# Two for the Price of One: Integrating Large Language Models to Learn Biophysical Interactions

Joseph D. Clark,[†,⊥] Tanner J. Dean,[‡,⊥] and Diwakar Shukla*[∗,‡,¶,§,‖]

†School of Molecular and Cellular Biology, University of Illinois at Urbana-Champaign, Urbana, IL 61801, USA

‡Center for Biophysics and Quantitative Biology, University of Illinois at Urbana-Champaign, Urbana, IL 61801, USA

¶Department of Chemical and Biomolecular Engineering, University of Illinois at Urbana-Champaign, Urbana, IL 61801, USA

§Department of Bioengineering, University of Illinois at Urbana-Champaign, Urbana, IL 61801, USA

‖Department of Chemistry, University of Illinois at Urbana-Chamapaign, Urbana, IL 61801, USA

⊥These authors contributed equally to this work.

E-mail: diwakar@illinois.edu

## Abstract

Deep learning models have become fundamental tools in drug design. In particular, large language models trained on biochemical sequences learn feature vectors that guide drug discovery through virtual screening. However, such models do not capture the molecular interactions important for binding affinity and specificity. Therefore, there

1

is a need to 'compose' representations from distinct biological modalities to effectively represent molecular complexes. We present an overview of the methods to combine molecular representations and propose that future work should balance computational efficiency and expressiveness. Specifically, we argue that improvements in both speed and accuracy are possible by learning to merge the representations from internal layers of domain specific biological language models. We demonstrate that 'composing' biochemical language models performs similar or better than standard methods representing molecular interactions despite having significantly fewer features. Finally, we discuss recent methods for interpreting and democratizing large language models that could aid the development of interaction aware foundation models for biology, as well as their shortcomings.

## Introduction

The vastness of chemical space severely limits experimental screening in drug design.[1] Advances in deep learning can help circumvent this issue by enabling large scale computational screening to identify potential drugs.[2,3] First, molecules are transformed into feature vectors (also called embeddings) which encode biochemical information (Figure 1A). In practice, embeddings can take the form of hand picked features,[4] topological encodings,[5] or the internal representations from large language models.[6] Machine learning (ML) models are then trained to predict molecular properties given embeddings as input. Popular ML models in drug discovery include random forests, support vector machines, and neural networks.[7–9] Accurate ML models enable efficient screening of molecular libraries to identify candidate molecules with desired properties. In the context of drug design, essential properties of interest include high binding affinity and specificity for protein target(s).[10,11] Therefore, ML models often predict properties of molecular complexes, such as protein-ligand, protein-protein, protein-peptide, or protein-nucleic acid interactions. However, molecular representations are typically unimodal in nature and lack explicit features describing intermolecular interactions (Figure 1B). Additionally, most molecular representations are derived from sequence alone,

2

which limits their ability to capture structural information important for binding. Finally, large language models containing up to billions of parameters produce high-dimensional embeddings with uninterpretable features that cannot be easily mapped to known biochemical concepts.

Failure to represent molecular interactions significantly reduces the performance of ML models trained on binding affinity and specificity. One solution is to join, merge, or 'compose' unimodal molecular representations to produce augmented embeddings of molecular complexes (Figure 1C). We review a variety of methods to merge standalone molecular representations into multimodal embeddings, and argue that future work should emphasize the fusion of molecular language models to balance a trade off between computational efficiency and representation ability. We emphasize recent work in natural language processing which seeks to 'compose' domain specific language models by learning to merge the representations from internal layers. We argue that these composition frameworks could become powerful tools to integrate information from different chemical modalities. We also take note of parameter efficient fine-tuning (PEFT) and language model interpretability methods that could advance the development of composed langauge models. Finally, training tasks such as pairwise contact prediction could help infuse structural information into multimodal embeddings, producing an emergent understanding of binding affinity and specificity.

## Concatenation Appends Standalone Representations

Specificity and binding affinity are properties of molecular complexes. Therefore, there is a need to merge features from distinct molecular modalities to effectively represent biological interactions. The simplest method to combine features from different modalities is to concatenate molecular representations. Let $h_a \in \mathbb{R}^m$ and $h_b \in \mathbb{R}^n$ be two vector representations for distinct molecules (e.g., a protein ligand pair). Concatenation appends the two representations to produce a new encoding $h_a h_b \in \mathbb{R}^{m+n}$ which contains features of both molecules (Figure 1D). The concatenated vector is used as the input to machine learning models that
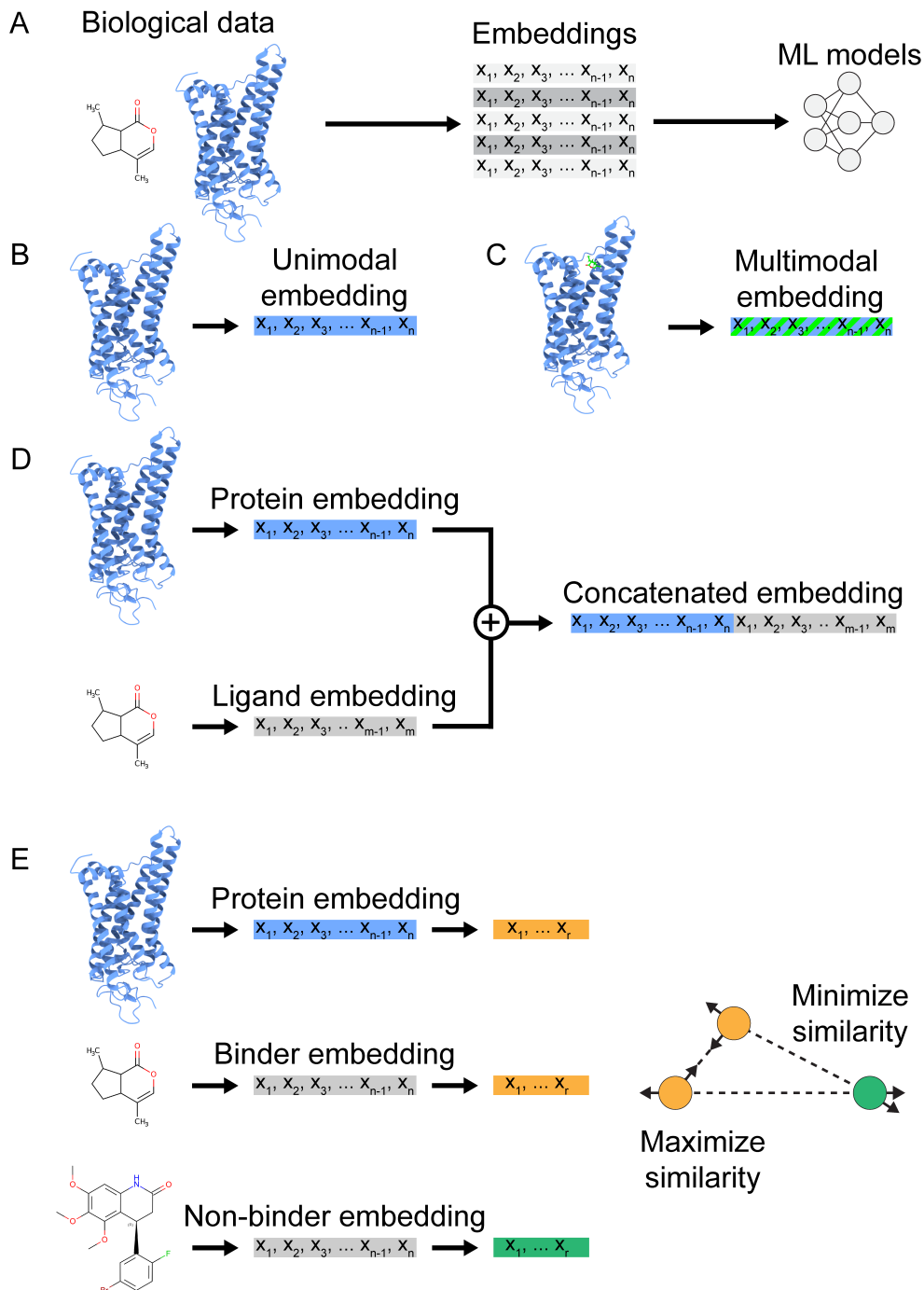
Figure 1: The workflow for computational screening and the development of multimodal embeddings. A) Sequence or structure data are mapped to feature vectors called embeddings and are used to train machine learning models to predict molecular properties. B) Unimodal embeddings contain information about a single molecule. C) Multimodal embeddings encode the features of molecular complexes. D) Concatenation appends two individual molecular embeddings to produce a larger multimodal representation. E) Contrastive learning projects two or more representations to a shared embedding space. The embeddings of interacting molecules are optimized to be closer in the latent space.

