

552 Appendix

553 A Experimental setting

554 **Data Descriptions for Different Stages of HTP.** (1) For the single-protein sequence level, we em-
555 ploy the state-of-the-art ESM-2 [25], which outperforms all tested single-sequence protein language
556 models across a wide range of structure prediction tasks and enables atomic resolution structure
557 prediction. It is trained on 86 billion amino acids across 250 million protein sequences spanning
558 evolutionary diversity. Specifically, ESM uses UniRef50, September 2021 version. The training
559 dataset was partitioned by randomly selecting 0.5% ($\approx 250,000$) sequences to form the validation
560 set. The training set has sequences removed via the procedure described in Hie et al. [48]. ESM-2
561 runs MMseqs search to obtain the query and target databases. All train sequences which match a
562 validation sequence with 50% sequence identity under this search are removed from the train set. The
563 details of the ESM series can be found in <https://github.com/facebookresearch/esm>.

564 (2) For the antibody sequence level, we use the Observed Antibody Space database (OAS) [32] and its
565 succeeding update [33] as the pretraining data. It currently contains over one billion sequences, from
566 over 80 different studies that cover diverse immune states, organisms, and individuals, which can
567 be downloaded from its official website at <https://opig.stats.ox.ac.uk/webapps/oas/>. We
568 upload the processed paired data in [https://pan.baidu.com/s/181B8g19Maf0nnNPIw83ZzA?](https://pan.baidu.com/s/181B8g19Maf0nnNPIw83ZzA?pwd=1212)
569 <https://pan.baidu.com/s/161gU8fso6rz6-QGfNoCoHQ?pwd=96uF> (password: 1212) as well as the unpaired data in [https://pan.baidu.com/s/](https://pan.baidu.com/s/161gU8fso6rz6-QGfNoCoHQ?pwd=96uF)
570 [161gU8fso6rz6-QGfNoCoHQ?pwd=96uF](https://pan.baidu.com/s/161gU8fso6rz6-QGfNoCoHQ?pwd=96uF) (password: 96uf).

571 (3) For the protein-protein complex structure level, we leverage the Database of Interacting Protein
572 Structures (DIPS) [35]. It is a large protein complex structure dataset than existing antibody-antigen
573 complex structure datasets and is mined from the Protein Data Bank [36]. We attain the database
574 from Atom3d in Zendo <https://zenodo.org/record/4911102>, which is a collection of both
575 novel and existing benchmark datasets spanning several key classes of biomolecules. Referring to
576 Atom3d, we split protein complexes by sequence identity at 30%, resulting in train/validation/test
577 sets with 87,303/31,050/15,268 instances.

578 (4) For the antibody-antigen complex structure level, we select all available antibody-antigen pro-
579 tein complexes from SAbDab [17] at [https://opig.stats.ox.ac.uk/webapps/newsabdab/](https://opig.stats.ox.ac.uk/webapps/newsabdab/sabdab/)
580 [sabdab/](https://opig.stats.ox.ac.uk/webapps/newsabdab/sabdab/), leading to a dataset containing 9,823 structures. CDRs are identified using the antibody
581 numbering program AbRSA [49]. Following the setting in [13], the chosen data points are divided
582 into training and test data based on their release date and CDR sequence identity. To be explicit,
583 the test split contains protein structures released after December 24, 2021, as well as structures with any
584 CDR similar to those released after this date with sequence identity higher than 50%. Antibodies in
585 the test set are further clustered with 50% CDR sequence identity to remove duplicates, finally lead-
586 ing to 21 antibody-antigen structures. The training and validation splits just include complexes not
587 involved during the curation of the test split. After that, we randomly split the remaining complexes
588 with a ratio of 90% and 10% into the training and validation sets.

589 **Implementation Details.** HPT is implemented in PyTorch and PyTorch Geometric packages. For
590 all four training stages, we leverage an Adam optimizer [50] with a weight decay of $1e-15$. All
591 experiments are run on multiple A100 GPUs, each with a memory storage of 80G.

592 (1) For ESM-2 in the single-protein sequence level training, we adopt a middle-size version, which
593 has a parameter number of 150M, 30 layers, and a hidden dimension of 640. Besides, we append the
594 ESM-2 with a three-layer perceptron to forecast the residue type for MLM.

595 (2) For the antibody sequence level training, we use a batch size of 2 to avoid the out-of-memory
596 error and 4 workers to load the data. The number of epochs is 100 and starting learning rate is $1e-5$.
597 Apart from that, we utilize a ReduceLRonPlateau scheduler with a factor of 0.6, patience of 5 epochs,
598 and a minimum learning rate of $1e-7$.

599 (3) For the protein-protein complex structure level training, we use a batch size of 32, 1000 epochs,
600 and 4 works to speed up data loading. The starting learning rate is $1e-4$, and a ReduceLRonPlateau
601 scheduler is utilized to adjust the learning rate automatically with a factor of 0.6 and patience of 3
602 epochs. We adopt a distance threshold of 8.0\AA to determine the connection between different graph
603 nodes (*i.e.*, the alpha carbon of each residue). As for the loss weight balance, we set $\lambda = 1$.

