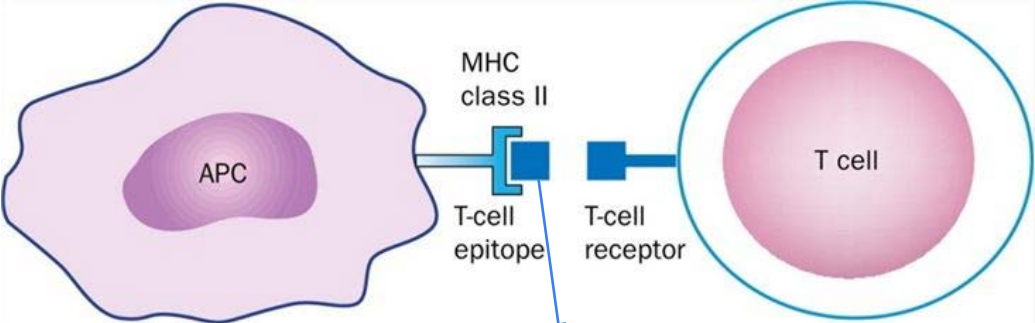




# PEPMatch: A Tool to Identify Short Peptide Sequence Matches in Large Sets of Proteins

Presented by: Daniel Marrama, Associate Bioinformatics Specialist I

# Peptide Sequence Matching



Search for matches



**Proteome**

\* Set of all proteins expressed by an organism

# Common Tools for Peptide Sequence Matching



diamond



MMseqs2

# Use Cases for Peptide Matching

## Curation for IEDB

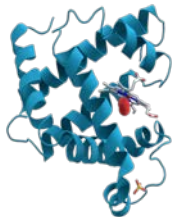
Reference  
Paper



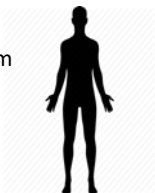
Epitope



Sourced  
from a  
protein



In an organism



## A large peptidome dataset improves HLA class I epitope prediction across most of the human population

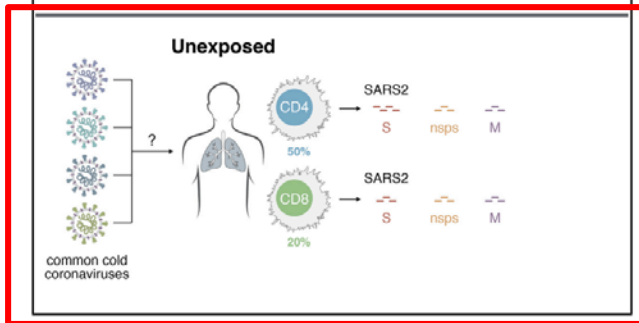
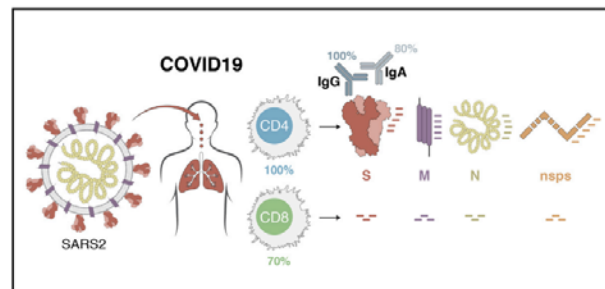
Siranush Sarkizova<sup>1,2,13</sup>, Susan Klaeger<sup>2,13</sup>, Phuong M. Le<sup>3</sup>, Letitia W. Li<sup>3</sup>, Giacomo Oliveira<sup>3</sup>, Hasmik Keshishian<sup>2</sup>, Christina R. Hartigan<sup>2</sup>, Wandi Zhang<sup>3</sup>, David A. Braun<sup>2,3,4,5</sup>, Keith L. Ligon<sup>2,4,6,7</sup>, Pavan Bachireddy<sup>2,3,5</sup>, Ioannis K. Zervantonakis<sup>8</sup>, Jennifer M. Rosenbluth<sup>8</sup>, Tamara Ouspenskaia<sup>2</sup>, Travis Law<sup>2</sup>, Sune Justesen<sup>9</sup>, Jonathan Stevens<sup>10</sup>, William J. Lane<sup>4,10</sup>, Thomas Eisenhaure<sup>2</sup>, Guang Lan Zhang<sup>3,4,11</sup>, Karl R. Clauser<sup>2</sup>, Nir Hacohen<sup>2,3,12\*</sup>, Steven A. Carr<sup>2\*</sup>, Catherine J. Wu<sup>2,3,4,5\*</sup> and Derin B. Keskin<sup>2,3,4,5,11\*</sup>

