

ProteinAligner: A Multi-modal Pretraining Framework for Protein Foundation Models

Li Zhang^{1,*}, Han Guo^{1,*}, Leah Schaffer², Young Su Ko³, Digvijay Singh⁴, Hamid Rahmani⁵, Danielle Grotjahn⁵, Elizabeth Villa^{4,6}, Michael Gilson⁷, Wei Wang^{3,8}, Trey Ideker^{2,9}, Eric Xing^{10,11}, and Pengtao Xie^{1,2,11} ✉

¹Department of Electrical and Computer Engineering, University of California San Diego, La Jolla, CA 92093, USA

²Department of Medicine, University of California San Diego, La Jolla, CA 92093, USA

³Department of Chemistry and Biochemistry, University of California San Diego, La Jolla, CA 92093, USA

⁴School of Biological Sciences, University of California San Diego, La Jolla, CA 92093, USA

⁵Department of Integrative Structural and Computational Biology, The Scripps Research Institute, La Jolla, CA 92037, USA

⁶Howard Hughes Medical Institute, University of California San Diego, La Jolla, CA 92093, USA

⁷Skaggs School of Pharmacy and Pharmaceutical Sciences, University of California San Diego, La Jolla, CA 92093, USA

⁸Department of Cellular and Molecular Medicine, University of California San Diego, La Jolla, CA 92093, USA

⁹Department of Bioengineering, University of California San Diego, La Jolla, CA 92093, USA

¹⁰Machine Learning Department, Carnegie Mellon University, Pittsburgh, PA 15213, USA

¹¹Mohamed bin Zayed University of Artificial Intelligence, Masdar City, Abu Dhabi, UAE

Protein foundation models, particularly protein language models, have demonstrated strong success in learning meaningful representations of proteins using transformer architectures pretrained on large-scale protein datasets with self-supervised learning. These representations have been highly effective for downstream tasks such as predicting protein functions and properties. However, most current protein foundation models focus on pretraining with amino acid sequences, often neglecting additional modalities like protein structures and related literature, both of which provide valuable insights. To address this gap, we propose a multi-modal pretraining approach that integrates three key modalities - protein sequences, structures, and literature text. In our framework, the protein sequence modality serves as the anchor, with the other two modalities aligned to it, enhancing the model's capacity to capture more comprehensive protein information. ProteinAligner outperformed state-of-the-art protein foundation models in predicting protein functions and properties across diverse downstream tasks.

Protein foundation model, multi-modal learning, protein function prediction, protein property prediction

Correspondence: p1xie@ucsd.edu

* Equal contribution

Introduction

Proteins play a fundamental role in virtually all biological processes. Understanding their functions and properties is central to advancing fields such as drug discovery (1), diagnostics (2), and biotechnology (3). Recent advances in artificial intelligence, particularly in transformer-based models (4), have led to the development of protein foundation models capable of learning rich representations from large-scale protein datasets (5–11). These models, particularly protein language models (PLMs) (5–7, 10, 11), have shown remarkable success in performing various downstream tasks such as protein function prediction (12, 13), property prediction (14, 15), structure prediction (16, 17), and protein design (18, 19).

Despite these successes, current PLMs predominantly focus on amino acid sequences while overlooking the wealth

of complementary information available in other modalities. Protein structures, for example, provide critical three-dimensional information that is essential for understanding how proteins fold and interact with other molecules, directly influencing their biological functions (20). The spatial arrangement of amino acids, which governs interactions such as binding affinities and functional sites, cannot be readily inferred from sequence data alone (21, 22), making the integration of structural data crucial for a more comprehensive understanding of protein behavior. Similarly, the vast amount of biological literature contains experimentally validated insights into protein mechanisms, behavior, and interactions that are often context-specific and difficult to infer from sequences or structures alone (23, 24). Literature captures critical information about post-translational modifications, protein dynamics in various environments, and interaction networks - details accumulated from years of experimental studies. By incorporating these additional modalities - protein structures and related literature - protein foundation models can move beyond sequence prediction to a more robust, context-aware understanding of protein biology. This multi-modal integration has the potential to greatly enhance the representational power of these models, enabling more accurate predictions of protein functions and behaviors in diverse biological scenarios. While some studies have explored the use of two modalities, such as protein sequence and text (24–26) or protein structure and sequence (8, 9, 27), the simultaneous integration of three modalities for pretraining protein foundation models remains underexplored.

To address this gap, we introduce ProteinAligner, a multi-modal pretraining framework that combines protein sequences, structures, and literature text. Our framework aligns these modalities with the protein sequence as the anchor, enabling the model to learn richer and more comprehensive representations of proteins. By integrating diverse protein-related data, ProteinAligner improves the model's ability to capture intricate biological phenomena, paving the way for more accurate predictions of protein

057 functions and properties. ProteinAligner utilizes three
058 specialized encoders - a sequence encoder, a structure
059 encoder, and a text encoder - to learn representations for
060 each modality. These distinct representations are projected
061 into a shared latent space, enabling direct comparison
062 across modalities. By employing a contrastive alignment
063 strategy (28), ProteinAligner uses protein sequences as
064 the anchor to align corresponding structures and textual
065 descriptions, encouraging similar representations for the
066 same protein and dissimilar representations for different
067 proteins. This approach not only maximizes data utiliza-
068 tion by allowing pretraining on incomplete modality data
069 but also captures the comprehensive biological insights
070 provided by each modality. Recently, another independent
071 study (29) also explored learning protein representations
072 across three modalities. Their work and ours were developed
073 independently, with approaches and experiments conceived
074 separately.

075
076 ProteinAligner demonstrated superior performance
077 compared to state-of-the-art baselines across various
078 downstream prediction tasks, including detecting type I
079 anti-CRISPR activities, predicting pathogenic missense
080 variants, predicting protein thermostability, identifying
081 potent bioactive peptides, and estimating the minimum
082 inhibitory concentration of antimicrobial peptides.

084 Results

085 **ProteinAligner overview.** ProteinAligner is a multi-modal
086 foundation model for protein representation learning,
087 integrating three distinct modalities: amino acid (AA)
088 sequences, 3D structures, and textual descriptions of pro-
089 teins. ProteinAligner contains three encoders - a protein
090 sequence encoder, a protein structure encoder, and a text
091 encoder - each dedicated to learning representations for its
092 corresponding modality (Fig. 1a). The protein sequence en-
093 coder is a protein language model that uses the transformer
094 architecture (4) to extract a representation for the input AA
095 sequence. It represents each AA as a token and employs
096 self-attention (4) to capture long-range dependencies among
097 AAs. The protein structure encoder utilizes Geometric
098 Vector Perceptron Graph Neural Network (GVP-GNN) (30)
099 layers for geometric representation learning of the input
100 protein structure, followed by transformer layers that capture
101 long-range interactions between atomic coordinates. The
102 text encoder employs a transformer architecture, utilizing
103 self-attention to capture long-range dependencies between
104 language tokens. Specifically, we employed ESM (5), a
105 leading protein language model, as the protein sequence
106 encoder, and ESM-IF1 (9) as the protein structure encoder.
107 ESM consists of 33 transformer layers and 650 million
108 parameters, pretrained on 65 million protein sequences.
109 ESM-IF1 features 20 layers and 124 million parameters,
110 pretrained on 12 million computed protein structures and
111 16,000 experimentally verified structures. The text encoder
112 includes eight transformer layers with a total of 78 million
113

parameters. ProteinAligner uses modality-specific linear
projection modules to map the representations extracted by
different encoders into a shared latent space with matching
dimensions, ensuring that representations from different
modalities are directly comparable. ProteinAligner consists
of 867 million model parameters in total.

ProteinAligner performs joint pretraining of the three
encoders by leveraging a modality alignment strategy, using
protein sequences as the anchor modality to align the other
two modalities (Fig. 1a). Specifically, given a protein
sequence and a protein structure, if they correspond to the
same protein, ProteinAligner encourages their represen-
tations to be similar, and dissimilar otherwise. The same
principle applies for protein text and protein sequences,
with representations aligned if they refer to the same protein
and separated if not. This alignment is accomplished by
minimizing contrastive losses (28, 31) defined on the rep-
resentations of sequence-structure pairs and sequence-text
pairs. ProteinAligner does not require all three modalities to
be present simultaneously for each protein in the pretraining
data. The alignment can be performed as long as the protein
sequence and at least one additional modality - either struc-
ture or text - are available. We chose protein sequences as
the anchor for alignment because they are the most prevalent
data modality in protein databases; nearly every protein has
an associated amino acid sequence, whereas information
on structures or textual descriptions is often incomplete.
By using sequences as the anchor, we can maximize data
utilization, ensuring the inclusion of as many proteins as
possible in the alignment process. With pretrained encoders
in place, they can be fine-tuned on task-specific data to
handle a variety of downstream tasks. During this process,
the encoders are integrated with task-specific modules,
creating models that are customized for specific prediction
tasks.

We curated a large-scale pretraining dataset for Pro-
teinAligner by integrating data from the UniProtKB/Swiss-
Prot (32) and RCSB PDB (33) databases. The dataset
consists of 290,480 proteins, each with an amino acid
sequence and a corresponding textual description. 133,726
of them are also associated with protein structures. In
total, the dataset contains 133,726 sequence-structure pairs
and 290,480 sequence-text pairs. The textual descriptions
provide information about the proteins' functions. Both the
structures and the functional descriptions were experimen-
tally validated and reviewed by domain experts. Fig. 1b
shows the distribution of protein taxonomy, functions, and
types in the dataset.

ProteinAligner detects type I anti-CRISPR activities.