predict molecular interactions and their structural or kinetic properties.

Numerous studies combine molecular representations through concatenation.[11–16] In Enzyme Substrate Prediction (ESP),[12] protein embeddings were obtained from the protein language model ESM-1b[17] and small molecule embeddings were encoded via Graph Neural Networks (GNNs). The concatenated embeddings were used to train a gradient boosting decision tree to predict enzyme specificity. DeepDTA[13] employed a similar method to predict the binding affinity of protein-ligand pairs. Two independent convolutional neural networks were trained to produce feature vectors for proteins and ligands, then a multi-layer perceptron was trained on the concatenated representations to regress binding affinity. Other studies have explored concatenation for other modalities, such as DeepLigand[14] which concatenated Major Histocompatability Complex (MHC) class I protein feature vectors with peptide sequence representations to predict specificity. In this case, amino acid sequences were represented as one-hot encodings and Blosum50 matrix scores, and a deep residual network was trained on the conactenated embeddings. The widespread use of concatenation is largely attributed to its role as a benchmark to compare more complex embedding strategies. Concatenation is a simple and parameter free strategy that supports the combination of arbitrary features from sources such as protein or chemical language models, hand crafted or biochemical features, and traditional molecular representations such as ECFPs. However, concatenated embeddings are usually high-dimensional, potentially leading to poorer performance on small data sets, overfitting, and longer training times for ML models that scale unfavorably with the number of input features. Lastly, appending the standalone representations of molecules does not explicitly encode their interactions which significantly limits the expressive power of concatenation.

## Contrastive Learning Aligns Independent Embedding Spaces

Contrastive learning encompasses a variety of related methods originally proposed for representation learning in computer vision.[18] Broadly, contrastive learning trains embeddings

by maximizing the similarity between encoded data points with shared properties (e.g., increasing the similarity of embeddings of ligands from the same group of agonists for a known protein). Importantly, separate encoders can be employed to perform multimodal contrastive learning in which similarity is optimized between embeddings of data points from distinct modalities (Figure 1E).[19] In the context of biology, contrastive learning models often maximize the agreement between learned embeddings of interacting molecules such as protein ligand pairs.[20] Specifically, a protein target and one or more ligands are featurized by encoders such as pretrained language models or fingerprinting methods. In most cases, the encoder models are not trained during contrastive learning. Each molecular representation is then projected to a shared latent space by learnable transformations such as linear projections or feed-forward layers. Finally, a loss is calculated based on the similarity of pairs, triplets, or batches of embeddings in the shared latent space. In this framework, the probability of binding is predicted by the similarity/distance of protein and ligand embeddings in the shared latent space. The shared embedding space is structured to encode binding specificity and is typically of lower dimensionality than the original embeddings. Protein targets are usually termed as anchors, while binding and non-binding molecules are denoted as positive and negative samples respectively. Cosine similarity is typically used as the distance metric, and there are multiple related contrastive losses.[21–23] The triplet loss minimizes the euclidean distance between one anchor and one positive while maximizing the distance between the anchor and a negative.[22] Other contrastive losses such as InfoNCE[24] and NT-Xent[25] stabilize training by expanding the triplet loss to include more negative samples.

A major advantage of contrastive learning is the inference speed of the trained model due to the inexpensive nature of the cosine similarity computation. Once a large number of ligand and protein embeddings are precomputed, all pairwise interactions can be efficiently calculated as a cosine similarity matrix. Therefore, contrastive learning has the potential to dramatically expedite computational screening. However, contrastive learning may perform poorly on out-of-distribution predictions involving unseen targets that were not present as

anchors in the shared embedding space.[20] Despite this, multiple studies have successfully used contrastive learning models to predict drug-target interactions.[20,26–29]

In ConPLex,[20] protein representations were extracted from the language model Prot-Bert,[30] and small molecules were represented as Morgan Fingerprints. The embeddings were projected to a shared space via a learnable fully connected layers, and the model was jointly trained via the triplet loss and binary classification of interacting protein-ligand pairs. The probability of a drug-target interaction was interpreted as the sigmoid of the cosine similarity between protein and ligand embeddings. Similarly, BALM[26] performed contrastive learning of protein-ligand binding affinity, but replaced ProtBert and Morgan Fingerprints with the language models ESM-2[6] and ChemBERTa-2[31] respectively. A parameter efficient fine-tuning (PEFT) method was used to update the encoders during contrastive training. PEFT methods add a small number of trainable parameters to an otherwise fixed language model to learn minor modifications to the embeddings. Using PEFT to update the encoder models during contrastive training markedly improved binding affinity prediction. We provide a more detailed discussion of PEFT in a later section.

DrugCLIP[27] predicted interacting protein-ligand pairs via contrastive learning utilizing structure-aware ligand and protein binding pocket encoders. In a later study, DrugCLIP was expanded to explore the scalability of contrastive virtual screening by predicting over 10 trillion protein-ligand interactions.[32] Finally, PepPrCLIP[33] is a contrastive learning framework tailored for peptide binder design. Binding peptide-protein pairs were each encoded by ESM-2, and a shared latent space was learned via binary classification based on the cosine similarity of projected embeddings. During inference, Gaussian noise was added to the ESM-2 embeddings of known peptide binders, and the perturbed embeddings were decoded into novel sequences. Candidate sequences were ranked by their predicted binding to select for high affinity binders. While most studies use contrastive learning for prediction and design of molecular interactions, other works align protein sequence embeddings with structure representations[34] or biophysical features.[35]

Contrastive learning models produce embeddings that implicitly capture binding specificity though distance in the latent space. The flexibility of contrastive learning supports the alignment of multiple modalities and fast inference speeds. However, such models still lack an explicit model of molecular interactions, suggesting a need for more sophisticated architectures to learn multimodal molecular representations.

