGen2	61.26	21.93	27.8802	18.21	28.97	17.0160	15.05
	ProtGPT2	57.44	26.67	18.8129	14.36	31.99	12.8339	14.75
Diffusion	TaxDiff(Our)	68.89	48.26	5.9075	26.60	46.02	4.5736	24.13

As shown in Table D1, TaxDiff outperforms competitive baseline methods across all metrics with a noticeable margin. Notably, the sequences generated by our model have exceeded those of the structure generation model RFdiffusion and FoldingDiff in terms of pLDDT. This result underscores the model’s efficacy in extracting structural information from protein sequences, validating the effectiveness of our sequence-based modeling approach. Moreover, in the structural alignment with AFDB and PDB datasets, TaxDiff significantly improves TM-score and RMSD, substantially outperforming other models, especially RMSD, which is less than half that of other models. Furthermore, the sequence consistency Fident also surpasses other models on both two datasets, showcasing the comprehensiveness and generalization capability of our model. Overall, the superior performance demonstrates TaxDiff’s enhanced ability to simulate protein sequence distributions and generate authentic and highly consistent protein sequences.

Appendix D.4 Taxonomic-Guided Sequences Generation

Table D2 Controllable generation on AFDB and PDB datasets. Metrics are calculated with 1,000 samples generated from each model. The sampling time was recorded on a single 4090.

Method	pLDDT \uparrow	Time(mins) \downarrow	AFDB Dataset			PDB Dataset		
			TM-score($\%$) \uparrow	RMSD \downarrow	Fident($\%$) \uparrow	TM-score($\%$) \uparrow	RMSD \downarrow	Fident($\%$) \uparrow
CARP	46.84 \pm 13.35	96.6	33.62	12.5753	12.5	32.38	10.8296	11.82
LRAR	47.33 \pm 14.26	79.31	26.94	17.9833	15.83	30.41	13.9929	16.33
EvoDiff	56.28 \pm 16.52	99.75	30.22	15.8664	16.17	33.02	12.3685	16.29
ESM-1b	67.91 \pm 11.59	37.4	37.13	13.6791	16.76	41.90	9.8445	17.27
TaxDiff(Our)	69.00 \pm 9.13	24.53	49.27	5.6518	25.02	48.80	4.8453	24.85

Under the taxonomic-guided, our objective is to perform controllable protein sequence generation to meet specific needs for proteins of particular species. This can be useful in realistic settings of protein and drug design where we are interested in discovering proteins with specific taxonomic preferences. For a valid comparison, we fine-tune a subset of representative networks with taxonomic-guided and same-label embedding architecture, enabling them to learn the distribution of taxonomic categories. Then, We set the conditional tax-id y to a fixed random variable representing the taxonomic group to which the protein belongs and generate 1,000 protein sequences to test the model’s controllable generation ability.

We report the numerical results in Table D2. As shown in the table, TaxDiff significantly outperforms baseline models on all the metrics, including the previous diffusion model named EvoDiff, running on the whole Uniref50. The results demonstrate that by taxonomic-guided, not only TaxDiff, almost all models acquired a higher capacity to incorporate given taxonomic information into the generation process. Furthermore, we compared the time to generate 1,000 protein sequences with these models. TaxDiff, fusing transformer into diffusion architectures, not only requires less inference time than models relying solely on either the Transformer architecture like ESM-1b or solely on the Diffusion architecture like Evodiff but also significantly outperforms other models in all metrics.

Appendix D.5 Sequences Global and Local Attention

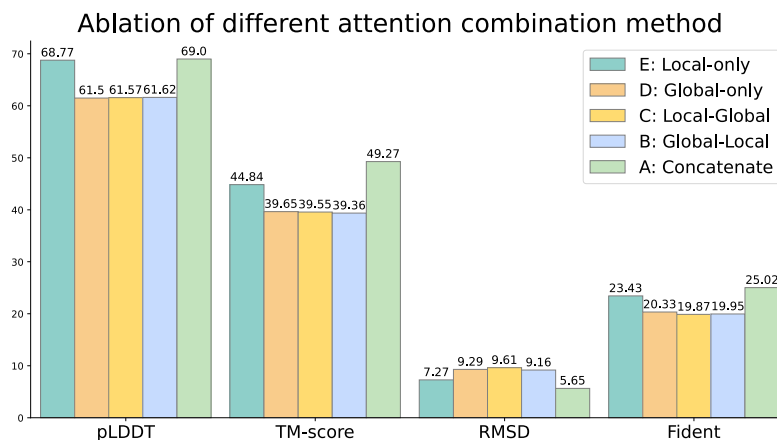


Figure D1 Ablation result of different attention combination on AFDB datasets. Metrics are calculated with 1,000 samples generated from each model.

