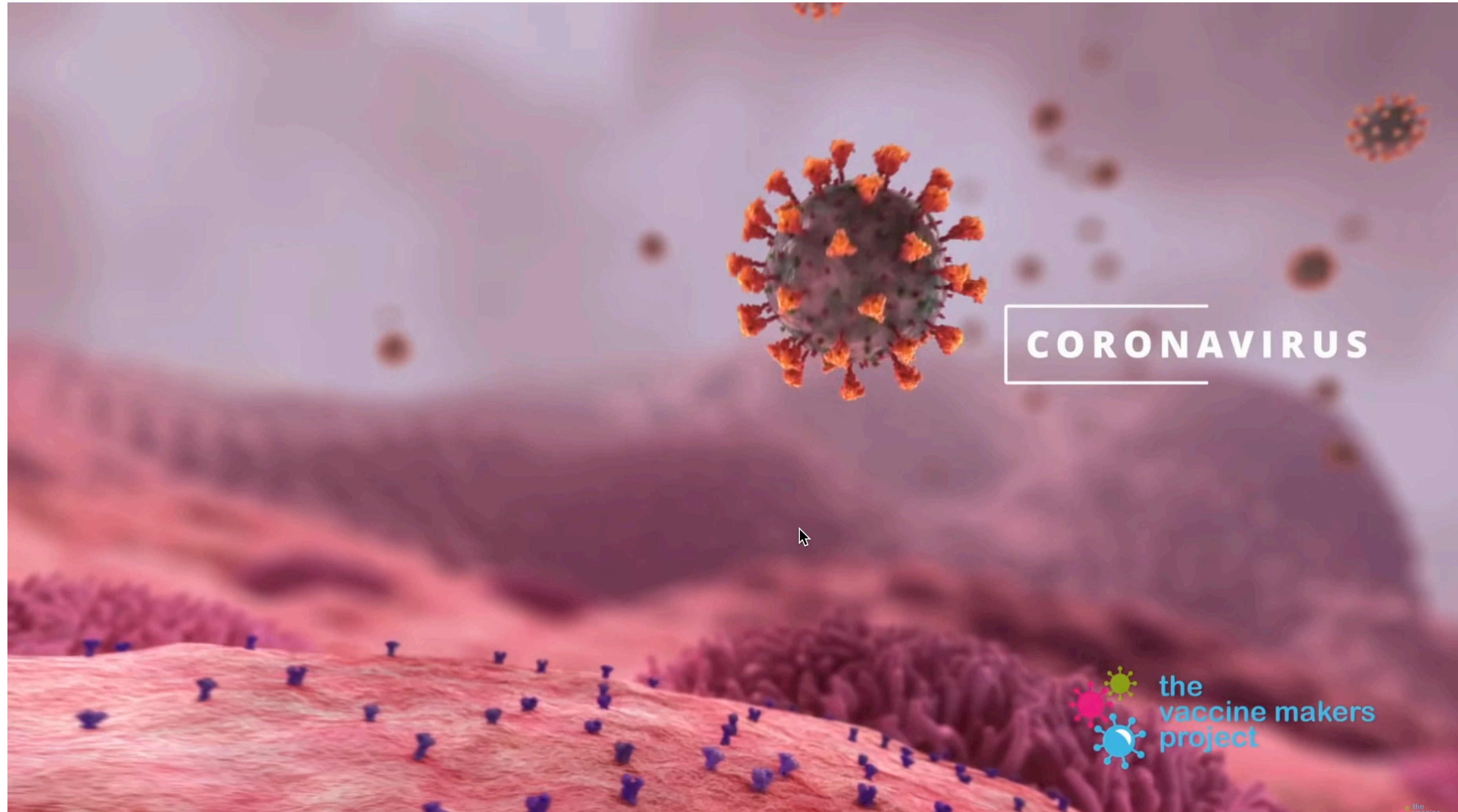


Designing protein-protein interactions with self-supervised geometric deep learning

Anton Bushuiev^{#*1}, Roman Bushuiev^{#1,4}, Anatolii Filkin¹, Petr Kouba^{1,3},
Marketa Gabrielova¹, Michal Gabriel¹, Jiri Sedlar¹,
Tomas Pluskal⁴, Jiri Damborsky^{2,3}, Stanislav Mazurenko^{2,3}, **Josef Sivic**¹

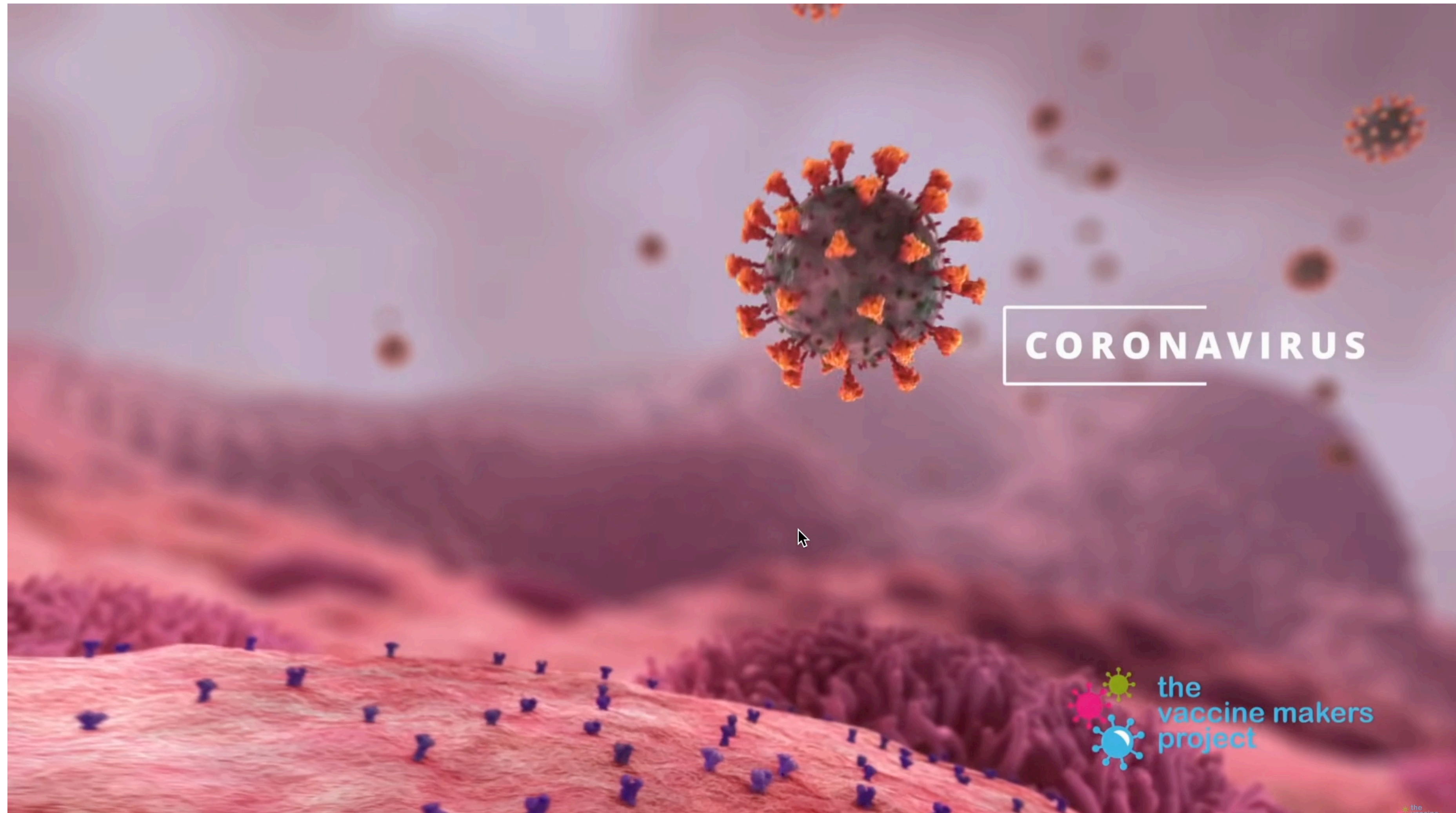


Protein—protein interactions



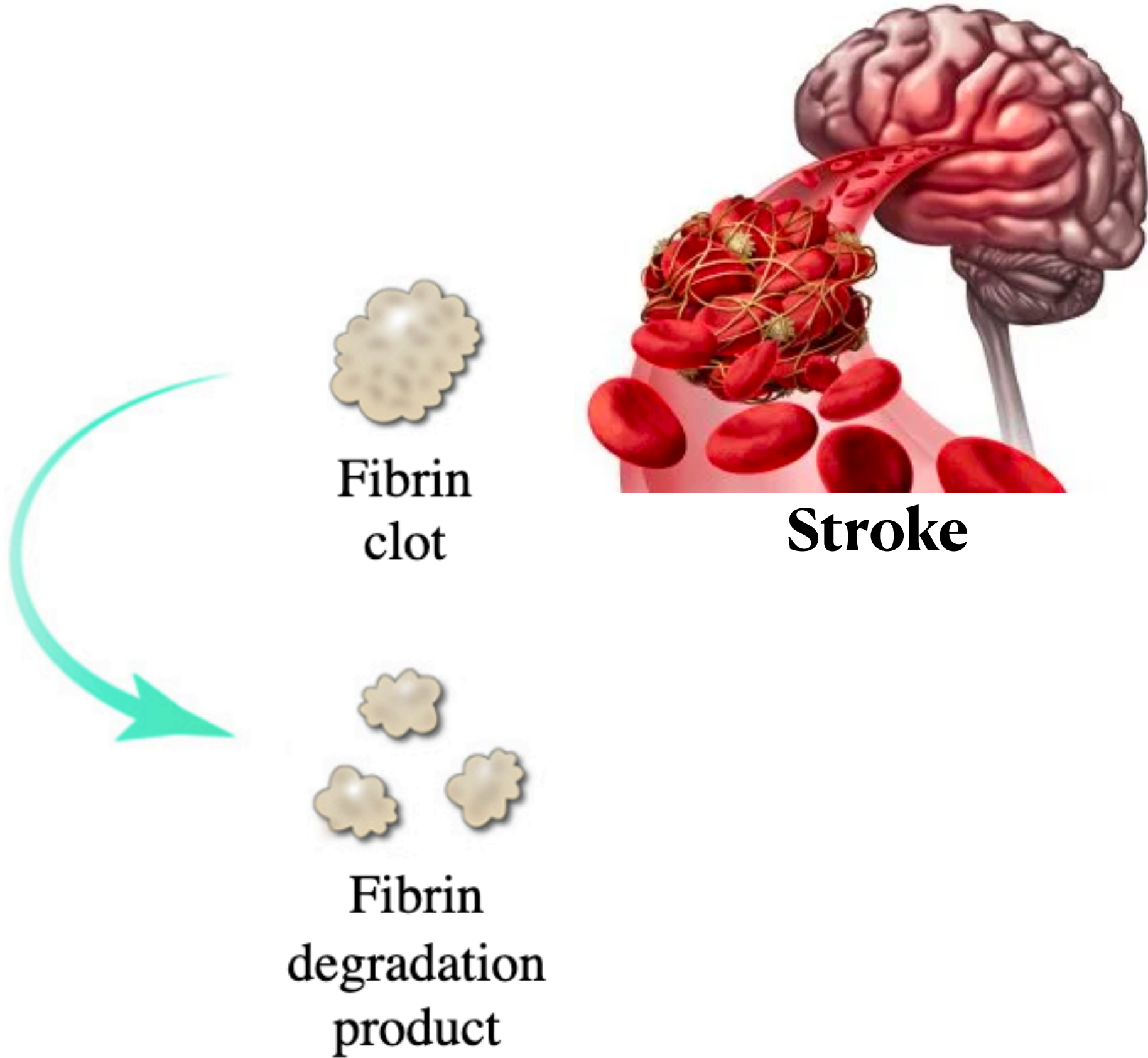
Directly linked to the development and treatment of **viruses, stroke, cancer, Alzheimer, ...**

Protein—protein interactions

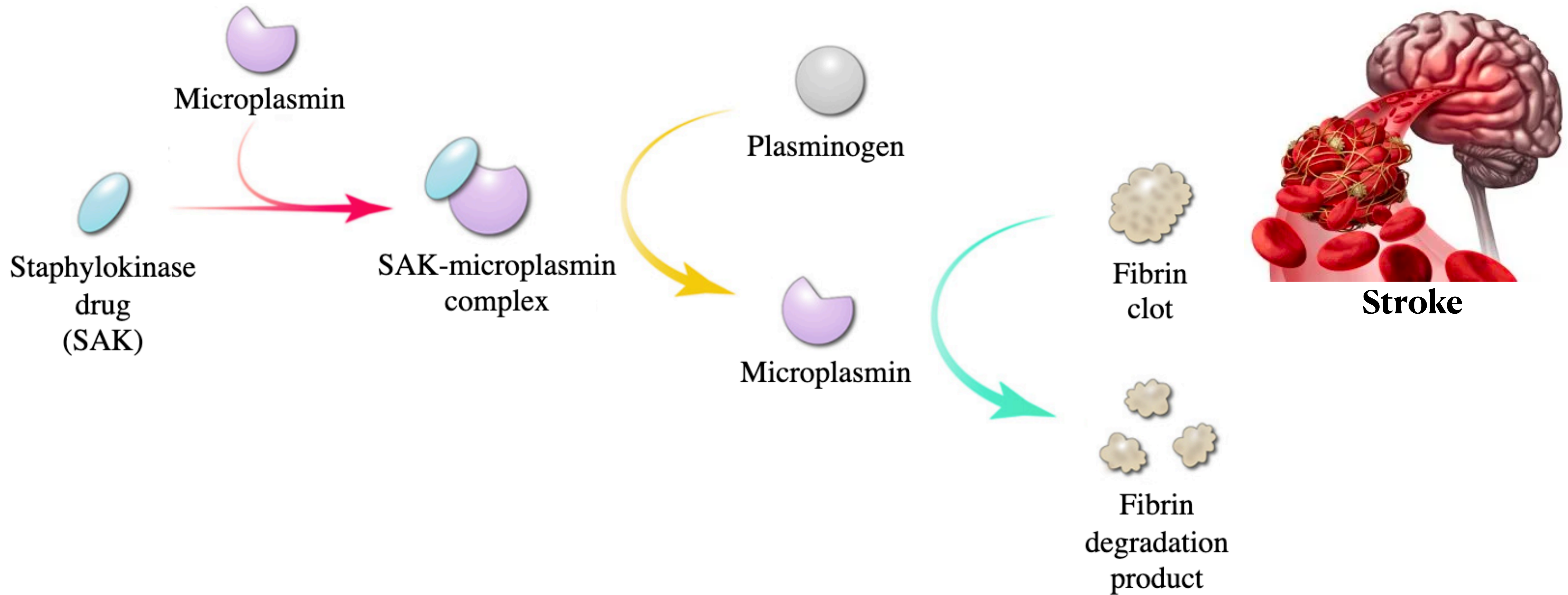


Directly linked to the development and treatment of **viruses, stroke, cancer, Alzheimer, ...**

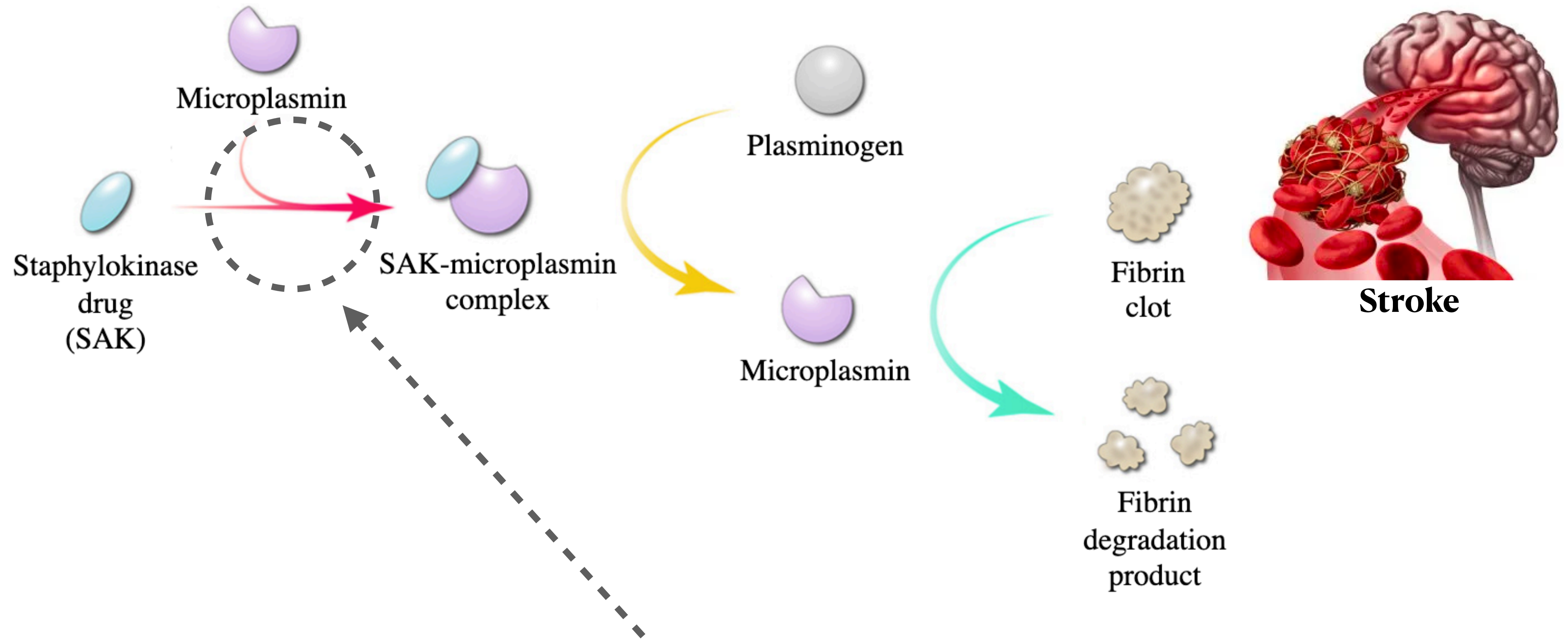
Staphylokinase: thrombolytic drug candidate