## Attention Learns Interactions via Dynamic Reweighting

In the early days of neural machine translation, state of the art methods consisted of recurrent neural network (RNN) encoder-decoder architectures. While these architectures performed well, a major problem was the compression of information into a fixed-length vector regardless of the length of the input sentence. These early models tended to perform poorly on sentences which extended beyond the maximum length in the training corpus. Proposed in 2014, the initial goal of the attention mechanism was to improve translation accuracy by allowing the model to focus on specific parts of the input sequence relevant to each output token.[36] Unlike prior methods, the attention mechanism enables the model to assign dynamic weights to each part of the input regardless of the sequence length. In the modern 'self-attention' mechanism, embeddings are duplicated and transformed into 'queries', 'keys', and 'values' by learned linear projections (Figure 2A). Queries and keys are multiplied and scaled to produce a square attention matrix describing the importance of each position in the sequence to each other position. The attention matrix is multiplied by the values to produce an updated representation in which the embedding for each position is a linear combination of all other embeddings in the sequence. This flexibility marked a turning point, allowing models to selectively prioritize information without relying on a rigid sequential structure like prior RNN architectures. Following in the success of the attention mechanism, the introduction of the Transformer model in 2017 revolutionized language models by incorporating self-attention as a core component.[37] This architecture enabled efficient handling of long sequences and facilitated scalability to larger datasets and models,

8

leading to the recent popularity of large language models (LLMs) like ChatGPT, Claude, and Llama.[38] Cross-attention mechanisms extend the concept of self-attention to connect information across different data modalities. Specifically, queries and key/value pairs come from distinct modalities allowing tokens from one modality (e.g., structure) to update tokens from another (e.g., residue embeddings), thus integrating complementary information (Figure 2B). In more recent years, this idea has been extended to biochemical research where information can be gained on protein-ligand tasks through cross attention of protein and ligand embeddings.

Following the success of natural language processing, many biological and chemical problems have been composed as sequence problems such as the sequence of nucleotides in DNA/RNA, sequence of amino acids in a protein, and even the SMILES string of small molecule ligands to take advantage of cross-attention's ability to connect information across sequence, structure, and biochemical features. This ability to connect information across different data has played a critical role in many recent applications of LLMs to biochemistry. For example, ChemGLaM used a cross-attention "interaction block" between the chemical language model MolFormer, and the protein language model ESM2 to improve compound-protein interaction (CPI) predictions.[6,39,40] These interaction-aware embeddings were then concatenated to ESM-2 embeddings and fed into a final fully-connected network for CPI predictions. ChemGLaM consistently outperformed previous state of the art methods for sequence only CPI predictions across four benchmarks. Additionally, recent work on lasso peptides and their corresponding cyclases has led to the development of LassoESM.[41] LassoESM uses cross attention between ESM2 and a fine-tuned language model for lasso peptides, enabling state-of-the-art performance in cyclase-peptide pair prediction along with other tasks such as the prediction of RNA polymerase inhibitory activity and non-cognate pairs of peptides and cyclases. Along with these prediction tasks, ablation studies conducted on LassoESM demonstrated the importance of cross attention over concatenation of embeddings for improved model performance. Similarly, the PepNN model leverages reciprocal
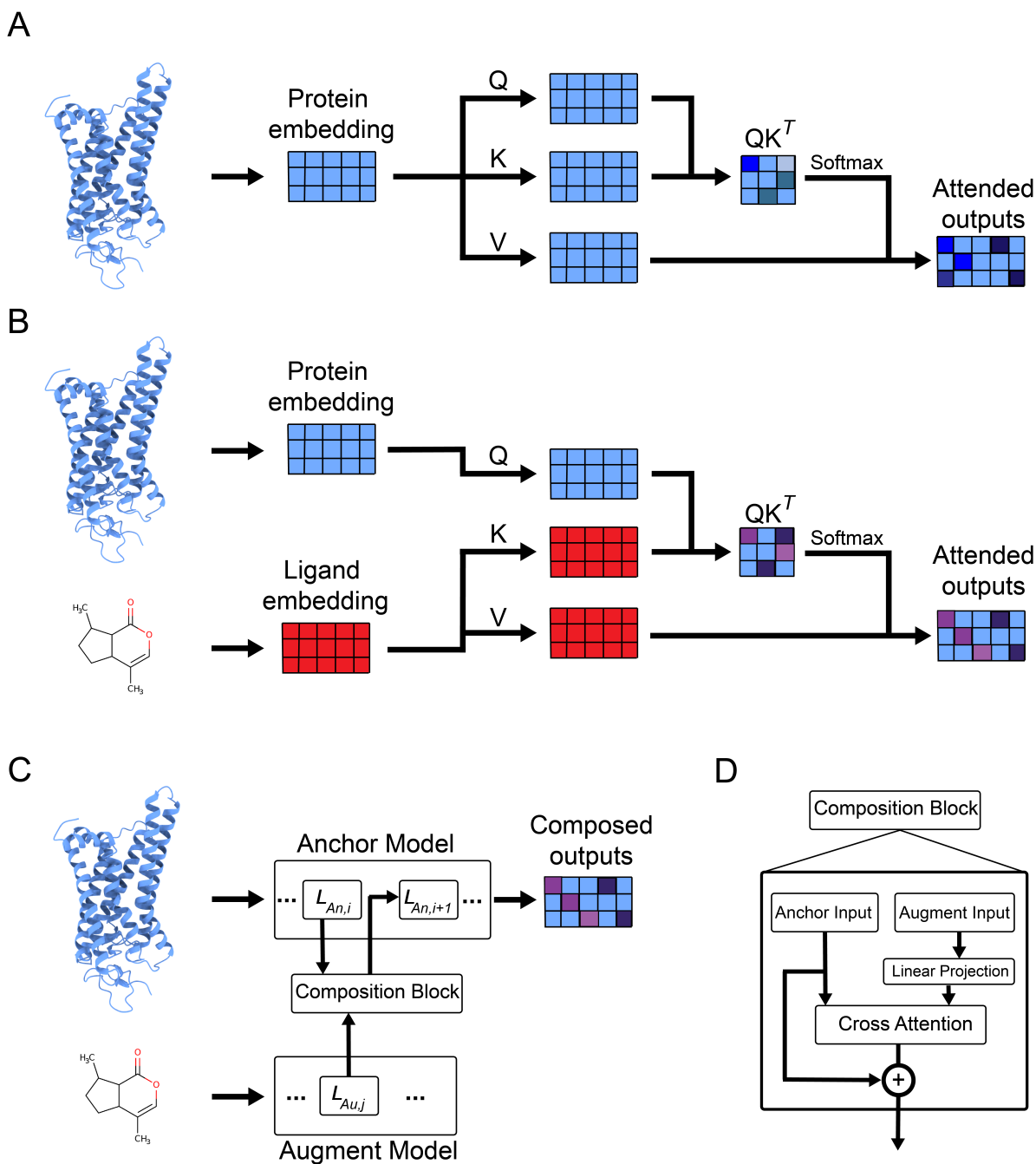
Figure 2: Schematic view of more complex methods of augmentation. A) The self-attention mechanism whereby an input tensor learns to place additional attention on specific elements of the sequence to dynamically focus on relevant parts of the input. B) The cross-attention mechanism, where a representative protein embedding acquires additional attention to key residues via a ligand embedding. C) Composition of two language models allows for the generation of ouptut embeddings of the anchor models dimensionality containing the relevant information of all augment models. D) The composition block of the CALM method uses a linear projection to match embedding dimensionalities prior to cross attention between embeddings and a skip connection into the next layer of the anchor model.

multi-head cross attention between peptide sequence and protein structure information to predict peptide binding sites.[42] The model was pre-trained on protein-protein interfaces from the Protein Data Bank[43] before being fine-tuned on protein-peptide complexes. Cross attention's ability to focus information between modalities enabled this architecture to achieve a strong performance of 0.88 ROC AUC on a per-residue binding score. Finally, CoNCISE[44] learned hierarchical discrete representations of chemical space and performed cross-attention with protein embeddings to achieve computationally efficient drug-target interaction prediction. Interestingly, cross attention mechanisms are agnostic to type of cross-modal information sharing within models, and multiple works develop strucutre-aware protein embeddings via cross attention between structure and sequence representations.[34,45]

Attention mechanisms, more specifically cross attention, revolutionized the ability to share information across modalities. The method has allowed for the integration of protein structure and ligand information among others, enhancing the prediction of interactions and binding affinities, which is crucial for drug discovery and protein-ligand docking studies. While these methods do improve feature information for cross-modal tasks, they do not address the problem of dimensionality that arises from sharing information between several embeddings leading to reduced model interpretability, a higher risk of overfitting to noise in the features, and increased computational complexity.