In this section, we investigate the impact of different attention mechanisms on model performance and further analyze the effect of patchify-size at the amino acid level. Specifically, we examine the influence of five different attention combinations in Figure C1 on the protein sequence generation capability. To standardize comparison, except for the experimental variables, we keep other parameters consistent and fix the random seed. This allows the five attention combinations to generate corresponding protein sequences under 1,000 fixed random taxonomic groups. The evaluation of the generated sequences is the same as controllable sequence generation.

The visual comparison in Figure D1 indicate that the parallel use of Global-Attention and Local-Attention (Method A) achieves the best performance across all metrics, significantly surpassing other combinations. However, it is noteworthy that the exclusive use of Local-Attention (Method E) secures the second-best result, exceeding the performance of what we initially anticipated as the sub-optimal (method B) and (method C). This suggests that, in contrast to natural language, where global context is often critical, protein sequences benefit more from focusing on local interactions using Local-Attention. It implies that protein sequences also contain short sentence like local semantic structures.

To further analyze the impact of different patchify-size on features representation at the amino acid level, we divide sequences by length into segments from 4 to 64, keeping other parameters constant, and use Method A to combine Global-Attention and Local-Attention. We report the numerical results in Table D3. The results demonstrate that dividing protein sequences into 16 amino acids per local patch provides significant advantages for protein structure modeling metrics, such as TM-score and RMSD. In contrast, using a larger patchify size, like dividing protein sequences into 64 amino acids per local patch, has a more substantial impact on improving sequence-related metrics like Fident.

Table D3 Contrast experiment of different patchify-size P in Local-Attention on AFDB and PDB datasets. Metrics are calculated with 1,000 samples generated from each model.

Patchify-size	pLDDT \uparrow	Time(mins) \downarrow	AFDB Dataset			PDB Dataset		
			TM-score \uparrow	RMSD \downarrow	Fident(%) \uparrow	TM-score \uparrow	RMSD \downarrow	Fident(%) \uparrow
4	67.56 \pm 10.35	20.57	46.70	6.8047	23.01	46.58	6.2253	20.57
8	68.88 \pm 9.55	22.94	45.57	6.5407	23.94	46.04	5.4404	22.94
16	69.01 \pm 9.03	24.85	48.88	5.4992	25.88	49.58	4.6262	24.95
32	65.25 \pm 12.12	20.72	42.82	8.094	21.71	43.27	6.8687	20.72
64	70.83 \pm 8.77	24.78	46.48	5.9851	26.16	45.88	5.3077	25.20

It is worth noting that the performance of *patchify-size* = 32 is inferior to that of 16 and 64. We attribute this performance disparity to several biologically relevant factors concerning protein structure and function, as outlined below:

Secondary Structure Formation [49]: Specific length sequences favor forming essential secondary structures such as α -helices and β -sheets, which are crucial for protein function and can significantly influence biological activity. A 32 amino acid length might not support these structures well, affecting performance.

Functional Domain Integrity: The function of proteins often relies on specific functional domains, which require a certain amino acid sequence length to maintain their structural and functional integrity. A partition length of 32 amino acids may inadvertently disrupt these critical functional domains, resulting in deactivation or diminished function.

Hydrogen and Disulfide Bond Formation [50]: Hydrogen and disulfide bonds are pivotal for the stability of protein tertiary structures. The length of 32 amino acid sequences may influence the formation of these bonds, especially at sequence partition boundaries.

Appendix D.6 Scaling Ability of TaxDiff

In order to verify the scaling ability of TaxDiff, we extend the original TaxDiff (layer=12) to layer=6 (Small) and layer=18(Large) in the DiT model scaling method in fellow. In addition, when the sample number was expanded to 1,000, the experimental results in table D4 showed that our TaxDiff(layer=12) was the best, especially in TM-score, RMSD and Fident.