Staphylokinase: thrombolytic drug candidate



Staphylokinase: thrombolytic drug candidate

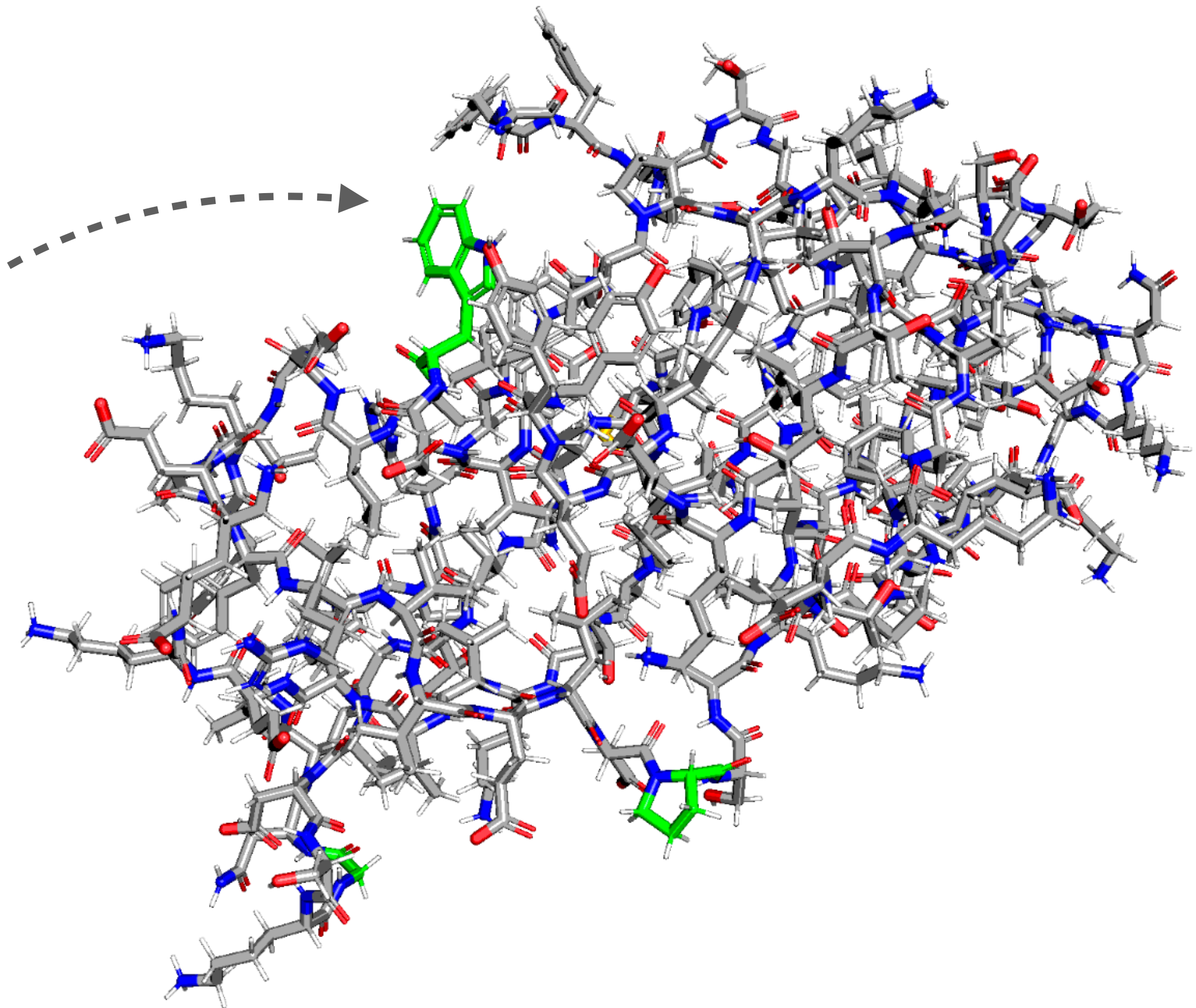
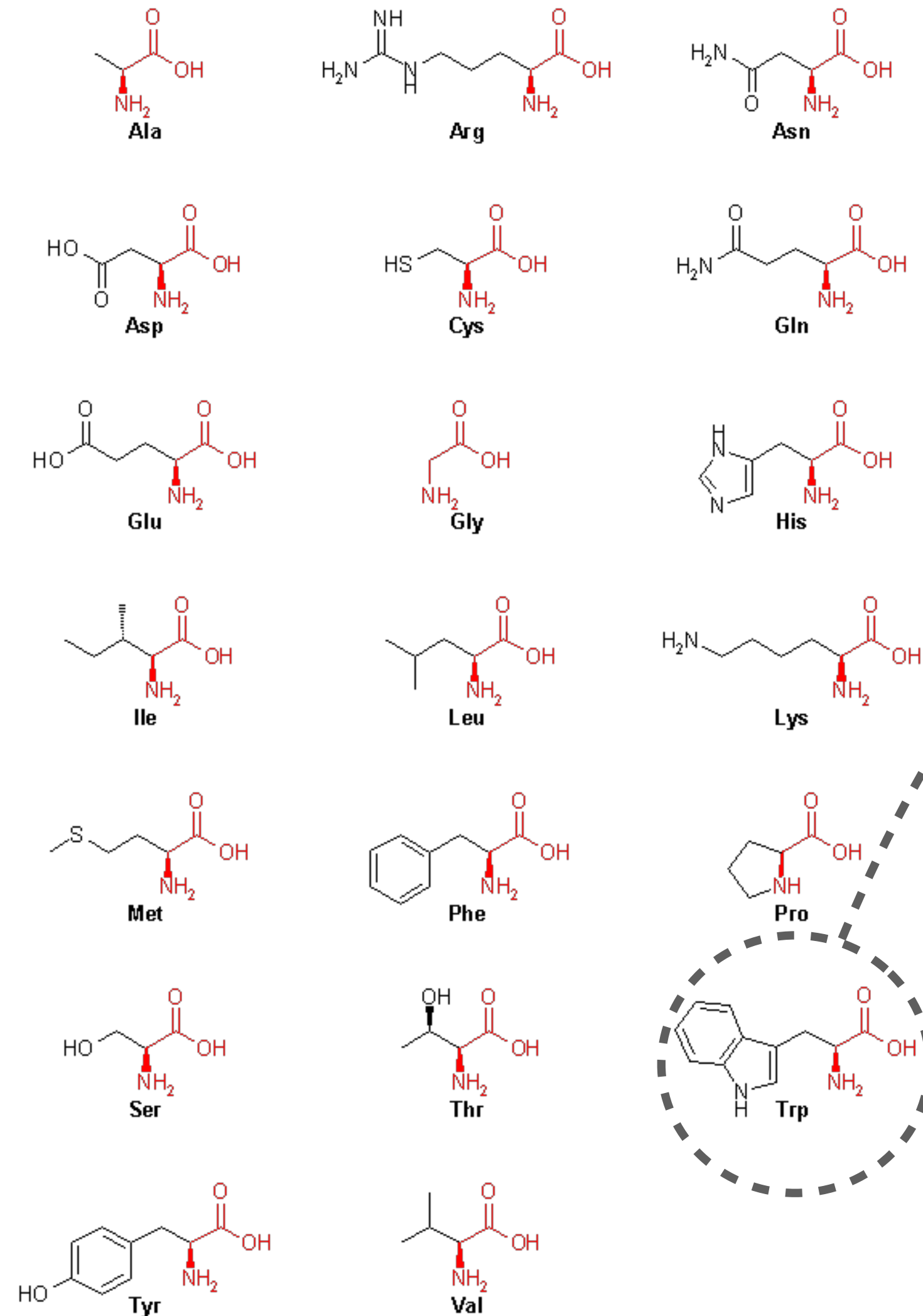


How to enhance the *binding affinity* of the interaction for effective thrombolysis?

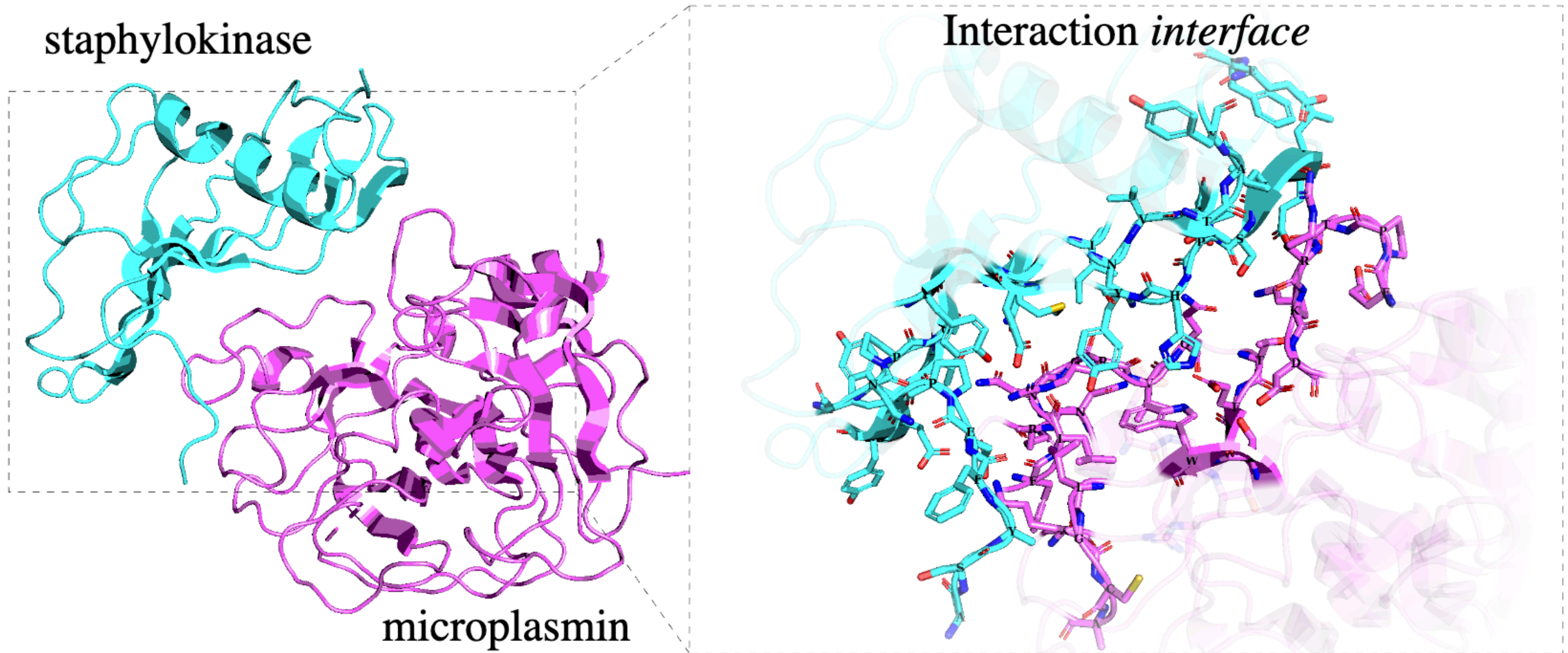
What is a protein?

20 amino acids (building block types)

Protein: a folded chain of amino acids

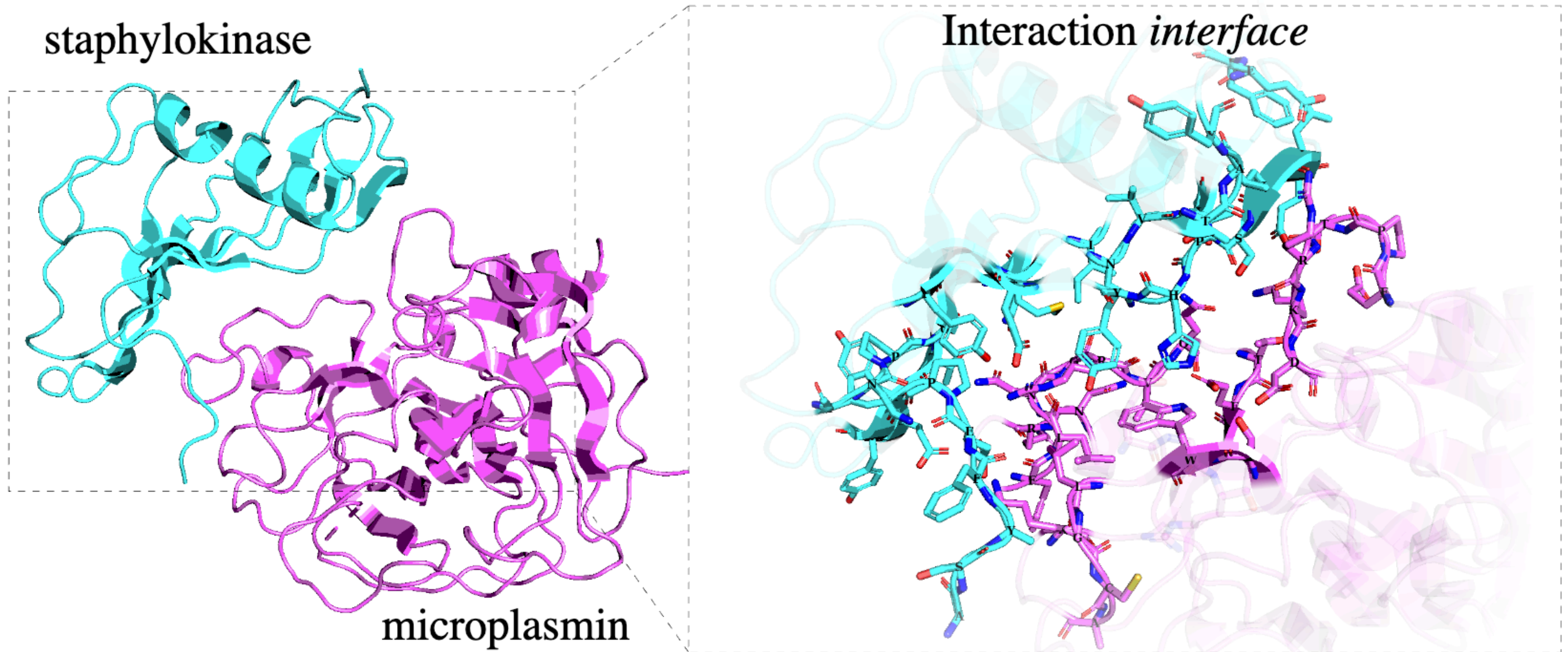


Staphylokinase—microplasmin interaction



How to enhance the *binding affinity* of the interaction?

Staphylokinase—microplasmin interaction



What amino acids of staphylokinase to mutate and how? 20^n combinations

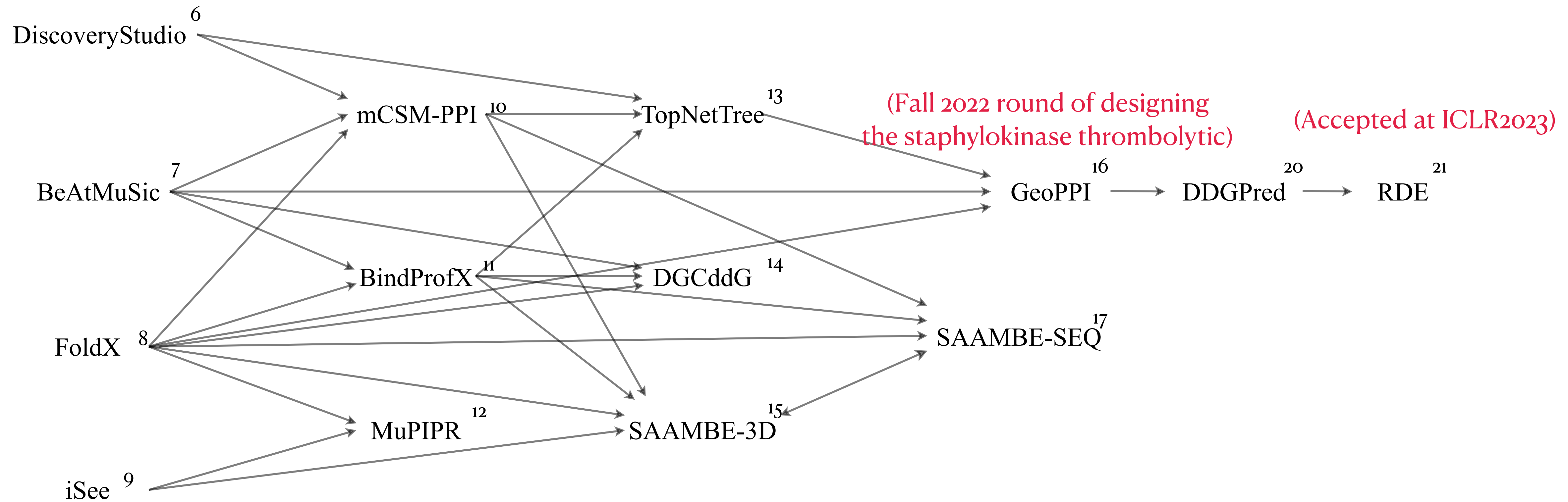
Standard approach: $\Delta\Delta G$ screening

1. Screen thousands or millions of mutations according to $\Delta\Delta G$ – binding energy change upon mutation ranging roughly in $[-12, 12]$
2. Select several best candidates (with lowest $\Delta\Delta G$) and test in a lab