We evaluated the effectiveness of ProteinAligner in detect-
ing type I anti-CRISPR (Acr) activities. Acr proteins are
produced by certain viruses, such as bacteriophages, or
mobile genetic elements to inhibit the type I CRISPR-Cas

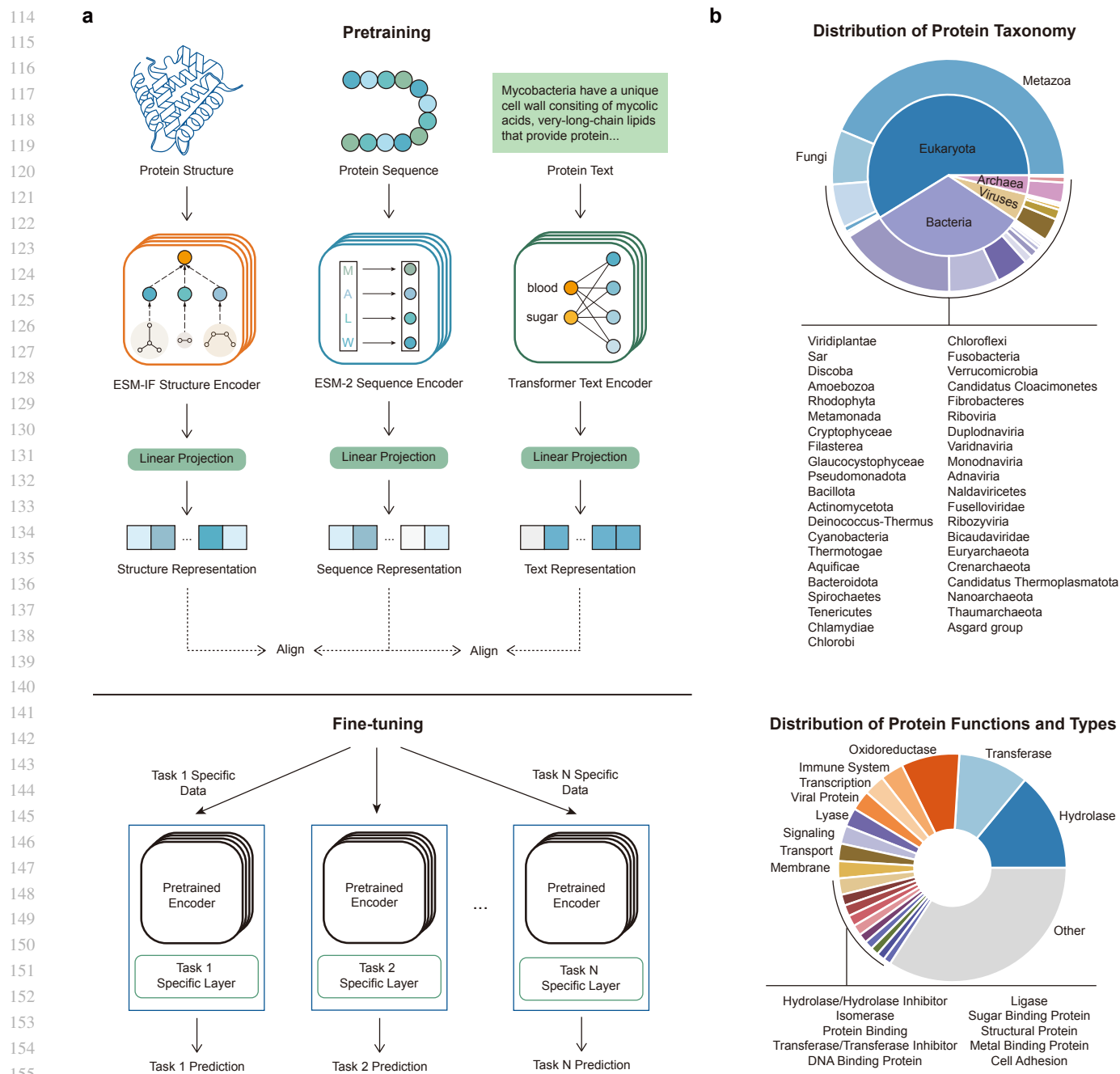
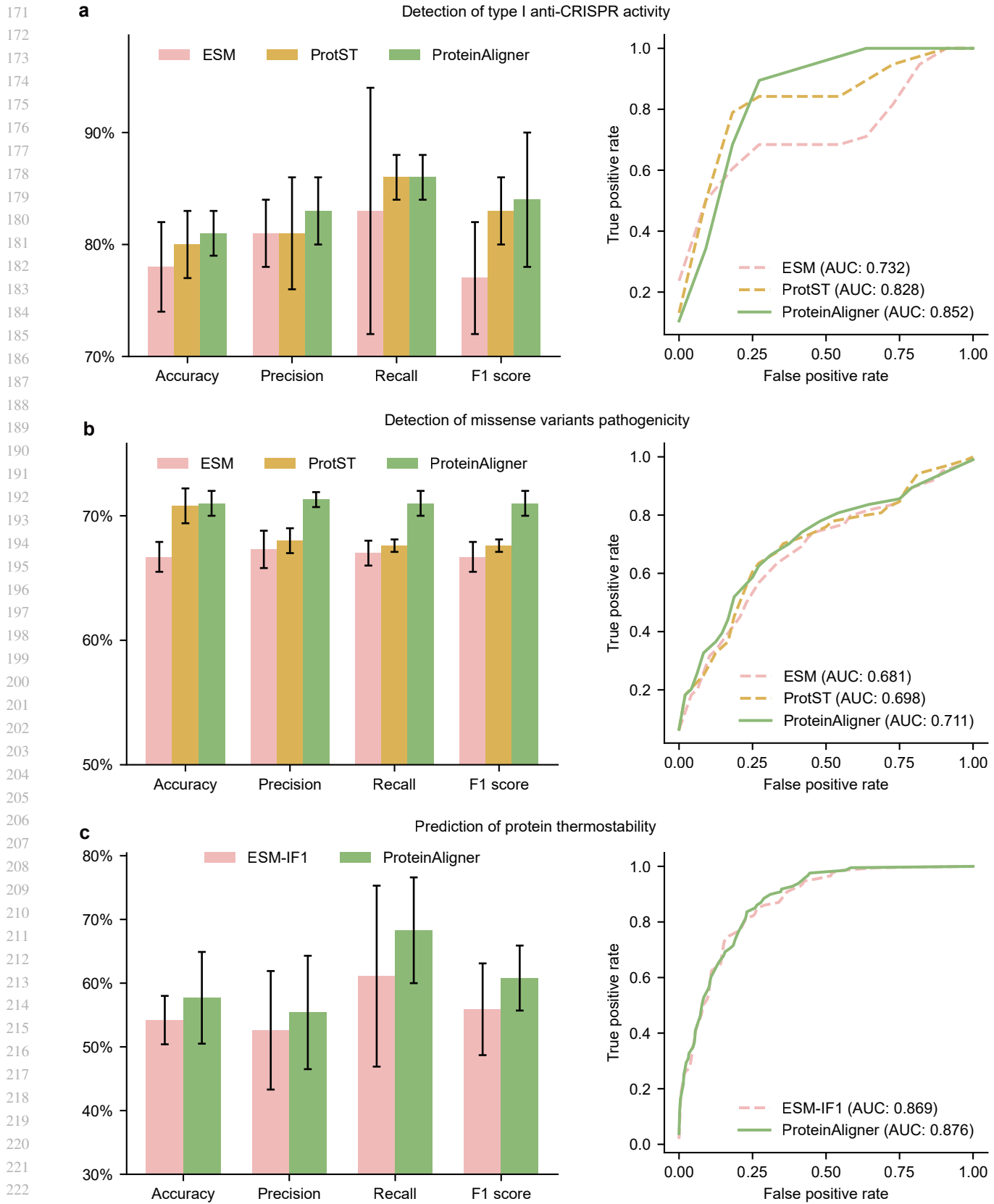


Fig. 1 | ProteinAligner facilitates multi-modal pretraining of protein foundation models by integrating diverse modalities including amino acid sequences, 3D structures, and textual data. **a**, ProteinAligner consists of three encoders: a protein sequence encoder based on ESM, a protein structure encoder based on ESM-IF1, and a transformer-based protein text encoder. These encoders learn representations for protein sequences, structures, and text, respectively. Modality-specific projection modules then transform these representations into a shared latent space, enabling direct comparison across modalities. Using protein sequences as the anchor, ProteinAligner aligns the other two modalities by minimizing contrastive losses. After pretraining, the encoders can be fine-tuned with task-specific data for various downstream applications. **b**, Our curated pretraining data for ProteinAligner spans a diverse range of proteins from various taxonomic groups, functions, and types. The upper chart displays the distribution of protein taxonomy, with the inner ring representing superkingdoms and the outer ring representing kingdoms. The lower chart illustrates the distribution of protein functions and types.



224 **Fig. 2 | ProteinAligner outperforms state-of-the-art protein foundation models in various downstream tasks. a-b**, ProteinAligner's sequence encoder outperforms
225 ESM and ProST in detecting type I anti-CRISPR activity (a) and pathogenic missense variants (b), achieving higher accuracy, F1 scores (which balance precision and
226 recall), and area under the ROC curve (AUC). c. ProteinAligner's structure encoder outperforms ESM-IF1 in predicting protein thermostability.
227

228 immune system in bacteria and archaea (34). The CRISPR-
229 Cas system functions as an adaptive immune mechanism
230 in these microorganisms, recognizing and cleaving foreign
231 DNA from viral invaders. In type I systems, which involve
232 multi-subunit Cas proteins, Acr proteins disrupt this defense
233 by preventing Cas proteins from binding to target DNA
234 or carrying out their cleavage functions. Understanding
235 and detecting these Acr activities is crucial for controlling
236 CRISPR-Cas systems in genetic engineering and leveraging
237 bacteriophages to combat antimicrobial resistance.

238
239 Given the amino acid sequences of an Arc protein and
240 a set of Cas proteins from a CRISPR-Cas system, we
241 employed ProteinAligner's pretrained sequence encoder
242 to extract representation vectors for each protein. These
243 vectors were then input into a convolutional neural network
244 (CNN) based classification module to predict whether the
245 Arc protein could inhibit the CRISPR-Cas system. We
246 utilized the Acr-CRISPR-Cas inhibition dataset (34) for
247 experiments, which comprises 227 pairs of Acr proteins
248 and CRISPR-Cas systems, including 132 experimentally
249 verified positive pairs (Acr inhibits CRISPR-Cas) and 95
250 negative pairs (Acr does not inhibit CRISPR-Cas). The
251 dataset was randomly split into training and test sets in
252 an 8:2 ratio. We compared ProteinAligner to ESM (5), a
253 protein language model pretrained on protein sequences,
254 and ProtST (24), which builds on ESM by incorporating
255 contrastive pretraining between protein sequences and
256 text. We used area under the ROC curve (AUC), accuracy,
257 precision, recall, and F1 score as evaluation metrics.