Table D4 Comparative experiments result TaxDiff layers on AFDB and PDB datasets. Metrics are calculated with 1,000 samples generated from each model.

Method	pLDDT↑	AFDB Dataset			PDB Dataset		
		TM-score↑	RMSD↓	Fident(%)↑	TM-score↑	RMSD↓	Fident(%)↑
Small	68.8507	43.82	6.7151	24.84	45.90	5.4525	24.36
Large	68.1949	45.62	6.4333	24.53	45.77	5.9830	23.70
TaxDiff	69.0000	49.27	5.6518	25.02	48.80	4.8453	24.85

Appendix D.7 Taxonomic Guidance Generation

To demonstrate the applicability of our model under different levels of taxonomic guidance, we evaluated the model’s ability to generate protein sequences guided by the three predominant domains within cellular organisms: Bacteria (63%), Eukaryota (33%), and Archaea (3%). The evaluation metrics were consistent with those used in the main experiments.

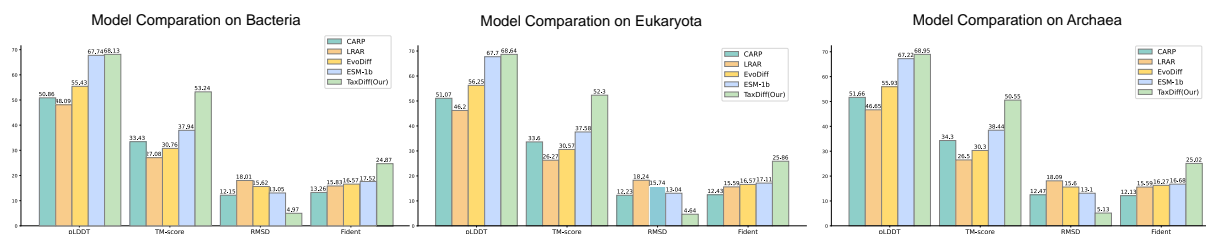


Figure D2 Taxonomic distribution and comparison of taxonomic-guided models. we compare the performance of different models under the condition of the second taxonomic level within cellular organisms.