Staphylokinase mutants

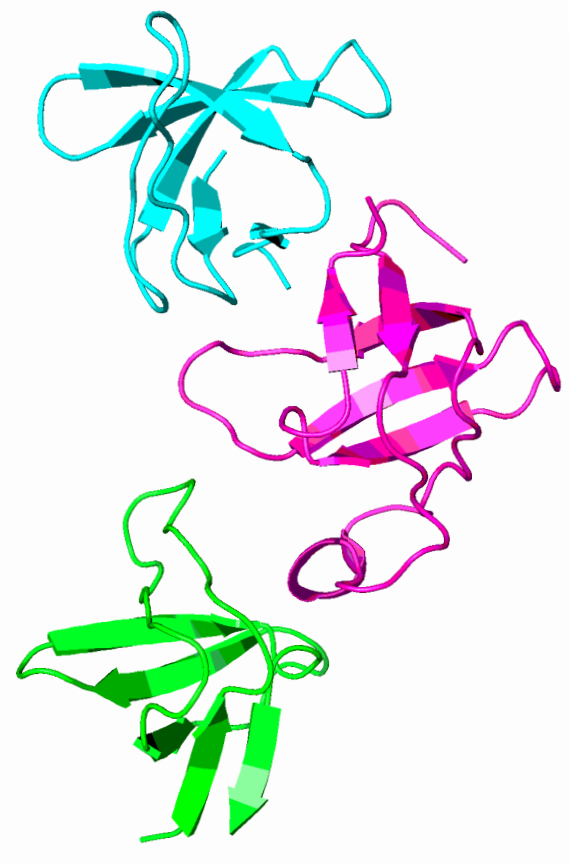
State of the art for predicting $\Delta\Delta G$



- Rely on **small data** (7K annotated mutations from SKEMPI2) → unstable, weak generalization
- Often require **mutant 3D structure** → slow
- **Weak evaluation protocol** → poor generalization

Labeled data (SKEMPI2)

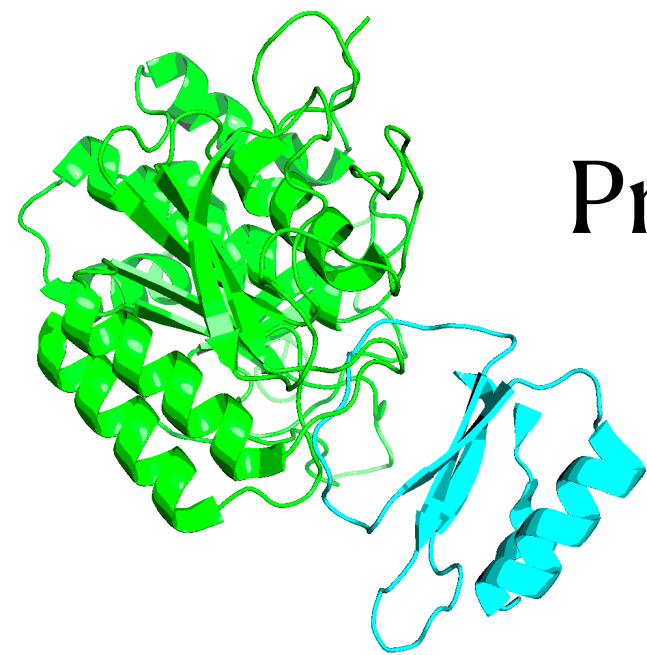
300 interactions, 7K mutations



Protein C, amino acid 21

Cys \rightarrow Val

$\Delta\Delta G = -0.025$



Protein I, amino acid 45

Leu \rightarrow Ser

$\Delta\Delta G = 1.17$

PPIRef: New large dataset of PPIs

Labeled data (SKEMPI2)

300 interactions, 7K mutations

Unlabeled data (Protein Data Bank)

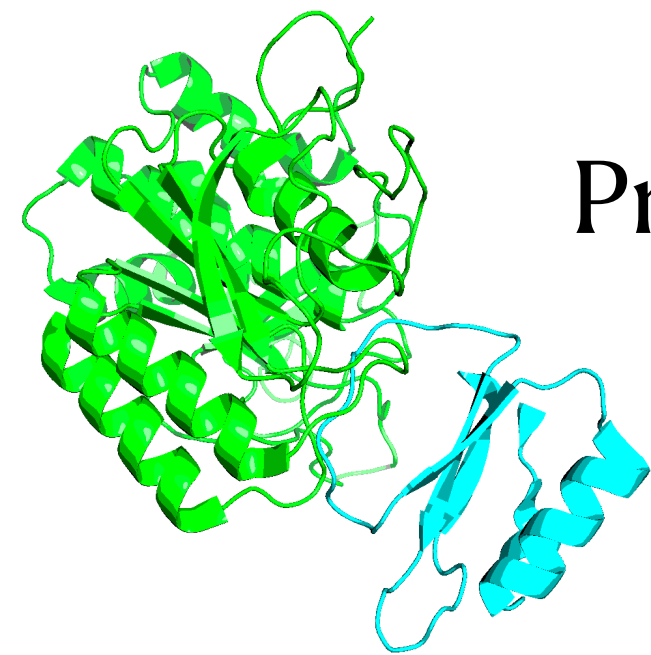
322K interactions in our **PPIRef**, **41K** in DIPS



Protein C, amino acid 21

Cys → Val

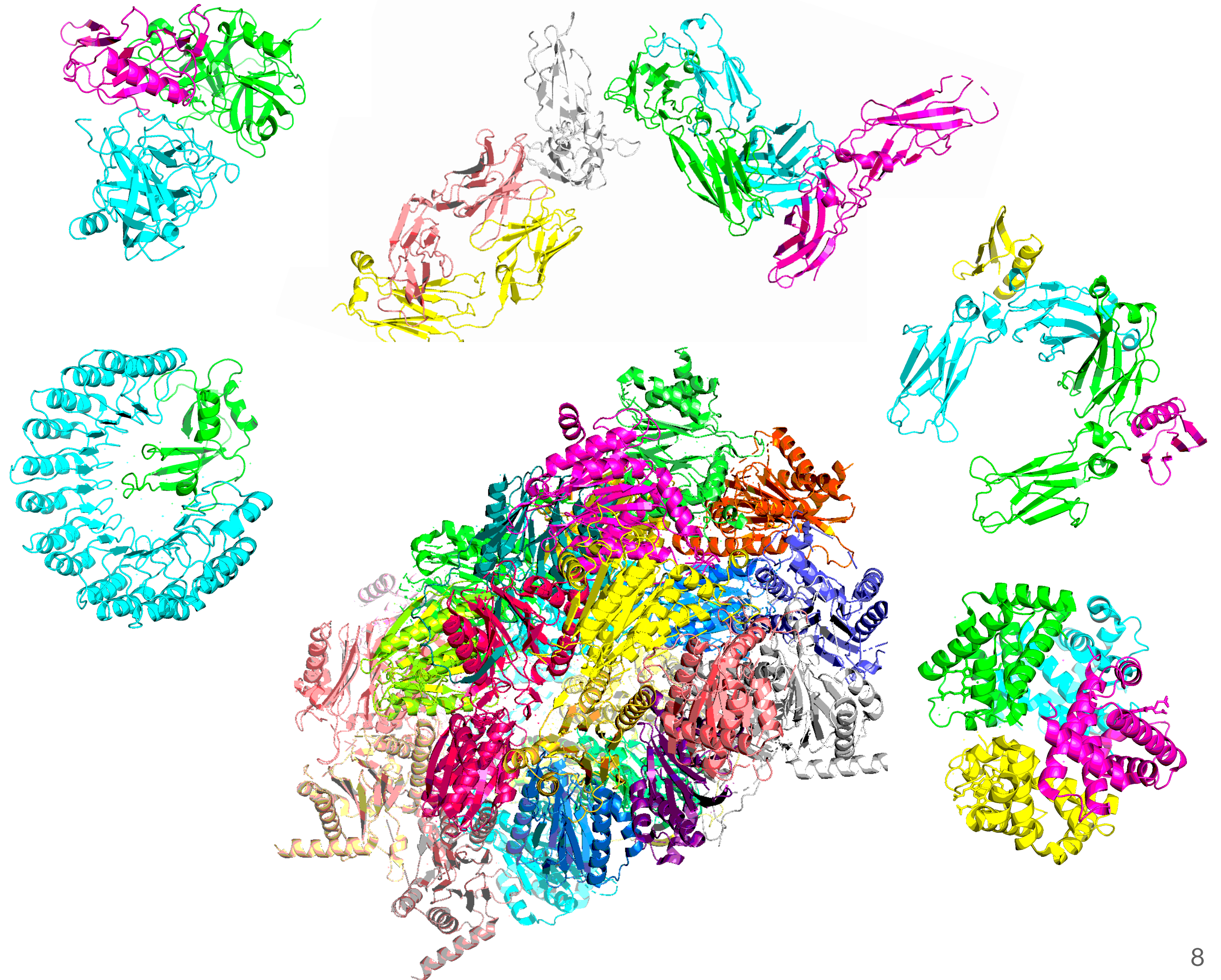
$\Delta\Delta G = -0.025$



Protein I, amino acid 45

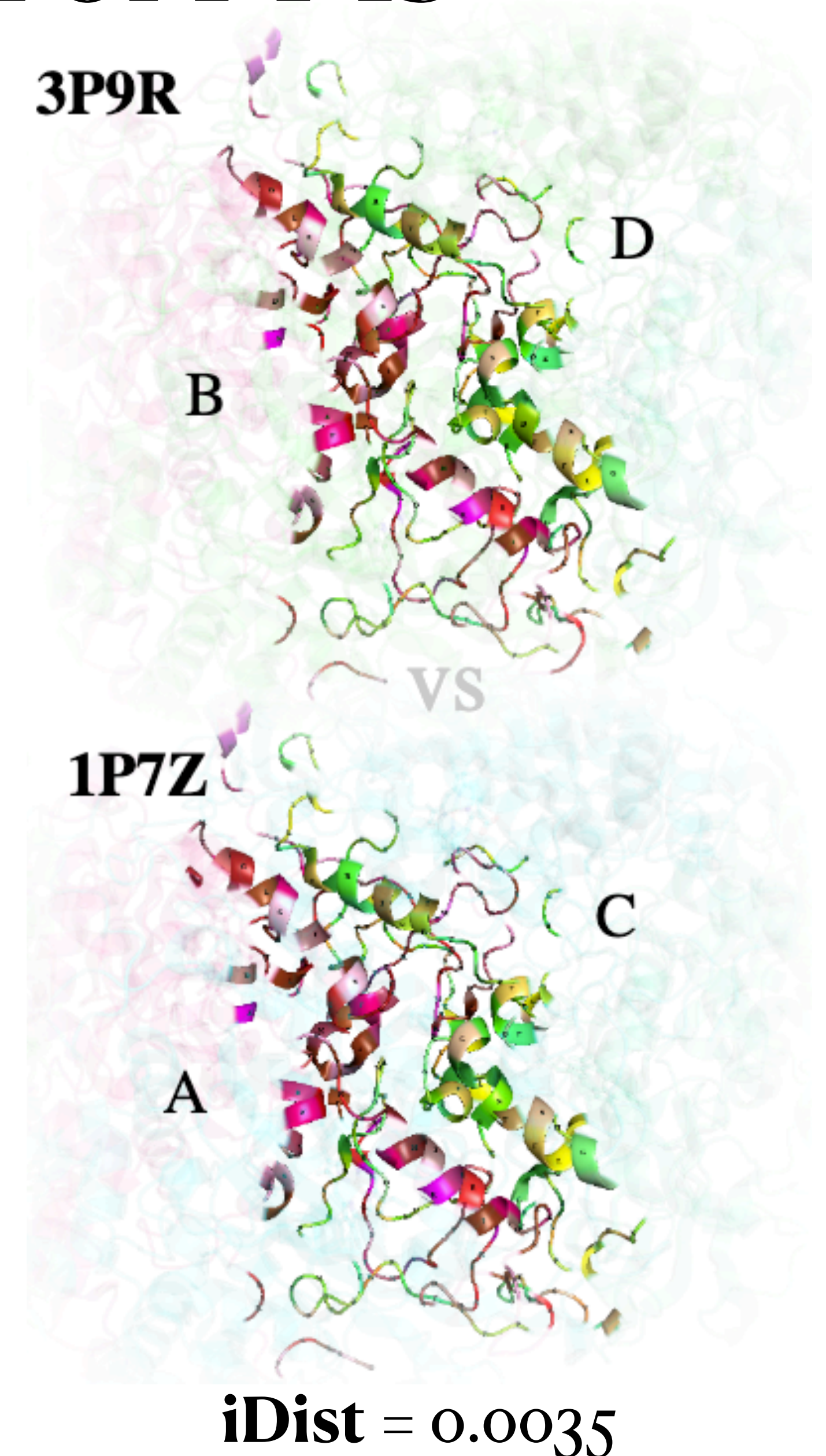
Leu → Ser

$\Delta\Delta G = 1.17$



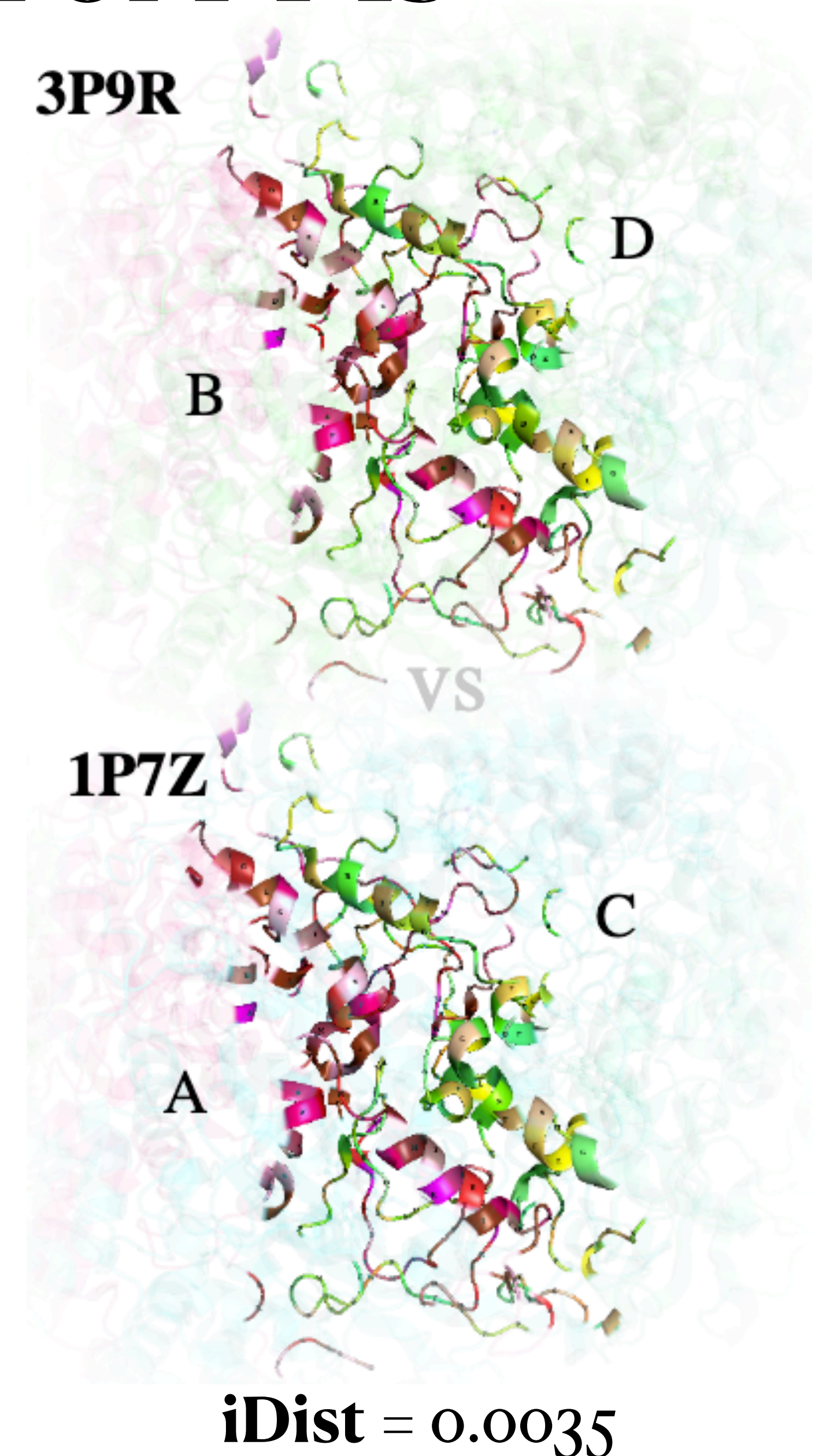
iDist: Scalable comparison of PPIs

- **iDist** accurately approximates iAlign¹⁸ (TM-score for PPIs)
(near-duplicate detection with 99% precision and 97% recall)
- **iDist** is ~500 times faster than iAlign