## Composition Merges Language Models to Integrate Domain-Specific Knowledge

Following attention, much focus has been placed on alternative methods to combine multi-modal information while avoiding the complexity of additional features. This area, known as composition, focuses on the modularity of language models which enables them to be modified and combined into larger systems for resolving more complex tasks. Composition methods merge language models in parameter space or data-flow space.[46] Composition in parameter space involves integrating the weights of multiple pretrained models with a shared

architecture. In the simplest case, model parameters are combined via linear interpolation or vector averaging.[47] Fisher-Weighted averaging[48] is a more rigorous method to average model weights while retaining maximal performance on distinct tasks. Finally, task-arithmetic[49] computes 'task vectors' defined as directions in parameter space representing the difference between pretrained and fine-tuned models. Combining the task vectors of multiple models enables multi-task capabilities. Composing models in the data flow space involves re-routing intermediate representations through distinct models.[46] For example, an intermediate language model embedding may pass through a layer from a different language model before being forwarded to the next layer.

A recent method building off of these ideas is Composition to Augment Language Models (CALM).[50] Unlike prior methods such as concatenation or attention, CALM merges two or more pretrained language models via a joint training task to achieve high performance on tasks requiring information from all composed models (Figure 2C). CALM includes a baseline model called an 'anchor' and one or more 'augment' models which infuse multimodal information by performing cross-attention between multiple internal layers. Composition blocks are composed of a linear projection layer to match the embedding dimensionalities, and a cross attention layer to combine information (Figure 2D). The attended output is then fed back into the anchor model via a residual connection (i.e., a vector addition). Composition blocks are interspaced throughout the encoder layers of a given anchor model and all augment models. During training on a joint task, the model learns to compose embeddings from distinct language models into a joint latent space while retaining the information of individual models and improving performance on tasks that leverage information from both models. Importantly, CALM retains the feature length of the anchor model while adding new information from multiple pretrained language models. Previous tasks demonstrated to work with this method have been limited to non-biological problems such as key-value arithmetic and low-resource language inclusivity. As biochemical language models continue to increase in complexity, such as ESM-3 with up to 98 billion parameters,[51] CALM-style

methods could enforce a balance of model performance and representation ability while retaining reduced feature scales for general applications. In particular, CALM frameworks applied to foundational chemical and biological language models could produce powerful multimodal encoders with minimal increase in computational cost.

## Performance Comparison of Multimodal Embeddings

We compared the performance of the four multimodal representation strategies and visualized their embeddings. We trained models to classify peptide binders of MHC Class I and II proteins using the data set from Motmaen *et. al.*[52] Protein and peptide sequences were first encoded by lightweight versions of ESM-2,[6] and 4 model architectures were trained to predict MHC specificity. A 2-layer multi-layer perceptron (MLP) was trained on concatenated peptide-protein embeddings. A contrastive learning model used two independent MLPs to project peptide/protein embeddings to a shared latent space with a cosine similarity loss for training. A cross attention model performed cross attention between the final-layer representations from each of the language models. Finally, a composition style model performed cross attention between the intermediate layers of the language models. As a baseline, we also trained 2 MLPs on unimodal peptide and protein embeddings respectively. We trained each model using 3 seeds and reported the average performance on a held-out validation set containing unseen peptide-protein pairs. For each seed, all models were trained for 100 epochs with a learning rate of $1 \times 10^{-3}$ and a batch size of 128.

Models trained solely on protein embeddings could not predict MHC specificity because a given protein may participate in multiple binding or non-binding pairs (Figure 3A). In contrast, models trained on peptide embeddings showed moderate predictive power since all MHC proteins share some degree of specificity. Concatenated peptide-protein embeddings and contrastive learning further improved classification performance. Cross attention and composition models performed the best, suggesting that these embeddings learned the most information rich features describing peptide-protein interactions. The multimodal embed-

dings of peptide-protein pairs unseen during training were then reduced to 2 dimensions with t-distributed stochastic neighbor embedding (t-SNE) for visualization. Concatenated embeddings failed to separate binding and non-binding peptide-protein pairs prior to model training (Figure 3B). The contrastive latent space grouped individual proteins into well-defined clusters based on shared specificity, and distributed peptide embeddings to ensure proximity to preferred targets (Figure 3C). Finally, cross attention and composition style models learned embeddings of peptide-protein complexes that clearly separated binding and non-binding pairs (Figure 3D, E). Of particular note, these models also clustered complexes based on the identity of the protein (individual string-like clusters in Figure 3E tend to contain similar proteins). In general, composition style models introduce minor perturbations to an original protein embedding based on the identity of a ligand, tailoring the protein embedding to context of a multimeric complex. We view composed embeddings as particularly powerful given their ability to simultaneously capture specificity while retaining valuable information about protein sequence identity.