Prediction of HLA epitopes is important for the development of cancer immunotherapies and vaccines. However, current prediction algorithms have limited predictive power, in part because they were not trained on high-quality epitope datasets covering a broad range of HLA alleles. To enable prediction of endogenous HLA class I-associated peptides across a large fraction of the human population, we used mass spectrometry to profile >185,000 peptides eluted from 95 HLA-A, -B, -C and -G mono-allelic cell lines. We identified canonical peptide motifs per HLA allele, unique and shared binding submotifs across alleles and distinct motifs associated with different peptide lengths. By integrating these data with transcript abundance and peptide processing, we developed HLAthena, providing allele-and-length-specific and pan-allele-pan-length prediction models for endogenous peptide presentation. These models predicted endogenous HLA class I-associated ligands with 1.5-fold improvement in positive predictive value compared with existing tools and correctly identified >75% of HLA-bound peptides that were observed experimentally in 11 patient-derived tumor cell lines.

# Use Cases for Peptide Matching

## Targets of T Cell Responses to SARS-CoV-2 Coronavirus in Humans with COVID-19 Disease and Unexposed Individuals

### Graphical Abstract



### Authors

Alba Grifoni, Daniela Weiskopf, Sydney I. Ramirez, ..., Davey M. Smith, Shane Crotty, Alessandro Sette

### Correspondence

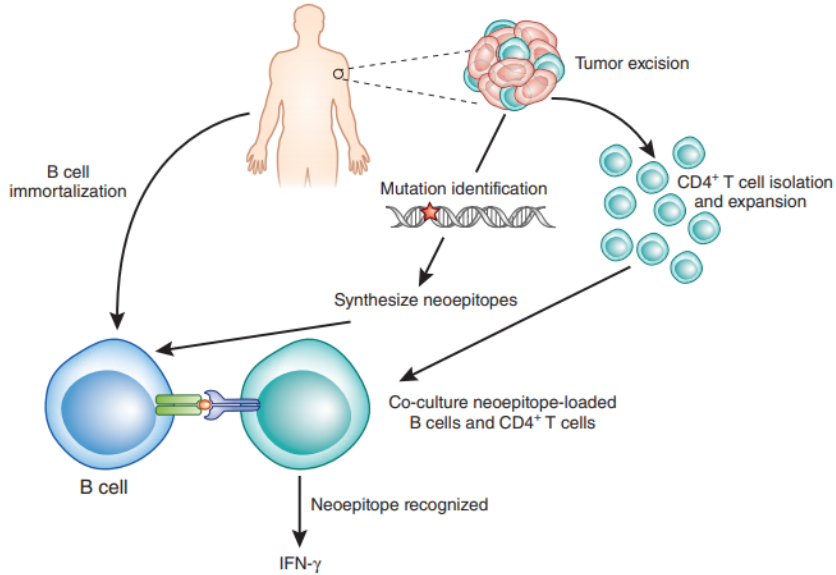
shane@lji.org (S.C.), alex@lji.org (A.S.)