iDist: Scalable comparison of PPIs

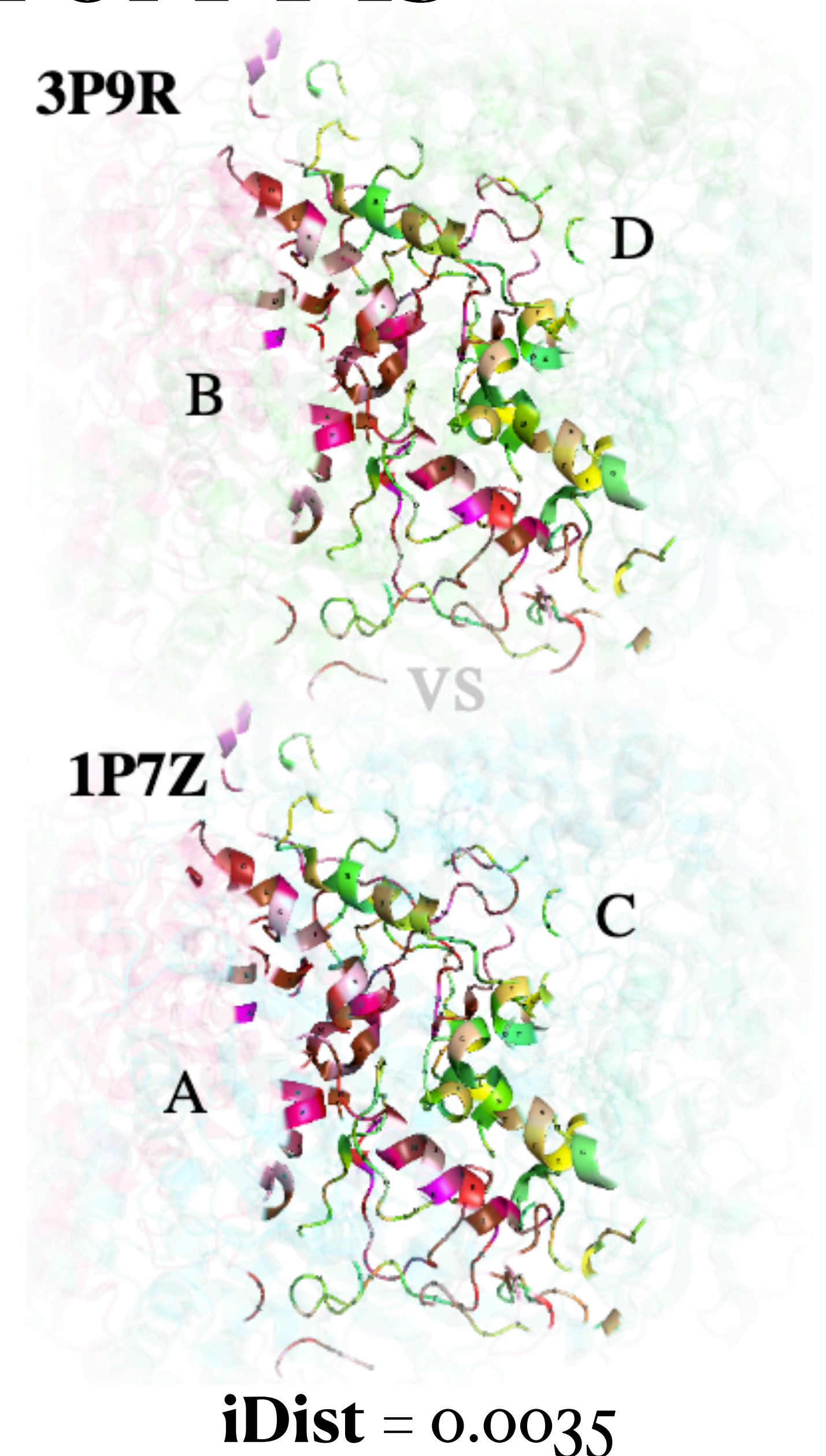
- **iDist** accurately approximates iAlign¹⁸ (TM-score for PPIs)
(near-duplicate detection with 99% precision and 97% recall)
- **iDist** is ~500 times faster than iAlign
- Available datasets are redundant and incomplete
(Some PPIs represented > 500 times, many missing)
- Existing train-test splits suffer from data leakage
(>53% of test sets have near-duplicates in train sets)



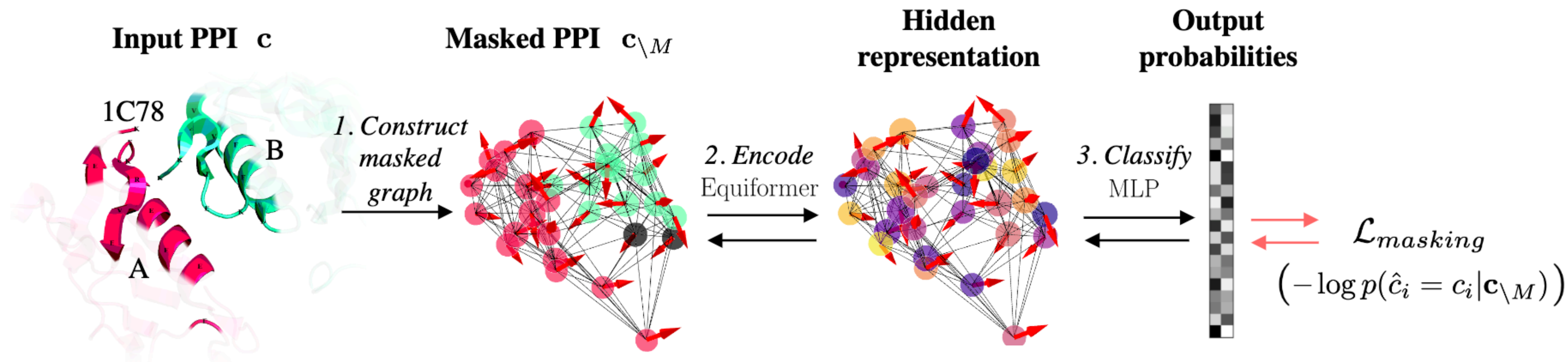
iDist: Scalable comparison of PPIs

- **iDist** accurately approximates iAlign¹⁸ (TM-score for PPIs)
(near-duplicate detection with 99% precision and 97% recall)
- **iDist** is ~500 times faster than iAlign
- Available datasets are redundant and incomplete
(Some PPIs represented > 500 times, many missing)
- Existing train-test splits suffer from data leakage
(>53% of test sets have near-duplicates in train sets)

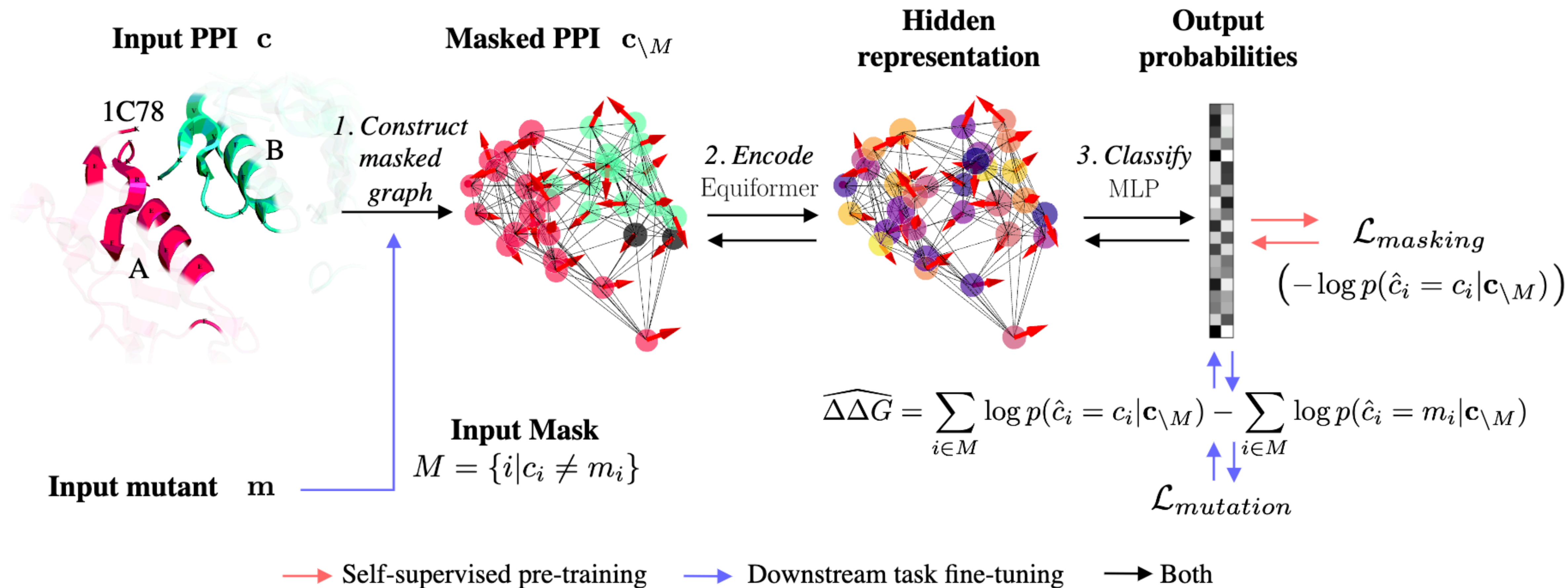
- **PPIRef**: non-redundant complete set of PPIs from PDB



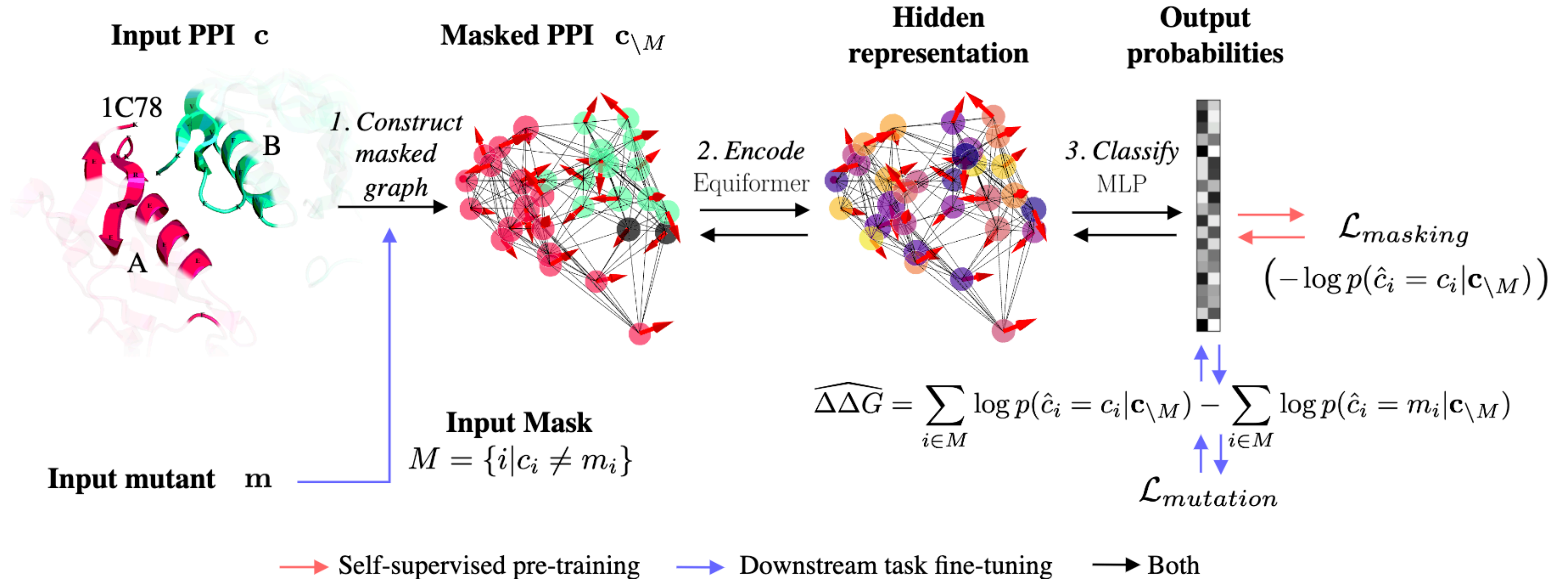
PPIformer



PPIformer



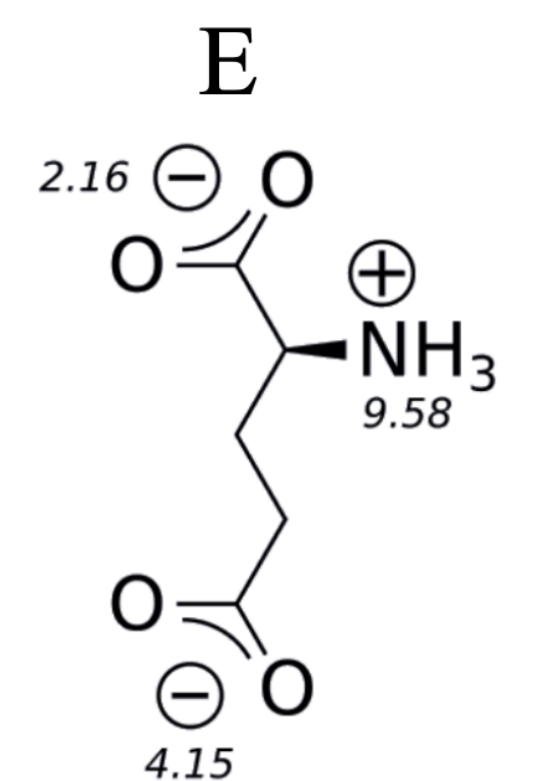
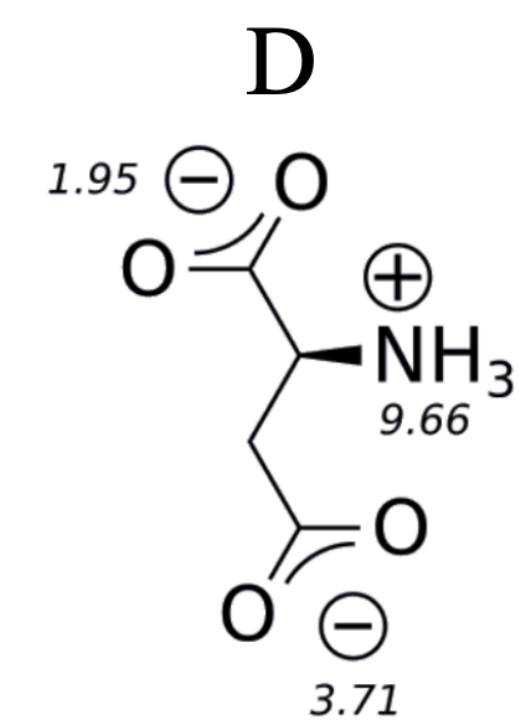
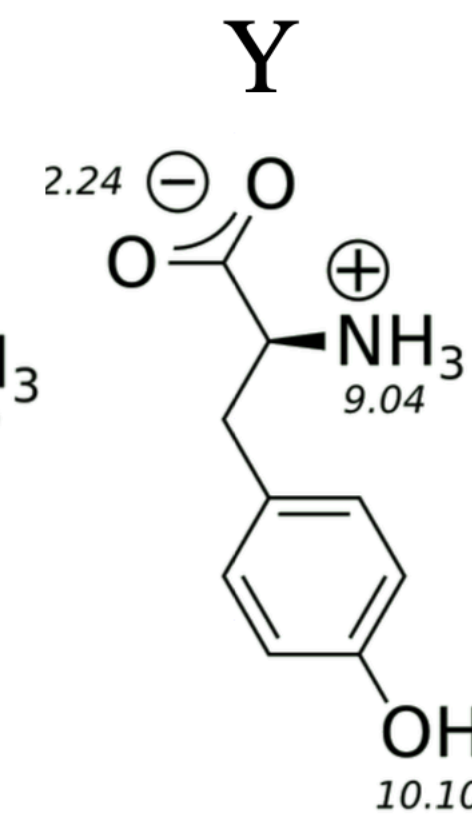
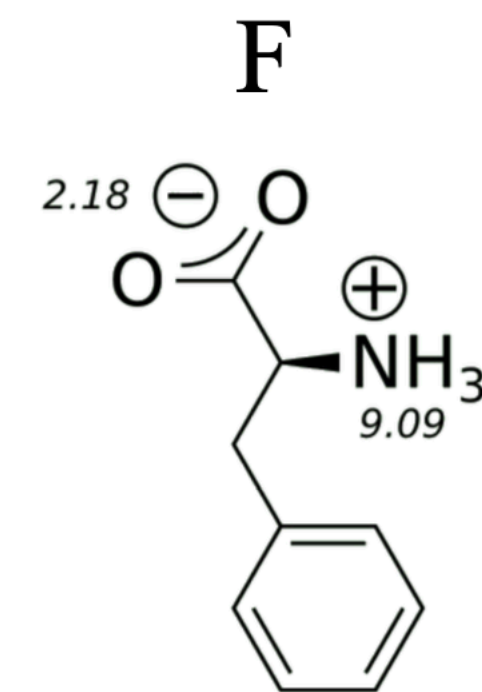
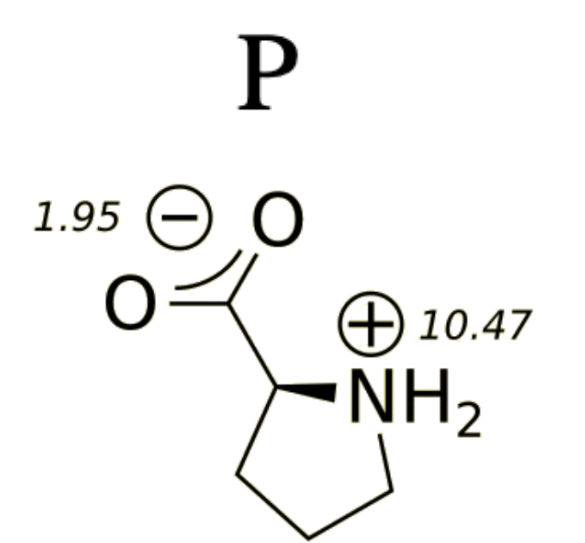
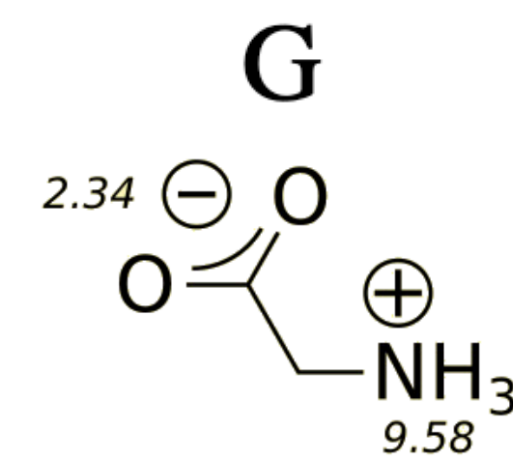
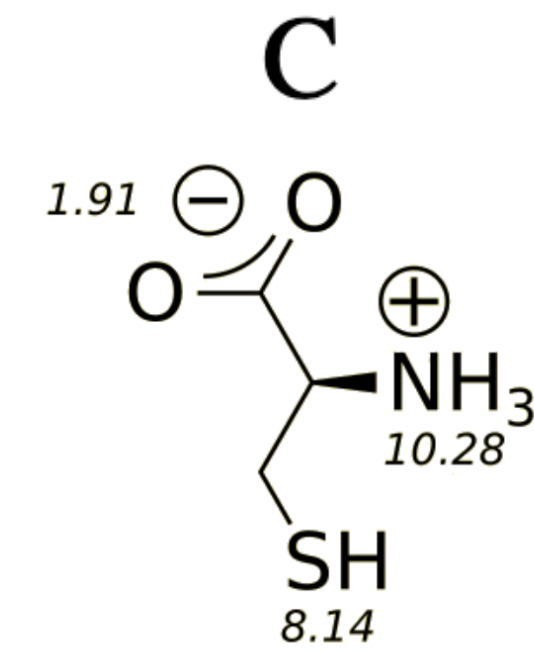
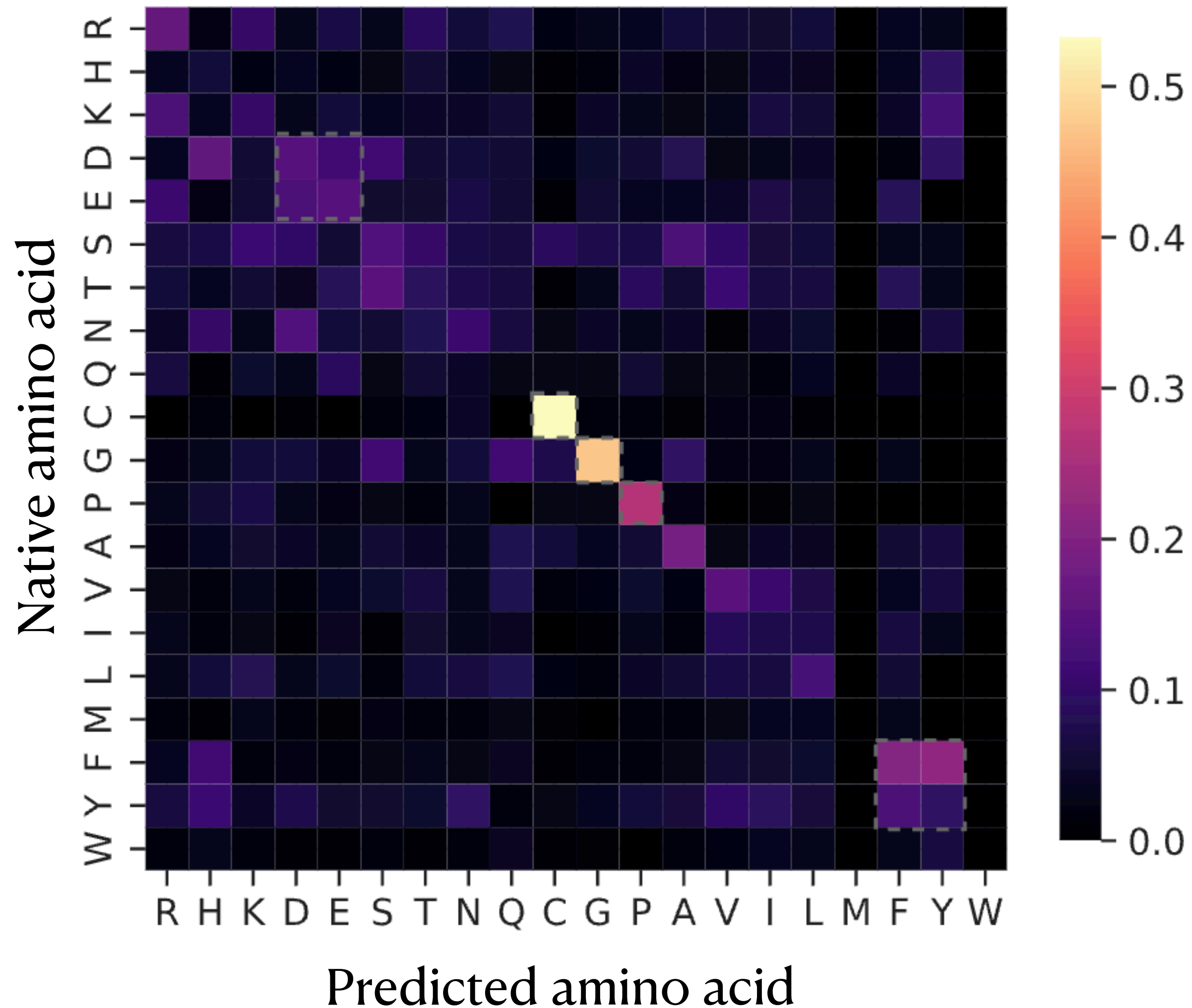
PPIformer