## Improving the Accessibility and Transparency of Composed Models

The favorable properties of composed embeddings do not address the uninterpretable nature of language models or the significant computational resources needed to train them. Unfortunately, accurate language model predictions are less impactful if they do not offer direct biological insight or are inaccessible without high performance compute resources. Fortunately, recent methods from the broader natural language processing field may alleviate these issues by promoting cheaper and more transparent language model training. Parameter efficient fine-tuning (PEFT) seeks to train/fine-tune language models while updating few or none of the original parameters.[53] Instead, PEFT methods introduce lightweight 'adapters' which learn to modify an existing language model's representations towards a distinct training objective. In Low-Rank Adaption (LoRA),[54] a language model's parameters are frozen and trainable adapters are added to the attention layers (Figure 4A). An adapter projects the
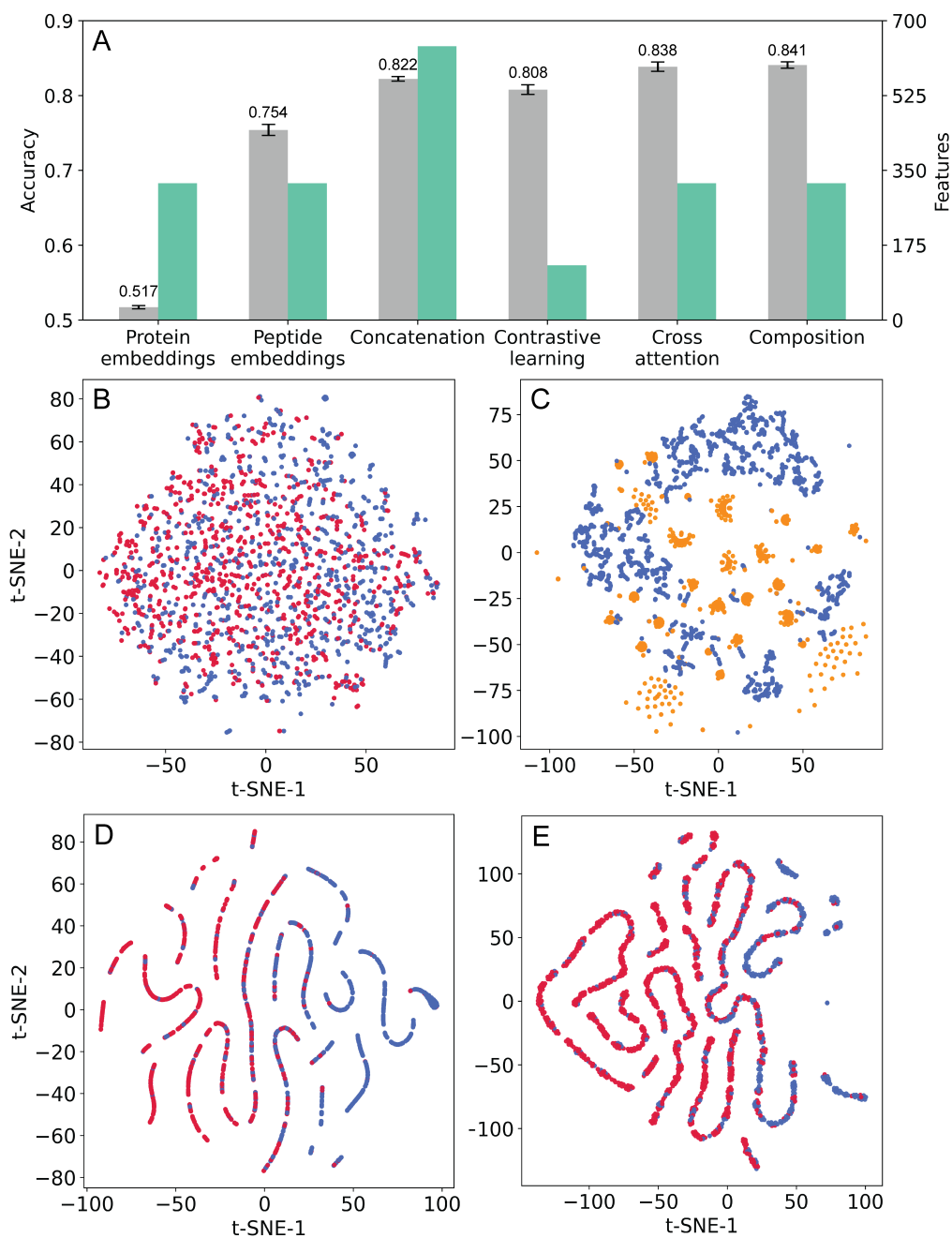
Figure 3: Comparison of multimodal embeddings on an MHC class I and II protein specificity prediction task. A) Cross attention and composed embeddings perform the best at peptide binding specificity prediction, with the embedding size shown on the right. B) Concatenated peptide-protein embeddings do not immediately distinguish between binding (red) and non-binding (blue) sequence pairs. C) Contrastive learning maps protein embeddings (orange) and peptide embeddings (blue) to a shared latent space. D) Cross attention and E) composition produce embedding spaces that cluster complexes by protein and separate binding/non-binding pairs.

attention layer's input to a low-rank representation before reshaping it to the original dimensionality. The adapter's output is then added to the output of the pretrained attention layer. The low-rank projection of the input can contain as few as one or two dimensions, leading the adapter to contain considerably fewer parameters than the pretrained layer. LoRA enables dramatic improvements in training efficiency and leads to minimal performance degradation on fine-tuning tasks.[54] This has lead to widespread use of LoRA including during training of protein language models.[55–57] Related PEFT methods find alternative modifications and decompositions of pretrained weights (see IA3[58] and DoRA[59]) or extend parameter efficient principles to language model pretraining (see GaLore[60]).

Language model interpretability remains an unsolved problem both in natural language processing and computational biology. Often, it is unclear what learned features enable models' powerful predictive performance. The high dimensionality and vast size of latent spaces make it impossible to parse interpretable features by directly analyzing molecular embeddings. Mechanistic interpretability provides insight into what a model has learned based on the assumption that language models represent more features than there are dimensions in their latent space.[61] In this method, a 'sparse-autoencoder' is trained to project language model embeddings to a substantially higher dimensional space (Figure 4B). The model is trained to reconstruct protein embeddings using a linear projection. Importantly, the high-dimensional projection is explicitly trained to be sparse (i.e., most values are zero). Given a sparse protein representation, the few active features are assumed to correspond to biochemical properties that describe the corresponding sequence. Mechanistic interpretability has identified features that are only active in the presence of specific protein properties (e.g. zinc fingers, kinase binding sites).[62] This presents and exciting opportunity to shed light on what biochemical language models have learned. Of particular note, comparing the features learned by biochemical language models and domain specific or composition models could provide valuable insight on the specific biological concepts that emerged during training.
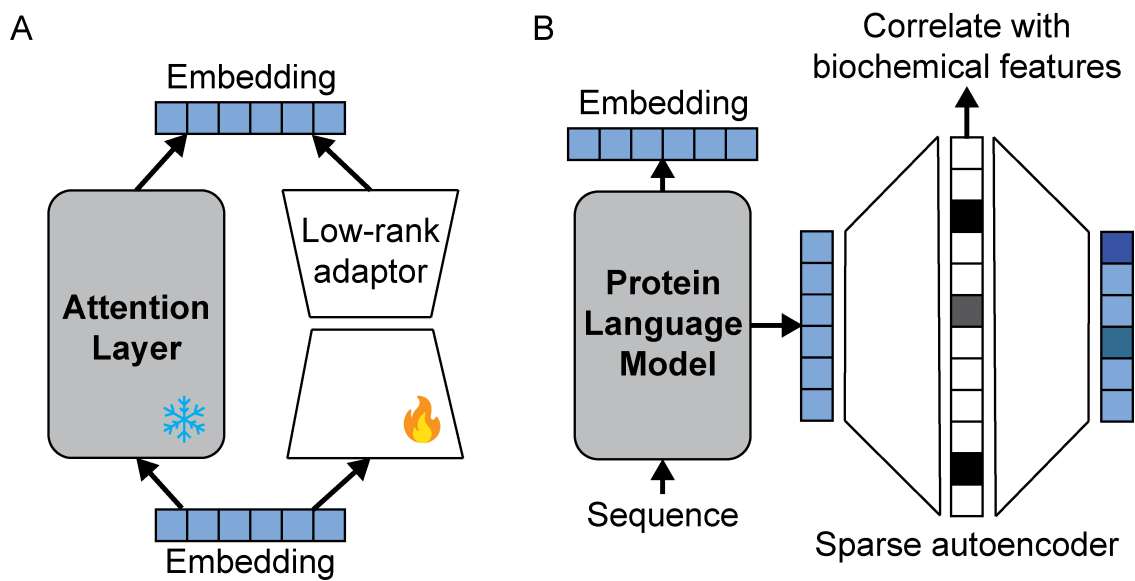
Figure 4: A schematic representation of large language model interpretability and parameter efficient fine-tuning (PEFT). A) Using sparse autoencoders, embeddings are projected to high dimensional, sparse representations which can be more easily interpreted. Sparse representations are used to reconstruct embeddings. B) In low rank adaptation (a representative PEFT method), language model representations are simultaneously modified by frozen pretrained weights, and 'adaptors' containing a small number of learnable parameters. The output of a layer is slightly modified by the adaptors, enabling efficient model training.