258
259 ProteinAligner demonstrated significantly better per-
260 formance across all metrics in detecting type I anti-CRISPR
261 activity compared to ESM (Fig. 2a). For instance,
262 ProteinAligner achieved an AUC of 0.852, substantially sur-
263 passing ESM's AUC of 0.732. Similarly, ProteinAligner's
264 F1 score of 0.84 was considerably higher than ESM's F1
265 score of 0.77. Moreover, ProteinAligner achieved superior
266 performance compared to ProtST (Fig. 2a).

267
268 **ProteinAligner predicts pathogenic missense vari-**
269 **ants.** Pathogenic missense variants refer to specific types
270 of genetic mutations where a single nucleotide change in a
271 DNA sequence results in the substitution of one amino acid
272 for another in the corresponding protein (35). This change
273 can disrupt the protein's normal function, potentially leading
274 to diseases or disorders. In the context of pathogenicity,
275 these variants are considered harmful because they alter
276 the protein's structure or function in a way that impairs
277 biological processes. Depending on the protein's role,
278 this can lead to a variety of outcomes, from minor effects
279 to severe genetic disorders, such as cystic fibrosis, sickle
280 cell disease, or certain forms of cancer. Identifying and
281 characterizing pathogenic missense variants is crucial in
282 genetic research and clinical diagnostics for understanding
283 inherited diseases and developing targeted treatments.

The inputs for this task are two protein sequences: the
wildtype sequence (before mutation) and the mutant
sequence (after mutation). We employed the sequence
encoder in ProteinAligner to extract representation vectors
for both proteins. These vectors were then concatenated
and passed through a multi-layer perceptron to predict
whether the mutant protein is pathogenic. We used 200
labeled examples from the VariPred (36) dataset, with 100
examples allocated for training and the remaining 100 for
testing. ProteinAligner outperformed both ESM and ProtST
(Fig. 2b). For example, it achieved an F1 score of 0.71
compared to 0.667 for ESM and 0.676 for ProtST.

ProteinAligner predicts protein thermostability. Protein
thermostability refers to a protein's ability to maintain
its structure and function when exposed to elevated temper-
atures (37). This characteristic is critical because proteins
typically lose their functional shape, or denature, at high
temperatures, rendering them ineffective. Thermostability
is an important factor in various biological processes and
industrial applications. For instance, enzymes with high
thermostability are essential in industries such as biotech-
nology and pharmaceuticals, where reactions often require
high temperatures for optimal efficiency. Predicting protein
thermostability allows researchers to design or engineer pro-
teins that can withstand challenging conditions, improving
their functionality and longevity. Additionally, thermostable
proteins are valuable in drug design, as they tend to have
better shelf lives and performance under physiological
conditions. Accurate predictions of thermostability are
crucial for advancing protein engineering and enhancing the
reliability of proteins in various applications.

Unlike the previous two tasks, this task takes the 3D
structures of proteins, specifically their atomic coordinates,
as input. The 3D structure of each protein was processed
through ProteinAligner's structure encoder, generating a
representation vector. This vector was then passed through
a multi-layer perceptron to predict the protein's thermosta-
bility class. We employed the HP-S²C5 dataset (38), which
comprises 1,040 proteins spanning five thermostability
classes: Hyperthermophilic (above 75°C), Thermophilic
(45-75°C), Mesophilic (25-45°C), Psychrophilic (5-25°C),
and Cryophilic (-20-5°C). 936 proteins were used for
training and 104 for testing. We compared ProteinAligner
with ESM-IF1 (9), a protein structure encoder pretrained
on both protein structures and sequences. ProteinAligner
remarkably outperformed ESM-IF1 (Fig. 2c), achieving an
F1 score of 0.608 compared to 0.559, and an accuracy of
0.577 compared to 0.542.

ProteinAligner identifies potent bioactive peptides.
Bioactive peptides are short chains of amino acids with
specific biological activities (39). They play critical roles
in regulating physiological processes, including immune

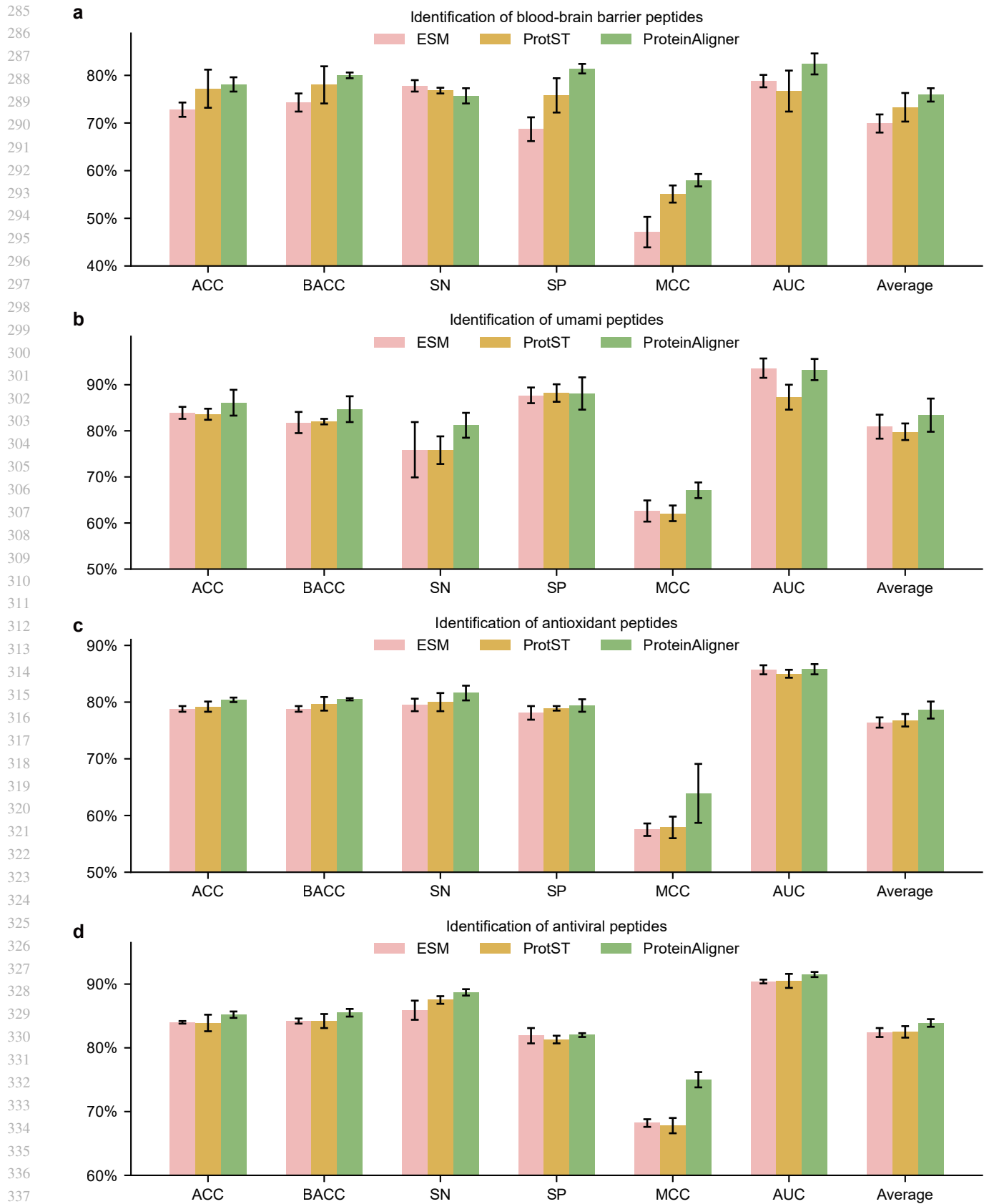


Fig. 3 | ProteinAligner surpasses ESM and ProtST in identifying potent bioactive peptides. It demonstrated superior performance in predicting blood-brain barrier penetration (a), umami taste induction (b), antioxidant activity (c), and antiviral properties (d). Model performance was evaluated using accuracy (ACC), balanced accuracy (BACC), sensitivity (SN), specificity (SP), Matthews correlation coefficient (MCC), and the area under the ROC curve (AUC).