- Leverages **big data** (millions of masked examples from PPIRef during pre-training)
- Very fast, requires a **single forward pass on the native 3D structure**
- Fine-tuned and evaluated on **non-leaking $\Delta\Delta G$ data** using **practically-important metrics**

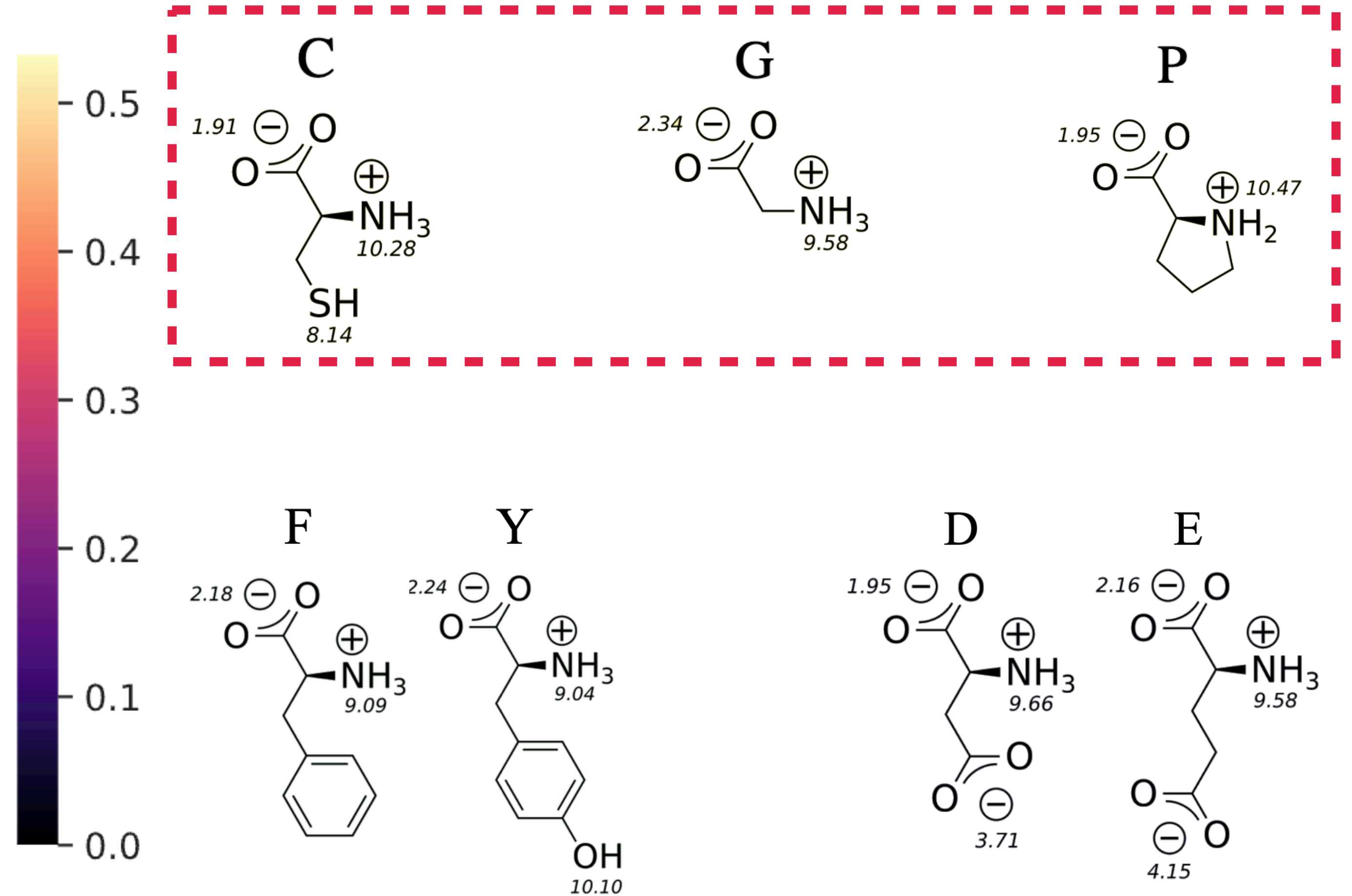
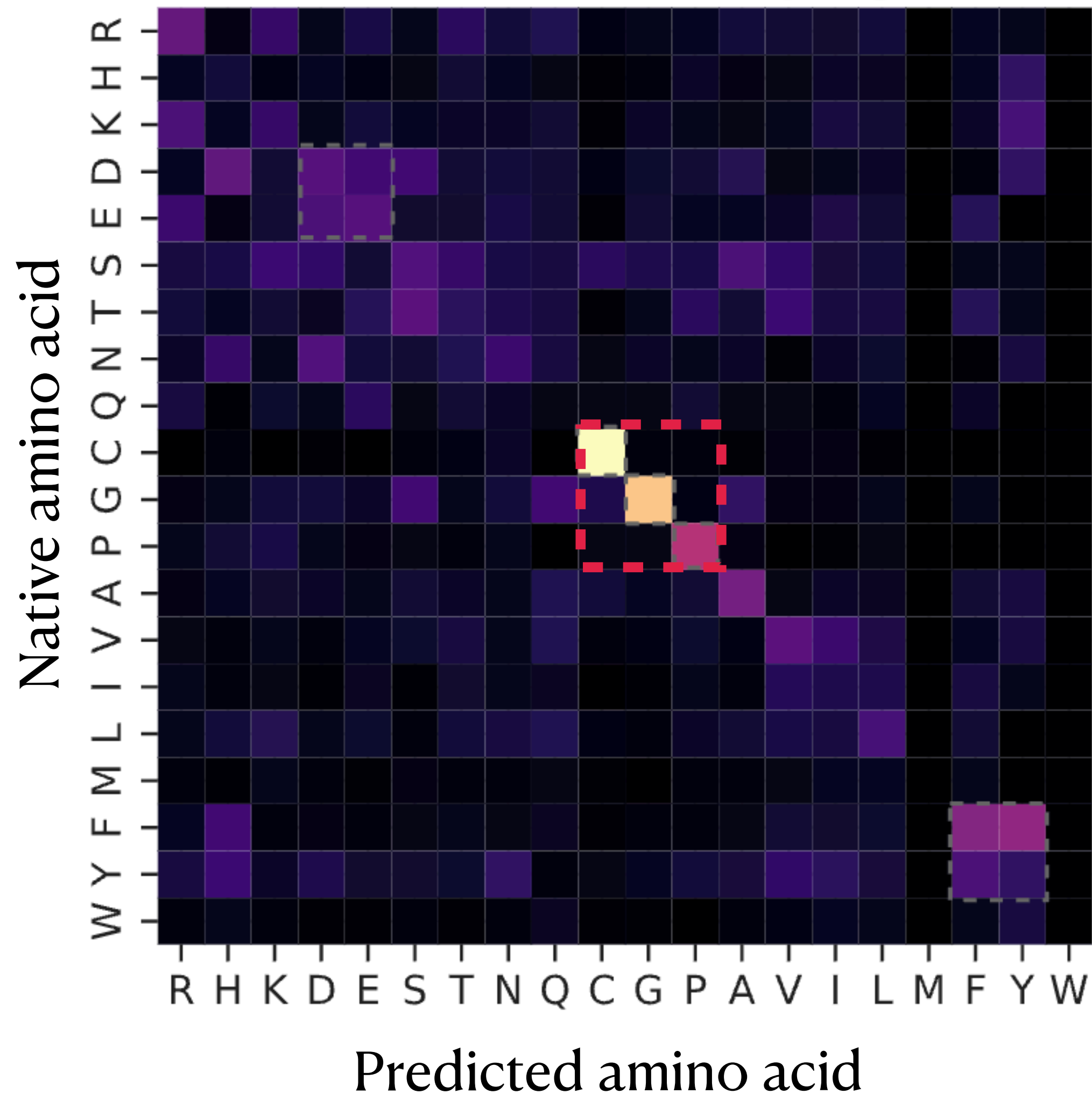
PPIformer captures biochemical principles

Confusion matrix



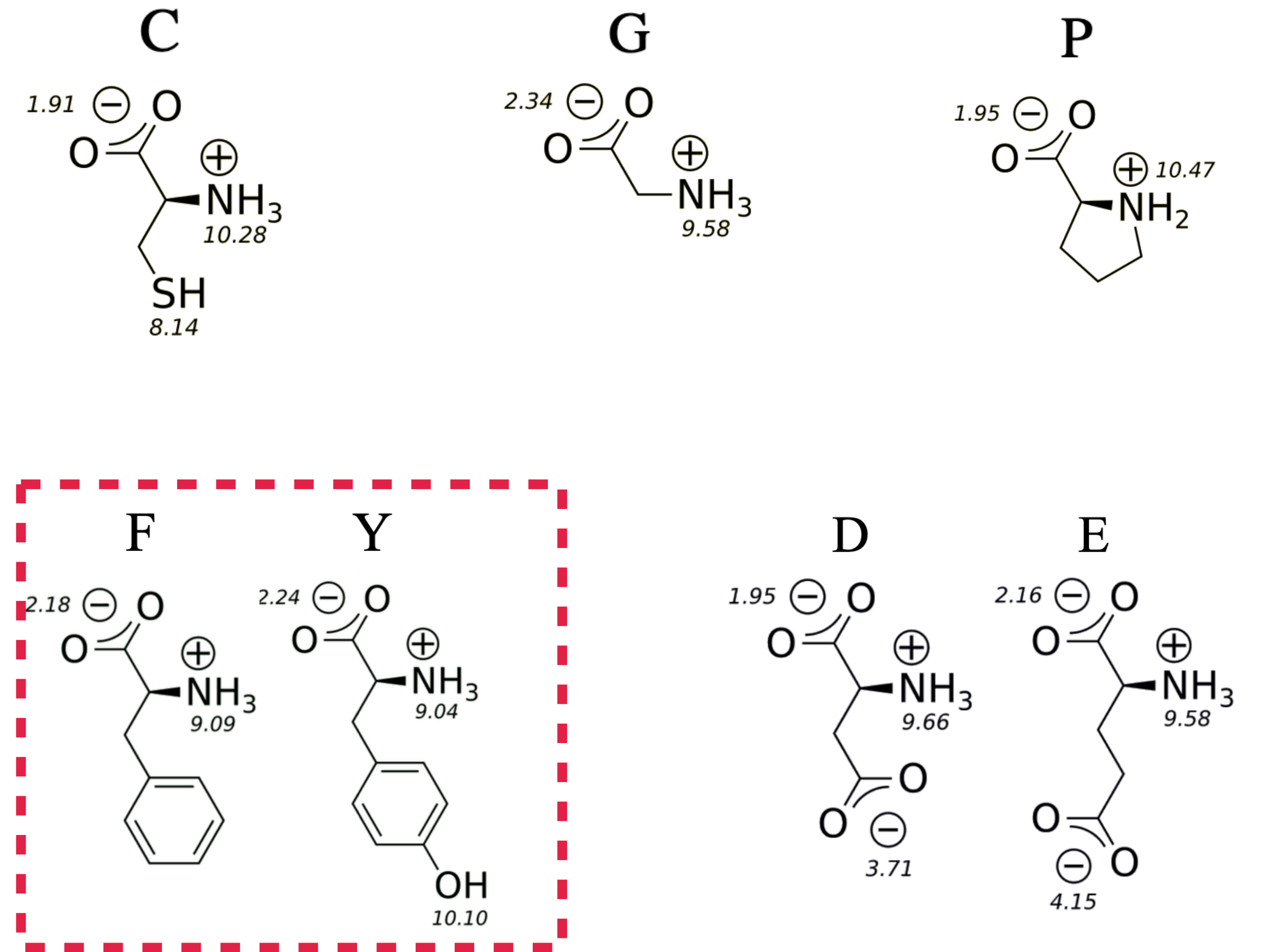
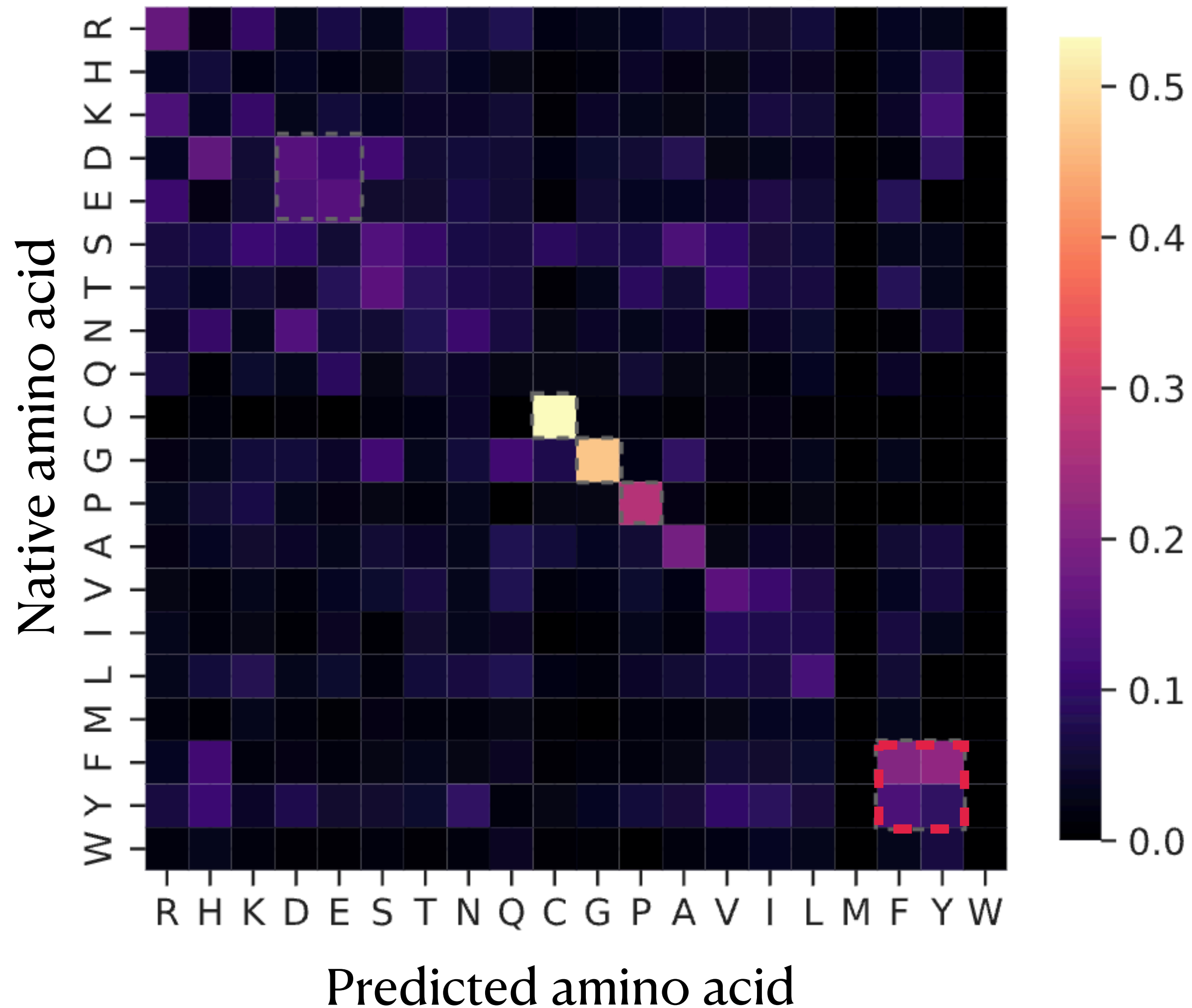
PPIformer captures biochemical principles