## Future Perspective

The development of methods to extract information from large language models—including contrastive learning, concatenation, attention mechanisms, and composition—has significantly enhanced their ability to process and generate complex language patterns. Each method offers unique advantages: contrastive learning improves model robustness through distinguishing relevant from irrelevant data; concatenation provides simplicity and flexibility by integrating diverse data streams; attention mechanisms enable efficient context capture and long-range dependencies; and composition promotes a structured and modular approach to knowledge representation.

Figure 5 presents a timeline of the major deep learning methods for merging multimodal embeddings. Contrastive learning and attention based methods have significantly grown in popularity, and methods that leverage multiple strategies have become more common. Moving forward, a key area for the development of language models lies in the integration of techniques to improve feature representations while evading the bottleneck of higher dimensions. Empirically, medium and small sized protein language models show competitive performance with larger models.[63] Recent work also suggests that protein language model representations can be compressed to significantly lower dimensionality with minimal loss of information needed for accurate structure prediction.[64,65] This calls into question the idea that larger models and more features are needed to develop high performing, multi-modal embeddings. Further, expanding the size of language models presents a challenge to their interpretability as the question of how to explain their learned features is still outstanding.[66] However, recent work suggests that protein language model features can be correlated with biological concepts through 'dictionary learning' methods called sparse auto encoders.[62] The features of sparse autoencoders can be correlated with the biochemical properties, offering insight into what the model has learned. The application of interpretability methods to biology could help discover the emergent features learned by composed models.

We propose that future work should prioritize using fewer, more salient features learned
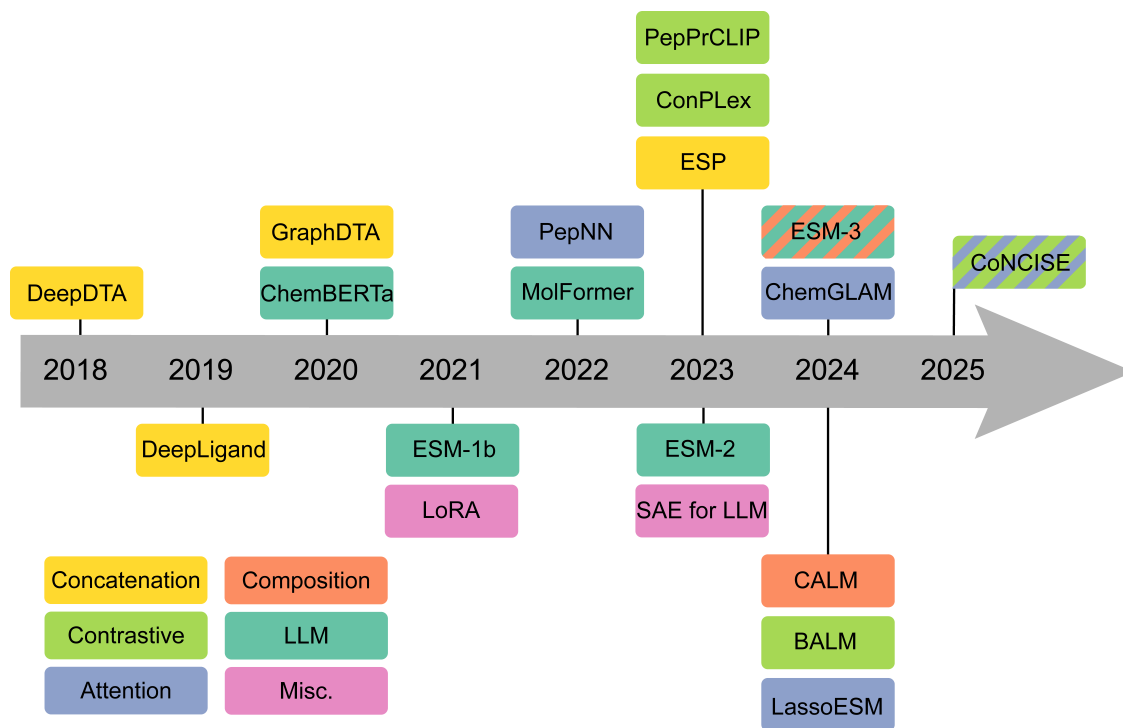
Figure 5: A timeline of deep learning methods to combine molecular representations. Methods are ordered by release date (e.g., as a preprint) rather than publication date. Attention and contrastive learning methods have grown in popularity along with hybrid methods.

from interaction aware language models based on composition inspired frameworks. In particular, training CALM-like frameworks on the features of molecular complexes such as contact maps, binding affinities, or specificity could produce powerful multimodal encoders with structure aware embeddings. Inspired by PEFT methods, CALM frameworks promote efficiency during training, as only small composition blocks must be trained instead of entire language models. The nature of the CALM framework, coupled with recent work on domain-specific biological language models[41,67–70] makes it a promising technique for developing augmented language models of molecular interactions.

# References

(1) Reymond, J.-L.; van Deursen, R.; Blum, L. C.; Ruddigkeit, L. Chemical space as a source for new drugs. *MedChemComm* **2010**, *1*, 30–38.

(2) Sadybekov, A. V.; Katritch, V. Computational approaches streamlining drug discovery. *Nature* **2023**, *616*, 673–685.

(3) Wu, H.; Liu, J.; Zhang, R.; Lu, Y.; Cui, G.; Cui, Z.; Ding, Y. A review of deep learning methods for ligand based drug virtual screening. *Fundamental Research* **2024**, *4*, 715–737.

(4) Wei, L.; Su, R.; Wang, B.; Li, X.; Zou, Q.; Gao, X. Integration of deep feature representations and handcrafted features to improve the prediction of N6-methyladenosine sites. *Neurocomputing* **2019**, *324*, 3–9.

(5) Rogers, D.; Hahn, M. Extended-Connectivity Fingerprints. *Journal of Chemical Information and Modeling* **2010**, *50*, 742–754.

(6) Lin, Z.; Akin, H.; Rao, R.; Hie, B.; Zhu, Z.; Lu, W.; Smetanin, N.; Verkuil, R.; Kabeli, O.; Shmueli, Y.; dos Santos Costa, A.; Fazel-Zarandi, M.; Sercu, T.; Candido, S.; Rives, A. Evolutionary-scale prediction of atomic-level protein structure with a language model. *Science* **2023**, *379*, 1123–1130.

(7) Rodríguez-Pérez, R.; Bajorath, J. Evolution of support vector machine and regression modeling in chemoinformatics and drug discovery. *Journal of Computer-Aided Molecular Design* **2022**, *36*, 355–362.

(8) Kapsiani, S.; Howlin, B. J. Random forest classification for predicting lifespan-extending chemical compounds. *Scientific reports* **2021**, *11*, 13812.

(9) Krishnan, S. R.; Bung, N.; Bulusu, G.; Roy, A. Accelerating de novo drug design against novel proteins using deep learning. *Journal of Chemical Information and Modeling* **2021**, *61*, 621–630.

(10) Puszkarska, A. M.; Taddese, B.; Revell, J.; Davies, G.; Field, J.; Hornigold, D. C.; Buchanan, A.; Vaughan, T. J.; Colwell, L. J. Machine learning designs new

GCGR/GLP-1R dual agonists with enhanced biological potency. *Nature Chemistry* **2024**, *16*, 1436–1444.

(11) Abbasi Mesrabadi, H.; Faez, K.; Pirgazi, J. Drug–target interaction prediction based on protein features, using wrapper feature selection. *Scientific Reports* **2023**, *13*, 3594.

(12) Kroll, A.; Ranjan, S.; Engqvist, M. K. M.; Lercher, M. J. A general model to predict small molecule substrates of enzymes based on machine and deep learning. *Nat. Commun.* **2023**, *14*, 2787.