### In Brief

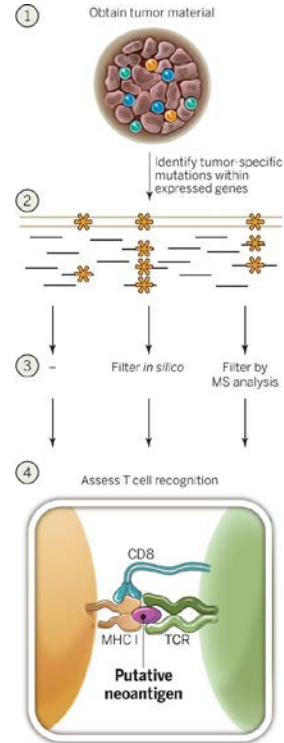
An analysis of immune cell responses to SARS-CoV-2 from recovered patients identifies the regions of the virus that is targeted and also reveals cross-reactivity with other common circulating coronaviruses

# Use Cases for Peptide Matching

## Neopeptide Similarity For Cancer Vaccine Candidates



Overwijk, W., et al. *Nat Med* (2015)



Schumacher TN, et al. *Science*. (2015)

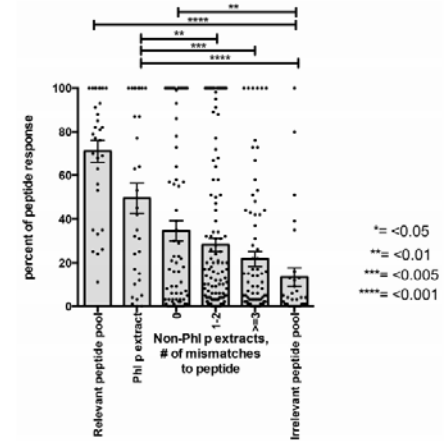
# Use Cases for Peptide Matching

## Conservation Analyses

- Milk allergens → do conserved peptides in human proteome drive stronger reactions?
- Fungus allergens → does conservation across allergen species determine immunogenicity?
- COVID-19 vaccination → does the spike protein encoded by the mRNA vaccines have similar enough homology that could lead to autoimmune myocarditis?

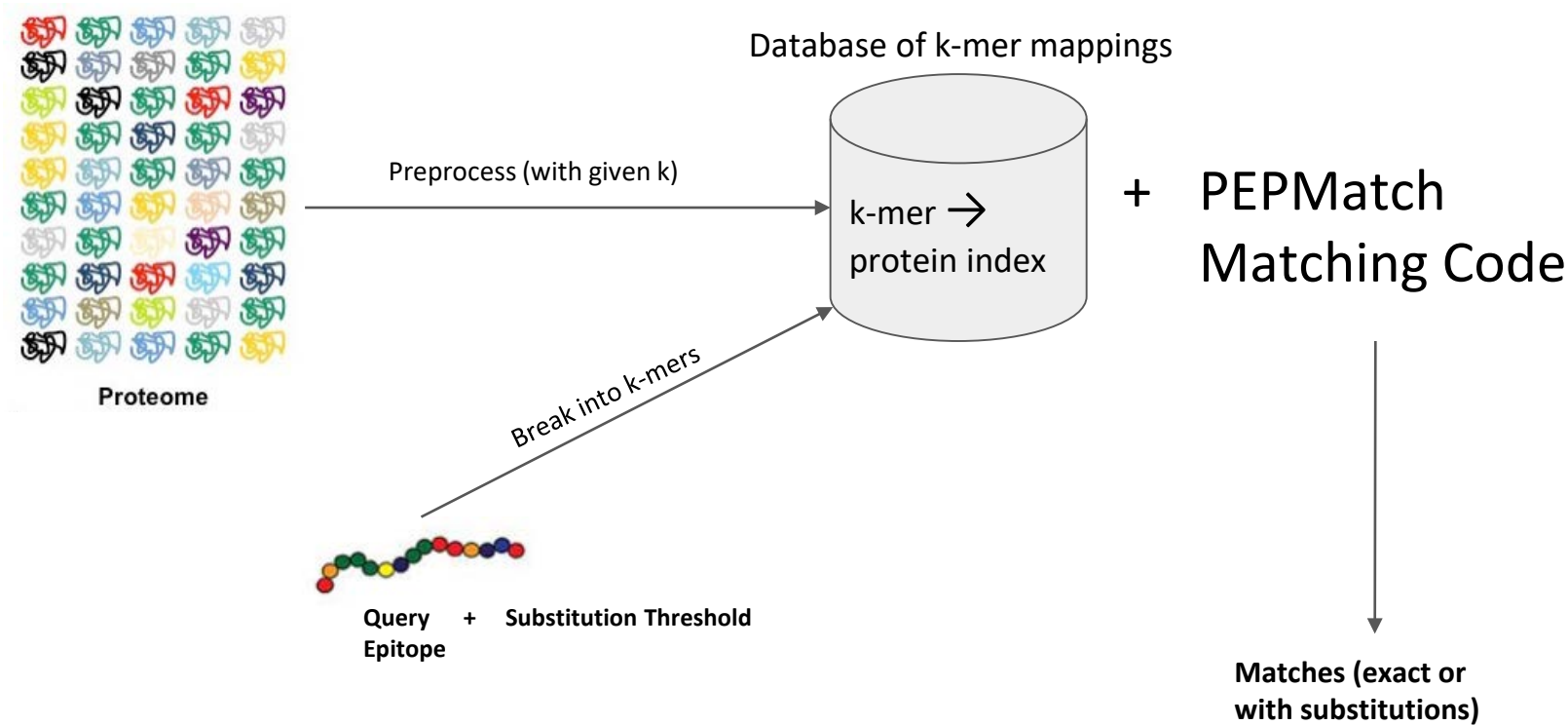
T-cell epitope conservation across allergen species is a major determinant of immunogenicity

