

PLMSearch: Protein language model
powers accurate and fast sequence
search for remote homology

Stephen Owesney



Overview

- Introduction to Homologous Protein Search
- Background Knowledge / Related Works
- Introduction to PLMSearch
- Methodology
- Results / Experiments
- Future Work
- Questions / Discussion



Homologous Protein Search

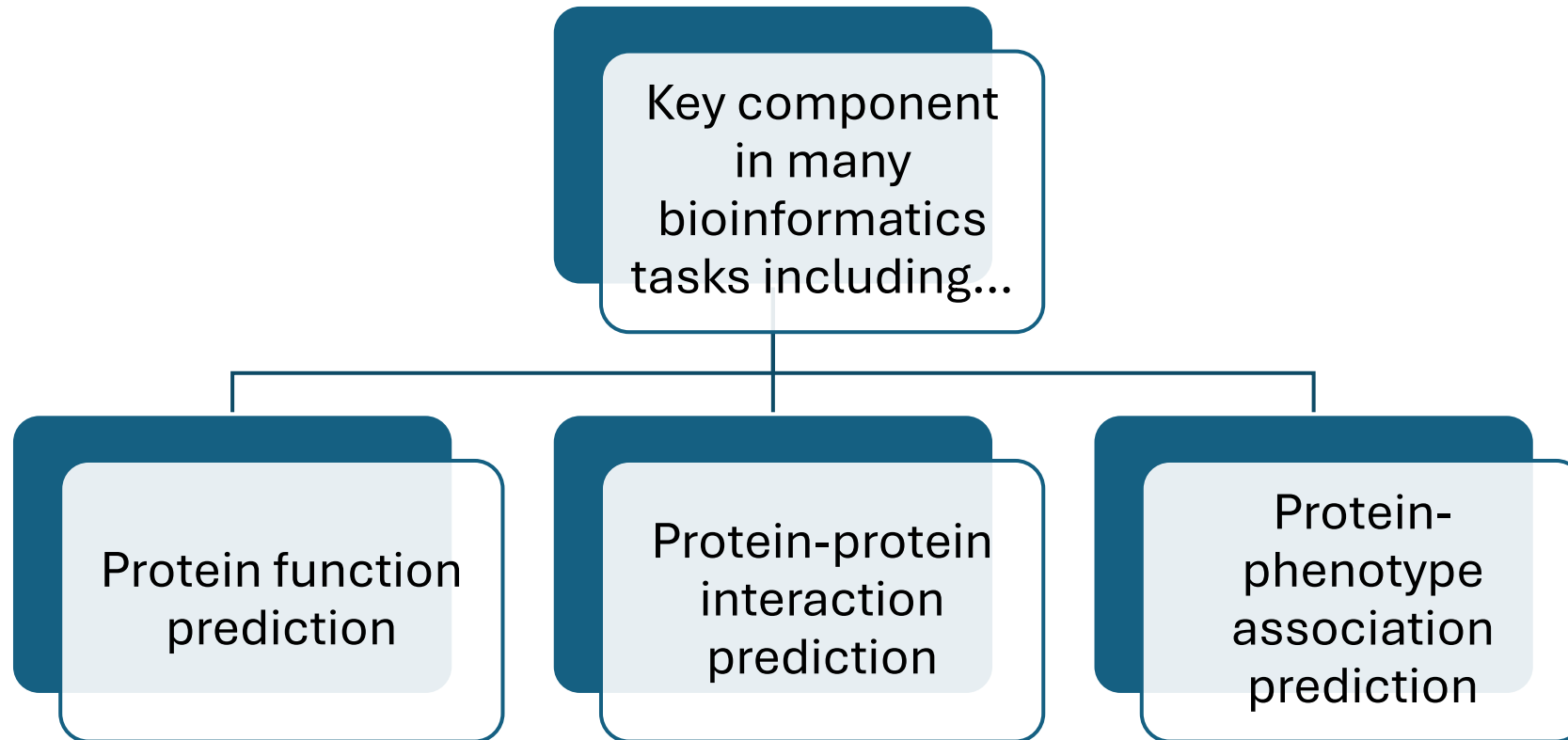
The goal of homologous protein search is to associate a query protein with homologous proteins



Homologous Protein Search

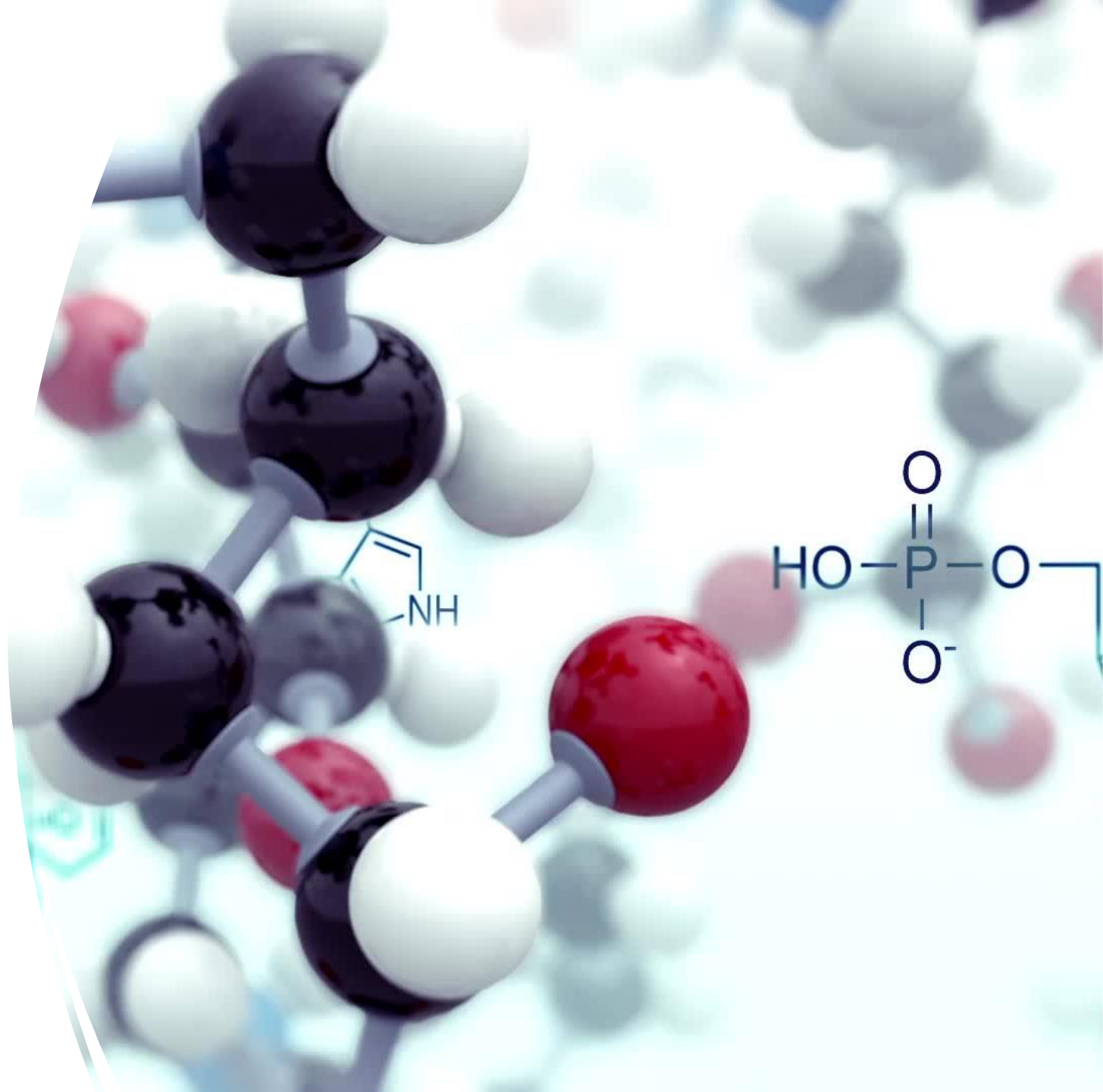
Two homologous proteins share evolutionary origins

Homologous Protein Search



Key Challenges for Homologous Protein Search

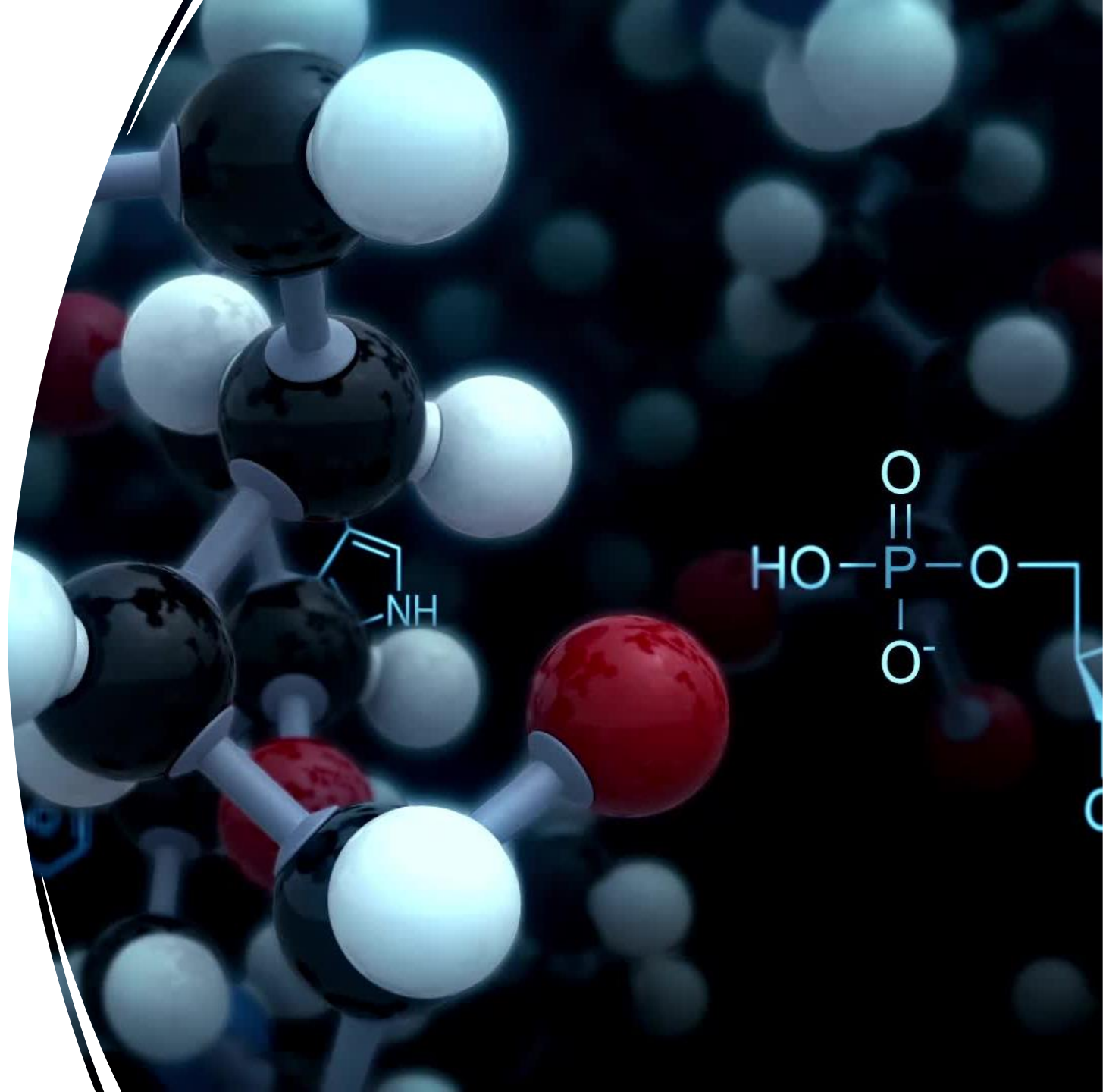
- 1) Computational Complexity
- 2) Noisy Datasets
- 3) Low Sequence Similarity



Homologous Protein Search

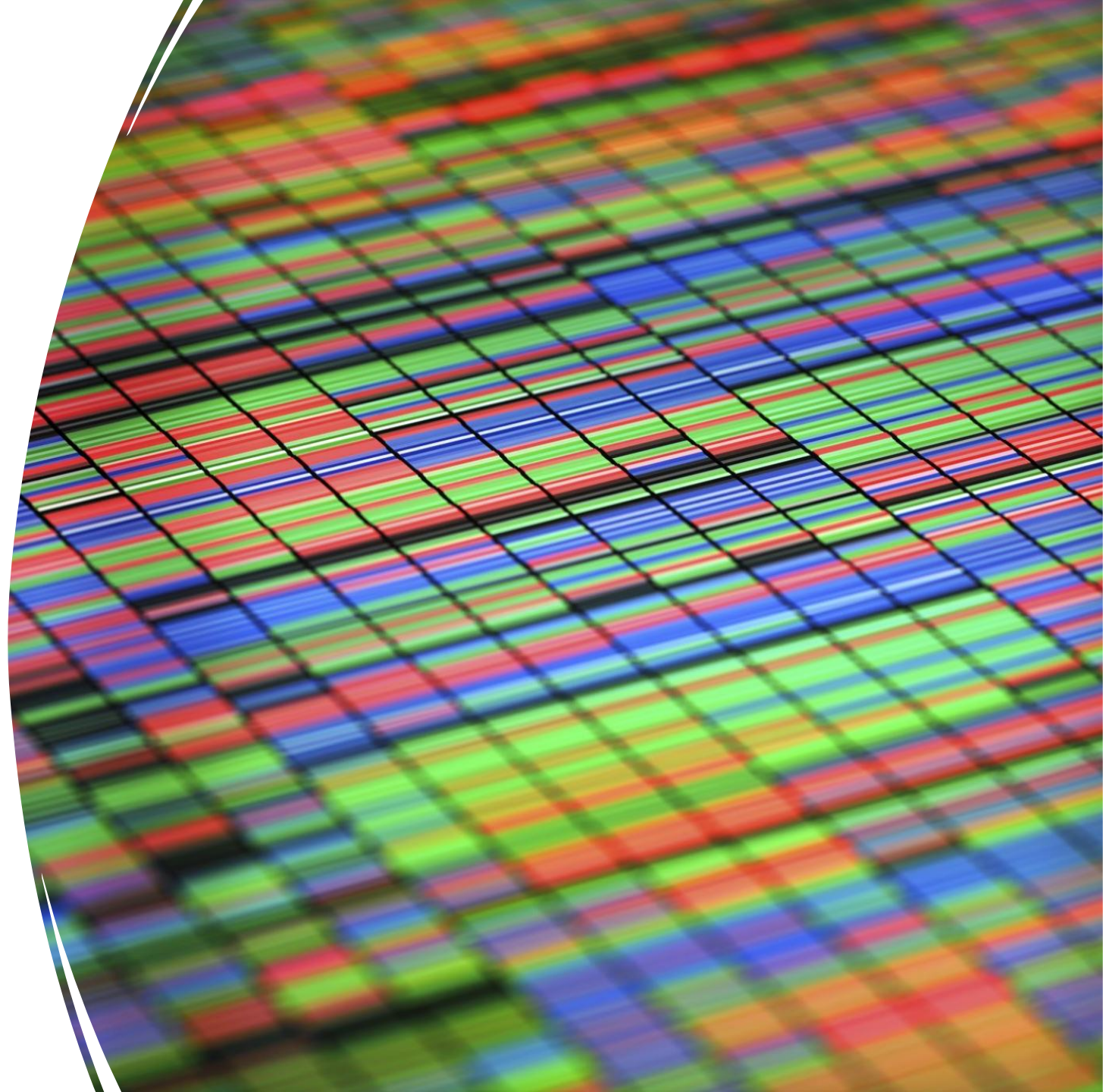
Homologous protein search can be roughly divided into two approaches or components