(13) Öztürk, H.; Özgür, A.; Ozkirimli, E. DeepDTA: deep drug–target binding affinity prediction. *Bioinformatics* **2018**, *34*, i821–i829.

(14) Zeng, H.; Gifford, D. K. DeepLigand: accurate prediction of MHC class I ligands using peptide embedding. *Bioinformatics* **2019**, *35*, i278–i283.

(15) Nguyen, T.; Le, H.; Quinn, T. P.; Nguyen, T.; Le, T. D.; Venkatesh, S. GraphDTA: predicting drug–target binding affinity with graph neural networks. *Bioinformatics* **2020**, *37*, 1140–1147.

(16) Jiang, M.; Li, Z.; Zhang, S.; Wang, S.; Wang, X.; Yuan, Q.; Wei, Z. Drug–target affinity prediction using graph neural network and contact maps. *RSC Advances* **2020**, *10*, 20701–20712.

(17) Rives, A.; Meier, J.; Sercu, T.; Goyal, S.; Lin, Z.; Liu, J.; Guo, D.; Ott, M.; Zitnick, C. L.; Ma, J.; Fergus, R. Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences. *Proceedings of the National Academy of Sciences* **2021**, *118*, e2016239118.

(18) Le-Khac, P. H.; Healy, G.; Smeaton, A. F. Contrastive Representation Learning: A Framework and Review. *IEEE Access* **2020**, *8*, 193907–193934.

(19) Yuan, X.; Lin, Z.; Kuen, J.; Zhang, J.; Wang, Y.; Maire, M.; Kale, A.; Faieta, B. Multimodal Contrastive Training for Visual Representation Learning. *arXiv* **2021**, 10.48550/ARXIV.2104.12836.

(20) Singh, R.; Sledzieski, S.; Bryson, B.; Cowen, L.; Berger, B. Contrastive learning in protein language space predicts interactions between drugs and protein targets. *Proceedings of the National Academy of Sciences* **2023**, *120*, e2220778120.

(21) Chopra, S.; Hadsell, R.; LeCun, Y. Learning a similarity metric discriminatively, with application to face verification. 2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05). 2005; pp 539–546 vol. 1.

(22) Schroff, F.; Kalenichenko, D.; Philbin, J. FaceNet: A Unified Embedding for Face Recognition and Clustering. *arXiv* **2015**, 10.48550/ARXIV.1503.03832.

(23) Song, H. O.; Xiang, Y.; Jegelka, S.; Savarese, S. Deep Metric Learning via Lifted Structured Feature Embedding. *arXiv* **2015**, 10.48550/ARXIV.1511.06452.

(24) Oord, A. v. d.; Li, Y.; Vinyals, O. Representation Learning with Contrastive Predictive Coding. *arXiv* **2018**, 10.48550/ARXIV.1807.03748.

(25) Chen, T.; Kornblith, S.; Norouzi, M.; Hinton, G. A Simple Framework for Contrastive Learning of Visual Representations. *arXiv* **2020**, 10.48550/ARXIV.2002.05709.

(26) Gorantla, R.; Gema, A. P.; Yang, I. X.; Serrano-Morras, A.; Suutari, B.; Jimenez, J. J.; Mey, A. S. J. S. Learning Binding Affinities via Fine-tuning of Protein and Ligand Language Models. *bioRxiv* **2024**, 10.1101/2024.11.01.621495.

(27) Gao, B.; Qiang, B.; Tan, H.; Ren, M.; Jia, Y.; Lu, M.; Liu, J.; Ma, W.; Lan, Y. Drug-CLIP: Contrastive Protein-Molecule Representation Learning for Virtual Screening. *arXiv* **2023**, 10.48550/ARXIV.2310.06367.

(28) Xu, H.; You, Y.; Shen, Y. Multi-Modal Contrastive Learning for Proteins by Combining Domain-Informed Views. *ICLR 2024 Workshop on Machine Learning for Genomics Explorations* **2024**, https://openreview.net/forum?id=xDcTugulVV.

(29) Du, Z.; Fu, W.; Guo, X.; Caragea, D.; Li, Y. FusionESP: Improved enzyme-substrate pair prediction by fusing protein and chemical knowledge. *bioRxiv* **2024**, 10.1101/2024.08.13.607829.

(30) Brandes, N.; Ofer, D.; Peleg, Y.; Rappoport, N.; Linial, M. ProteinBERT: a universal deep-learning model of protein sequence and function. *Bioinformatics* **2022**, *38*, 2102–2110.

(31) Ahmad, W.; Simon, E.; Chithrananda, S.; Grand, G.; Ramsundar, B. ChemBERTa-2: Towards Chemical Foundation Models. *arXiv* **2022**, 10.48550/ARXIV.2209.01712.

(32) Jia, Y. et al. Deep contrastive learning enables genome-wide virtual screening. *arXiv* **2024**, 10.1101/2024.09.02.610777.

(33) Bhat, S. et al. De novo design of peptide binders to conformationally diverse targets with contrastive language modeling. *Science Advances* **2025**, *11*, eadr8638.

(34) Wang, D.; Pourmirzaei, M.; Abbas, U. L.; Zeng, S.; Manshour, N.; Esmaili, F.; Poudel, B.; Jiang, Y.; Shao, Q.; Chen, J.; Xu, D. S-PLM: Structure-Aware Protein Language Model via Contrastive Learning Between Sequence and Structure. *Advanced Science* **2024**, *12*, 2404212.

(35) Peng, Y.; Wu, J.; Sun, Y.; Zhang, Y.; Wang, Q.; Shao, S. Contrastive-learning of language embedding and biological features for cross modality encoding and effector prediction. *Nature Communications* **2025**, *16*, 1299.

(36) Bahdanau, D.; Cho, K.; Bengio, Y. Neural Machine Translation by Jointly Learning to Align and Translate. *arXiv* **2014**, 10.48550/ARXIV.1409.0473.

(37) Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A. N.; Kaiser, L.; Polosukhin, I. Attention Is All You Need. *arXiv* **2017**, 10.48550/ARXIV.1706.03762.

(38) Touvron, H.; Lavril, T.; Izacard, G.; Martinet, X.; Lachaux, M.-A.; Lacroix, T.; Rozière, B.; Goyal, N.; Hambro, E.; Azhar, F.; Rodriguez, A.; Joulin, A.; Grave, E.; Lample, G. LLaMA: Open and Efficient Foundation Language Models. *arXiv* **2023**, 10.48550/ARXIV.2302.13971.

(39) Koyama, T.; Tsumura, H.; Matsumoto, S.; Okita, R.; Kojima, R.; Okuno, Y. ChemGLaM: Chemical-Genomics Language Models for Compound-Protein Interaction Prediction. *bioRxiv* **2024**, 10.1101/2024.02.13.580100.

(40) Ross, J.; Belgodere, B.; Chenthamarakshan, V.; Padhi, I.; Mroueh, Y.; Das, P. Large-scale chemical language representations capture molecular structure and properties. *Nature Machine Intelligence* **2022**, *4*, 1256–1264.

(41) Mi, X.; Barrett, S. E.; Mitchell, D. A.; Shukla, D. LassoESM: A tailored language model for enhanced lasso peptide property prediction. *bioRxiv* **2024**, 10.1101/2024.10.25.620295.

(42) Abdin, O.; Nim, S.; Wen, H.; Kim, P. M. PepNN: a deep attention model for the identification of peptide binding sites. *Communications Biology* **2022**, *5*, 503.