Luise Westernberg<sup>1</sup>, Véronique Schulten<sup>1</sup>, Jason A Greenbaum<sup>1</sup>, Sara Natali<sup>3</sup>, Victoria Tripple<sup>1</sup>, Denise M. McKinney<sup>1</sup>, April Frazier<sup>1</sup>, Heidi Hofer<sup>2</sup>, Michael Wallner<sup>2</sup>, Federica Sallusto<sup>3,4</sup>, Alessandro Sette<sup>1</sup>, and Bjoern Peters<sup>1</sup>



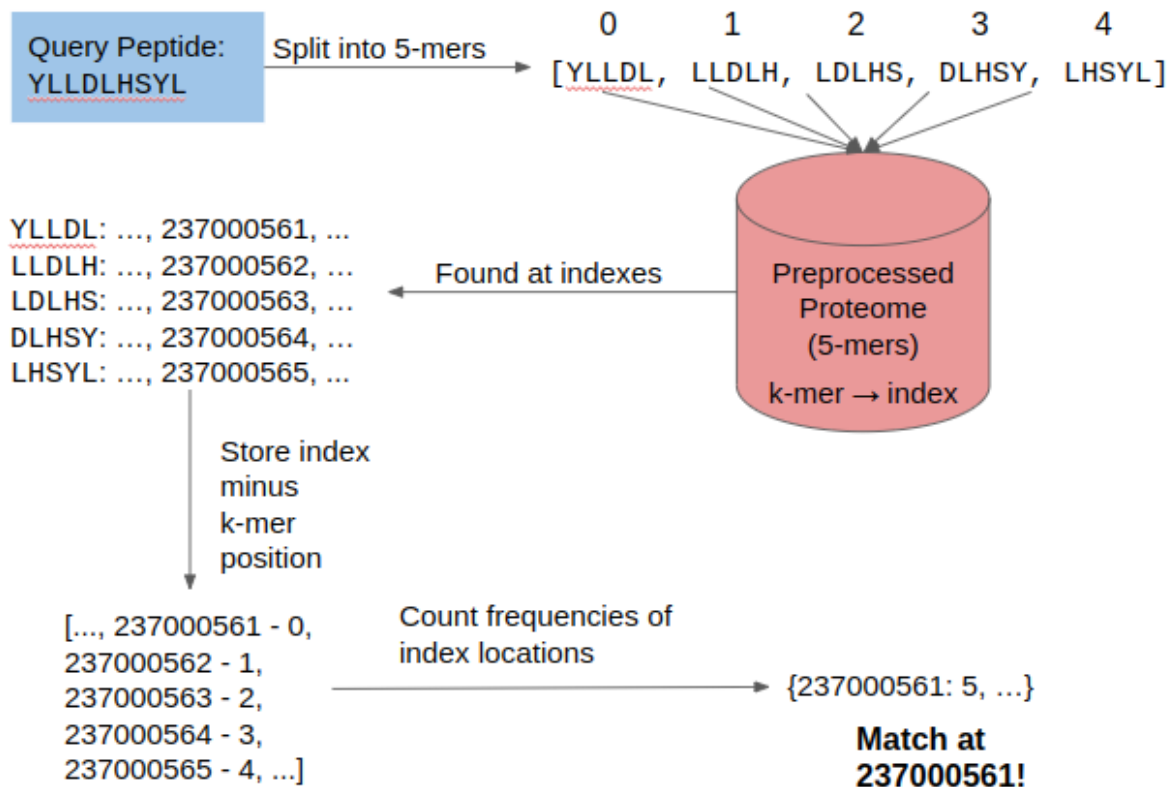
Westernberg, Luise et al. Journal of Allergy and Clinical Immunology (2016)