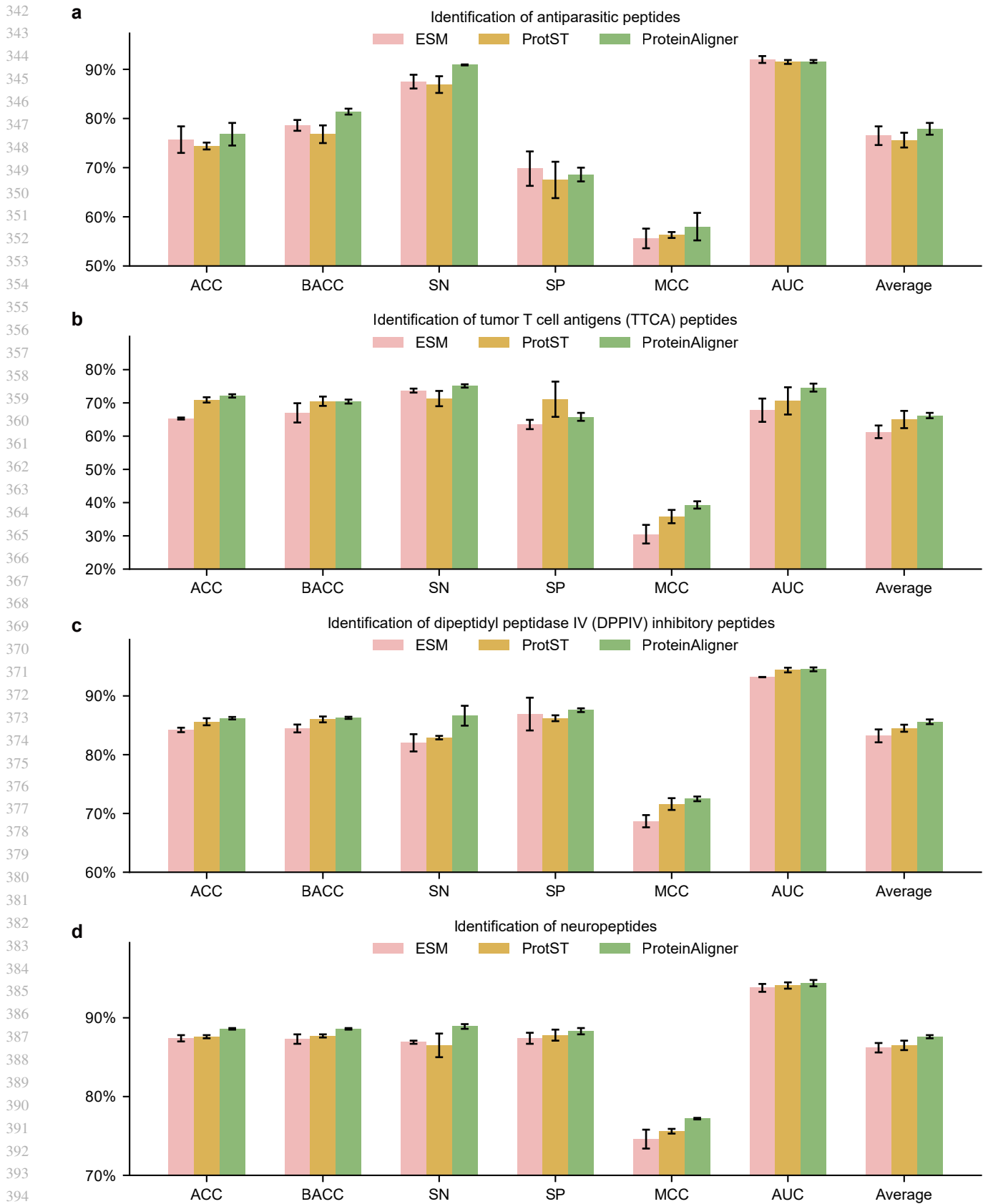


Fig. 4 | ProteinAligner outperforms ESM and ProtST in identifying potent bioactive peptides. It achieved superior performance in predicting antiparasitic effects (a), T-cell immune response induction (b), inhibition of dipeptidyl peptidase IV (DPP-IV) (c), and modulation of brain activity (d). Model performance was assessed using accuracy (ACC), balanced accuracy (BACC), sensitivity (SN), specificity (SP), Matthews correlation coefficient (MCC), and the area under the ROC curve (AUC).

399 function, metabolism, and cardiovascular health. Identifying
400 bioactive peptides is important because they offer significant
401 potential for developing new therapeutic agents and func-
402 tional foods. These peptides can serve as natural, targeted
403 treatments with fewer side effects compared to traditional
404 drugs, and their discovery can lead to advancements in both
405 medical applications and nutrition, benefiting public health
406 and disease prevention efforts.

407
408 Given the amino acid sequence of a peptide, we em-
409 ployed ProteinAligner’s protein sequence encoder to extract
410 a representation vector, which was subsequently input into
411 a convolutional neural network (CNN)-based classification
412 head to predict whether the peptide has a specific bioactivity.
413 We examined eight distinct bioactivities, including blood-
414 brain barrier penetration (40), umami taste induction (41),
415 antioxidant activity (42), antiviral properties (43), antipar-
416 asitic effects (44), T-cell immune response induction (45),
417 inhibition of dipeptidyl peptidase IV (DPP-IV) (46), and
418 modulation of brain activity (47). Given that a peptide can
419 exhibit multiple bioactivities concurrently, we approached
420 each bioactivity prediction as a binary classification task,
421 avoiding the use of a multi-class model that would assign
422 the peptide to a single category. Separate datasets were
423 used for each bioactivity (Methods). The evaluation metrics
424 for this task included accuracy (ACC), balanced accuracy
425 (BACC) (48), sensitivity (SN), specificity (SP), Matthews
426 correlation coefficient (MCC) (49), and the area under the
427 ROC curve (AUC).

428
429 ProteinAligner consistently surpassed ESM and ProtST
430 across all eight evaluated tasks (Figs. 3 and 4). For example,
431 in predicting blood-brain barrier penetration, ProteinAligner
432 achieved an AUC of 0.824, outperforming both ESM
433 (0.788) and ProtST (0.767). Similarly, in predicting T-cell
434 immune response induction, ProteinAligner reached an
435 AUC of 0.746, exceeding the performance of ESM (0.678)
436 and ProtST (0.706).

437
438 **ProteinAligner predicts the minimum inhibitory**
439 **concentration (MIC) of antimicrobial peptides.** Antimi-
440 crobial peptides (AMPs) are short chains of amino acids
441 that serve as a crucial part of the innate immune response in
442 many organisms, exhibiting broad-spectrum activity against
443 bacteria, viruses, fungi, and even cancer cells (50). They
444 function by disrupting microbial membranes, leading to cell
445 death, and are considered potential alternatives to conven-
446 tional antibiotics, especially in the face of rising antibiotic
447 resistance. The minimum inhibitory concentration (MIC) is
448 the lowest concentration of an antimicrobial agent, such as
449 an AMP, that prevents visible microbial growth. Accurately
450 predicting the MIC values of AMPs is essential as it allows
451 for the optimization of peptide design for therapeutic use,
452 minimizes potential toxicity, and helps in the early-stage
453 screening of effective peptides before in vitro or in vivo
454 testing. This predictive capability is vital for accelerating
455

the development of AMPs as a novel class of antimicrobial
agents in clinical applications.

Given the amino acid sequence of a peptide, we ap-
plied ProteinAligner’s protein sequence encoder to extract a
representation vector, which was then fed into a multi-layer
perceptron-based regression module to predict the MIC of
the peptide against a specific pathogen. We focused on
Escherichia coli (*E. coli*), a gram-negative bacterium. We
utilized the dataset from (51), comprising 3,695 training
and 924 testing examples. Mean squared error was used
as the evaluation metric. ProteinAligner achieved lower
prediction error compared to ESM and ProtST (Extended
Data Fig. 1).

Discussion

ProteinAligner demonstrates superior performance com-
pared to models like ESM and ProtST by employing a
multi-modal approach that integrates sequences, structures,
and textual information, offering a more comprehensive
understanding of proteins (Figs. 2a, 2b, 3, 4, and Ex-
tended Data Fig.1). While ESM focuses on learning from
amino acid sequences, ProteinAligner incorporates protein
structures, which provide critical insights into folding,
stability, and interaction dynamics that sequences alone
cannot reveal. Additionally, by drawing from experimental
literature, ProteinAligner captures contextual information
like functional annotations and post-translational modi-
fications, bridging gaps in sequence and structure-based
models. The multi-modal alignment in a shared latent space
enables ProteinAligner to generate representations that are
not only sequence-based but also grounded in functional and
structural knowledge, leading to more accurate predictions.
Compared to ProtST, which integrates sequence and text
data, ProteinAligner adds the essential dimension of protein
structure, crucial for understanding spatial arrangements and
physical properties. This enriched representation enables the
model to better predict protein functions influenced by struc-
tural conformation, outperforming ProtST’s dual-modality
approach. Additionally, ProteinAligner surpasses ESM-IF1
(Fig. 2c), which combines sequences and structures, by in-
corporating textual descriptions that provide further context
from experimental literature. ProteinAligner’s contrastive
alignment strategy ensures that the representations from
sequences, structures, and text are effectively integrated,
enhancing its predictive capacity across different tasks.
This holistic approach equips ProteinAligner to excel in
both sequence- and structure-based predictions, offering
a more nuanced understanding of protein properties and
delivering superior performance across a range of biological
applications.

ProteinAligner’s ability to perform a wide range of
prediction tasks presents promising applications across
biology, drug discovery, and medicine. In drug discovery,
ProteinAligner’s accurate identification of bioactive pep-

456 tides, such as DPP-IV inhibitors (Fig. 4c), is particularly
457 relevant for developing treatments for metabolic disorders
458 like diabetes. Its capability to predict antimicrobial peptide
459 properties, such as minimum inhibitory concentration
460 (MIC) (Extended Data Fig. 1), is critical for advancing
461 new antimicrobial therapies, especially in addressing the
462 challenge of antibiotic-resistant pathogens. This has impor-
463 tant implications for the global fight against antimicrobial
464 resistance. The model's ability to detect type I anti-CRISPR
465 activities supports the design of more efficient and precise
466 CRISPR-based tools for both research and therapeutic
467 applications (Fig. 2a). Anti-CRISPR systems could be
468 used to enhance the safety of gene editing by mitigating
469 off-target effects or enabling reversible gene modifications.
470 In precision medicine, the prediction of pathogenic missense
471 variants aids in the early detection and diagnosis of genetic
472 disorders. By identifying harmful mutations that may lead
473 to diseases, ProteinAligner can contribute to personalized
474 treatment strategies, improving patient outcomes (Fig. 2b).
475 Additionally, ProteinAligner's accurate prediction of protein
476 thermostability is vital for protein engineering, biopharma-
477 ceutical development, and industrial biotechnology, where
478 stable proteins are necessary for drug formulations and
479 biocatalysts (Fig. 2c). Overall, ProteinAligner's diverse
480 prediction capabilities position it as a valuable tool that
481 can accelerate innovation in multiple fields, enabling faster
482 therapeutic discoveries, more precise gene-editing tools,
483 and advancements in personalized medicine and protein
484 engineering.

486 Despite the advantages of ProteinAligner, the model
487 has several limitations. One of the key challenges is the
488 dependency on high-quality structural and textual data,
489 which is not always available for all proteins. While
490 ProteinAligner can perform pretraining even when only
491 sequences and one additional modality (either structure or
492 text) are present, the absence of full multi-modal data for
493 many proteins can limit the model's ability to learn com-
494 prehensive representations. Additionally, ProteinAligner's
495 reliance on contrastive loss for modality alignment may
496 not fully capture subtle biological nuances in cases where
497 sequence, structure, and text data are not perfectly aligned.
498 Another limitation is the computational cost associated with
499 training multi-modal models, especially when dealing with
500 large-scale protein datasets that involve high-dimensional
501 structural information and large text corpora. Finally,
502 while ProteinAligner improves upon previous models by
503 integrating structure, sequence, and text, it still does not
504 account for other potentially informative modalities, such
505 as protein-protein interactions or functional annotations
506 from various databases, which could further enhance its
507 predictive capabilities.