- 1) sequence search
- 2) structure search



Sequence Search Approach

Goal: Compare amino acid sequence of query to targets in database



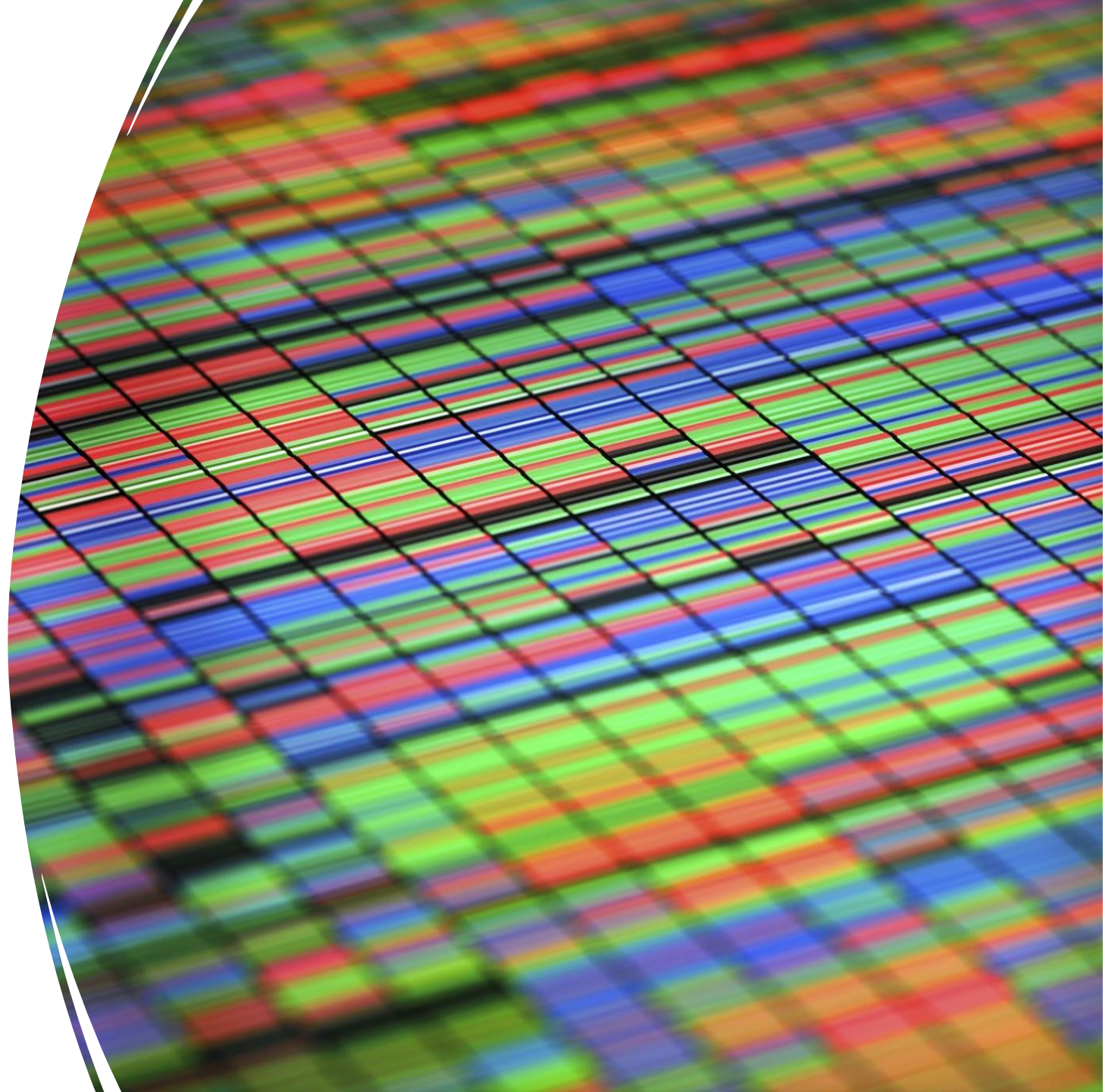
Sequence Search Approach

Strengths

- 1) Low cost
- 2) Data abundance
- 3) Well established

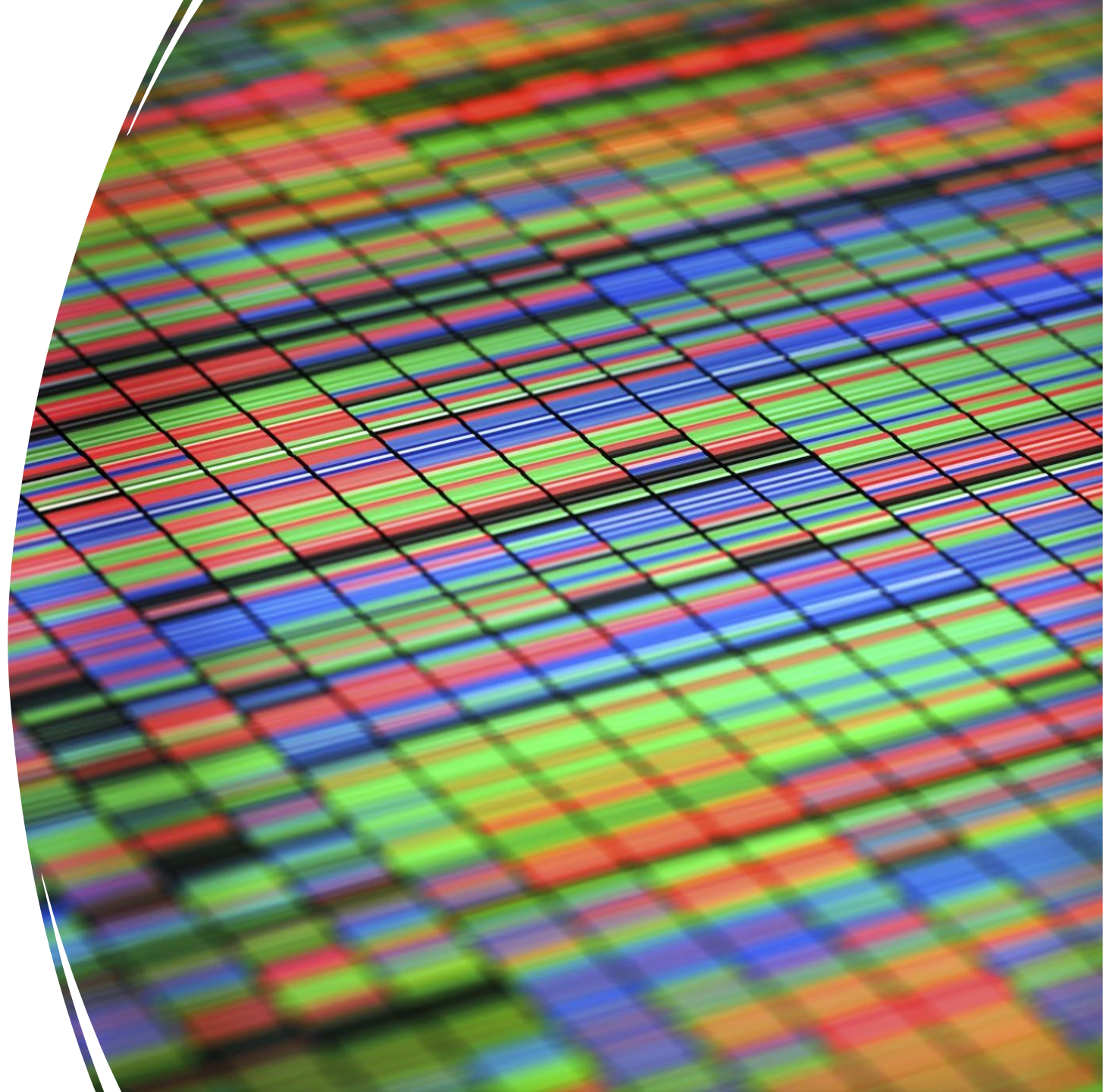
Challenges

- 1) Sequence divergence
- 2) Poor remote performance



Sequence Search Approach

- MMseqs2
- BLASTp
- Diamond





Structure Search Approach

- **Goal:** Perform the comparison of the query and target proteins leveraging 3D structure of the proteins