Confusion matrix



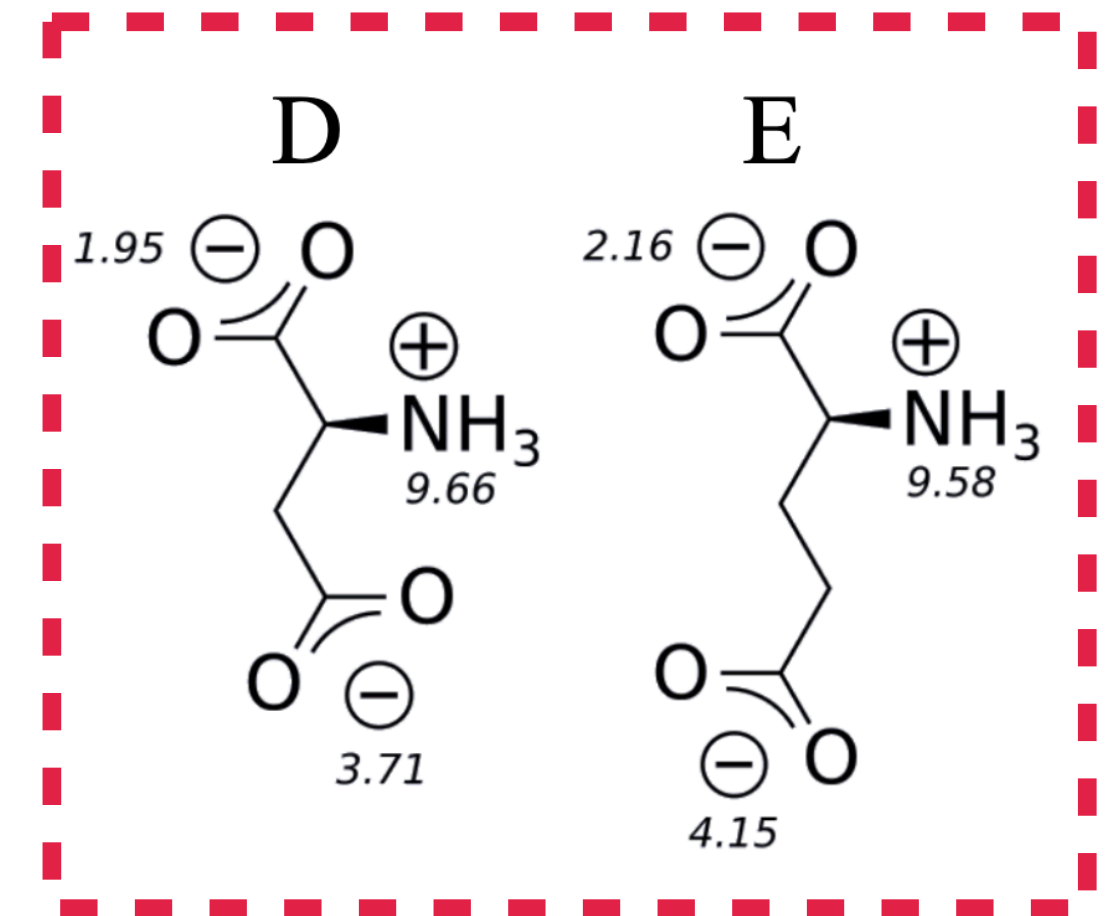
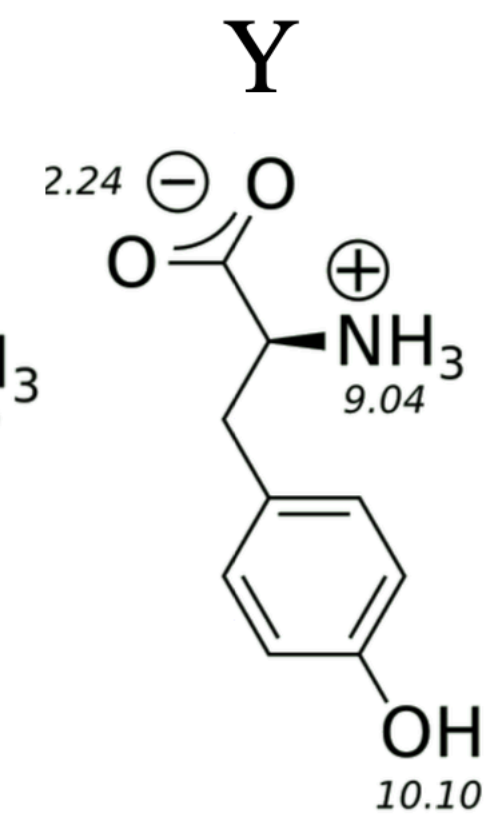
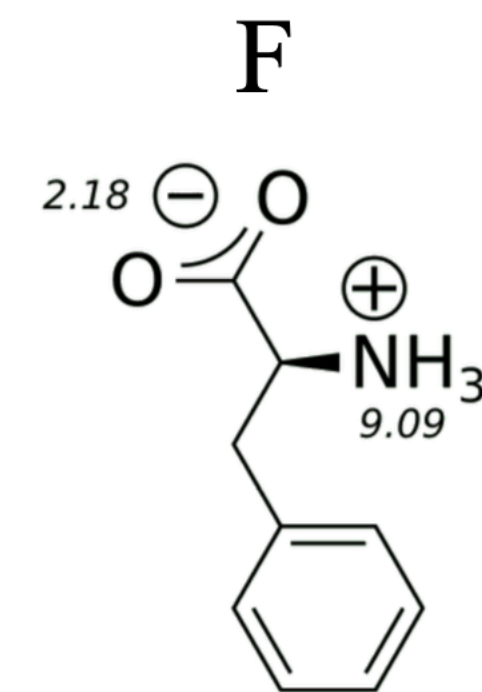
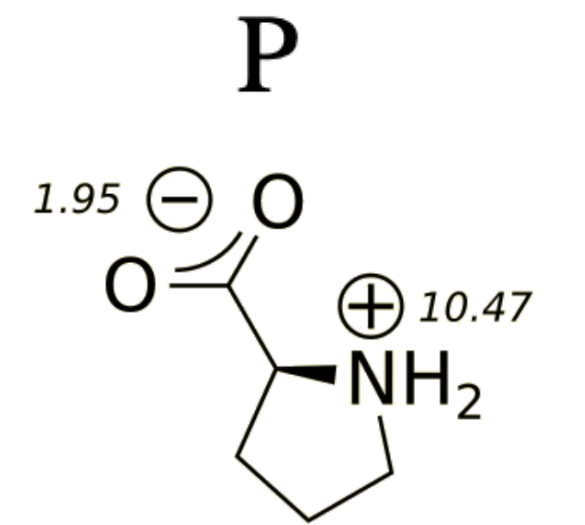
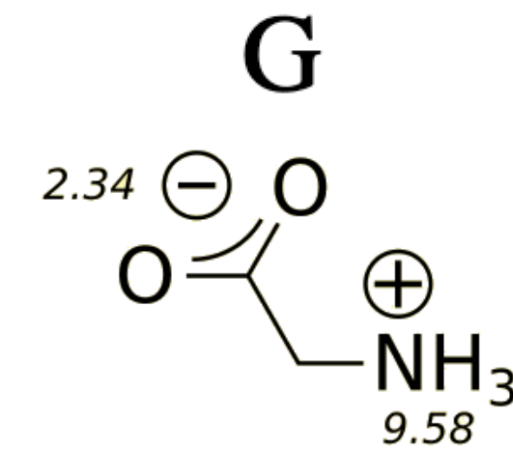
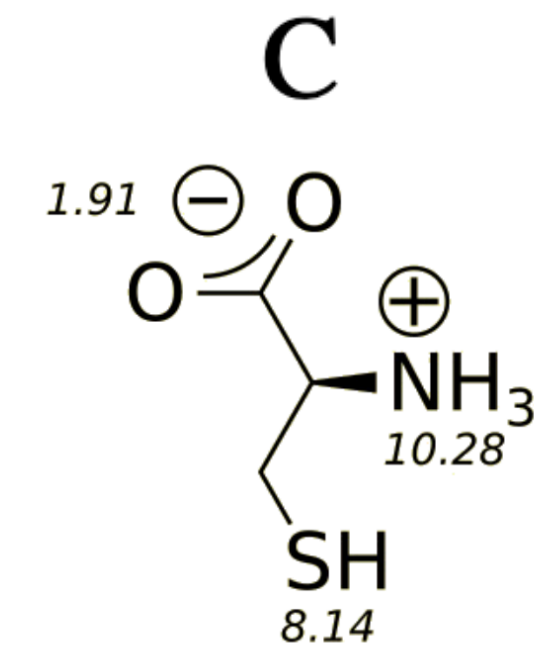
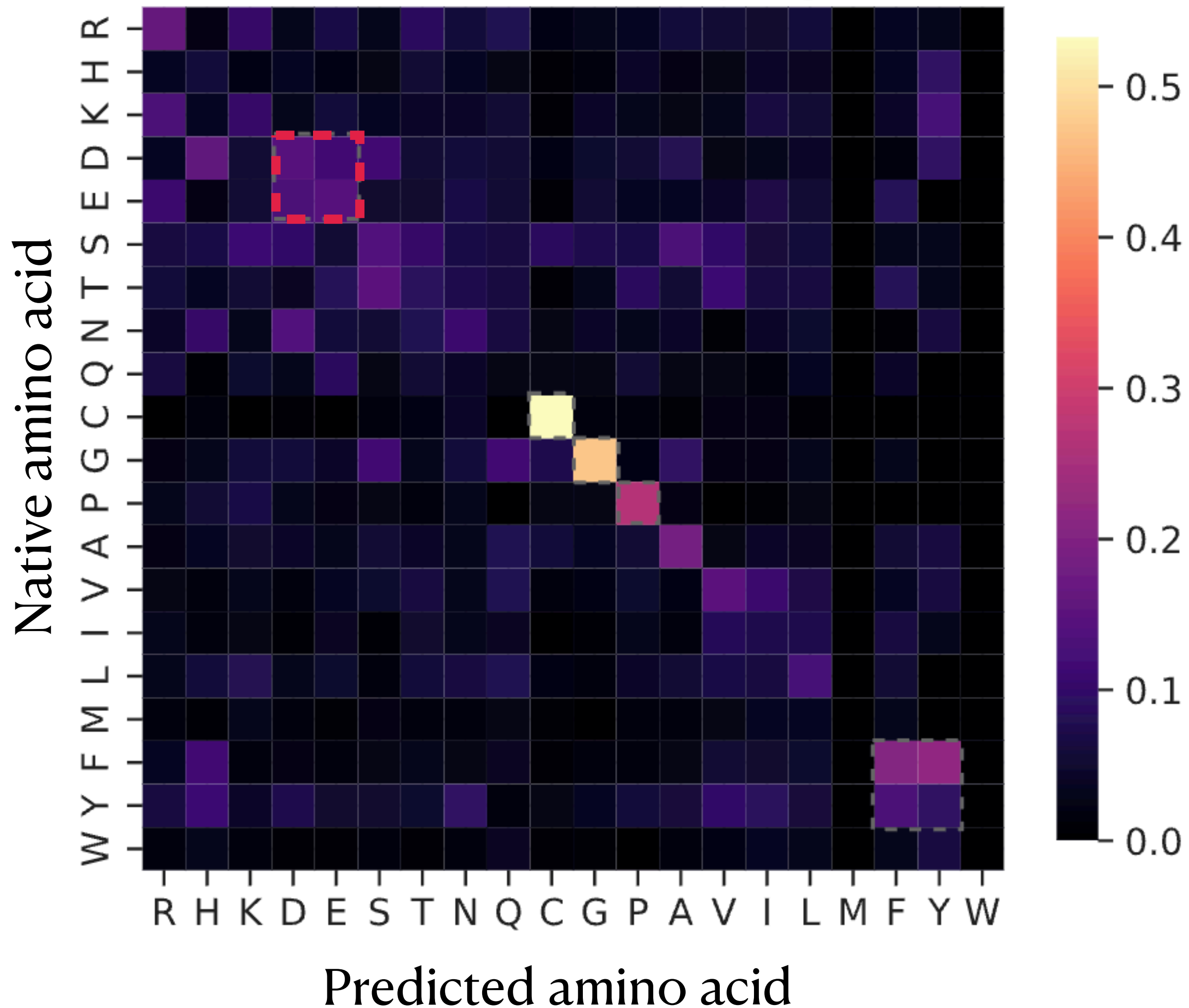
PPIformer captures biochemical principles