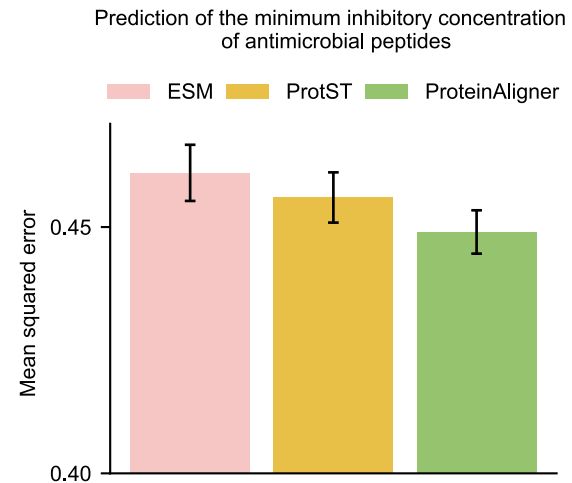
509 Future work on ProteinAligner could focus on several
510 key directions to further enhance its performance and
511 applicability. One promising area is the incorporation of

additional modalities, such as protein-protein interaction
networks and post-translational modifications. These addi-
tional data sources could provide deeper insights into protein
behavior and interactions, leading to even more robust and
comprehensive protein representations. Another direction
for future work is to improve the model's ability to handle
incomplete or noisy data by developing more sophisticated
alignment strategies that better tolerate inconsistencies
between modalities. Enhancing the interpretability of
ProteinAligner's predictions is also a critical area for future
research, which could involve incorporating explainability
techniques to make the model's decision-making process
more transparent, particularly in cases where sequence,
structure, and text data converge. Lastly, expanding Pro-
teinAligner's applications beyond protein function and
property prediction - such as protein design and structure
prediction - could broaden its impact across a wide range of
biological and biomedical challenges.

References

1. James A Wells and Christopher L McCleendon. Reaching for high-hanging fruit in drug discovery at protein-protein interfaces. *Nature*, 450(7172):1001–1009, 2007.
2. Carl AK Borrebaeck. Precision diagnostics: moving towards protein biomarker signatures of clinical utility in cancer. *Nature Reviews Cancer*, 17(3):199–204, 2017.
3. Irene Nobeli, Angelo D Favia, and Janet M Thornton. Protein promiscuity and its implications for biotechnology. *Nature biotechnology*, 27(2):157–167, 2009.
4. Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
5. Alexander Rives, Joshua Meier, Tom Sercu, Siddharth Goyal, Zeming Lin, Jason Liu, Demi Guo, Myle Ott, C Lawrence Zitnick, Jerry Ma, et al. Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences. *Proceedings of the National Academy of Sciences*, 118(15):e2016239118, 2021.
6. Ahmed Elnaggar, Michael Heinzinger, Christian Dallago, Ghalia Rehawi, Yu Wang, Llion Jones, Tom Gibbs, Tamas Feher, Christoph Angerer, Martin Steinegger, et al. ProtTrans: Toward understanding the language of life through self-supervised learning. *IEEE transactions on pattern analysis and machine intelligence*, 44(10):7112–7127, 2021.
7. Tristan Bepler and Bonnie Berger. Learning the protein language: Evolution, structure, and function. *Cell systems*, 12(6):654–669, 2021.
8. 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.
9. Chloe Hsu, Robert Verkuil, Jason Liu, Zeming Lin, Brian Hie, Tom Sercu, Adam Lerer, and Alexander Rives. Learning inverse folding from millions of predicted structures. In *International Conference on Machine Learning*, pages 8946–8970. PMLR, 2022.
10. Nadav Brandes, Dan Ofer, Yam Peleg, Nadav Rappoport, and Michal Linial. ProteinBERT: a universal deep-learning model of protein sequence and function. *Bioinformatics*, 38(8):2102–2110, 2022.
11. Bo Chen, Xingyi Cheng, Pan Li, Yangli-ao Geng, Jing Gong, Shen Li, Zhilei Bei, Xu Tan, Boyan Wang, Xin Zeng, et al. xtrimgplm: unified 100b-scale pre-trained transformer for deciphering the language of protein. *arXiv preprint arXiv:2401.06199*, 2024.
12. Serbulent Unsal, Heval Atas, Muammer Albayrak, Kemal Turhan, Aybar C Acar, and Tunca Doğan. Learning functional properties of proteins with language models. *Nature Machine Intelligence*, 4(3):227–245, 2022.
13. Tianhao Yu, Haiyang Cui, Jianan Canal Li, Yunan Luo, Guangde Jiang, and Huimin Zhao. Enzyme function prediction using contrastive learning. *Science*, 379(6639):1358–1363, 2023.
14. Zachary N Flamholz, Steven J Biller, and Libusha Kelly. Large language models improve annotation of prokaryotic viral proteins. *Nature Microbiology*, 9(2):537–549, 2024.
15. Felix Teufel, José Juan Almagro Armenteros, Alexander Rosenberg Johansen, Magnús Halldór Gíslason, Silas Irby Pihl, Konstantinos D Tsirigos, Ole Winther, Søren Brunak, Gunnar von Heijne, and Henrik Nielsen. Signalp 6.0 predicts all five types of signal peptides using protein language models. *Nature biotechnology*, 40(7):1023–1025, 2022.
16. 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.
17. Ratul Chowdhury, Nazim Bouatta, Surojit Biswas, Christina Floristean, Anant Kharkar, Koushik Roy, Charlotte Rochereau, Gustaf Ahdrizt, Joanna Zhang, George M Church, et al. Single-sequence protein structure prediction using a language model and deep learning. *Nature Biotechnology*, 40(11):1617–1623, 2022.
18. Ali Madani, Ben Krause, Eric R Greene, Subu Subramanian, Benjamin P Mohr, James M

- 513 Holton, Jose Luis Olmos, Caiming Xiong, Zachary Z Sun, Richard Socher, et al. Large
514 language models generate functional protein sequences across diverse families. *Nature*
515 *Biotechnology*, 41(8):1099–1106, 2023.
- 516 19. Noelia Ferruz, Steffen Schmidt, and Birte Höcker. Protgpt2 is a deep unsupervised
517 language model for protein design. *Nature communications*, 13(1):4348, 2022.
- 518 20. Gregory A Petsko and Dagmar Ringe. *Protein structure and function*. New Science Press,
519 2004.
- 520 21. Qiangfeng Cliff Zhang, Donald Petrey, Lei Deng, Li Qiang, Yu Shi, Chan Aye Thu, Brygida
521 Bisikirska, Celine Lefebvre, Domenico Accili, Tony Hunter, et al. Structure-based prediction
522 of protein–protein interactions on a genome-wide scale. *Nature*, 490(7421):556–560, 2012.
- 523 22. Roberto Mosca, Arnaud Céol, and Patrick Aloy. Interactome3d: adding structural details to
524 protein networks. *Nature methods*, 10(1):47–53, 2013.
- 525 23. Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So,
526 and Jaewoo Kang. Biobert: a pre-trained biomedical language representation model for
527 biomedical text mining. *Bioinformatics*, 36(4):1234–1240, 2020.
- 528 24. Minghao Xu, Xinyu Yuan, Santiago Miret, and Jian Tang. Protst: Multi-modality learning of
529 protein sequences and biomedical texts. In *International Conference on Machine Learning*,
530 pages 38749–38767. PMLR, 2023.
- 531 25. Ningyu Zhang, Zhen Bi, Xiaozhuan Liang, Siyuan Cheng, Haosen Hong, Shumin Deng,
532 Jiazhang Lian, Qiang Zhang, and Huajun Chen. Ontoprotein: Protein pretraining with
533 gene ontology embedding. *International Conference on Learning Representations*, 2022.
- 534 26. Kevin E Wu, Howard Chang, and James Zou. Proteinclip: enhancing protein language
535 models with natural language. *bioRxiv*, pages 2024–05, 2024.
- 536 27. Varun R Shanker, Theodora UJ Bruun, Brian L Hie, and Peter S Kim. Unsupervised evolution
537 of protein and antibody complexes with a structure-informed language model. *Science*,
538 385(6704):46–53, 2024.
- 539 28. Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for
540 unsupervised visual representation learning. In *2020 IEEE/CVF Conference on Computer*
541 *Vision and Pattern Recognition (CVPR)*, pages 9726–9735. IEEE Computer Society, 2020.
- 542 29. Jin Su, Xibin Zhou, Xuting Zhang, and Fajie Yuan. Protrek: Navigating the protein universe
543 through tri-modal contrastive learning. *bioRxiv*, pages 2024–05, 2024.
- 544 30. Bowen Jing, Stephan Eismann, Patricia Suriana, Raphael John Lamarre Townshend, and
545 Ron Dror. Learning from protein structure with geometric vector perceptrons. In *International*
546 *Conference on Learning Representations*, 2020.
- 547 31. Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive
548 predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- 549 32. UniProt Consortium. Uniprot: a worldwide hub of protein knowledge. *Nucleic acids re-*
550 *search*, 47(D1):D506–D515, 2019.
- 551 33. Stephen K Burley, Charmi Bhikadiya, Chunxiao Bi, Sebastian Bittrich, Henry Chao,
552 Li Chen, Paul A Craig, Gregg V Crichlow, Kenneth Dalenberg, Jose M Duarte, et al. Rcsb
553 protein data bank (rcsb.org): delivery of experimentally-determined pdb structures along-
554 side one million computed structure models of proteins from artificial intelligence/machine
555 learning. *Nucleic acids research*, 51(D1):D488–D508, 2023.
- 556 34. Moein Hasani, Chantel N. Trost, Nolen Timmerman, and Lingling Jin. Acctransact: Pre-
557 trained protein transformer models for the detection of type i anti-crispr activities. In *Pro-*
558 *ceedings of The 14th ACM Conference on Bioinformatics, Computational Biology, and*
559 *Health Informatics (ACM-BCB)*, page 6, Houston, TX, USA, 2023. ACM.
- 560 35. Jun Cheng, Guido Novati, Joshua Pan, Clare Bycroft, Akvilė Žemgulytė, Taylor Apple-
561 baum, Alexander Pritzel, Lai Hong Wong, Michal Zieliński, Tobias Sargeant, et al. Accu-
562 rate proteome-wide missense variant effect prediction with alphamissense. *Science*, 381
563 (6664):eadg7492, 2023.
- 564 36. Weining Lin, Jude Wells, Zeyuan Wang, Christine Orengo, and Andrew CR Martin. Enhanc-
565 ing missense variant pathogenicity prediction with protein language models using varipred.
566 *Scientific Reports*, 14(1):8136, 2024.
- 567 37. H Pezeshgi Modarres, MR Mofrad, and AJRA Sanati-Nezhad. Protein thermostability en-
568 gineering. *RSC advances*, 6(116):115252–115270, 2016.
- 569 38. Tianlong Chen, Chengyue Gong, Daniel Jesus Diaz, Xuxi Chen, Jordan Tyler Wells, Qiang
570 Liu, Zhangyang Wang, Andrew Ellington, Alexandros Dimakis, and Adam Klivans. Hotpro-
571 tein: A novel framework for protein thermostability prediction and editing. In *International*
572 *Conference on Learning Representations (ICLR)*, 2023.
- 573 39. Ali Adem Bahar and Dacheng Ren. Antimicrobial peptides. *Pharmaceuticals*, 6(12):1543–
574 1575, 2013.
- 575 40. Ruyi Dai, Wei Zhang, Wending Tang, Evelien Wynendaele, Qizhi Zhu, Yannan Bin, Bart
576 De Spiegeleer, and Junfeng Xia. Bbppred: sequence-based prediction of blood-brain bar-
577 rier peptides with feature representation learning and logistic regression. *Journal of Chem-*
578 *ical Information and Modeling*, 61(1):525–534, 2021.
- 579 41. Yin Zhang, Chandrasekar Venkatasamy, Zhongli Pan, Wenlong Liu, and Liming Zhao. Novel
580 umami ingredients: Umami peptides and their taste. *Journal of food science*, 82(1):16–23,
581 2017.
- 582 42. Tang-Bin Zou, Tai-Ping He, Hua-Bin Li, Huan-Wen Tang, and En-Qin Xia. The structure-
583 activity relationship of the antioxidant peptides from natural proteins. *Molecules*, 21(1):72,
584 2016.
- 585 43. Liana Costa Pereira Vilas Boas, Marcelo Lattarulo Campos, Rhayfa Lorraine Araujo
586 Berlanda, Natan de Carvalho Neves, and Octávio Luiz Franco. Antiviral peptides as promis-
587 ing therapeutic drugs. *Cellular and Molecular Life Sciences*, 76:3525–3542, 2019.
- 588 44. Amram Mor. Multifunctional host defense peptides: antiparasitic activities. *The FEBS*
589 *journal*, 276(22):6474–6482, 2009.
- 590 45. Phasit Charoenkwan, Chanin Nantasenam, Md Mehedi Hasan, and Watshara Shoombatong.
591 itca-hybrid: Improved and robust identification of tumor t cell antigens by utilizing
592 hybrid feature representation. *Analytical biochemistry*, 599:113747, 2020.
- 593 46. Hanne B Rasmussen, Sven Branner, Finn C Wiberg, and Nicolai Wagtmann. Crystal struc-
594 ture of human dipeptidyl peptidase iv/cd26 in complex with a substrate analog. *Nature*
595 *structural biology*, 10(1):19–25, 2003.
- 596 47. Yannan Bin, Wei Zhang, Wending Tang, Ruyi Dai, Menglu Li, Qizhi Zhu, and Junfeng