Structure Search Approach

Strengths

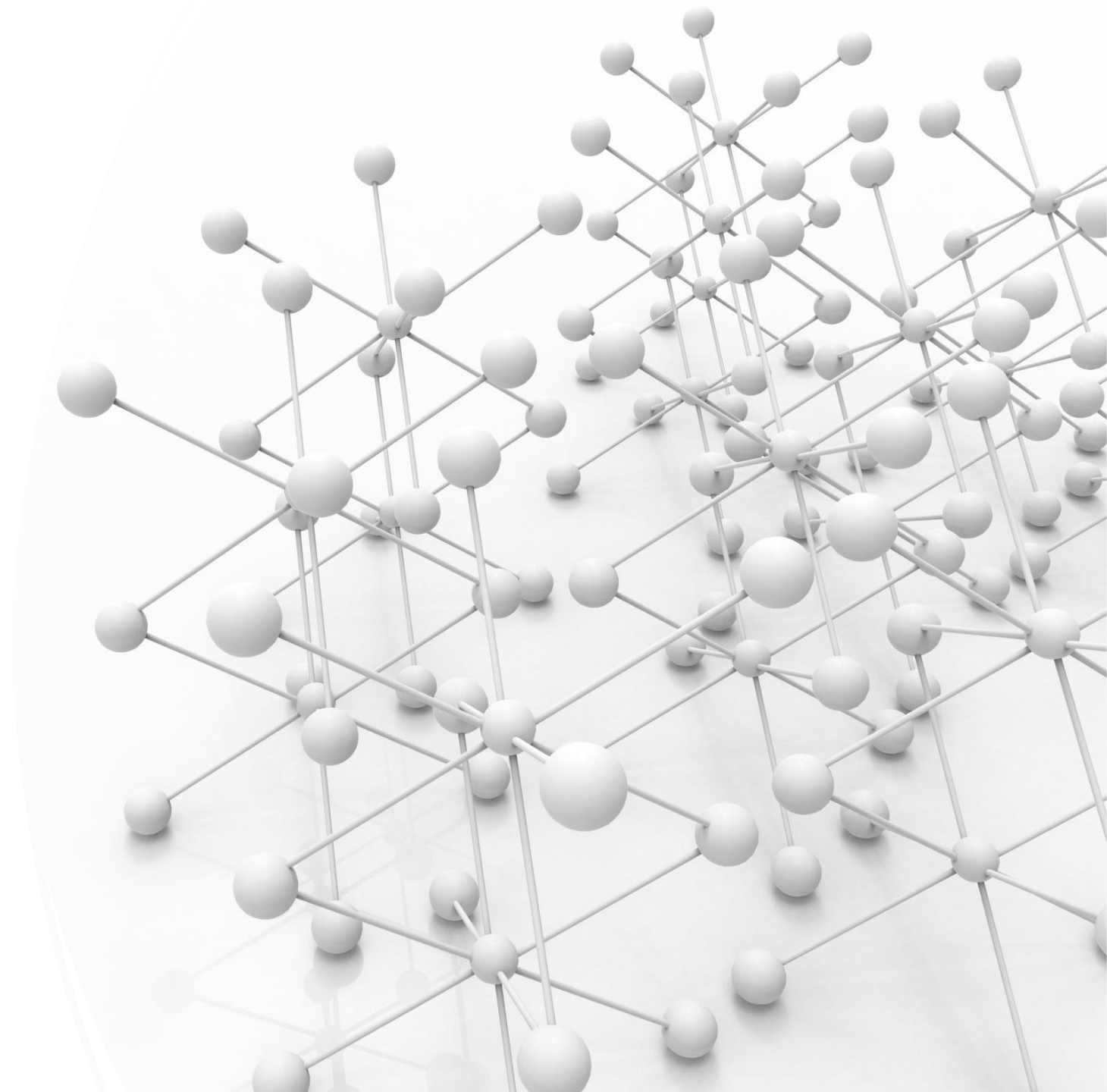
- 1) Better remote performance

Challenges

- 1) High cost
- 2) Data scarcity

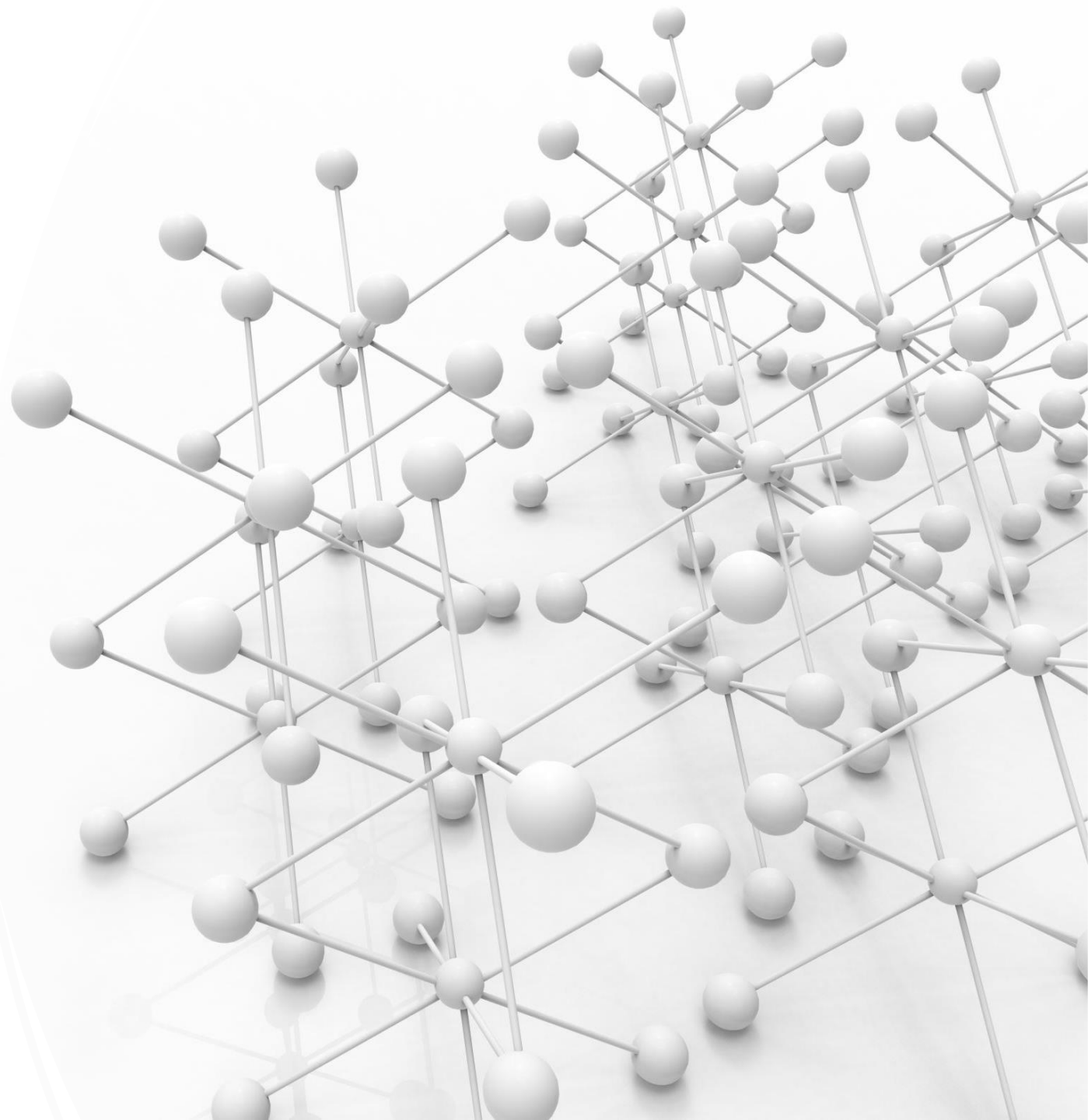
Structure Search Approach

Protein structure search methods can be further divided into 3 main approaches



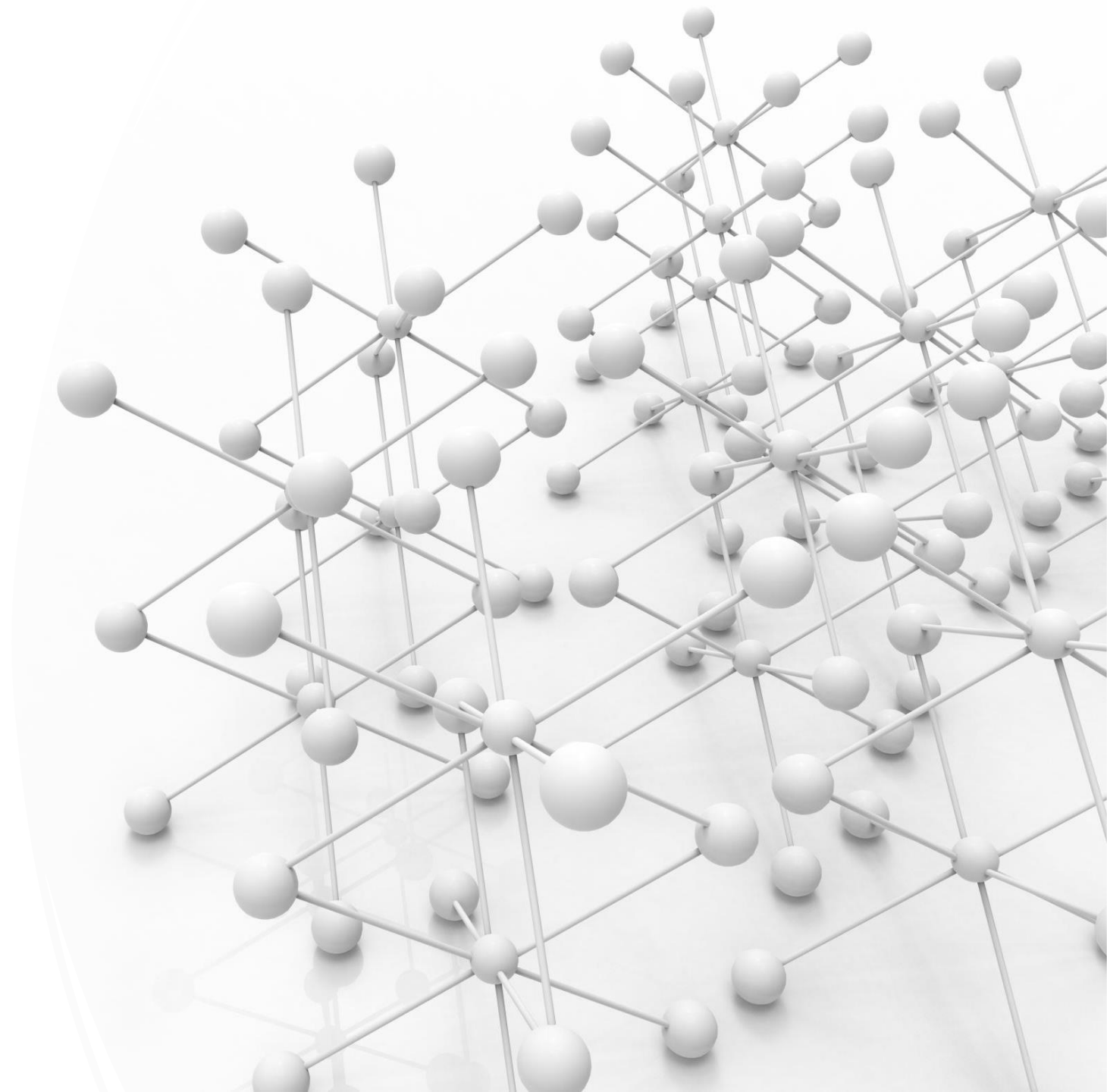
Contact/Distance Map-Based

- Map_align
- EigenTHREADER
- DiscoVER



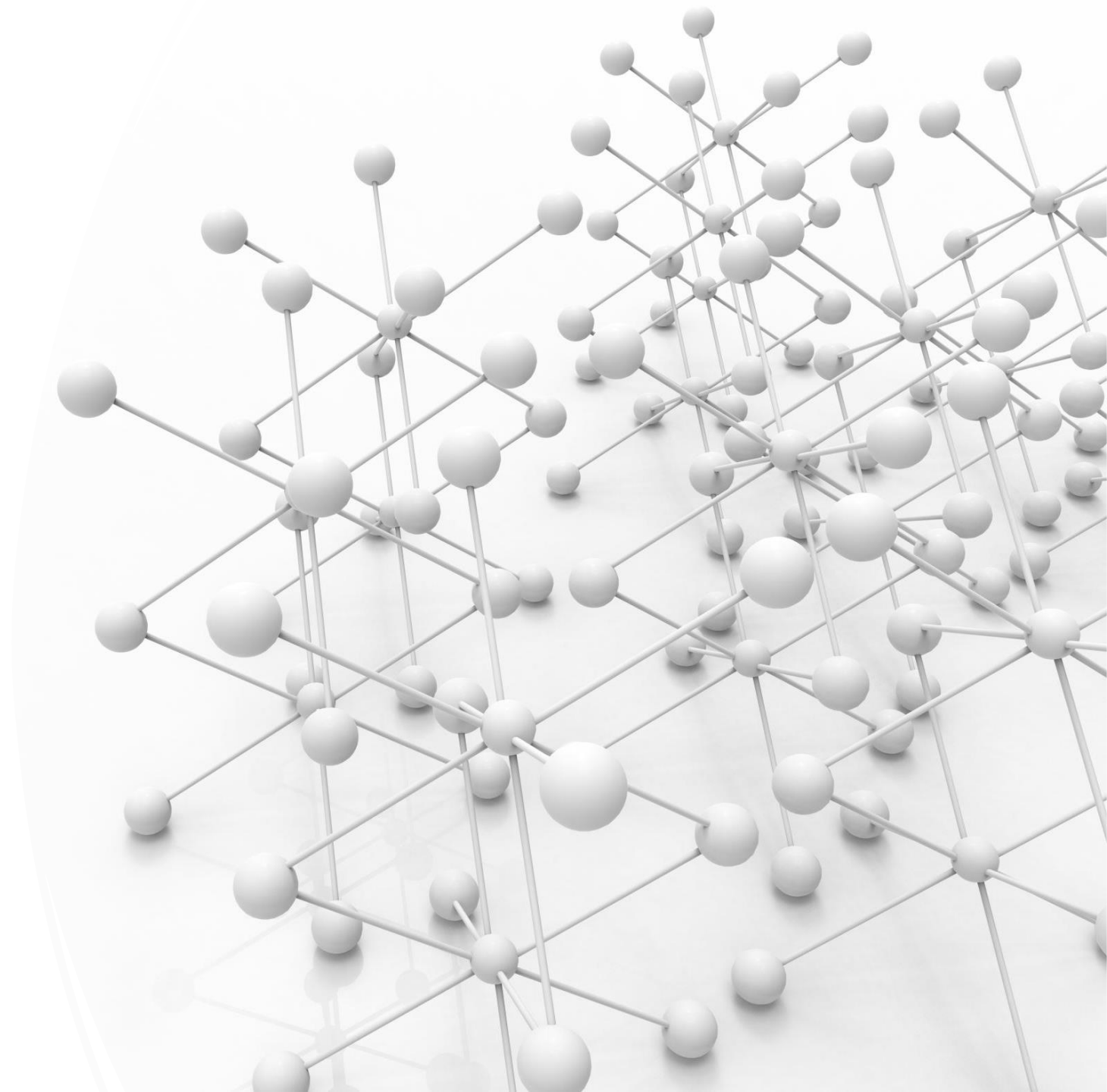
Structural Alphabet-Based

- 3D-BLAST-SW
- CLE-SW
- Foldseek
- Foldseek-TM



Structural Alignment-Based

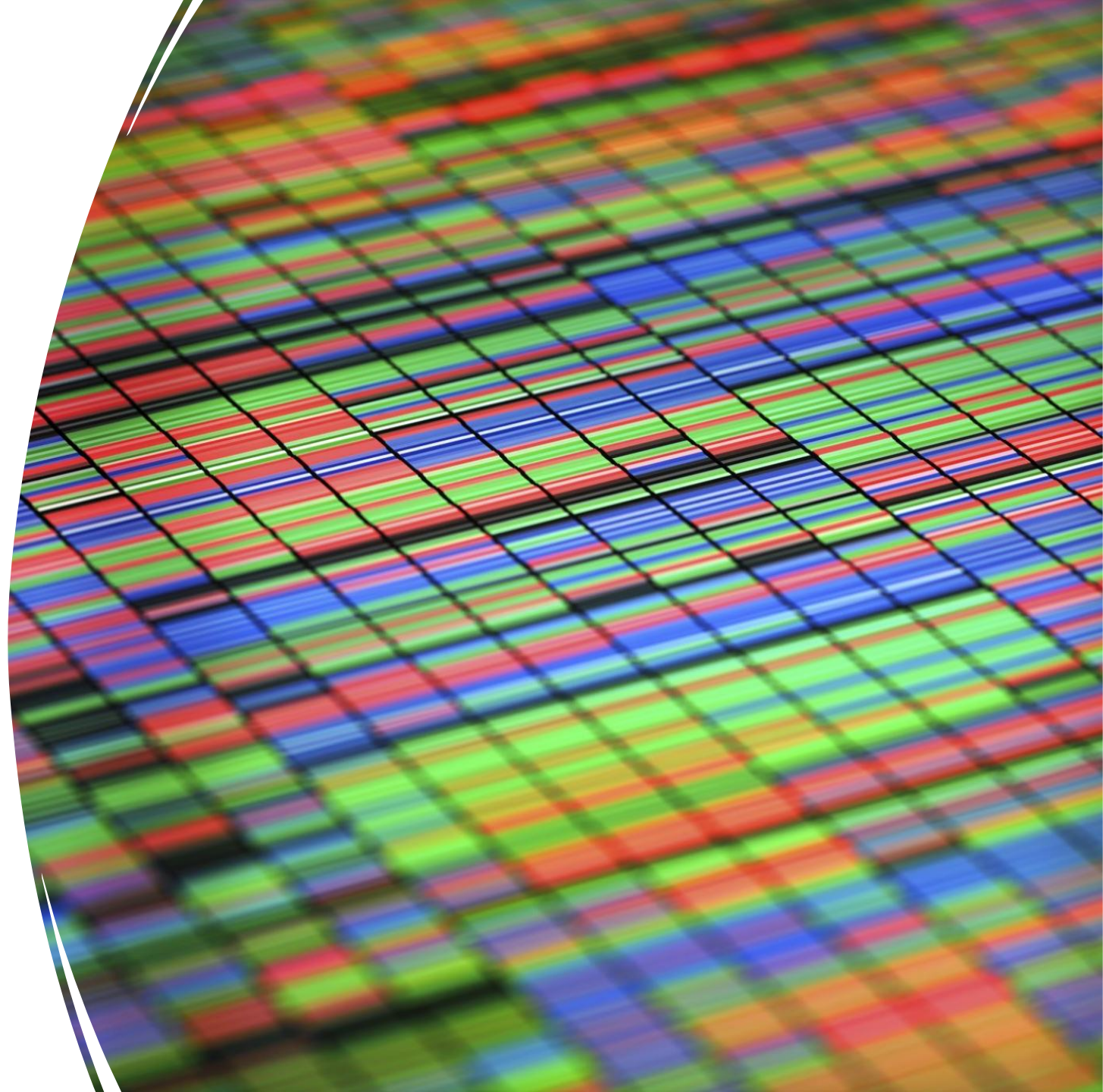
- CE
- Dali
- TM-align

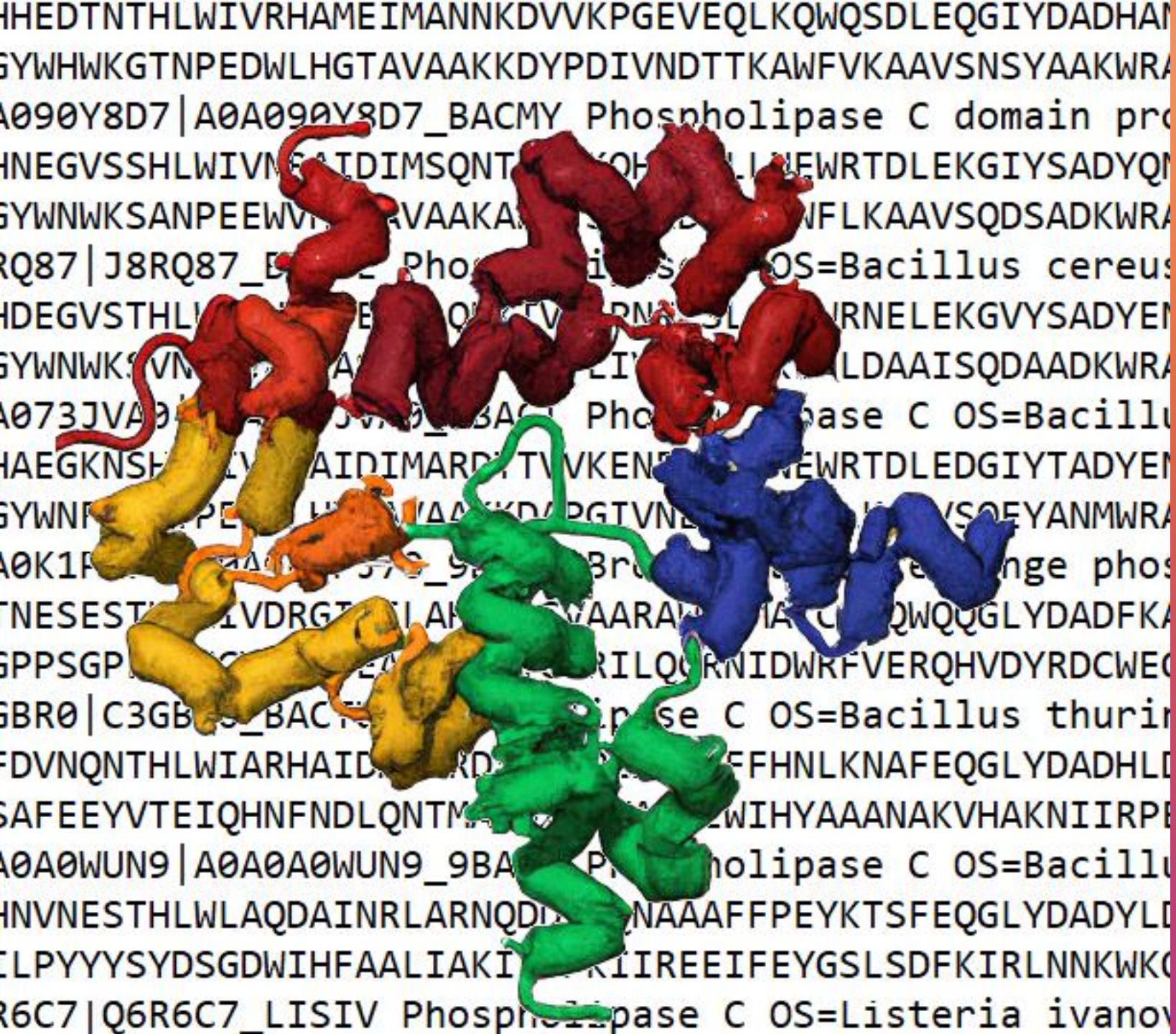


An Old Fix: HMM

Hidden Markov Models
(HMMs)

- HMMER
- HHsearch
- HHblits





The New Wave: Protein Language Models

(PLM)

Goal: Learn contextual representations of proteins by treating sequences as language

DTNTHLWIVRHAMEIMANNKDVVKPGEVEQLKQWQSDLE
JHWKGTNPEDWLHGTAVAAKKDYPDIVNDTTKAWFVKA
0Y8D7 | A0A090Y8D7_BACMY Phospholipase C
GVSSHLWIVMFAIDIMSQNT...KOH...LI...NEWRTDLE