References

- 1 John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Židek, Anna Potapenko, et al. Highly accurate protein structure prediction with alphafold. *Nature*, 596(7873):583–589, 2021.
- 2 Zhenqiao Song and Lei Li. Importance weighted expectation-maximization for protein sequence design. *arXiv preprint arXiv:2305.00386*, 2023.
- 3 Zaixiang Zheng, Yifan Deng, Dongyu Xue, Yi Zhou, Fei Ye, and Quanquan Gu. Structure-informed language models are protein designers. *bioRxiv*, pages 2023–02, 2023.
- 4 Liuzhenghao Lv, Zongying Lin, Hao Li, Yuyang Liu, Jiayi Cui, Calvin Yu-Chian Chen, Li Yuan, and Yonghong Tian. Prollama: A protein large language model for multi-task protein language processing. *arXiv e-prints*, pages arXiv–2402, 2024.
- 5 Donatas Repecka, Vykintas Jauniskis, Laurynas Karpus, Elzbieta Rembeza, Irmantas Rokaitis, Jan Zrimec, Simona Poviloniene, Audrius Laurynenas, Sandra Viknander, Wissam Abuajwa, et al. Expanding functional protein sequence spaces using generative adversarial networks. *Nature Machine Intelligence*, 3(4):324–333, 2021.
- 6 Richard J Fox, S Christopher Davis, Emily C Mundorff, Lisa M Newman, Vesna Gavrilovic, Steven K Ma, Loleta M Chung, Charlene Ching, Sarena Tam, Sheela Muley, et al. Improving catalytic function by prosar-driven enzyme evolution. *Nature biotechnology*, 25(3):338–344, 2007.
- 7 Philip A Romero and Frances H Arnold. Exploring protein fitness landscapes by directed evolution. *Nature reviews Molecular cell biology*, 10(12):866–876, 2009.
- 8 Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, 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*, 2022:500902, 2022.
- 9 Dominic Phillips, Hans-Christof Gasser, Sebestyén Kamp, Aleksander Pałkowski, Lukasz Rabalski, Diego A Oyarzún, Ajitha Rajan, and Javier Antonio Alfaro. Generating immune-aware sars-cov-2 spike proteins for universal vaccine design. In *Workshop on Healthcare AI and COVID-19*, pages 100–116. PMLR, 2022.
- 10 Surojit Biswas, Grigory Khimulya, Ethan C Alley, Kevin M Esvelt, and George M Church. Low-n protein engineering with data-efficient deep learning. *Nature methods*, 18(4):389–396, 2021.
- 11 Zheng Tan, Yan Li, Xin Wu, Ziyang Zhang, Weimei Shi, Shiqing Yang, and Wanli Zhang. De novo creation of fluorescent molecules via adversarial generative modeling. *RSC advances*, 13(2):1031–1040, 2023.
- 12 Sarah Alamdari, Nitya Thakkar, Rianne van den Berg, Alex Xijie Lu, Nicolo Fusi, Ava Pardis Amini, and Kevin K Yang. Protein generation with evolutionary diffusion: sequence is all you need. *bioRxiv*, pages 2023–09, 2023.
- 13 Alex Hawkins-Hooker, Florence Depardieu, Sebastien Baur, Guillaume Couairon, Arthur Chen, and David Bikard. Generating functional protein variants with variational autoencoders. *PLoS computational biology*, 17(2):e1008736, 2021.
- 14 Kevin E Wu, Kevin K Yang, Rianne van den Berg, James Y Zou, Alex X Lu, and Ava P Amini. Protein structure generation via folding diffusion. *arXiv preprint arXiv:2209.15611*, 2022.
- 15 Joseph L Watson, David Juergens, Nathaniel R Bennett, Brian L Trippe, Jason Yim, Helen E Eisenach, Woody Ahern, Andrew J Borst, Robert J Ragotte, Lukas F Milles, et al. De novo design of protein structure and function with rfdiffusion. *Nature*, 620(7976):1089–1100, 2023.
- 16 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020.
- 17 Jung-Eun Shin, Adam J Riesselman, Aaron W Kollasch, Conor McMahon, Elana Simon, Chris Sander, Aashish Manglik, Andrew C Kruse, and Debora S Marks. Protein design and variant prediction using autoregressive generative models. *Nature communications*, 12(1):2403, 2021.
- 18 Po-Ssu Huang, Scott E Boyken, and David Baker. The coming of age of de novo protein design. *Nature*, 537(7620):320–327, 2016.
- 19 Erik Nijkamp, Jeffrey A Ruffolo, Eli N Weinstein, Nikhil Naik, and Ali Madani. Progen2: exploring the boundaries of protein language models. *Cell Systems*, 14(11):968–978, 2023.
- 20 Uniprot: the universal protein knowledgebase in 2023. *Nucleic Acids Research*, 51(D1):D523–D531, 2023.
- 21 Ruidong Wu, Fan Ding, Rui Wang, Rui Shen, Xiwen Zhang, Shitong Luo, Chenpeng Su, Zuofan Wu, Qi Xie, Bonnie Berger, Jianzhu Ma, and Jian Peng. High-resolution de novo structure prediction from primary sequence. *bioRxiv*, 2022.
- 22 Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International conference on machine learning*, pages 2256–2265. PMLR, 2015.
- 23 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- 24 Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952*, 2023.
- 25 Chenyang Si, Ziqi Huang, Yuming Jiang, and Ziwei Liu. Free: Free lunch in diffusion u-net. *arXiv preprint arXiv:2309.11497*, 2023.
- 26 William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4195–4205, 2023.
- 27 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- 28 Ethan Perez, Florian Strub, Harm De Vries, Vincent Dumoulin, and Aaron Courville. Film: Visual reasoning with a general conditioning layer. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
- 29 Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3836–3847, 2023.
- 30 Ling Yang, Zhilong Zhang, Yang Song, Shenda Hong, Runsheng Xu, Yue Zhao, Wentao Zhang, Bin Cui, and Ming-Hsuan Yang. Diffusion models: A comprehensive survey of methods and applications. *ACM Computing Surveys*, 56(4):1–39, 2023.
- 31 Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022.
- 32 Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021.
- 33 James J Luby. Taxonomic classification and brief history. In *Apples: botany, production and uses*, pages 1–14. Cabi Publishing Wallingford UK, 2003.
- 34 David L Wheeler, Tanya Barrett, Dennis A Benson, Stephen H Bryant, Kathi Canese, Vyacheslav Chetvernin, Deanna M Church, Michael DiCuccio, Ron Edgar, Scott Federhen, et al. Database resources of the national center for biotechnology information. *Nucleic acids research*, 35(suppl.1):D5–D12, 2007.
- 35 Yaghoob Safdari, Safar Farajnia, Mohammad Asgharzadeh, and Masoumeh Khalili. Antibody humanization methods—a review and update. *Biotechnology and Genetic Engineering Reviews*, 29(2):175–186, 2013.
- 36 Andy Chi-Lung Lee, Janelle Louise Harris, Kum Kum Khanna, and Ji-Hong Hong. A comprehensive review on current advances in peptide drug development and design. *International journal of molecular sciences*, 20(10):2383, 2019.
- 37 Michel van Kempen, Stephanie S Kim, Charlotte Tumescheit, Milot Mirdita, Jeongjae Lee, Cameron LM Gilchrist, Johannes Söding, and Martin Steinegger. Fast and accurate protein structure search with foldseek. *Nature Biotechnology*, pages 1–4, 2023.
- 38 Priya Goyal, Piotr Dollár, Ross Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, Andrew Tulloch, Yangqing Jia, and Kaiming He. Accurate, large minibatch sgd: Training imagenet in 1 hour. *arXiv preprint arXiv:1706.02677*, 2017.
- 39 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.