Confusion matrix



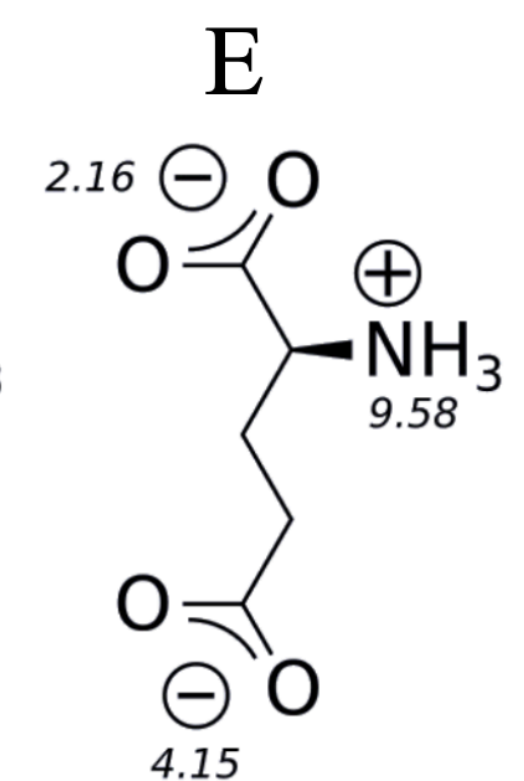
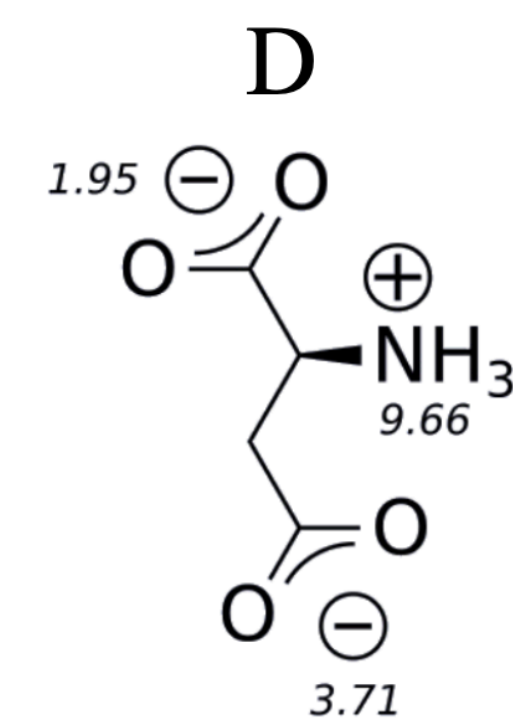
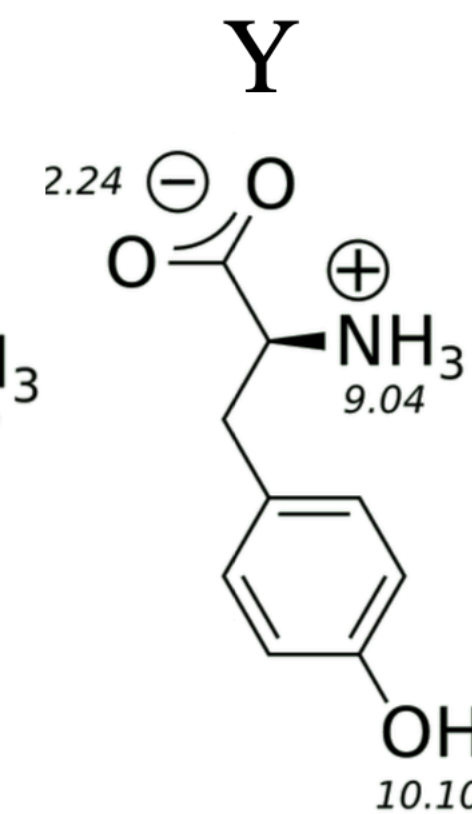
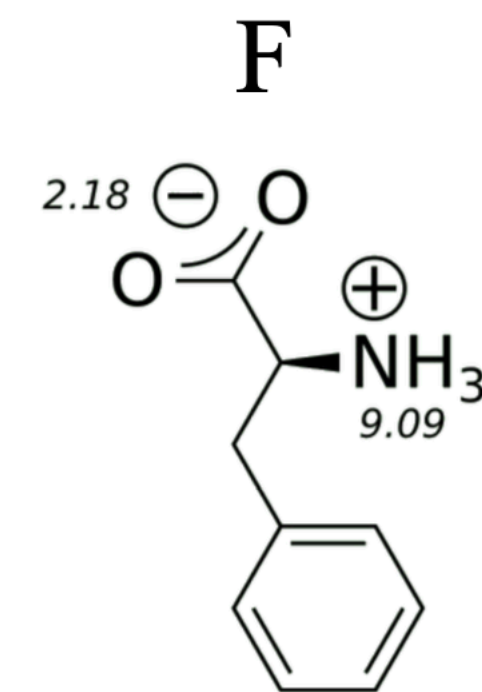
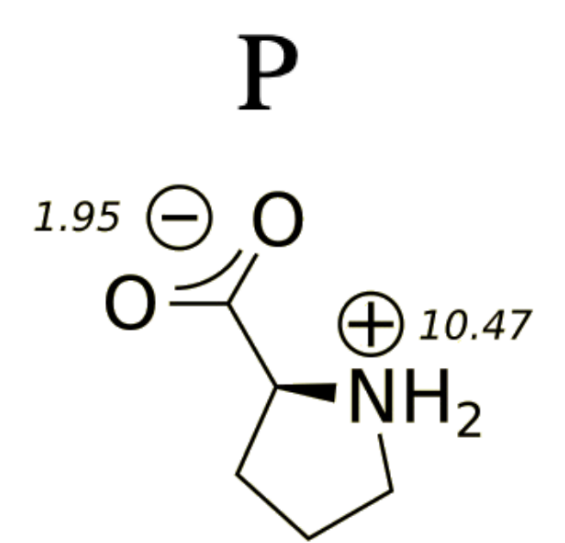
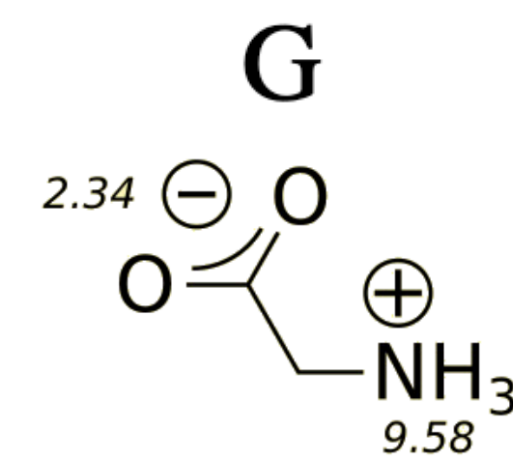
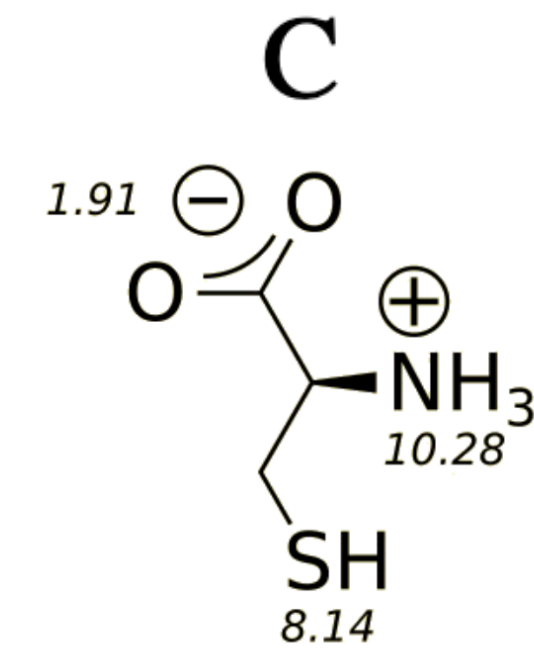
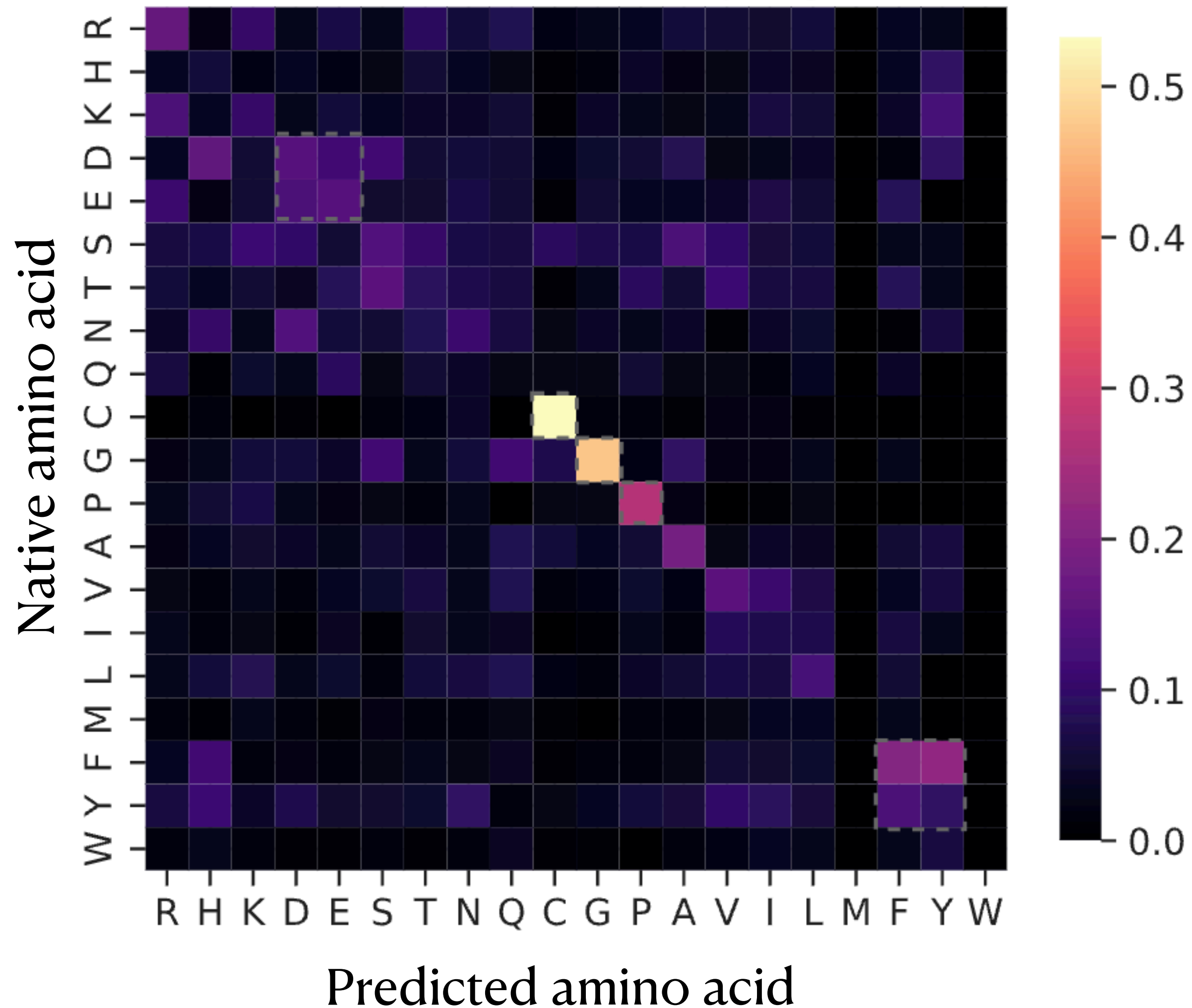
PPIformer captures biochemical principles

Confusion matrix



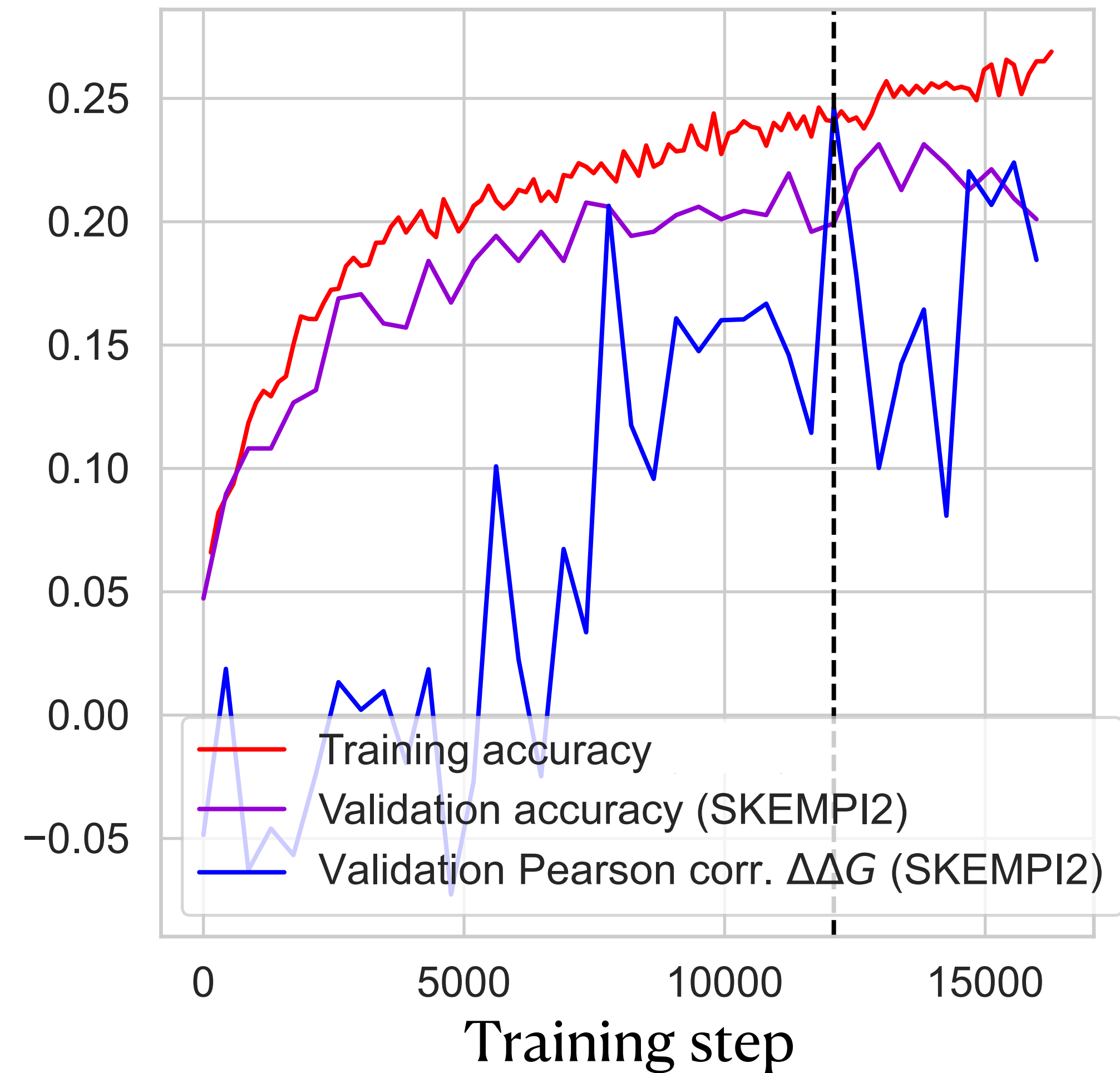
PPIformer captures biochemical principles

Confusion matrix



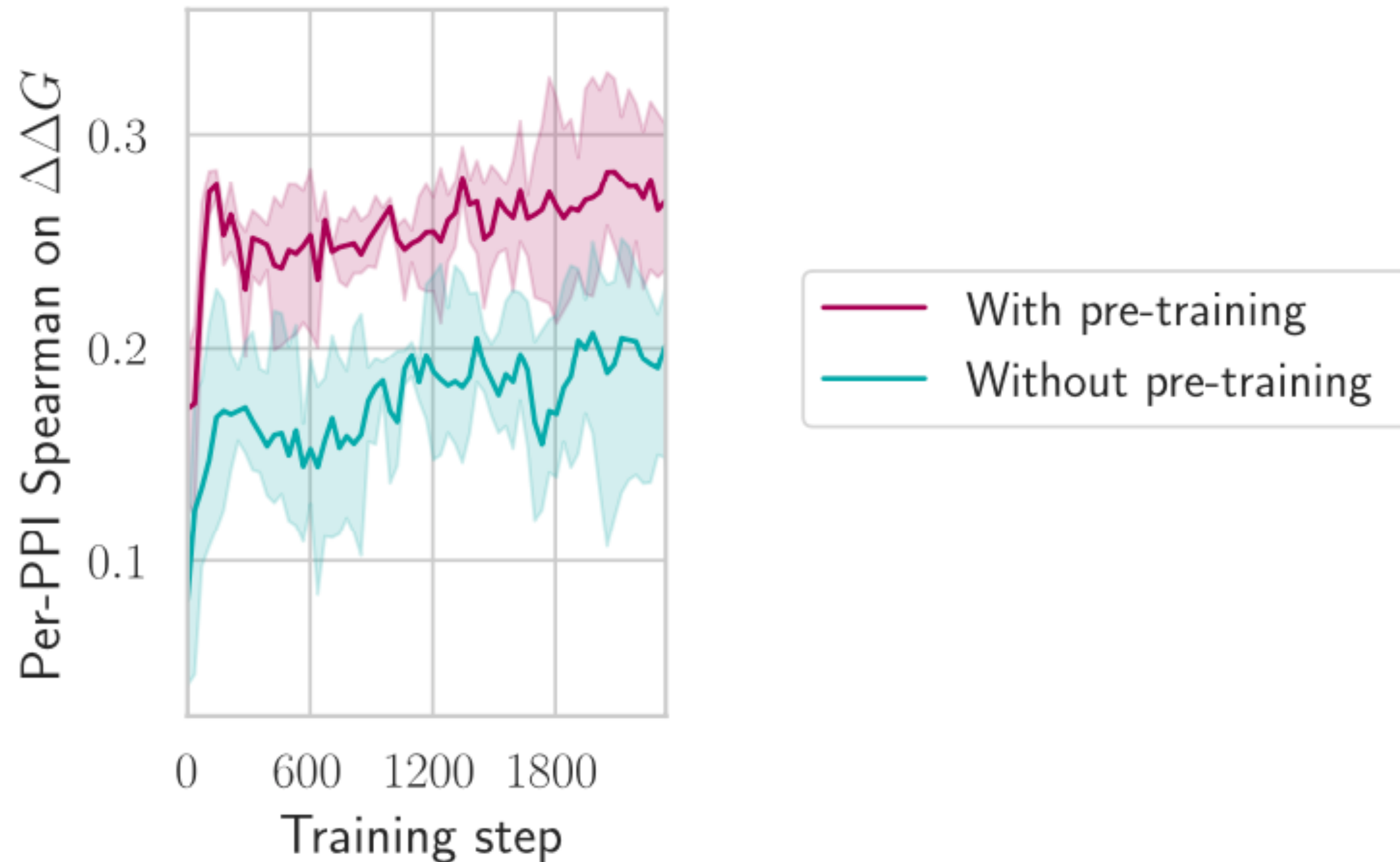
Emergence of mutation scoring capabilities

Self-supervised pre-training



Pre-training is crucial for fine-tuning

Supervised $\Delta\Delta G$ fine-tuning

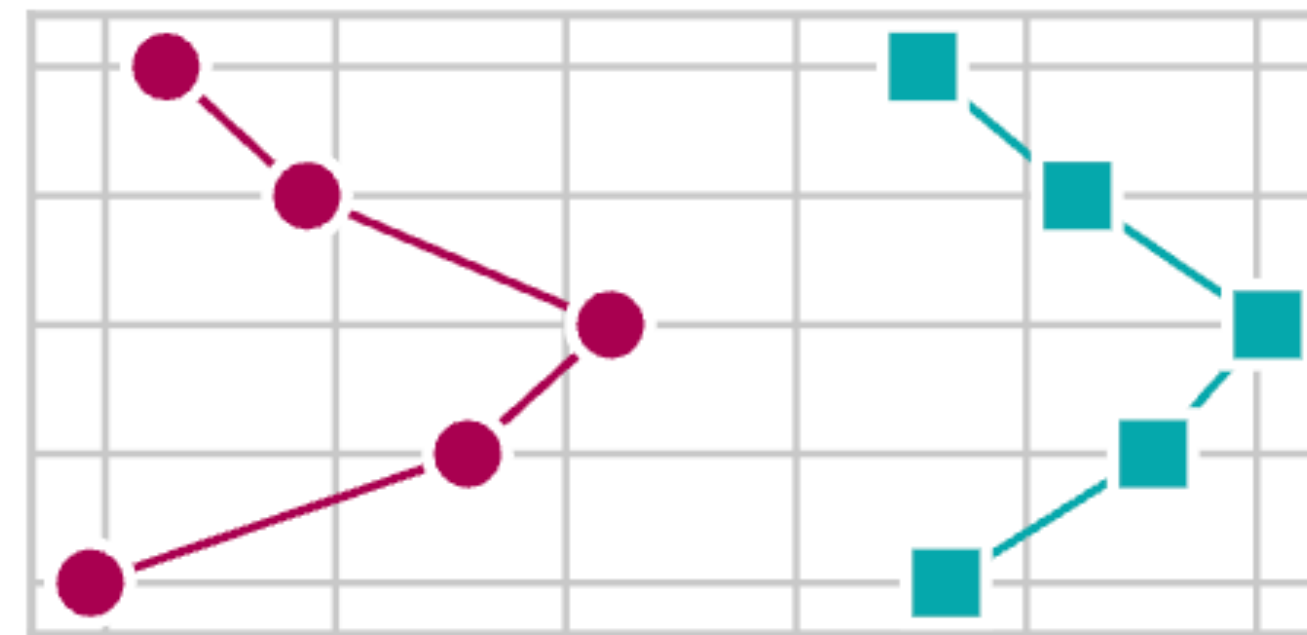


Pre-training from PPIRef is crucial

Pre-training

- Precision on negative $\Delta\Delta G$ (zero-shot)
- Per-PPI Spearman on $\Delta\Delta G$ (zero-shot)