JNWKSANPEEWM...RVAAKA...S...NFLKAAV
37 | J8RQ87_E...Pho...is...OS=Bacil
GVSTHL...E...Q...TV...RNA...SL...RNELE
JNWK...VN...A...LI...LDAAI
3JVA...BAC...Pho...ase C
GKNSH...AIDIMARD...TV...KENE...NEWRTDLE
JNF...PE...H...A...K...D...P...G...TV...N...
1F...44...57...9...Bro...e
SEST...IVDRGT...I...A...VAARA...IA...C...QWQ
SGP...C...E...RILOQ...RNIDWRFVERQ
0 | C3GB..._BAC...ip...se C OS=Bacil
VNQNTHLWIARHAID...RD...FFHNLKNAFE
EEYVTEIQHNFNDLQNTM...WIHYAAANAK
0WUN9 | A0A0A0WUN9_9BA...Pho...lipase C
VNESTHLWLAQDAINRLARNQDI...NAAAFFPEYKTSFE
YYYSYDSGDWIHFAALIAKT...RIIREEIFEYGSLSDF
7 | Q6R6C7_LISIV Phospholipase C OS=Liste



Protein Language Models

- ESMs
- ProtTrans
- ProtENN
- CATHe
- DEDAL
- DeepBLAST
- pLM-BLAST



PLMSearch

PLMSearch is designed to improve the **sensitivity** of sequence-based protein searches while maintaining their **efficiency** and **universality**.

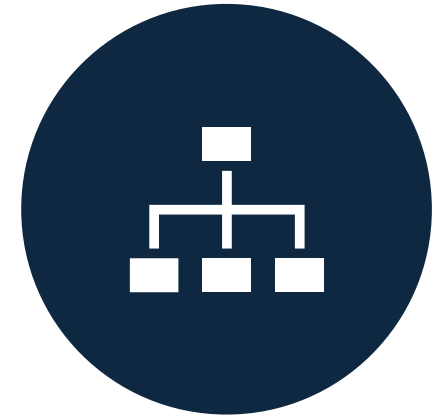
PLMSearch High Level



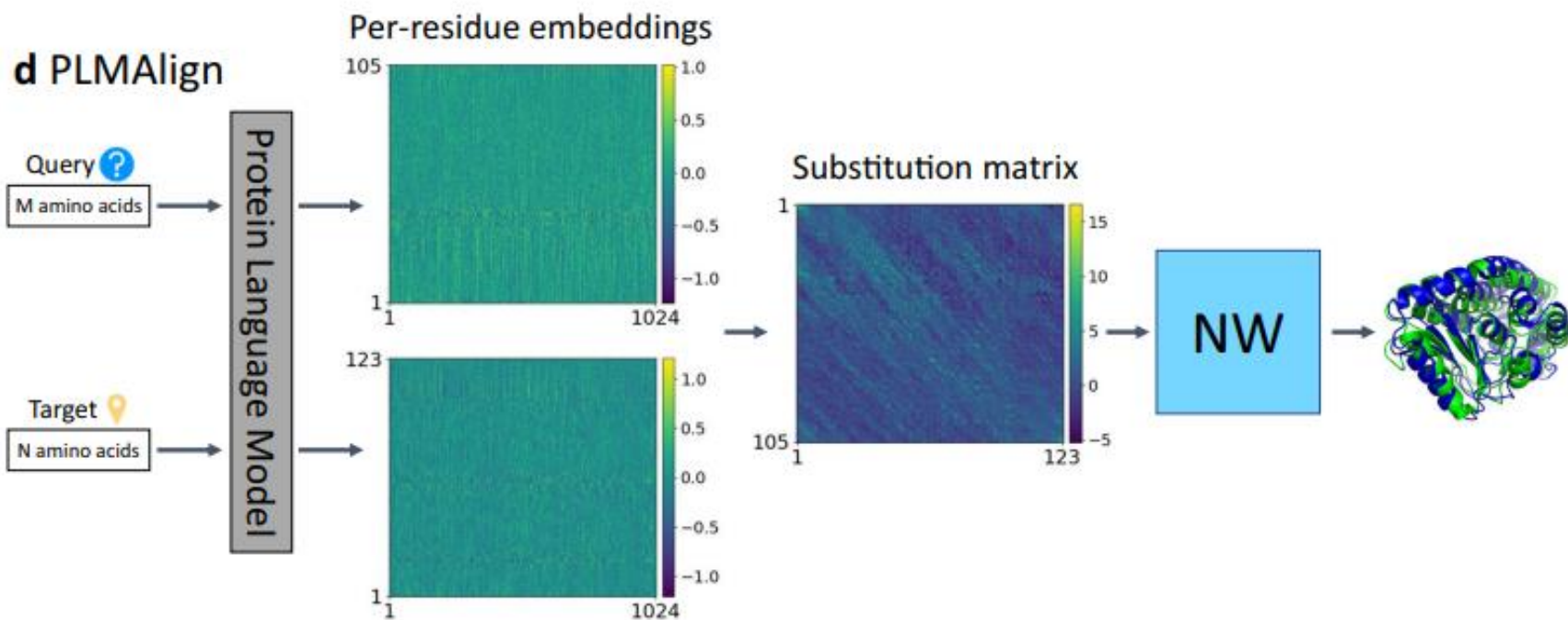
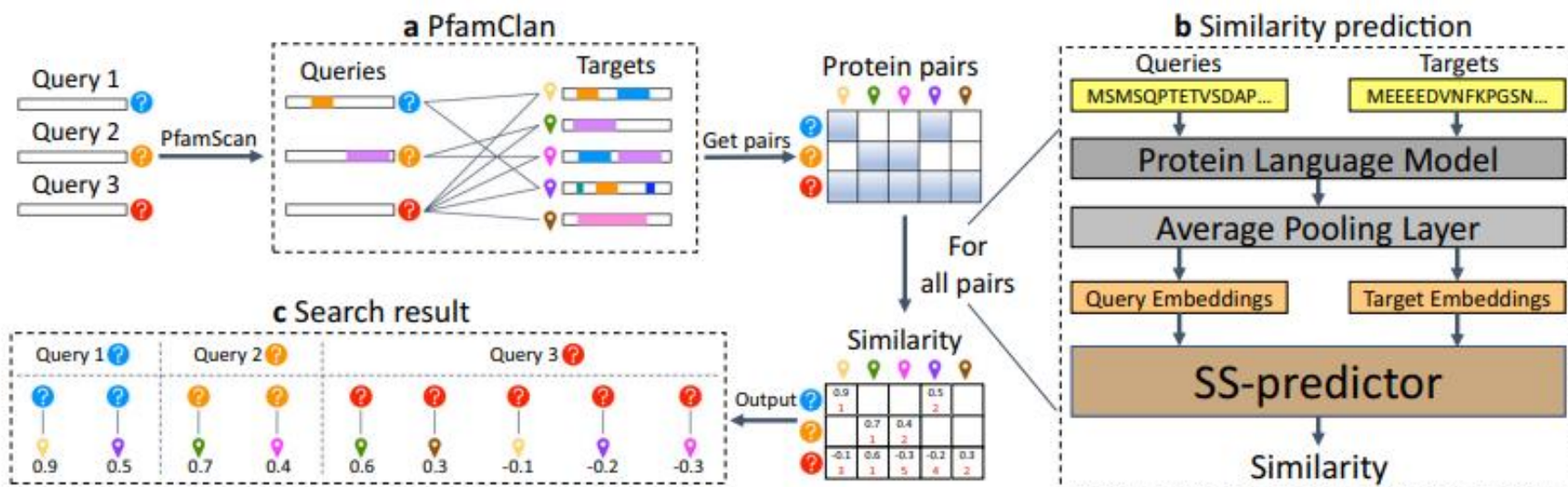
1) PFAMCLAN
FILTERING