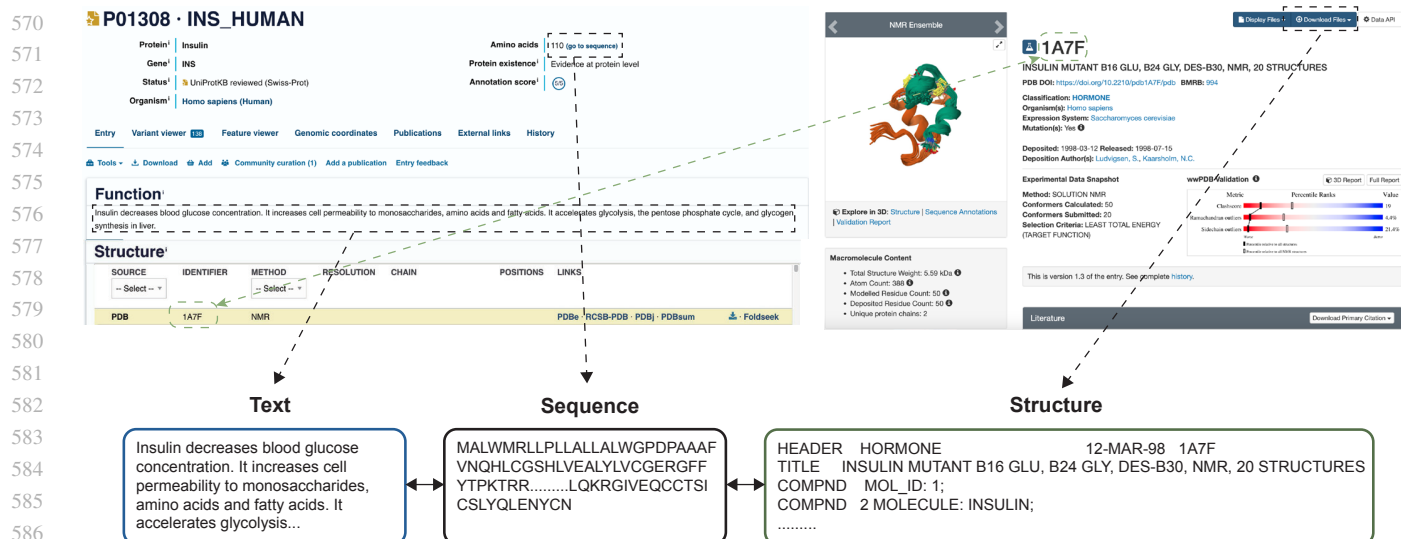
Extended Data Fig. 1 | ProteinAligner surpasses ESM and ProtST in predicting the minimum inhibitory concentration (MIC) of antimicrobial peptides, achieving a lower mean squared error.

- Xia. Prediction of neuropeptides from sequence information using ensemble classifier and hybrid features. *Journal of proteome research*, 19(9):3732–3740, 2020.
48. Haibo He and Edward A. Garcia. Learning from imbalanced data. *IEEE Transactions on Knowledge and Data Engineering*, 21(9):1263–1284, 2009.
49. Brian W. Matthews. Comparison of the predicted and observed secondary structure of t4 phage lysozyme. *Biochimica et Biophysica Acta (BBA) - Protein Structure*, 405(2):442–451, 1975.
50. Amir Pandi, David Adam, Amir Zare, Van Tuan Trinh, Stefan L Schaefer, Marie Burt, Björn Klabunde, Elizaveta Bobkova, Manish Kushwaha, Yeganeh Forouhijabbari, et al. Cell-free biosynthesis combined with deep learning accelerates de novo-development of antimicrobial peptides. *Nature Communications*, 14(1):7197, 2023.
51. Alba Ledesma-Fernandez, Susana Velasco-Lozano, Javier Santiago-Arcos, Fernando López-Gallego, and Aitziber L Cortajarena. Engineered repeat proteins as scaffolds to assemble multi-enzyme systems for efficient cell-free biosynthesis. *Nature Communications*, 14(1):2587, 2023.

Method

Dataset preprocessing. The pretraining data for ProteinAligner was sourced from the UniProtKB/Swiss-Prot (32) and RCSB PDB (33) databases. UniProtKB/Swiss-Prot is a well-curated repository containing high-quality protein sequences across a wide variety of organisms, along with detailed annotations on protein functions and properties. We utilized version UniProt 2023_02, which was released on May 2, 2023. The RCSB PDB database offers a comprehensive collection of experimentally determined 3D protein structures, derived from methods such as X-ray crystallography, nuclear magnetic resonance (NMR) spectroscopy, and cryo-electron microscopy (cryo-EM). It includes protein structures from a wide range of proteins, such as enzymes, receptors, and antibodies, originating from diverse organisms.

From these databases, we collected sequence-text pairs and sequence-structure pairs (Extended Data Fig. 2). The sequence-text pairs were sourced from UniProtKB/Swiss-Prot. We first obtained a collection of protein entries from Swiss-Prot that included textual descriptions of their functions, by filtering for entries where the `commentType` field



Extended Data Fig. 2 | An illustration of constructing sequence-text and sequence-structure pairs from the UniProtKB/Swiss-Prot and RCSB PDB databases. Sequence-text pairs were generated by linking each protein's amino acid sequence with its functional description, both available on the same UniProtKB/Swiss-Prot webpage. Sequence-structure pairs were created by matching the protein sequences from UniProtKB/Swiss-Prot with their corresponding structures in RCSB PDB, using PDB IDs.

was set to 'Function'. We then retrieved the corresponding sequence for each protein in this collection. Specifically, we accessed the UniProt ID from the `primaryAccession` field and used it to retrieve the corresponding protein FASTA file from the UniProt website, which contains the protein's sequence. We downloaded all available PDB files from the May 2, 2023 dataset release (33), comprising 200,734 experimentally determined protein structures. We then employed the UniProt ID mapping tool¹ to link the structures in the PDB files to their corresponding amino acid sequences in the FASTA files. To address memory constraints during pretraining, we excluded protein sequences longer than 300 residues, yielding 133,726 sequence-structure pairs and 290,480 sequence-text pairs.