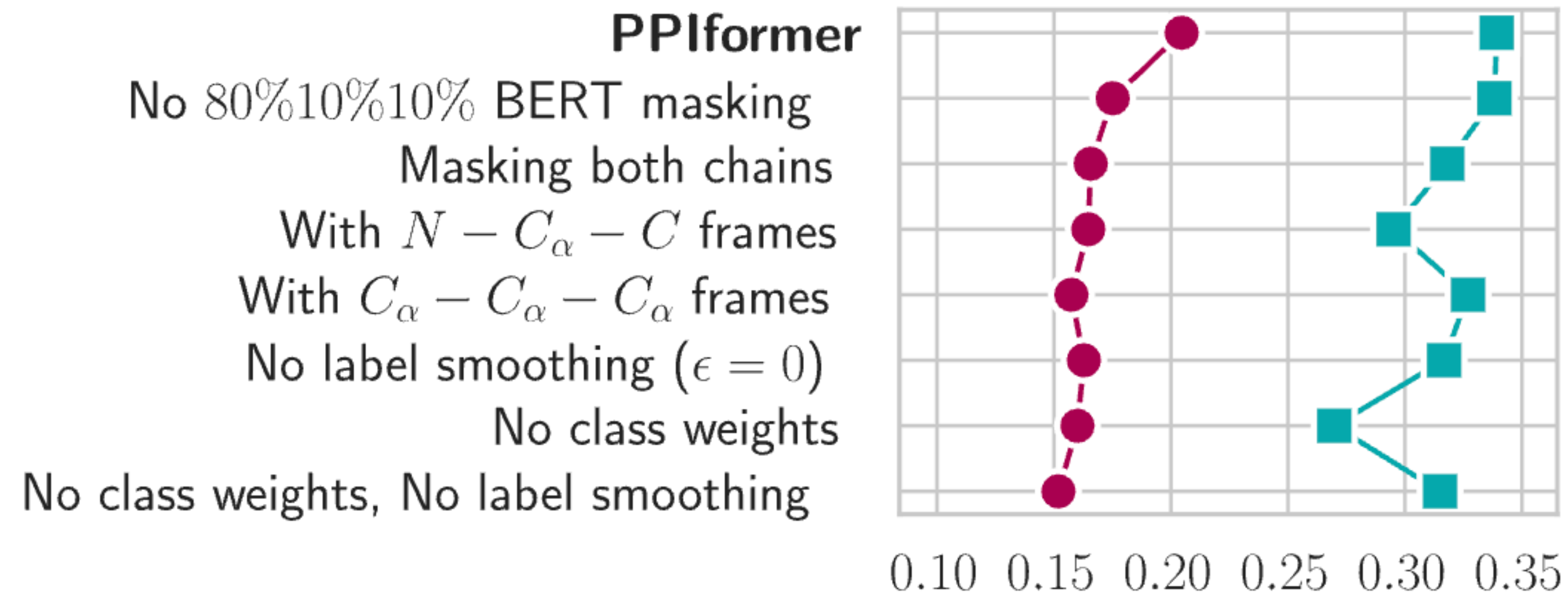
DIPS/DIPS-Plus (deduplicated; 8K)
DIPS/DIPS-Plus (40K)
PPIRef50K (filtered, deduplicated)
PPIRef300K (filtered)
PPIRef800K (raw)



Other key ingredients

Pre-training

- Precision on negative $\Delta\Delta G$ (zero-shot)
- Per-PPI Spearman on $\Delta\Delta G$ (zero-shot)



Comparison with the state of the art: 5 independent PPIs from SKEMPI

Category	Method	Spearman \uparrow	Pearson \uparrow	Precision \uparrow	Recall \uparrow	ROC AUC \uparrow	MAE \downarrow	RMSE \downarrow
Force field simulations	FLEX DDG* ²²	0.54	0.57	0.63	0.62	0.84	1.60	2.00
Machine learning	MSA TRANSFORMER ²³	<u>0.37</u>	<u>0.45</u>	0.51	0.38	<u>0.76</u>	5.99	6.77
	ESM-IF ²⁴	0.32	0.31	0.36	0.28	0.69	1.84	2.11
	RDE-NET. ²¹	0.24	0.30	<u>0.54</u>	0.65	0.67	<u>1.70</u>	<u>2.02</u>
	PPIFORMER (OURS)	0.42	0.46	0.58	<u>0.61</u>	0.77	1.64	1.94

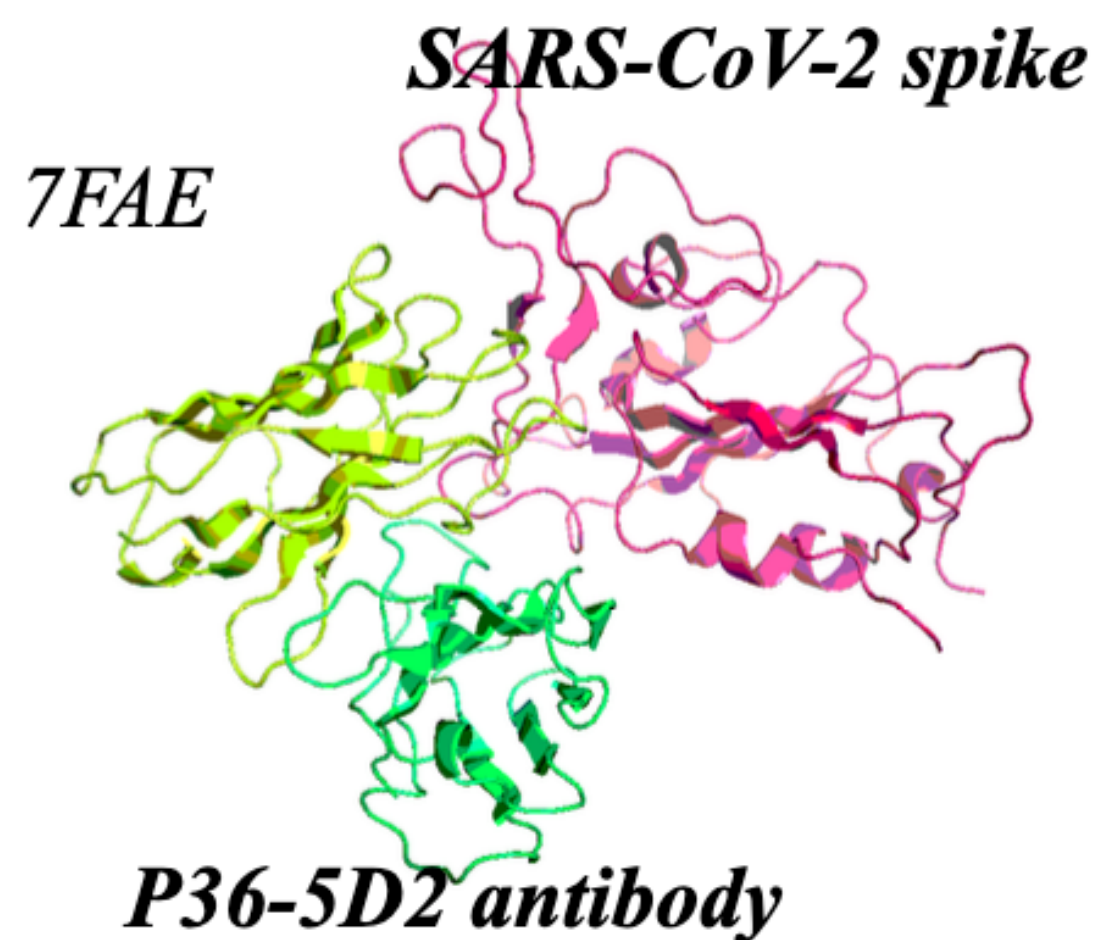
* A single prediction requires ~1 CPU hour (5 orders of magnitude slower than other methods).



Comparison with the state of the art: COVID

Method	known favorable mutations					precision		
	TH31W ↓	AH53F ↓	NH57L ↓	RH103M ↓	LH104F ↓	P@1 ↑	P@5% ↑	P@10% ↑
MSA TRANSFORMER	56.88	42.11	63.56	49.19	18.83	<u>0.00</u>	0.00	0.00
ROSETTA	10.73	76.72	93.93	13.56	6.88	<u>0.00</u>	0.00	2.04
FOLDX	5.67**	68.22	2.63	12.35	29.96	<u>0.00</u>	<u>4.00</u>	<u>4.08</u>
DDGPRED	2.02	14.17	24.49	4.05	6.48	<u>0.00</u>	8.00	6.12
END-TO-END	11.34	16.60	8.30	52.43	80.36	<u>0.00</u>	0.00	2.04
MIF-NET.	27.94	66.19	8.50	17.21	36.23	<u>0.00</u>	0.00	2.04
ESM-IF	49.39	17.61	17.00	51.42	48.58	<u>0.00</u>	0.00	0.00
RDE-NET.	1.62	2.02	20.65	61.54	5.47	<u>0.00</u>	8.00	6.12
PPIFORMER (OURS)	18.02	0.20	7.69	21.46	10.93	100	<u>4.00</u>	<u>4.08</u>

** Mutations that are in top 10% of predictions are in bold.

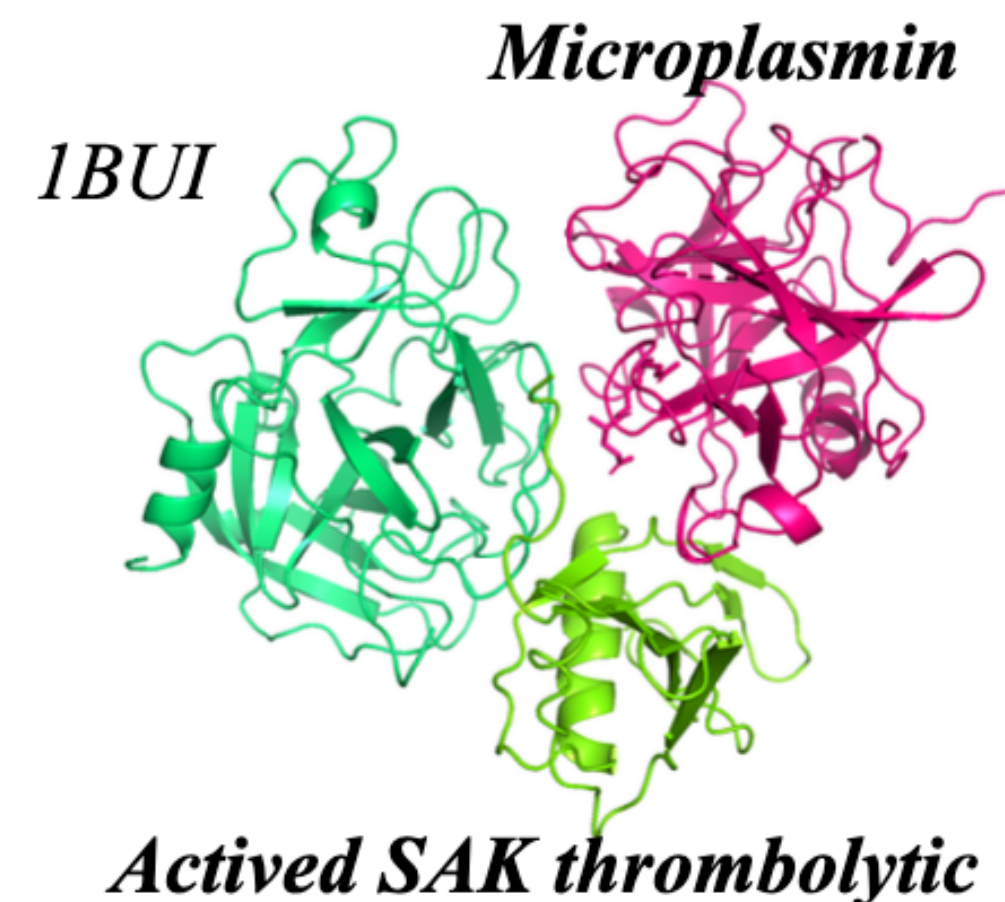


Comparison with the state of the art: stroke

known favorable mutations

precision

Method	Mutations with $\geq 2\times$ activity enhancement						Activity enhancement		
	KC130A ↓	KC130T ↓	KC130T ↓	KC135A ↓	KC135R ↓	KC135R ↓	P@1 ↑	P@5% ↑	P@10% ↑
MSA TRANSFORMER	52.50	32.50	55.00	40.0	70.00	78.75	100	<u>50.00</u>	37.50
ESM-IF	45.00	33.75	46.25	25.0	42.50	58.75	<u>0.00</u>	0.00	25.00
RDE-NET.	51.25	33.75	22.50	15.00	27.50	5.00	<u>0.00</u>	<u>50.00</u>	<u>62.50</u>
PPIFORMER (OURS)	66.25	15.00	2.50	52.50	33.75	1.25	100	75.00	87.50



LEARNING TO DESIGN PROTEIN–PROTEIN INTERACTIONS WITH ENHANCED GENERALIZATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Discovering mutations enhancing protein–protein interactions (PPIs) is critical for advancing biomedical research and developing improved therapeutics. While machine learning approaches have substantially advanced the field, they often struggle to generalize beyond training data in practical scenarios. The contributions of this work are three-fold. First, we construct PPIRef, the largest and non-redundant dataset of 3D protein–protein interactions, enabling effective large-scale learning. Second, we leverage PPIRef to pre-train PPIformer, a new SE(3)-equivariant model generalizing across diverse protein-binder variants. We fine-tune PPIformer to predict effects of mutations on protein–protein interactions via a thermodynamically motivated adjustment of the pre-training loss function. Finally, we demonstrate the enhanced generalization of our new PPIformer approach by outperforming other state-of-the-art methods on the new non-leaking splits of the standard labeled PPI mutational data and independent case studies optimizing a human antibody against SARS-CoV-2 and increasing staphylokinase thrombolytic activity.

PPIformer Github



anton.bushuiev@cvut.cz

References

1. Beckstette et al., Fast index based algorithms and software for matching position specific scoring matrices, BMC bioinformatics, 2006
2. Sumbalova et al., Hotspot wizard 3.0: web server for automated design of mutations and smart libraries based on sequence input information, Nucleic acids research, 2018
3. Meier et al., Language models enable zero-shot prediction of the effects of mutations on protein function, Advances in Neural Information Processing Systems, 2021
4. Shroff et al., Discovery of novel gain-of-function mutations guided by structure-based deep learning, ACS synthetic biology, 2020
5. Dauparas et al., Robust deep learning-based protein sequence design using ProteinMPNN, Science, 2022
6. Biota et al., Discovery studio modeling environment, Dassault Systemes, 2017
7. Dehouck et al., Beatmusic: prediction of changes in protein-protein binding affinity on mutations, Nucleic acids research, 2013
8. Delgado et al., Foldx 5.0: working with rna, small molecules and a new graphical interface, Bioinformatics, 2019
9. Geng et al., isee: Interface structure, evolution, and energy-based machine learning predictor of binding affinity changes upon mutations, Proteins: Structure, Function and Bioinformatics, 2019
10. Rodrigues et al., mmcsmpi: predicting the effects of multiple point mutations on protein-protein interactions, Nucleic Acids Research, 2021
11. Xiong et al., Bindprofx: assessing mutation-induced binding affinity change by protein interface profiles with pseudo-counts, Journal of molecular biology, 2017
12. Zhou et al., Mutation effect estimation on protein-protein interactions using deep contextualized representation learning 2020, NAR genomics and bioinformatics, 2020
13. Wang et al., A topology-based network tree for the prediction of protein-protein binding affinity changes following mutation, Nature Machine Intelligence, 2020
14. Jiang et al., Dgcddg: Deep graph convolution for predicting protein-protein binding affinity changes upon mutations, IEEE/ACM Transactions on Computational Biology and Bioinformatics, 2023
15. Pahari et al., Saambe-3d: predicting effect of mutations on protein-protein interactions, International journal of molecular sciences, 2020
16. Liu et al., Deep geometric representations for modeling effects of mutations on protein-protein binding affinity, PLoS computational biology, 2021
17. Li et al., Saambe-seq: a sequence-based method for predicting mutation effect on protein-protein binding affinity, Bioinformatics, 2021
18. Gao and Skolnick, ialign: a method for the structural comparison of protein-protein interfaces, Bioinformatics, 2010
19. Townshend et al., End-to-end learning on 3d protein structure for interface prediction, Advances in Neural Information Processing Systems, 2019
20. Shan, Sisi, et al. "Deep learning guided optimization of human antibody against SARS-CoV-2 variants with broad neutralization." Proceedings of the National Academy of Sciences 119.11 (2022): e2122954119.
21. Luo, Shitong, et al. "Rotamer Density Estimator is an Unsupervised Learner of the Effect of Mutations on Protein-Protein Interaction." bioRxiv (2023): 2023-02.
22. Barlow, Kyle A., et al. "Flex ddG: Rosetta ensemble-based estimation of changes in protein-protein binding affinity upon mutation." The Journal of Physical Chemistry B 122.21 (2018): 5389-5399.
23. Rao, Roshan M., et al. "MSA transformer." International Conference on Machine Learning. PMLR, 2021.
24. Hsu, Chloe, et al. "Learning inverse folding from millions of predicted structures." International Conference on Machine Learning. PMLR, 2022.