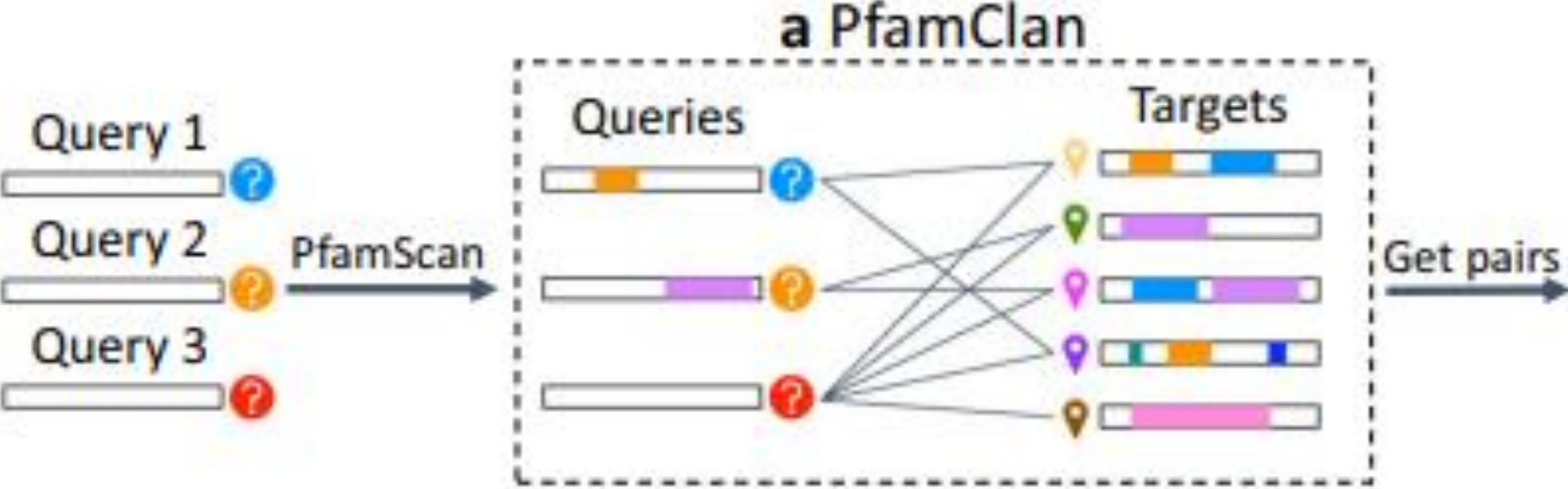
2) SS-PREDICTOR



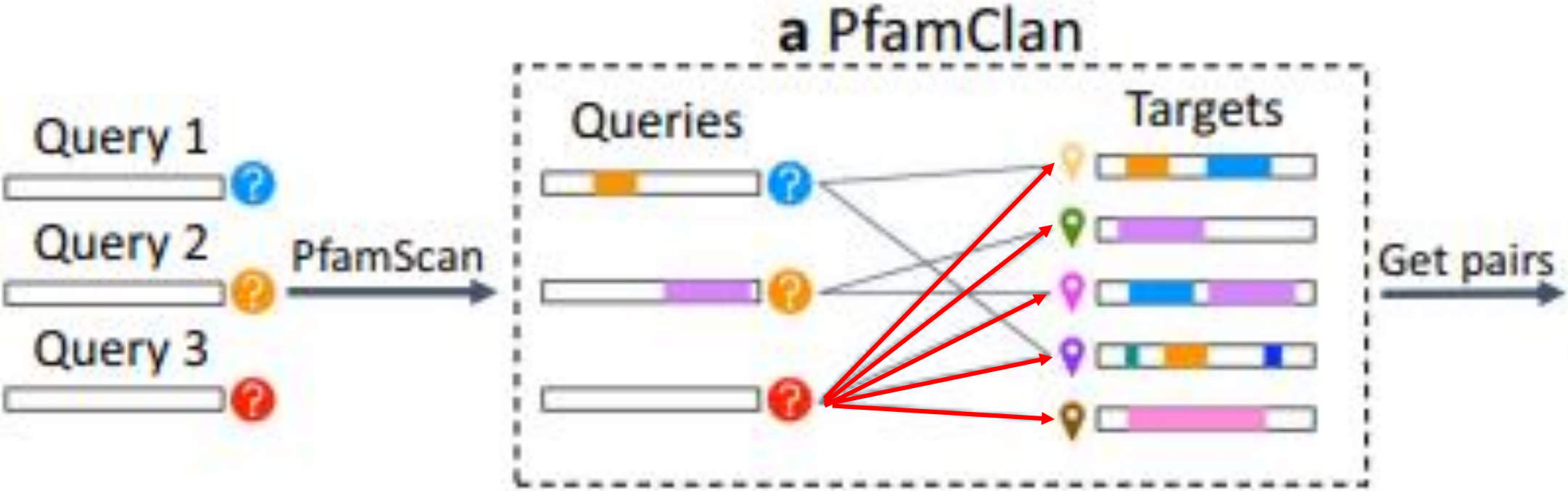
3) RANKING /
PLMALIGN



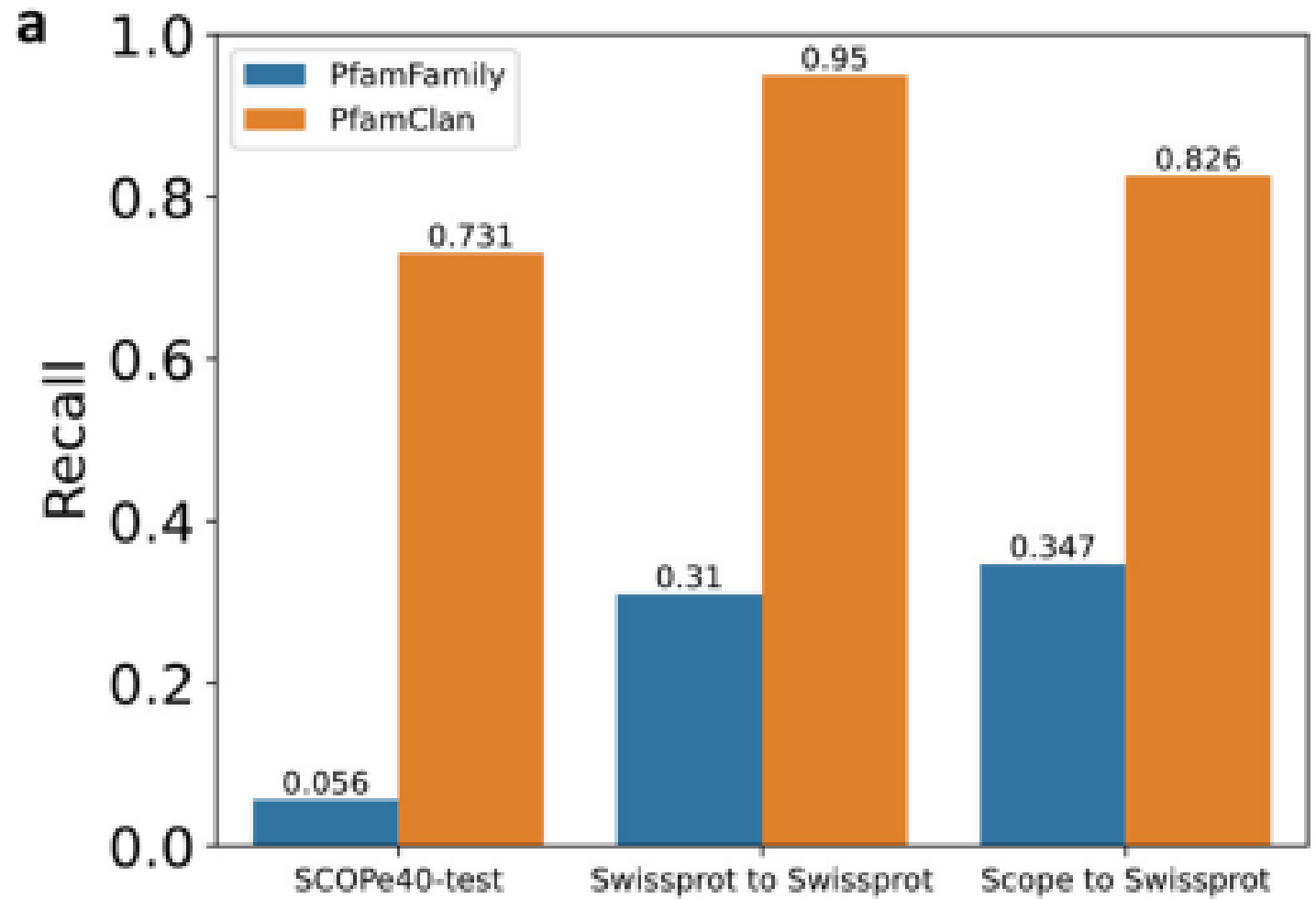
PfamClan

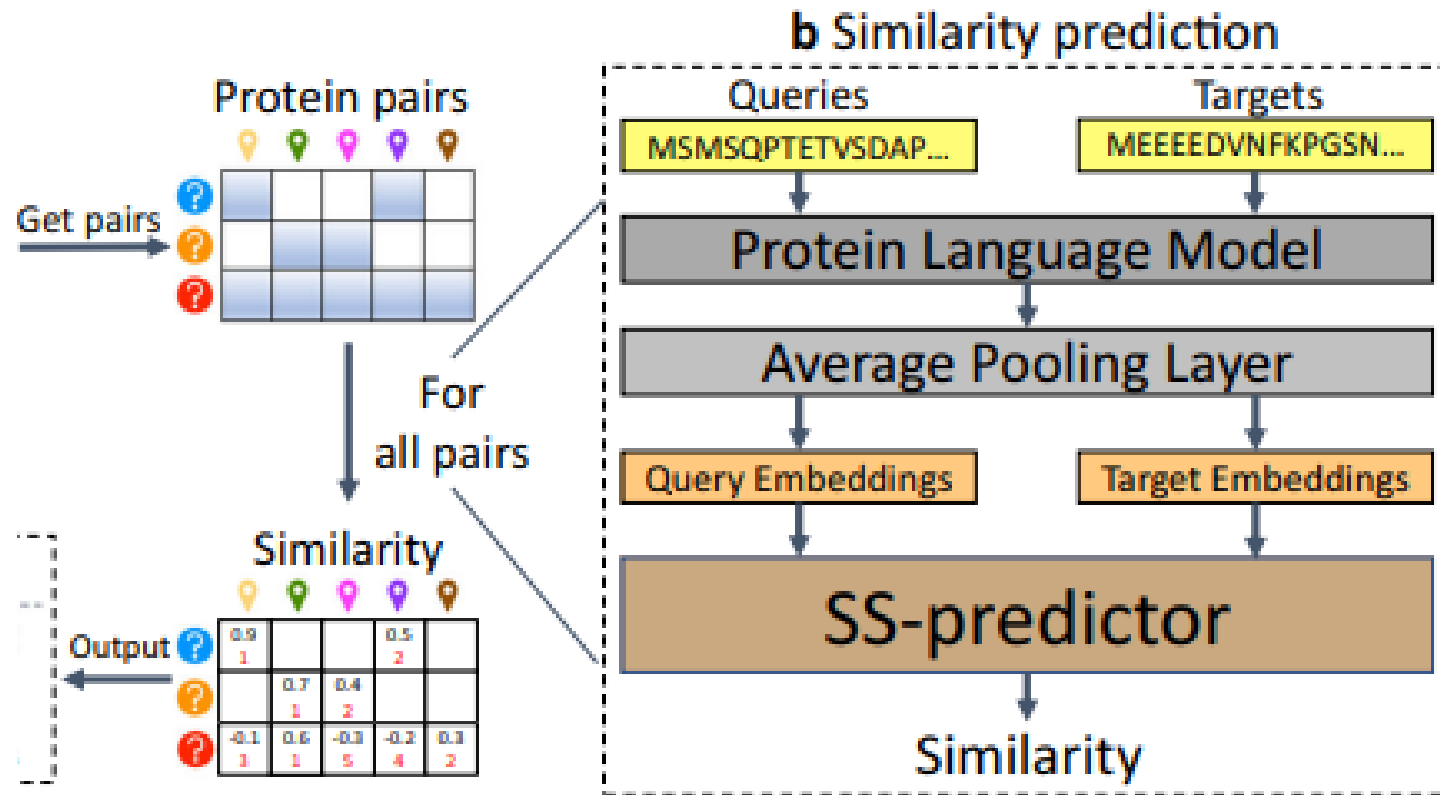


PfamClan

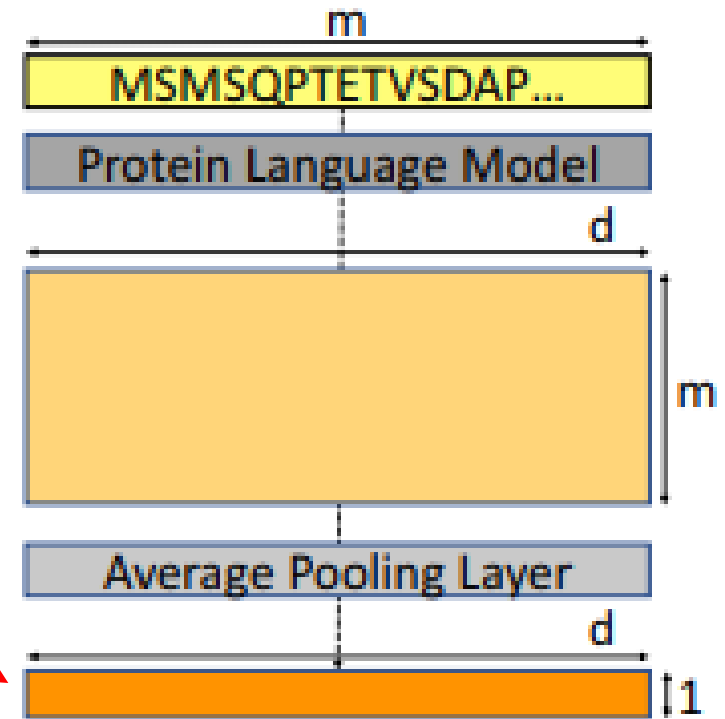
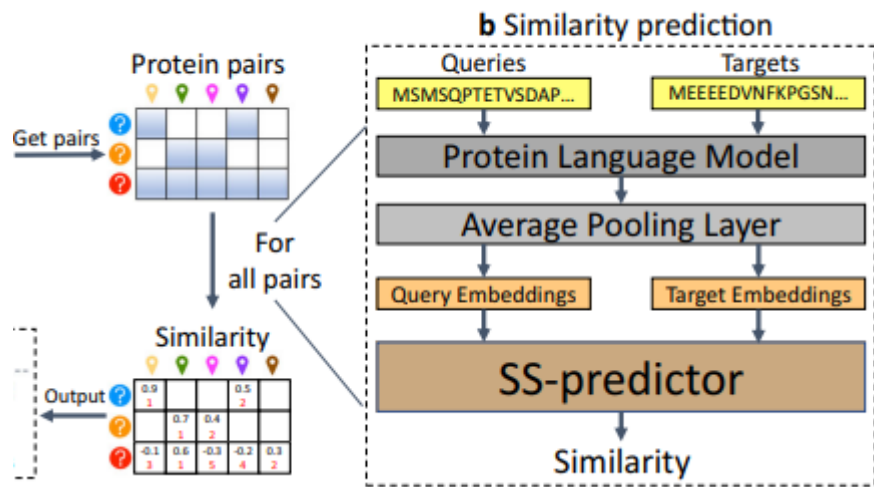