604 (4) For the antibody-antigen complex structure level training, we also adopt the distance threshold of
 605 8.0Å to build the graph connection. For random initialization of CDR coordinates, we use a noise of
 606 $\epsilon = 0.1$. As for the other important hyperparameters, we use a grid search mechanism to find the
 607 optimal combination. Notably, the geometric neural networks used in the third and fourth levels are
 608 matched to each other. If we alter the setting of GGNNs in the antibody-antigen complex structure
 609 level training, we need to retrain it in the protein-protein complex structure level first. The entire
 hyperparameter search space is depicted in Table 4.

Table 4: Hyperparameters setup for HTP.

Hyperparameters Search Space	Symbol	Value
Training Setup		
Epochs	–	[100, 500, 1000]
Batch size	–	[32, 64, 128]
Learning rate	–	[1e-4, 5e-5, 1e-6, 1e-7]
Warmup	–	[Yes, No]
Warmup epochs	–	[10, 20]
Loss Balance weight for Coordinates and Residue Types	λ	[0.1, 0.3, 0.5, 0.7]
GNN Architecture		
Dropout rate	–	[0.1, 0.2]
Number of GNN layers	L	[2, 4, 6]
Tanh activation function	–	[Yes, No]
Coordinate Normalization	–	[Yes, No]
The hidden dimension of node representations	–	[320, 640]
The hidden dimension of edge representations	–	[16, 32, 64]

610

611 **Reproduction of Baselines.** Concerning the implementation of several baseline methods, we
 612 use the official repositories for conditional RefineGNN ([https://github.com/wengong-jin/](https://github.com/wengong-jin/RefineGNN/)
 613 [RefineGNN/](https://github.com/wengong-jin/abdockgen)), HERN (<https://github.com/wengong-jin/abdockgen>), DiffAb (<https://github.com/luost26/diffab>). To reproduce the performance of existing antibody-specific
 614 pretrained PLMs, we download the code from <https://github.com/alchemab/antiberta> for
 615 AntiBERTa and <https://github.com/oxpig/AbLang> for AbLang. In our comparison, we directly
 616 use their pretrained residue features as the input for GGNNs without any fine-tuning.
 617

618 **Code Availability.** All relevant Python code to reproduce the results in our paper is stored in
 619 GitHub repository at <https://github.com/smiles724/HTP>.

620 B Additional Results

621 B.1 Ablation Study

622 We investigate the effectiveness and necessity of each component of our HTP. As shown in Table 5,
 623 the removal of protein-protein complex structure level induces performance detriment, where RMSD
 624 increases from 2.06 to 2.49. Moreover, we implement a variant of HTP by replacing features obtained
 625 by pretrained PLMs with learnable embedding features, whose performance is worse than HTP. To
 626 be concise, AAR declines from 40.98 to 25.31, and RMSD rises from 2.06 to 2.65. In summary, our
 627 HTP brings significant relative improvements of 78.56% in AAR, 41.97% in RMSD, and 2.94% in
 628 TM-Score. This phenomenon strongly supports the superiority of our approach over existing naive
 629 co-design algorithms that are trained only on antibody-specific structure data.

630 C Related Work

631 **Antibody Design.** The majority of old-school computational approaches for antibody design
 632 are based on sampling algorithms over hand-crafted and statistical energy functions to iteratively
 633 modify protein sequences and structures [7, 9, 10, 51, 52]. These physics-based algorithms are
 634 computationally expensive and prone to be stuck in local energy minimum, which triggers the

Table 5: Effects of each module, where SPS stands for the single-protein sequence level, PPCS denotes the protein-protein complex structure level, and AS represents the antibody sequence level. The last row computes the relative improvements of HTP over the primitive baseline without any protein data augmentation.

	SPS	AS	PPCS	SAbDab (CDR-H3)		
				AAR (%) \uparrow	RMSD \downarrow	TM-Score
1	\times	\times	\times	22.95 ± 0.5	3.55 ± 0.01	0.9146 ± 0.003
2	\checkmark	\times	\times	33.87 ± 0.8	2.77 ± 0.04	0.9450 ± 0.006
3	\checkmark	\checkmark	\times	38.42 ± 1.6	2.49 ± 0.03	0.9538 ± 0.004
4	\times	\times	\checkmark	25.31 ± 0.7	2.95 ± 0.02	0.9391 ± 0.005
5	\checkmark	\checkmark	\checkmark	40.98 ± 1.5	2.06 ± 0.03	0.9621 ± 0.005
Imp.	–	–	–	78.56%	41.97%	2.94%

635 adaptation of deep learning in this sub-field. The initial researchers [4, 5, 53, 54] use pure PLMs to
 636 generate protein sequences but disregard the available antigen structures.