- 40 Kevin K Yang, Nicolo Fusi, and Alex X Lu. Convolutions are competitive with transformers for protein sequence pretraining. *bioRxiv*, pages 2022–05, 2022.
- 41 Roshan M Rao, Jason Liu, Robert Verkuil, Joshua Meier, John Canny, Pieter Abbeel, Tom Sercu, and Alexander Rives. Msa transformer. In *International Conference on Machine Learning*, pages 8844–8856. PMLR, 2021.
- 42 Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Nikita Smetanin, Robert Verkuil, Ori Kabeli, Yaniv Shmueli, et al. Evolutionary-scale prediction of atomic-level protein structure with a language model. *Science*, 379(6637):1123–1130, 2023.
- 43 Noelia Ferruz, Steffen Schmidt, and Birte Höcker. Protgpt2 is a deep unsupervised language model for protein design. *Nature communications*, 13(1):4348, 2022.
- 44 Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- 45 Brian L Hie, Varun R Shanker, Duo Xu, Theodora UJ Bruun, Payton A Weidenbacher, Shaogeng Tang, Wesley Wu, John E Pak, and Peter S Kim. Efficient evolution of human antibodies from general protein language models. *Nature Biotechnology*, 42(2):275–283, 2024.
- 46 Mihaly Varadi, Stephen Anyango, Mandar Deshpande, Sreenath Nair, Cindy Natassia, Galabina Yordanova, David Yuan, Oana Stroe, Gemma Wood, Agata Laydon, et al. Alphafold protein structure database: massively expanding the structural coverage of protein-sequence space with high-accuracy models. *Nucleic acids research*, 50(D1):D439–D444, 2022.
- 47 Yang Zhang and Jeffrey Skolnick. Scoring function for automated assessment of protein structure template quality. *Proteins: Structure, Function, and Bioinformatics*, 57(4):702–710, 2004.
- 48 Helen M Berman, Tammy Battistuz, Talapady N Bhat, Wolfgang F Bluhm, Philip E Bourne, Kyle Burkhardt, Zukang Feng, Gary L Gilliland, Lisa Iype, Shri Jain, et al. The protein data bank. *Acta Crystallographica Section D: Biological Crystallography*, 58(6):899–907, 2002.
- 49 Walter Pirovano and Jaap Heringa. Protein secondary structure prediction. *Data Mining Techniques for the Life Sciences*, pages 327–348, 2010.
- 50 Thomas E Creighton. Disulfide bond formation in proteins. In *Methods in enzymology*, volume 107, pages 305–329. Elsevier, 1984.