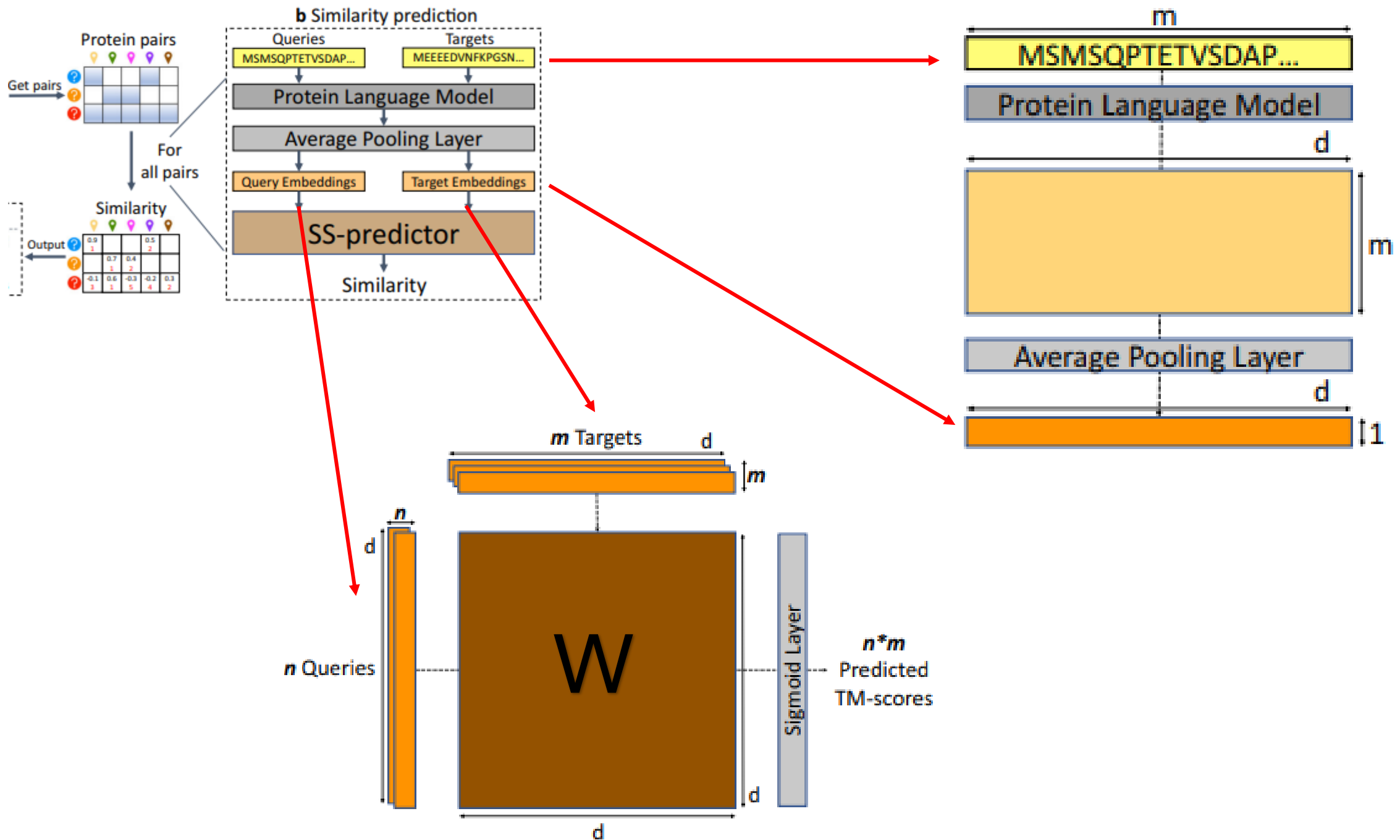
Pfam Clans vs Pfam Families



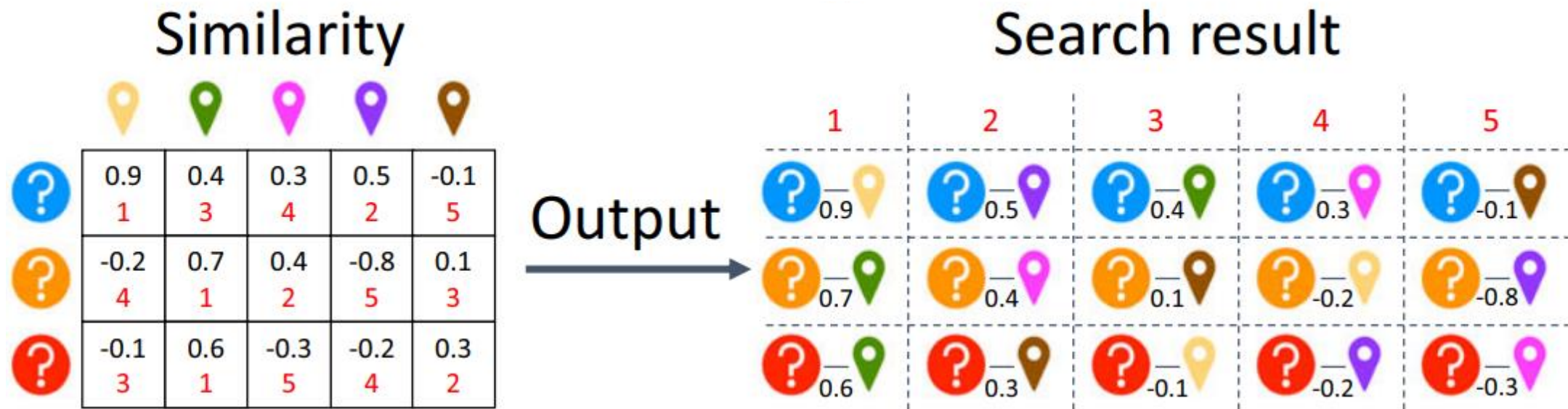


Structural Similarity Predictor (SS-predictor)