637 To circumvent this, Jin et al. [11] introduce RefineGNN, the first co-design architecture that aims
 638 to neutralize SARS-CoV-2. Later, HERN [12] is proposed as a more general version for paratope
 639 docking and design, which opens the door to produce antibodies given arbitrary antigen structures.
 640 Subsequent efforts are spent in either modifying the generative style or utilizing more advanced
 641 deep learning architectures such as diffusion denoise probabilistic models (DDPMs). For example,
 642 DiffAb [13] achieves atomic-resolution antibody design with SO(3)-equivariance, while MEAN [15]
 643 corrects the autoregressive manner with a full-shot one to prevent low efficiency and accumulated
 644 errors during inference.

645 **Protein Sequence Modeling.** Sequence-based protein representation learning is mainly inspired
 646 by the field of natural language processing. A large body of early works concentrates on modeling
 647 individual protein families [55], solving problems like functional nanobody design [5]. Its success,
 648 then, motivates the prospective trend to model large-scale databases of protein sequences by means
 649 of unsupervised learning. This line of study targets capturing the biochemical and co-evolutionary
 650 knowledge that underlies a large-scale protein sequence corpus by self-supervised pertaining. Thanks
 651 to them, a number of pertaining objects have been explored such as the next amino acid predic-
 652 tion [4, 26], masked language modeling (MLM) [55, 23], pairwise MLM [56], contrastive predictive
 653 coding [57], conditional generation [58], and position-specific scoring matrix prediction [59]. In
 654 addition, another line [60, 61] is based on multiple sequence alignment (MSA), leveraging sequences
 655 within a protein family to seize the conserved and variable regions of homologous sequences. Notably,
 656 some schemes for protein sequence modeling also seek to incorporate structural information in either
 657 the pretraining stage [62, 63] or the finetuning stage [64].

658 The improvements in model scale and architecture are also crucial to the recent achievement of
 659 PLMs. Explicitly, Rao et al. [55] evaluate various PLMs in a panel of benchmarks and discover that
 660 multi-head attention outpaces the Potts model in contact prediction, even if using a single sequence
 661 for inference. Concurrently, Vig et al. [65] observe that specific attention heads of pretrained
 662 Transformers have straight correlations with protein contact. Others [26] investigate a variety of
 663 Transformer variants [66] and demonstrate that large Transformers can procure state-of-the-art
 664 features across diverse tasks. Apart from that, the latest ESM-2 [25] trains the largest PLM with 15B
 665 parameters and shows that as models are scaled, they learn information enabling the protein structure
 666 prediction at the resolution of individual atoms.

667 **Protein Structure Learning.** With the rapid advance of geometric deep learning, it has been
 668 increasingly attractive and challenging to represent and reason about structures of macromolecules
 669 in the 3D space. For the sake of encoding spatial information in protein structures including bond
 670 lengths and dihedral angles, numerous 3D geometric neural networks such as 3DCNN [67–70] or
 671 GNNs [45, 46, 71, 72] have been invented. They excel at capturing complex interactions between
 672 sets of amino acids [73] and attain pivotal Euclidean geometry, *e.g.*, E(3) or SE(3)-equivariance and
 673 symmetry.

674 However, compared to protein sequences in databases like UniProt [74] or Pfam [75], the known
675 structures in the PDB are scarce and hard to obtain. Therefore, it becomes an urgent need to develop
676 structure-based mechanisms to efficiently learn protein representations with much less pretraining
677 data. For instance, Hermosilla and Ropinski [18] use contrastive learning in terms of molecular
678 substructures to help models understand protein structure similarity and functionality. Moreover, Chen
679 et al. [76] propose a self-supervised framework that predicts angles and inter-residue distances.
680 Additionally, Guo et al. [77] present a coordinate denoising score matching method. Wu et al. [19] put
681 forward a novel prompt-based denoising conformation generative pretraining method based on the
682 trajectories of molecular dynamics simulations. A recent attempt [34] makes a combination of both
683 contrastive learning and self-prediction with more intriguing augmentation functions. Despite this
684 progress, all of them are dealing with single-protein structures. No preceding studies have considered
685 structure-based pretraining in the circumstance of multiple proteins. That is, how to pretrain on
686 protein-protein complex, or more specifically, the antibody-antigen complex, remains unexplored.

687 **D Limitations and Future Work**

688 In spite of the promising progress of our HTP, there is still some space left for future explorations. First,
689 more abundant databases can be exploited in our framework. For instance, AntiBodies Chemically
690 Defined (ABCD) [78] is a large-sized antibody sequence database that can be used to enhance the
691 capacity of protein language models at the second level. We do not use it in our work because our
692 request for this database has not been approved by the authors so far. Secondly, we fix the language
693 models during the last two levels of training (*i.e.*, levels that need complex structure prediction) for
694 simplicity and use them as the node feature initializer. It might be beneficial if both PLM and the
695 geometric encoder are tuned.