(43) Berman, H. M. The Protein Data Bank. *Nucleic Acids Research* **2000**, *28*, 235–242.

(44) Erden, M.; Devkota, K.; Varghese, L.; Cowen, L.; Singh, R. Learning a CoNCISE language for small-molecule binding. *bioRxiv* **2025**, 10.1101/2025.01.08.632039.

(45) Tan, Y.; Li, M.; Zhou, B.; Zhong, B.; Zheng, L.; Tan, P.; Zhou, Z.; Yu, H.; Fan, G.; Hong, L. Simple, Efficient, and Scalable Structure-Aware Adapter Boosts Protein Language Models. *Journal of Chemical Information and Modeling* **2024**, *64*, 6338–6349.

(46) Akiba, T.; Shing, M.; Tang, Y.; Sun, Q.; Ha, D. Evolutionary optimization of model merging recipes. *Nature Machine Intelligence* **2025**, *7*, 195–204.

(47) Wortsman, M.; Ilharco, G.; Gadre, S. Y.; Roelofs, R.; Gontijo-Lopes, R.; Morcos, A. S.; Namkoong, H.; Farhadi, A.; Carmon, Y.; Kornblith, S.; Schmidt, L. Model soups: averaging weights of multiple fine-tuned models improves accuracy without increasing inference time. *arXiv* **2022**, 10.48550/ARXIV.2203.05482.

(48) Matena, M.; Raffel, C. Merging Models with Fisher-Weighted Averaging. *arXiv* **2021**, 10.48550/ARXIV.2111.09832.

(49) Ilharco, G.; Ribeiro, M. T.; Wortsman, M.; Gururangan, S.; Schmidt, L.; Hajishirzi, H.; Farhadi, A. Editing Models with Task Arithmetic. *arXiv* **2022**, 10.48550/ARXIV.2212.04089.

(50) Bansal, R.; Samanta, B.; Dalmia, S.; Gupta, N.; Vashishth, S.; Ganapathy, S.; Bapna, A.; Jain, P.; Talukdar, P. LLM Augmented LLMs: Expanding Capabilities through Composition. *arXiv* **2024**, 10.48550/ARXIV.2401.02412.

(51) Hayes, T. et al. Simulating 500 million years of evolution with a language model. *Science* **2025**, *387*, 850–858.

(52) Motmaen, A.; Dauparas, J.; Baek, M.; Abedi, M. H.; Baker, D.; Bradley, P. Peptide-binding specificity prediction using fine-tuned protein structure prediction networks. *Proceedings of the National Academy of Sciences* **2023**, *120*, e2216697120.

(53) Xu, L.; Xie, H.; Qin, S.-Z. J.; Tao, X.; Wang, F. L. Parameter-Efficient Fine-Tuning Methods for Pretrained Language Models: A Critical Review and Assessment. *arXiv* **2023**, 10.48550/ARXIV.2312.12148.

(54) Hu, E. J.; Shen, Y.; Wallis, P.; Allen-Zhu, Z.; Li, Y.; Wang, S.; Wang, L.;

Chen, W. LoRA: Low-Rank Adaptation of Large Language Models. *arXiv* **2021**, 10.48550/ARXIV.2106.09685.

(55) Sledzieski, S.; Kshirsagar, M.; Baek, M.; Dodhia, R.; Lavista Ferres, J.; Berger, B. Democratizing protein language models with parameter-efficient fine-tuning. *Proceedings of the National Academy of Sciences* **2024**, *121*, e2405840121.

(56) Zhou, Z.; Zhang, L.; Yu, Y.; Wu, B.; Li, M.; Hong, L.; Tan, P. Enhancing efficiency of protein language models with minimal wet-lab data through few-shot learning. *Nature Communications* **2024**, *15*, 5566.

(57) Schmirler, R.; Heinzinger, M.; Rost, B. Fine-tuning protein language models boosts predictions across diverse tasks. *Nature Communications* **2024**, *15*, 7407.

(58) Liu, H.; Tam, D.; Muqeeth, M.; Mohta, J.; Huang, T.; Bansal, M.; Raffel, C. Few-Shot Parameter-Efficient Fine-Tuning is Better and Cheaper than In-Context Learning. *arXiv* **2022**, 10.48550/ARXIV.2205.05638.

(59) Liu, S.-Y.; Wang, C.-Y.; Yin, H.; Molchanov, P.; Wang, Y.-C. F.; Cheng, K.-T.; Chen, M.-H. DoRA: Weight-Decomposed Low-Rank Adaptation. *arXiv* **2024**, 10.48550/ARXIV.2402.09353.

(60) Zhao, J.; Zhang, Z.; Chen, B.; Wang, Z.; Anandkumar, A.; Tian, Y. GaLore: Memory-Efficient LLM Training by Gradient Low-Rank Projection. *arXiv* **2024**, 10.48550/ARXIV.2403.03507.

(61) Bricken, T. et al. Towards Monosemanticity: Decomposing Language Models With Dictionary Learning. *Transformer Circuits Thread* **2023**, https://transformer-circuits.pub/2023/monosemantic-features/index.html.

(62) Simon, E.; Zou, J. InterPLM: Discovering Interpretable Features in Protein Language Models via Sparse Autoencoders. *bioRxiv* **2024**, 10.1101/2024.11.14.623630.

(63) Vieira, L. C.; Handojo, M. L.; Wilke, C. O. Scaling down for efficiency: Medium-sized protein language models perform well at transfer learning on realistic datasets. *bioRxiv* **2024**, 10.1101/2024.11.22.624936.

(64) Lu, A. X.; Yan, W.; Yang, K. K.; Gligorijevic, V.; Cho, K.; Abbeel, P.; Bonneau, R.; Frey, N. Tokenized and Continuous Embedding Compressions of Protein Sequence and Structure. *bioRxiv* **2024**, 10.1101/2024.08.06.606920.

(65) Lu, A. X.; Yan, W.; Robinson, S. A.; Yang, K. K.; Gligorijevic, V.; Cho, K.; Bonneau, R.; Abbeel, P.; Frey, N. Generating All-Atom Protein Structure from Sequence-Only Training Data. *bioRxiv* **2024**, 10.1101/2024.12.02.626353.

(66) Singh, C.; Inala, J. P.; Galley, M.; Caruana, R.; Gao, J. Rethinking Interpretability in the Era of Large Language Models. *arXiv* **2024**, 10.48550/ARXIV.2402.01761.

(67) Zeng, W.; Dou, Y.; Pan, L.; Xu, L.; Peng, S. Improving prediction performance of general protein language model by domain-adaptive pretraining on DNA-binding protein. *Nature Communications* **2024**, *15*, 7838.

(68) Wang, Y.; Lv, H.; Teo, Q. W.; Lei, R.; Gopal, A. B.; Ouyang, W. O.; Yeung, Y.-H.; Tan, T. J.; Choi, D.; Shen, I. R.; Chen, X.; Graham, C. S.; Wu, N. C. An explainable language model for antibody specificity prediction using curated influenza hemagglutinin antibodies. *Immunity* **2024**, *57*, 2453–2465.e7.

(69) Clark, J. D.; Mi, X.; Mitchell, D. A.; Shukla, D. Substrate prediction for RiPP biosynthetic enzymes via masked language modeling and transfer learning. *Digital Discovery* **2024**, *4*, 343–354.

(70) Vincoff, S.; Goel, S.; Kholina, K.; Pulugurta, R.; Vure, P.; Chatterjee, P. FusOn-pLM: a fusion oncoprotein-specific language model via adjusted rate masking. *Nature Communications* **2025**, *16*, 1436.