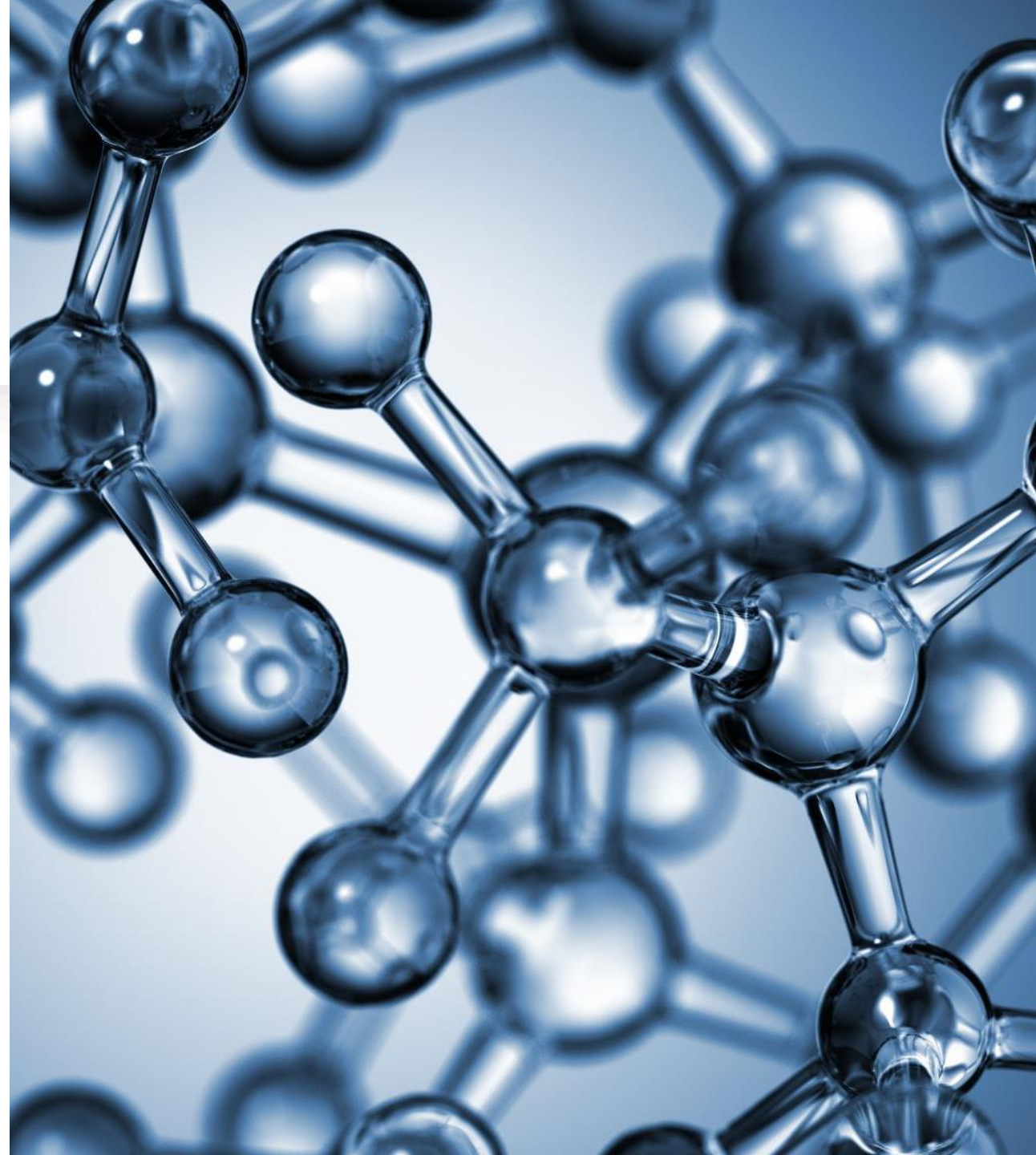


SS-predictor: Final Output

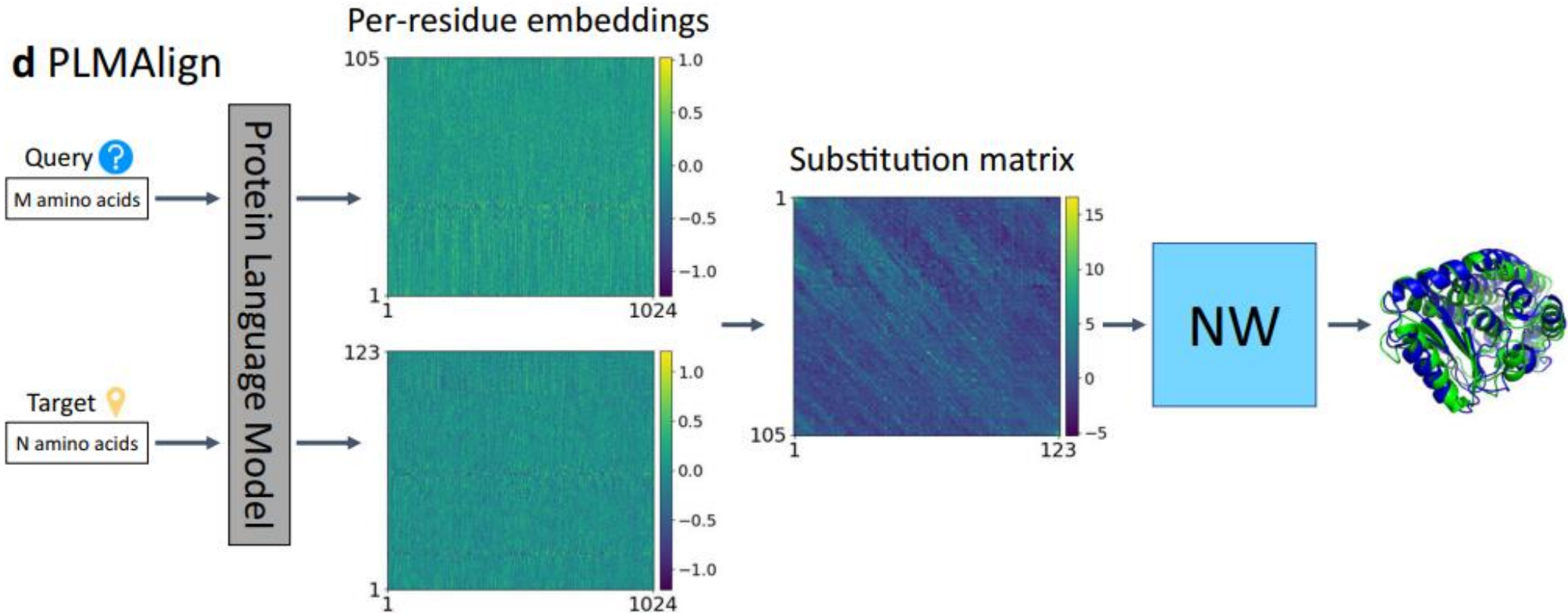


SS-predictor: Key Advantages

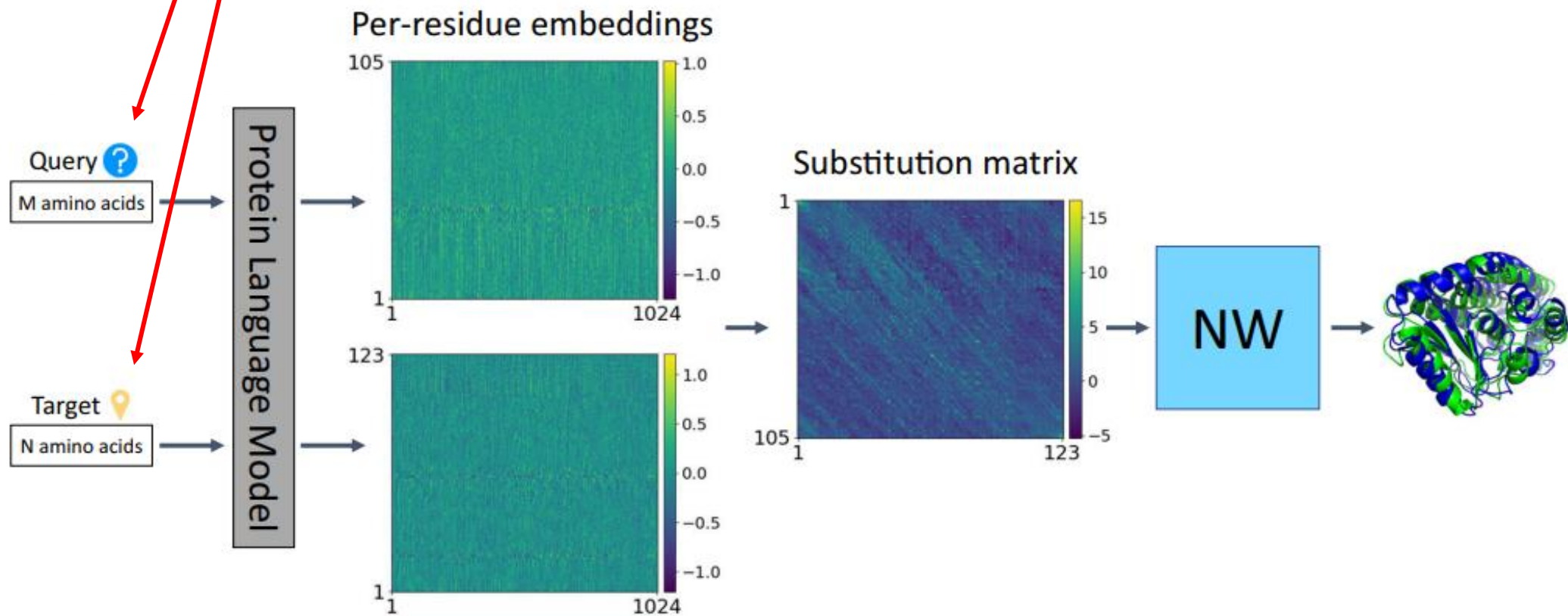
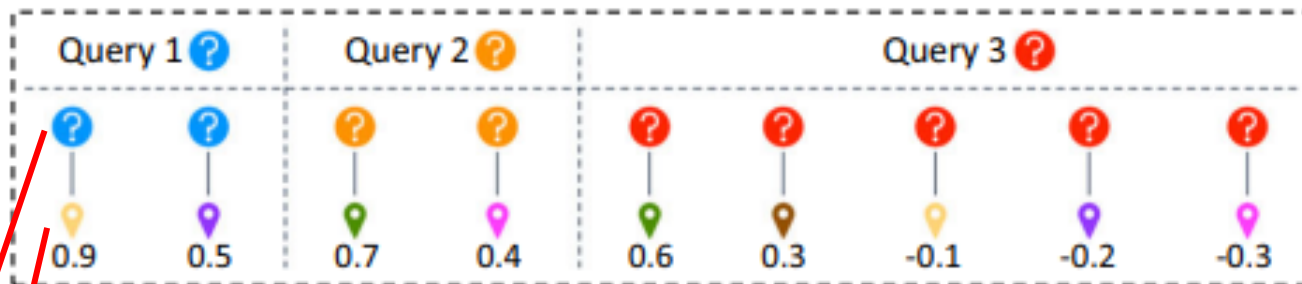
- 1) **Efficiency:** Avoids the computational cost of structure-based alignment for large datasets
- 2) **Scalability:** Handles the increasing size of protein clusters effectively through sequence-based predictions
- 3) **Precision:** Combines TM-score and COS similarity to handle both low and high sequence identity cases
- 4) **Global Similarity Focus:** Optimized for detecting **global structural similarity**, making it suitable for applications requiring a complete fold detection