# PEPMatch: Tool Overview

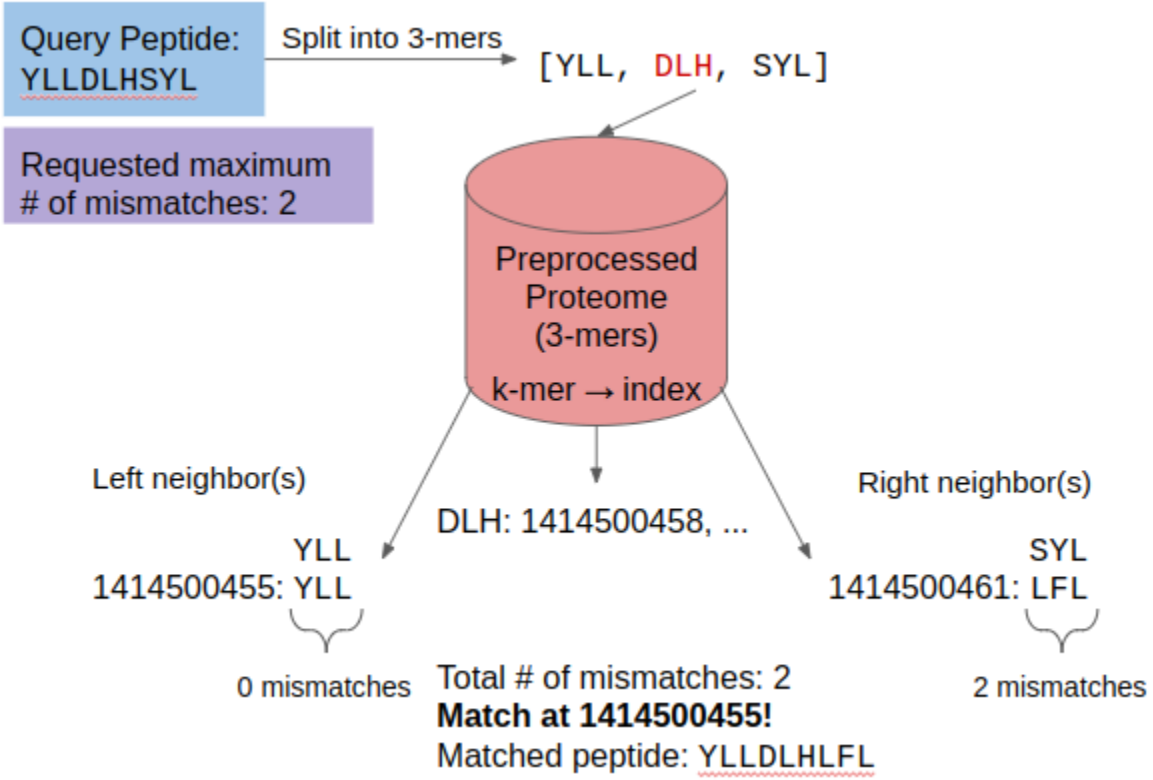




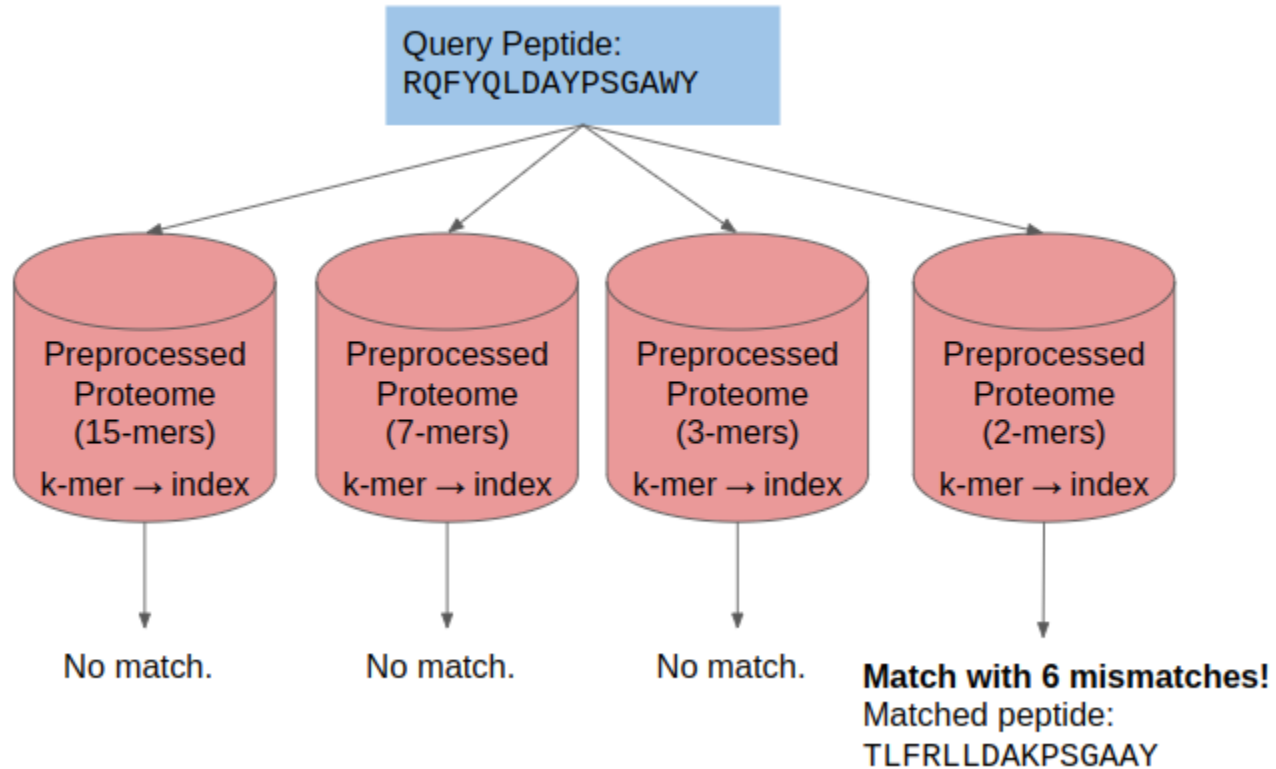
# PEPMatch: Exact Match Search



# PEPMatch: Searching with Substitutions (Mismatches)



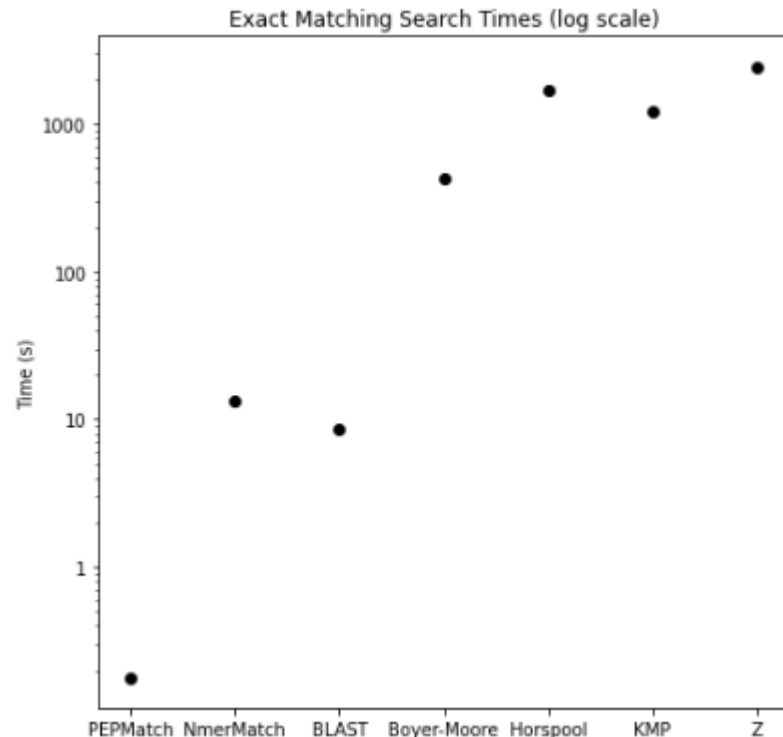
# PEPMatch: Best Match Search



# Benchmarking: Exact Matching Search

Task: Find 1,000 9-mers in the human proteome.

Algorithm/Tool	Proteome Preprocessing Time (s)	Query Preprocessing Time (s)	Searching Time (s)	Total Time (s)	Accuracy (%)