Encoders in ProteinAligner. ProteinAligner utilizes ESM (16) to learn representations for protein sequences. ESM, a protein language model, was pretrained on 65 million protein sequences from UniRef50 (52) by predicting masked amino acids. The model features 33 transformer layers and an embedding dimension of 1280, allowing it to effectively capture the complexities inherent in protein sequences. To encode protein structures, ProteinAligner employs ESM-IF1 (9), a model trained to address the inverse folding problem - predicting the amino acid sequence from a protein's backbone atom coordinates. ESM-IF1 comprises an encoder and a decoder, where the encoder extracts a representation vector from the input structure, which is then fed into the decoder to generate the corresponding sequence. ProteinAligner utilizes only the encoder from ESM-IF1, omitting the decoder component. The encoder is composed of four Geometric Vector Perceptron Graph Neural Network

¹<https://www.uniprot.org/id-mapping>

(GVP-GNN) (30) layers for geometric feature extraction, followed by eight transformer encoder layers to capture long-range interactions between these features. ESM-IF1 was trained on 12 million AlphaFold2 (8) computed protein structures and 16,000 experimentally verified structures, along with their associated sequences from the UniRef50 dataset (52). The text encoder is a transformer model comprising eight layers and a total of 78 million parameters.

ProteinAligner pretraining. Given a protein sequence S and a protein structure R , we employ the sequence encoder $E_s(\cdot)$ and structure encoder $E_r(\cdot)$ to extract representation vectors $\mathbf{s} = E_s(S)$ and $\mathbf{r} = E_r(R)$ for S and R , respectively. To ensure that the representations are similar when S and R belong to the same protein, and dissimilar when they do not, we minimize the InfoNCE (31) contrastive learning loss:

$$\mathcal{L}_{s,r} = -\log \frac{\exp(\mathbf{s}_i^\top \mathbf{r}_i / \tau)}{\exp(\mathbf{s}_i^\top \mathbf{r}_i / \tau) + \sum_{j \neq i} \exp(\mathbf{s}_i^\top \mathbf{r}_j / \tau)}, \quad (1)$$

where \mathbf{s}_i and \mathbf{r}_i represent the sequence and structure representations of the same protein i , while \mathbf{s}_i and \mathbf{r}_j represent the representations of different proteins i and j . This loss function encourages the alignment of \mathbf{s}_i and \mathbf{r}_i and discourages the similarity between \mathbf{s}_i and \mathbf{r}_j . The temperature parameter τ controls the sharpness of the softmax distribution.

Similarly, given a protein sequence S and a protein text description T , we use the sequence encoder $E_s(\cdot)$ and the text encoder $E_t(\cdot)$ to obtain representation vectors $\mathbf{s} = E_s(S)$ and $\mathbf{t} = E_t(T)$ for S and T , respectively. To ensure that the representations are similar when S and T correspond to the same protein, and dissimilar when they

do not, we minimize the following InfoNCE contrastive learning loss:

$$\mathcal{L}_{s,t} = -\log \frac{\exp(\mathbf{s}_i^\top \mathbf{t}_i / \tau)}{\exp(\mathbf{s}_i^\top \mathbf{t}_i / \tau) + \sum_{j \neq i} \exp(\mathbf{s}_i^\top \mathbf{t}_j / \tau)}. \quad (2)$$

During pretraining, we minimized the sum of the two loss functions with equal weights. The temperature parameter τ was configured to 0.07. Pretraining was carried out over 20 epochs with total batch size of 80 using 40 A100 GPUs. We optimized the model weights using the AdamW optimizer (53), with an initial learning rate of 5×10^{-6} , weight decay of 1×10^{-4} , and betas of (0.9, 0.95). The learning rate was dynamically adjusted throughout pretraining via the CosineAnnealingLR scheduler (54).

Type I anti-CRISPR activity detection. The overall model architecture for this task is illustrated in Extended Data Fig. 3a. The sequence encoder pretrained by ProteinAligner was fine-tuned using data specific to this task. The classification module was based on a CNN which was composed of two 1D convolutional layers, each with a stride of 1 and a kernel size of 7. The first convolutional layer takes an input of 1280 channels and outputs 4 channels. The second convolutional layer maintains the same input and output dimensions. Batch normalization (55) and ReLU activation (56) are applied after each convolutional layer. Following the convolutional layers, two fully connected layers, each with a hidden size of 4, are employed for final classification.

During training, we optimized model weights using the Adam (57) optimizer with an initial learning rate of 3×10^{-3} , over a maximum of 250 epochs with a batch size of 32. To prevent overfitting, we employed early stopping when the decrease in training loss fell below 0.005 and applied weight decay, starting at 0.01 and gradually reducing to 0.001. We also performed a hyperparameter sweep on the dropout rate (58), exploring values between 0.3 and 0.5. Additionally, we implemented a learning rate reduction strategy, decreasing the rate by a factor of 0.9 if validation performance did not improve for 10 consecutive epochs.

The metrics used to evaluate model performance in this task include accuracy, precision, recall rate, and F1 score. These metrics are calculated based on the number of true positives (TP), false positives (FP), false negatives (FN), and true negatives (TN), and are defined as follows:

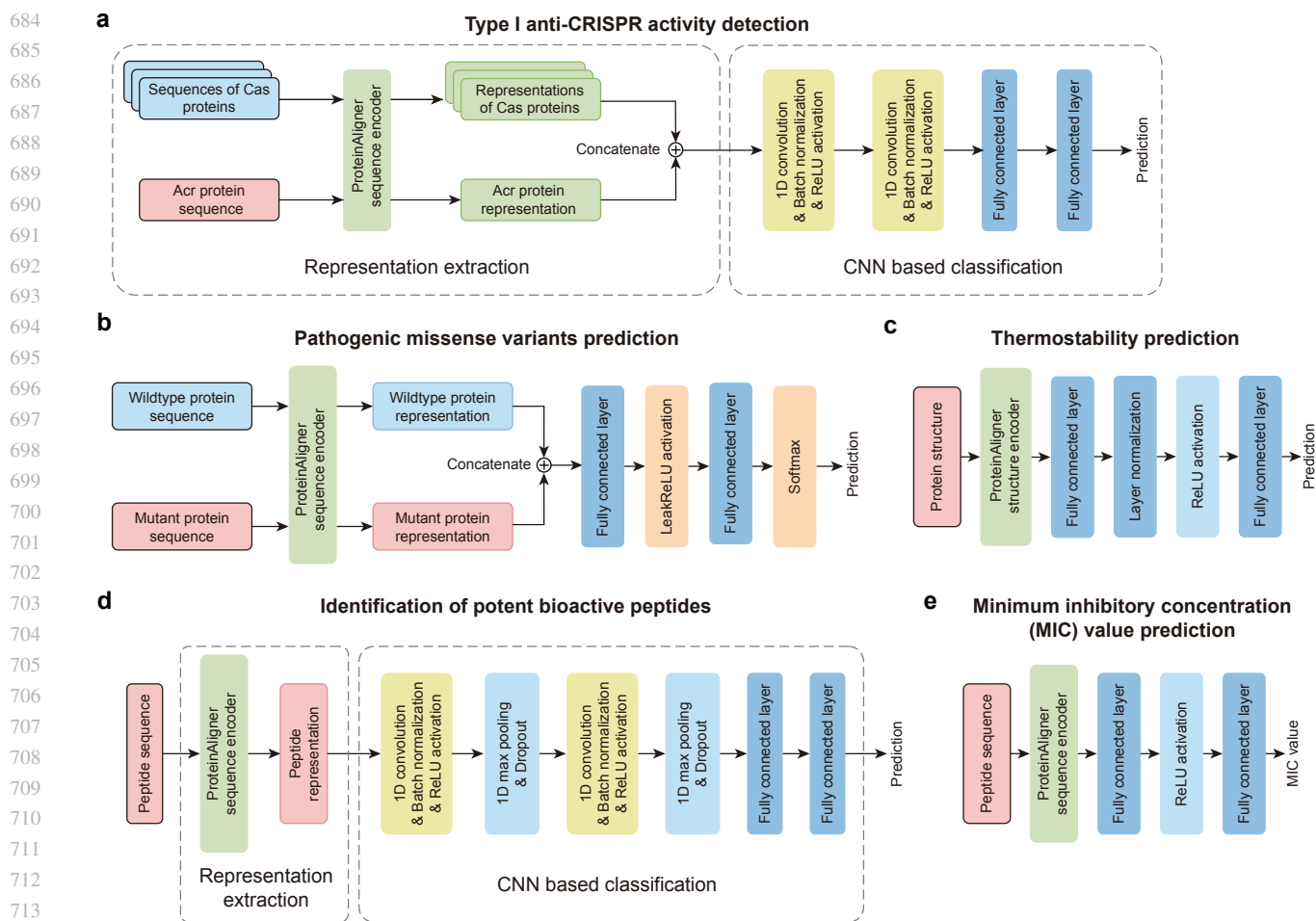
$$\begin{aligned} \text{Accuracy} &= \frac{TP + TN}{TP + TN + FP + FN}, \\ \text{Precision} &= \frac{TP}{TP + FP}, \\ \text{Recall} &= \frac{TP}{TP + FN}, \\ \text{F1 score} &= 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}. \end{aligned} \quad (3)$$

Pathogenic missense variants prediction. The overall model architecture is depicted in Extended Data Fig. 3b. For this task, the sequence encoder pretrained by ProteinAligner was fine-tuned. The classification module was based on a multi-layer perceptron, which comprises a fully connected layer with a hidden state size of 1280, a dropout layer with a probability of 0.5, a leaky ReLU activation (59), and a second fully connected layer with a softmax activation function. We employed the Adam optimizer with a learning rate of 1×10^{-4} and a weight decay of 1×10^{-3} , training the model for a maximum of 200 epochs with a batch size of 32. To mitigate overfitting, we applied an early stopping strategy: if validation performance did not improve over 10 consecutive epochs, training was halted, and the model checkpoint with the best validation performance was retained as the final model. The model's performance was evaluated using accuracy, precision, recall, and F1 score, as defined in Equation 3.

Thermostability prediction. The overall model architecture is depicted in Extended Data Fig. 3c. The structure encoder, pretrained using ProteinAligner, was fine-tuned for this task. The classification module, implemented as a multi-layer perceptron (MLP), includes a fully connected layer with a hidden dimension of 128, followed by layer normalization (60), a ReLU activation, and a second fully connected layer to produce the classification logits. For model optimization, we employed the AdamW (53) optimizer with a learning rate of 2×10^{-2} , a batch size of 4, and a weight decay of 5×10^{-2} . A one-cycle learning rate scheduler (61) was applied, and the model was trained for 8 epochs. Performance was evaluated based on accuracy, precision, recall, and F1 score, as defined in Equation 3.