PLMAlign



c Search result





Evaluation/Results

Results: Sensitivity Benchmarks

This paper's work



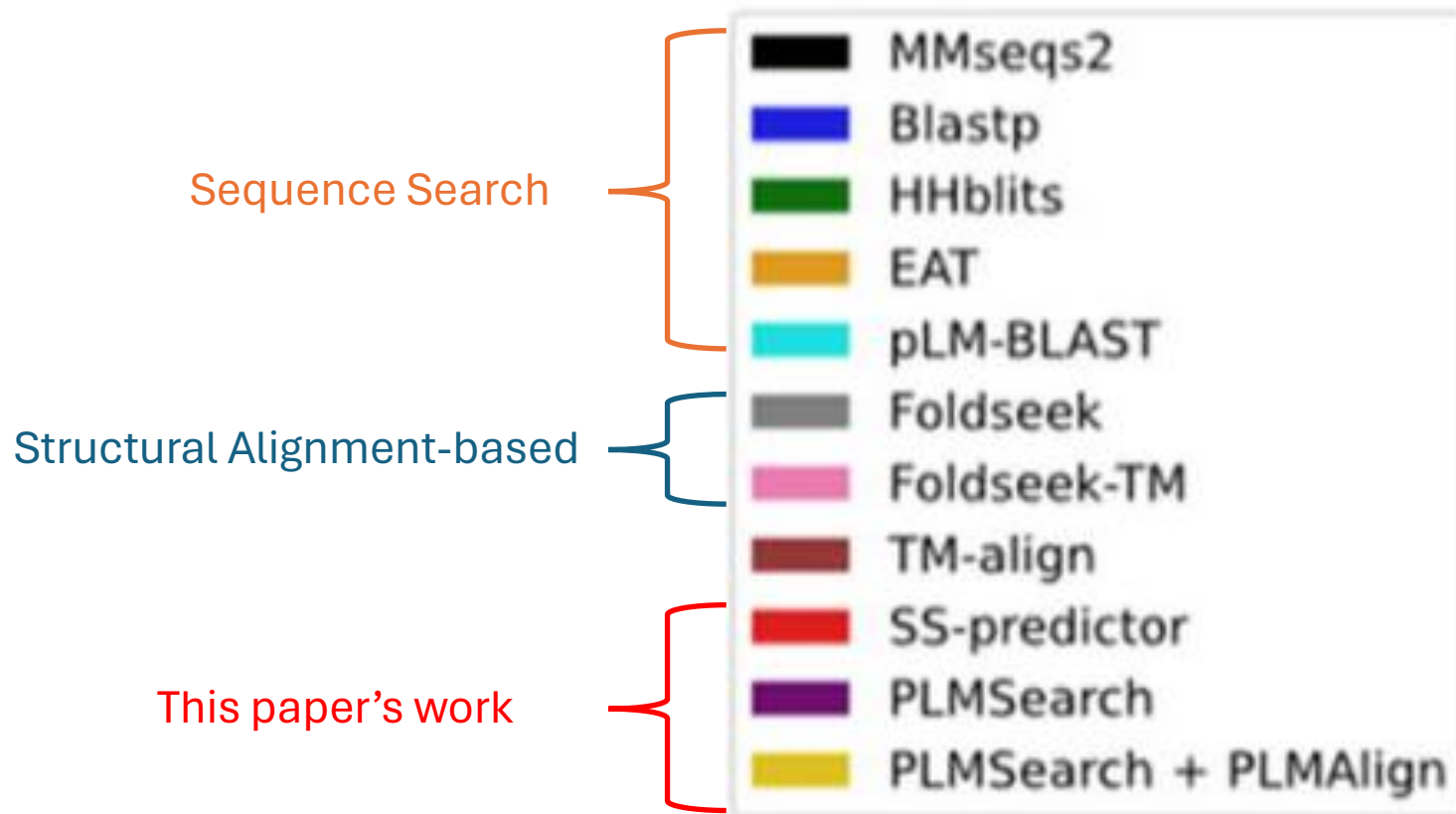
Results: Sensitivity Benchmarks

Sequence Search

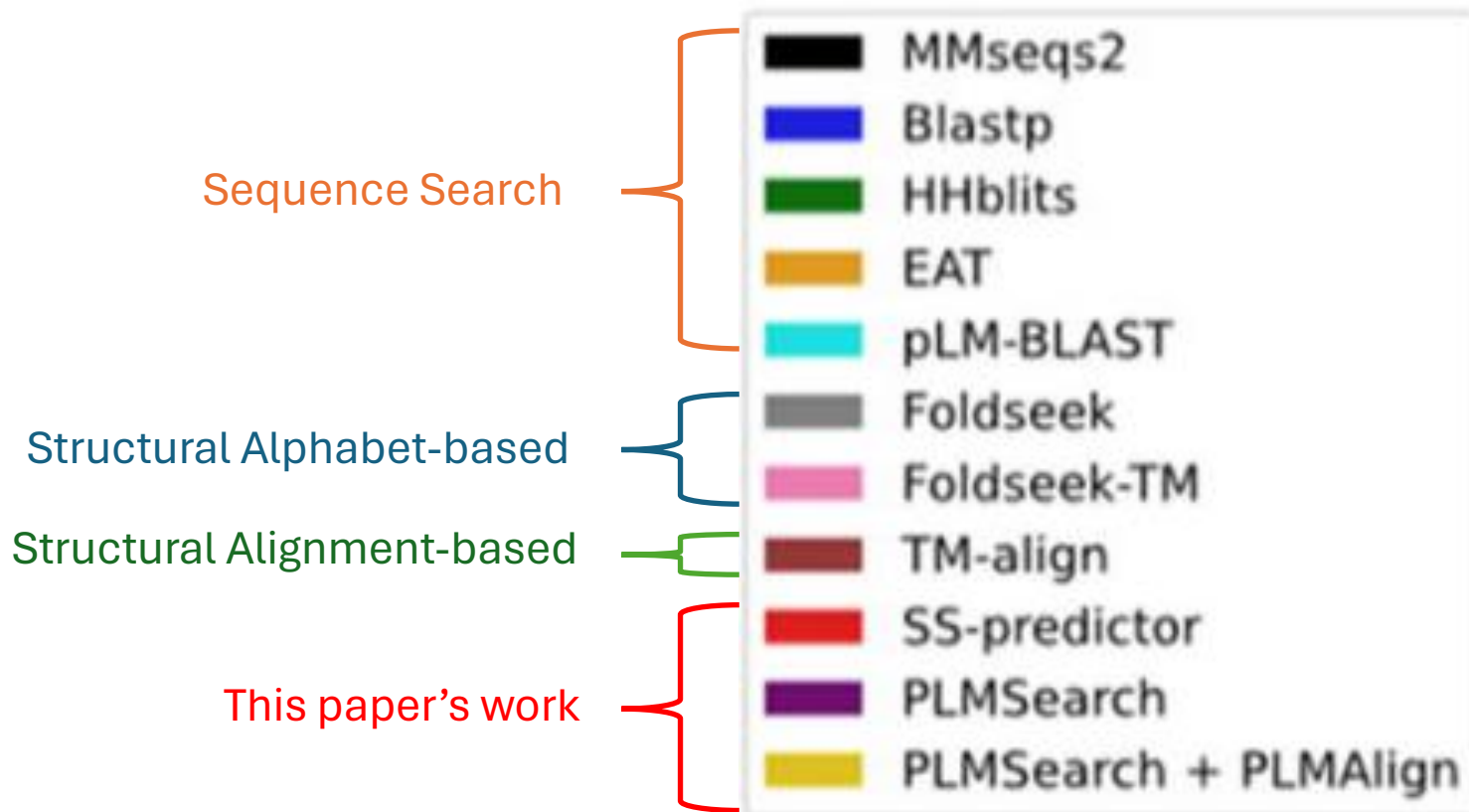
This paper's work



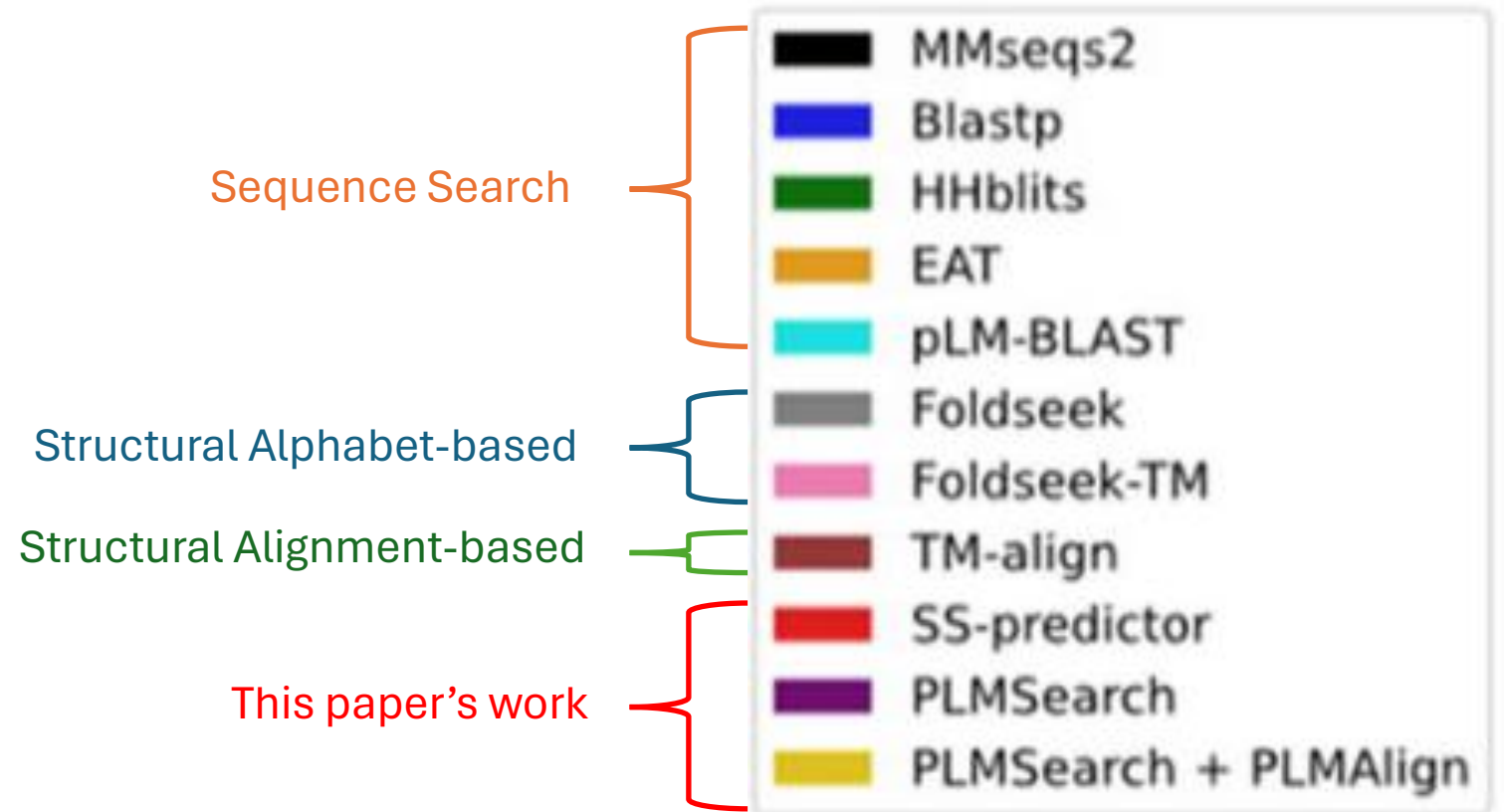
Results: Sensitivity Benchmarks



Results: Sensitivity Benchmarks

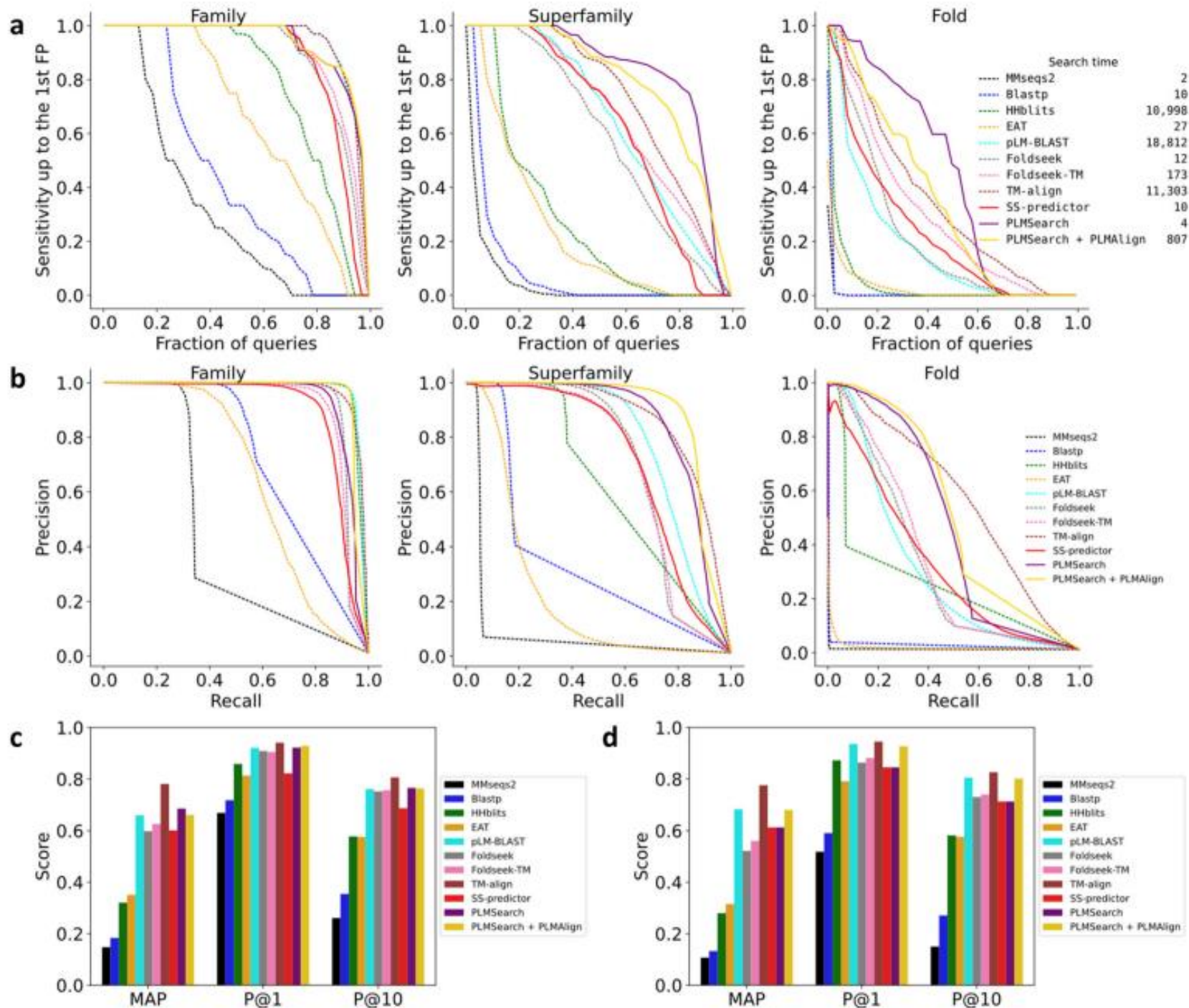


Results: Sensitivity Benchmarks

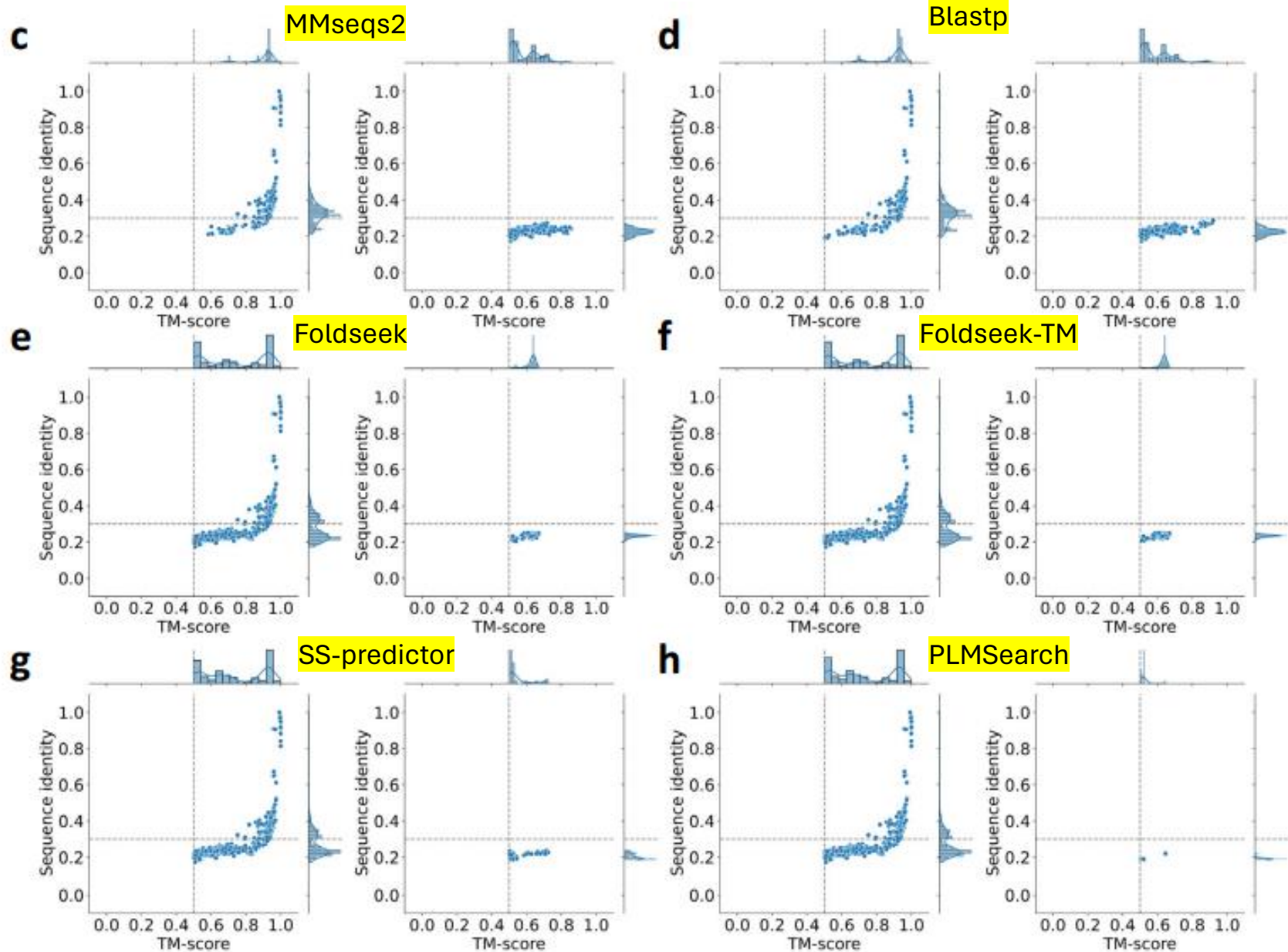


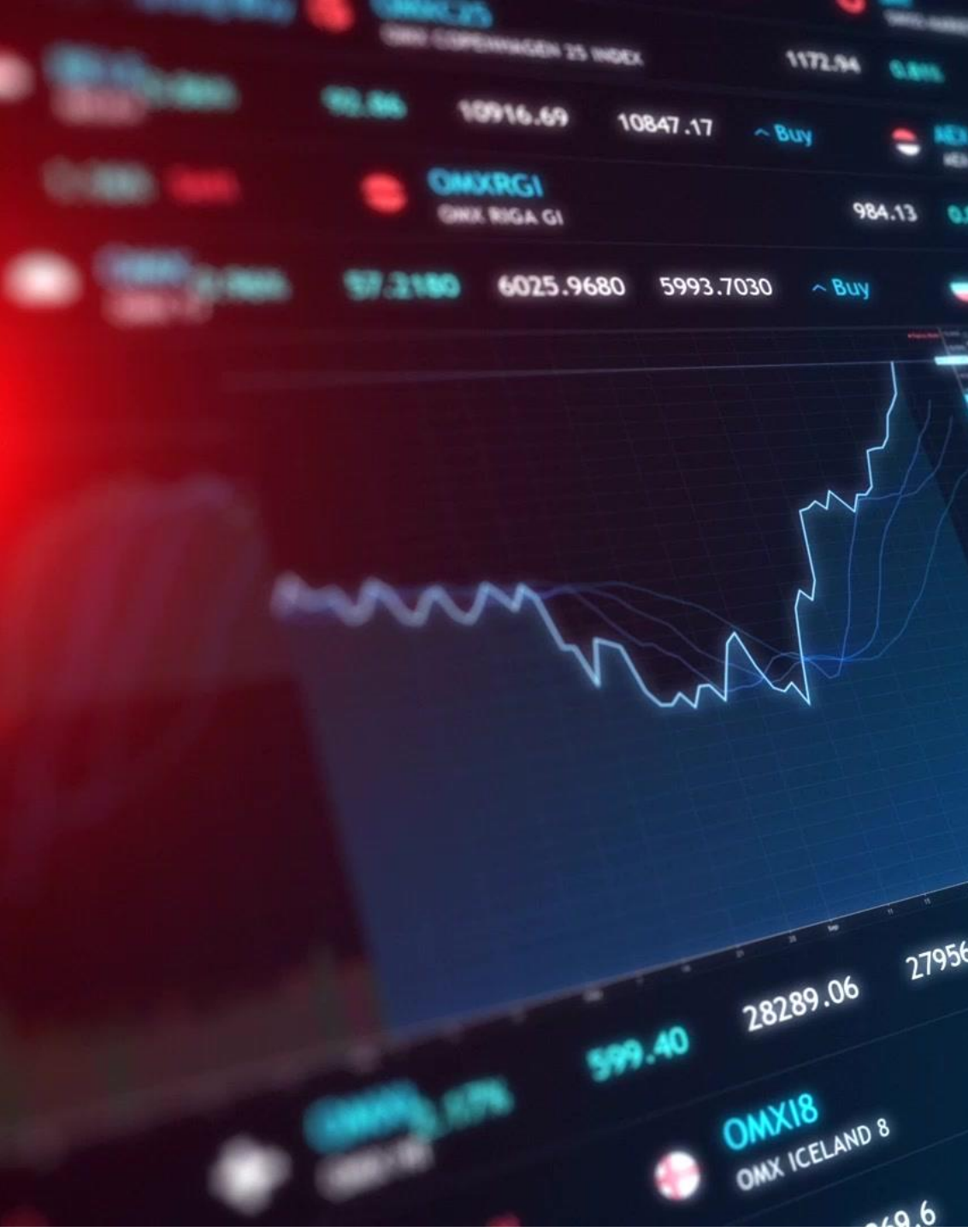
- SCOPe40-test dataset (2207 proteins)
- All-versus-all search test is performed
- 4,870,849 total query-target pairs

All-vs-All Search Test on SCOPe40-test



All-vs-All Search Test on SCOPe40-test





Results Summary

- PLMSearch is comparable to S.O.T.A. structural search methods
- PLMSearch overcame low sensitivity cap on sequence methods
- Improvements are concentrated in remote homology pairs
- PLMSearch is one of the fastest search methods
- Residue embedding-based alignment methods are limited by size of target dataset



Future Works

Efficient Embedding Storage





Thank You!





Questions?