PEPMatch	36.6	N/A	0.18	36.8	100
NmerMatch	28.7	0.003	13.2	41.9	100
BLAST	0.53	N/A	8.48	9.01	99.4
Horspool	N/A	N/A	423.7	423.7	100
Boyer-Moore	N/A	N/A	1700.6	1700.6	100
Knuth-Morris-Pratt (KMP)	N/A	N/A	1204.2	1204.2	100
Z	N/A	N/A	2448.4	2448.4	100



PEPMatch outperformed BLAST and other string searching algorithms.

\* excluded DIAMOND and MMseqs2 due

# Benchmarking: Searching with Substitutions (Mismatches)

Task: Find 628 peptides of various lengths (8-15) in all betacoronavirus proteomes with 2 substitutions or less.

Tool	Proteome Preprocessing Time (s)	Query Preprocessing Time (s)	Searching Time (s)	Total Time (s)	Accuracy (%)
PEPMatch	40.7	N/A	51.9	92.6	100
NmerMatch	280.7	0.003	25.2	305.9	100
BLAST	0.61	N/A	156.2	156.8	73.4
DIAMOND	0.25	N/A	6.38	6.63	6.34
MMseqs2	3.76	N/A	1.14	4.90	7.34

PEPMatch and NmerMatch outperformed BLAST in speed and accuracy. They outperformed DIAMOND and MMseqs2 in accuracy by a lot. NmerMatch takes longer to preprocess.