Identification of potent bioactive peptides. The overall model architecture is illustrated in Extended Data Fig. 3d. The sequence encoder, pretrained by ProteinAligner, was fine-tuned for each of the eight tasks. The classification module employs a convolutional neural network (CNN) with six layers, structured as follows: a 1D convolutional layer (kernel size = 3, stride = 1, padding = 2), followed by BatchNorm and ReLU activation; a max pooling layer (kernel size = 2, padding = 1) and a dropout layer with a probability of 0.15; another 1D convolutional layer (kernel size = 3, stride = 1, padding = 2), followed by BatchNorm and ReLU activation; a max pooling layer (kernel size = 2, padding = 1) and a dropout layer with a probability of 0.15; a fully connected layer with a hidden state size of 64, followed by ReLU activation and a dropout layer (probability = 0.15); and finally, a fully connected binary classification layer with sigmoid activation.

In the blood-brain barrier peptide (BBP) prediction task, the objective is to classify whether a peptide can penetrate the blood-brain barrier (i.e., BBP) (40). We employed the BBPpred dataset (40), consisting of 100



Extended Data Fig. 3 | Model architectures used in downstream tasks. **a**, Model architecture used in type I anti-CRISPR activity detection. **b**, Model architecture used in pathogenic missense variants prediction. **c**, Model architecture used in thermostability prediction. **d**, Model architecture used for identifying potent bioactive peptides. **e**, Model architecture used for predicting minimum inhibitory concentration values.

BBPs and 100 non-BBPs for training, and 19 BBPs and 19 non-BBPs for testing. In the umami peptide prediction task, the goal is to determine whether a peptide elicits an umami taste (41). For this task, we used the iUmami-SCM dataset (62), with a training set of 112 umami peptides and 241 non-umami peptides, and a test set of 28 umami peptides and 61 non-umami peptides. In the antioxidant peptide prediction task, the aim is to classify peptides based on their antioxidant properties (42). We used the AnOxPePred dataset (63), containing 582 antioxidative peptides and 241 non-antioxidative peptides for training, with a test set comprising 28 antioxidative peptides and 61 non-antioxidative peptides. For the antiviral peptide prediction task, the objective is to predict whether a peptide has antiviral activity (preventative and therapeutic against viral infections) (43). We utilized the ABPDiscover dataset (64), which includes 2321 antiviral peptides and 2321 non-antiviral peptides for training, and 623 antiviral peptides and 623 non-antiviral peptides for testing. In the antiparasitic peptide prediction task, the goal is to identify peptides with antiparasitic activity (44). Using the PredAPP

dataset (65), we trained on 255 antiparasitic peptides and 255 non-antiparasitic peptides, and tested on 46 antiparasitic peptides and 46 non-antiparasitic peptides. The tumor T cell antigen prediction task aims to classify peptides capable of inducing a T-cell immune response (45). We used the iTTCA-Hybrid dataset (45), including 470 antigenic peptides and 318 non-antigenic peptides for training, with 122 antigenic peptides and 75 non-antigenic peptides for testing. For the dipeptidyl peptidase IV (DPP-IV) inhibitory peptide prediction task, the goal is to identify peptides that inhibit DPP-IV activity (66). We used the iDPPIV-SCM dataset (66), containing 532 inhibitory peptides and 532 non-inhibitory peptides for training, and 133 inhibitory peptides and 133 non-inhibitory peptides for testing. Lastly, in the neuropeptide (NP) prediction task, the aim is to classify peptides as neuropeptides or non-neuropeptides (47). We used the PredNeuroP dataset (47), containing 1940 neuropeptides and 1940 non-neuropeptides for training, and 485 neuropeptides and 485 non-neuropeptides for testing.

During training, we optimized the model using stochastic

741 gradient descent (SGD) with a learning rate of 1×10^{-2} ,
742 momentum of 0.5, and no weight decay, over 200 epochs.
743 Additionally, we applied step decay for learning rate
744 adjustment and utilized early stopping based on validation
745 accuracy, halting training if no improvement was observed
746 for 40 consecutive epochs. Model performance was assessed
747 using several metrics, including accuracy (ACC), balanced
748 accuracy (BACC) (48), sensitivity (SN), specificity (SP),
749 Matthews correlation coefficient (MCC) (49), and area
750 under the ROC curve (AUC). These metrics were derived
751 from the counts of true positives (TP), false positives (FP),
752 false negatives (FN), and true negatives (TN), and are
753 defined by the following equations:

$$\begin{aligned} \text{ACC} &= \frac{TP + TN}{TP + TN + FP + FN}, \\ \text{SN} &= \frac{TP}{TP + FN}, \\ \text{SP} &= \frac{TN}{TN + FP}, \\ \text{BACC} &= 0.5 \times \text{SN} + 0.5 \times \text{SP}, \\ \text{MCC} &= \frac{(TP \times TN) - (FN \times FP)}{\sqrt{(TP + FN)(TN + FP)(TP + FP)(TN + FN)}}. \end{aligned} \quad (4)$$

754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000

Minimum inhibitory concentration (MIC) value prediction. The model architecture is illustrated in Extended Data Fig. 3e. The sequence encoder, pretrained with ProteinAligner, was fine-tuned to address this task. The classification module is a multi-layer perceptron (MLP) consisting of two fully connected layers, with a hidden size of 256 and a ReLU activation function. We employed the Adam optimizer with an initial learning rate of 1×10^{-4} , training the model for 200 epochs. Throughout the training process, the learning rate was dynamically adjusted at each epoch using the LambdaLR (67) scheduler. To assess the model's performance, we used mean squared error (MSE) as the evaluation metric, defined as:

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2, \quad (5)$$

where n denotes the number of examples, \hat{y}_i represents the model's prediction for the i -th example, and y_i is the corresponding ground-truth value.

Data availability

The FASTA and PDB datasets are publicly available in UniProtKB Swiss-Prot and RCSB PDB, respectively. The FASTA and PDB entries for protein sequences and structures in the pretraining data, along with their textual descriptions, are available at <https://drive.google.com/file/d/1Ff28ajdyUDQl9JtSNQcUz-6Nbz7FXdEA/view?usp=sharing>. All the data used in downstream tasks is also publicly available. The data used in type I anti-CRISPR

activity detection is available at AcrTransAct. The data for predicting the pathogenicity of missense variants is available at VariPred. The data used in thermostability prediction is available at HotProtein. The data used in peptide bioactivity prediction is available at UniDL4BioPep. The data for predicting the minimum inhibitory concentration (MIC) of antimicrobial peptides is available at DeepAMP.

Code availability

The source code for ProteinAligner pretraining, along with the pretrained checkpoints, is available at <https://github.com/Alexiland/ProteinAligner>. Additionally, links to the code for downstream tasks can be found in the README file.

References

- Baris E Suzek, Yuqi Wang, Hongzhan Huang, Peter B McGarvey, Cathy H Wu, and UniProt Consortium. Uniref clusters: a comprehensive and scalable alternative for improving sequence similarity searches. *Bioinformatics*, 31(6):926–932, 2015.
- Ilya Loshchilov, Frank Hutter, et al. Fixing weight decay regularization in adam. *arXiv preprint arXiv:1711.05101*, 5, 2017.
- Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. In *International Conference on Learning Representations*, 2022.
- Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International Conference on Machine Learning*, pages 448–456. PMLR, 2015.
- Vinod Nair and Geoffrey E Hinton. Rectified linear units improve restricted boltzmann machines. In *International Conference on Machine Learning*, pages 807–814, 2010.
- Diederik P Kingma. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15(56):1929–1958, 2014.
- Andrew L. Maas, Awni Y. Hannun, and Andrew Y. Ng. Rectifier nonlinearities improve neural network acoustic models. In *International Conference on Machine Learning*, pages 3–8, 2013.
- Jingjing Xu, Xu Sun, Zhiyuan Zhang, Guangxiang Zhao, and Junyang Lin. Understanding and improving layer normalization. *Advances in neural information processing systems*, 32, 2019.
- Leslie N. Smith. Super-convergence: Very fast training of neural networks using large learning rates. In *International Conference on Machine Learning (ICML)*, pages 1–10, 2018.
- Phasit Charoenkwan, Janchai Yana, Chanin Nantasenamat, Md Mehedi Hasan, and Watshara Shoombuatong. iumami-scm: a novel sequence-based predictor for prediction and analysis of umami peptides using a scoring card method with propensity scores of dipeptides. *Journal of Chemical Information and Modeling*, 60(12):6666–6678, 2020.
- Tobias Hegelund Olsen, Betül Yesiltas, Frederikke Isa Marin, Margarita Pertseva, Pedro J Garcia-Moreno, Simon Gregersen, Michael Toft Overgaard, Charlotte Jacobsen, Ole Lund, Egon Bech Hansen, et al. Anoxpepred: using deep learning for the prediction of antioxidative properties of peptides. *Scientific reports*, 10(1):21471, 2020.
- Sergio A Pinacho-Castellanos, César R García-Jacas, Michael K Gilson, and Carlos A Brizuela. Alignment-free antimicrobial peptide predictors: improving performance by a thorough analysis of the largest available data set. *Journal of Chemical Information and Modeling*, 61(6):3141–3157, 2021.
- Wei Zhang, Enhua Xia, Ruyi Dai, Wending Tang, Yannan Bin, and Junfeng Xia. Predapp: predicting anti-parasitic peptides with undersampling and ensemble approaches. *Interdisciplinary Sciences: Computational Life Sciences*, pages 1–11, 2022.
- Phasit Charoenkwan, Sakawrat Kanthawong, Chanin Nantasenamat, Md Mehedi Hasan, and Watshara Shoombuatong. idppiv-scm: a sequence-based predictor for identifying and analyzing dipeptidyl peptidase iv (dpp-iv) inhibitory peptides using a scoring card method. *Journal of proteome research*, 19(10):4125–4136, 2020.
- Leslie N. Smith. Cyclical learning rates for training neural networks. In *2017 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 464–472. IEEE, 2017.