# Benchmarking: Best Match Search

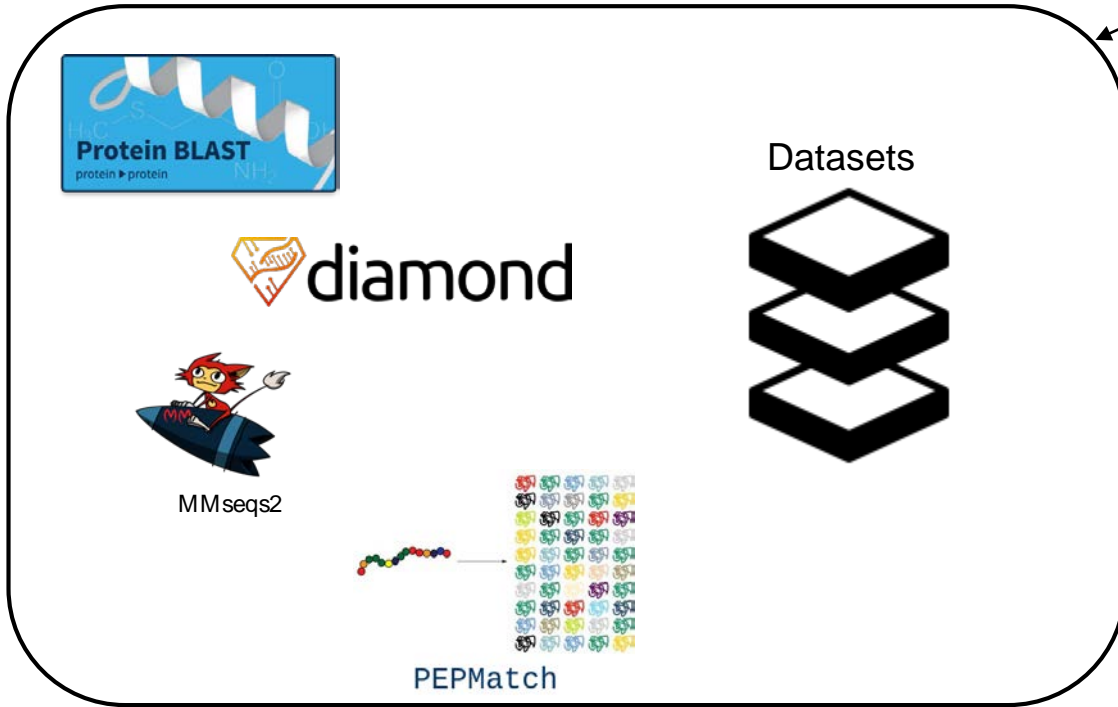
Task: Find the best match of 111 15-mers in the human proteome. Some matches may have 7-8 mismatches.

Tool	Proteome Preprocessing Time (s)	Query Preprocessing Time (s)	Searching Time (s)	Total Time (s)	Accuracy (%)
PEPMatch	45.2	N/A	155.4	200.6	100
NmerMatch	52.02	0.001	356.1	410.1	100
BLAST	0.56	N/A	18.21	18.77	24.3
DIAMOND	0.18	N/A	4.24	4.42	12.4
MMseqs2	2.98	N/A	0.64	3.62	15.0

PEPMatch outperformed NmerMatch in speed and BLAST, DIAMOND and MMseqs2 in accuracy.

# Benchmarking Framework

Your Tool  
(any language)



Expected output, time logging and accuracy.

<https://github.com/IEDB/PEPMatch>

# Using PEPMatch

## Command Line

```
[dan@foxp3 ~]$ pip install pepmatch  
> pepmatch-preprocess -p 9606.fasta  
> -k 5 -f sql
```

```
[dan@foxp3 ~]$ pepmatch-match -q peptides.fasta  
> -p 9606 -m 0 -k 5
```

## Python Module

```
1 #!/usr/bin/env python3  
2  
3 from pepmatch import Preprocessor, Matcher  
4  
5 Preprocessor('human_proteome.fasta', k=5, preprocess_format='sql').preprocess()  
6 Matcher('peptides.fasta', 'human_proteome', k=5, max_mismatches=0).match()
```

**Soon to be on  
the IEDB!**



IMMUNE EPITOPE DATABASE  
AND ANALYSIS RESOURCE



Thank you  
for listening!