



# Petasearch: Fast, approximate comparison of huge sequence datasets

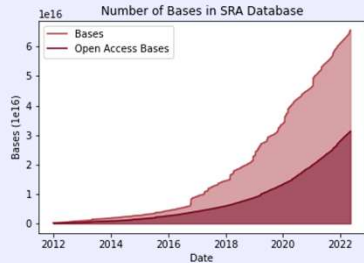
Minghang Li<sup>\*,1</sup>, Milot Mirdita<sup>\*,2</sup>, Jonas Hügél<sup>3</sup>, Johannes Söding<sup>2</sup> and Martin Steinegger<sup>1,4,5</sup>

<sup>\*</sup>Contributed equally, <sup>1</sup>School of Biological Sciences, Seoul National University, Seoul, South Korea, <sup>2</sup>Quantitative and Computational Biology, Max Planck Institute for Multidisciplinary Sciences, Göttingen, Germany, <sup>3</sup>Department of Medical Informatics, The University Medical Center Göttingen, Göttingen, Germany, <sup>4</sup>Institute of Molecular Biology and Genetics, Seoul National University, Seoul, South Korea, <sup>5</sup>Artificial Intelligence Institute, Seoul National University, Seoul, South Korea,



## Abstract

The Sequence Read Archive (SRA) currently holds over 60 petabases and is growing rapidly in size.



This vast amount of data represents a treasure trove for medicine and biotechnology. Bloom-filter and sketching based approaches to find k-mer seeds were proposed to accelerate searches, however they offer only limited sensitivity.

We developed **petasearch** to enable fast and sensitive searching for proteins extracted from huge databases. Its algorithm contains three stages:

1. We pre-process the database sequences to extract k-mers, sort them and store them in a highly compressed k-mer index.
2. We extract query k-mers, add similar k-mers and find matches between query and database k-mers. We remove matches with only a single shared k-mer to reject non-homologous sequences early. To maximize throughput, we exploit the caching and prefetch infrastructure of modern CPUs, advanced Linux IO techniques, as well as the enormous read bandwidth of NVMe-SSDs.
3. We compute SIMD-accelerated banded Smith-Waterman alignments between sequences of high-scoring k-mer matches.

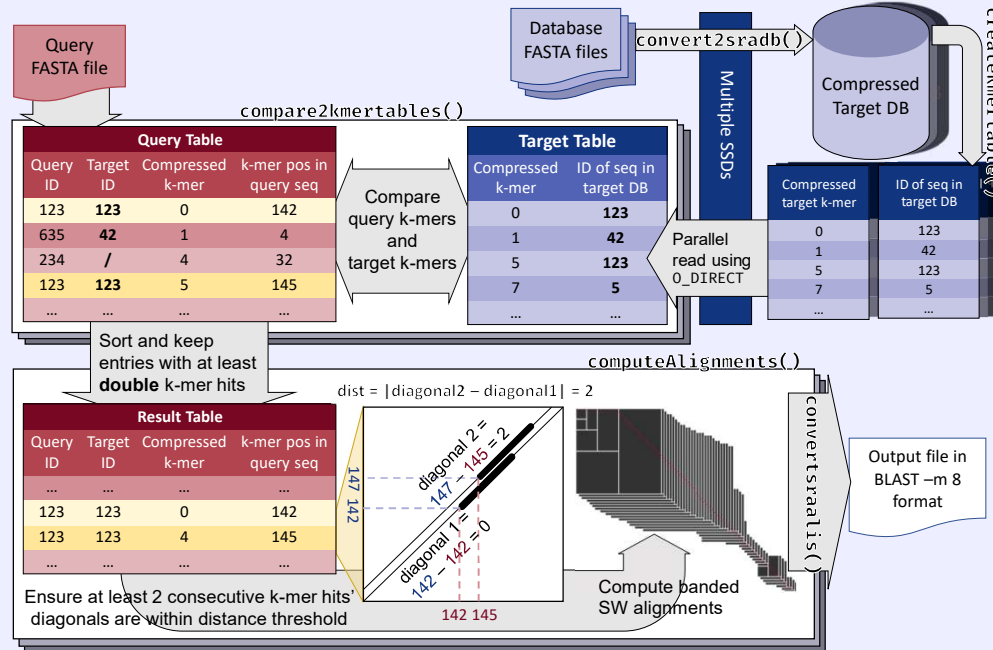
We show that **petasearch** is **190x** faster while querying a **9.3TB** dataset on **21** NVMe-SSDs. At much accelerated search speeds, **petasearch** matches state-of-the-art algorithms on sensitivity down to sequence identities of **60%**. On a SCOP25 benchmark we show that **petasearch**'s profile search detects sequence homology down to at least **40%** sequence identity.

## Code Availability

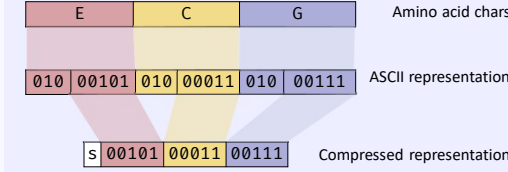


**petasearch** is a GPLV3-licensed open-source software available at <https://github.com/steineggerlab/petasearch>

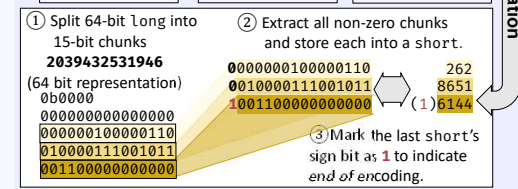
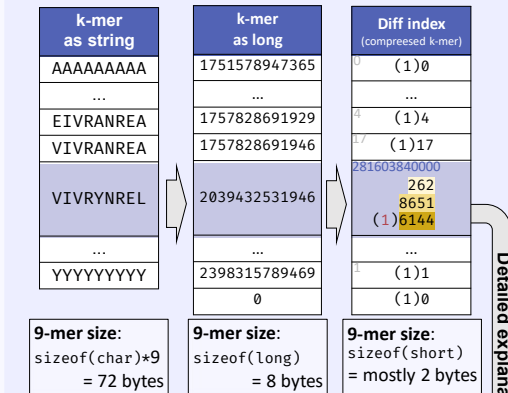
## Architecture Overview



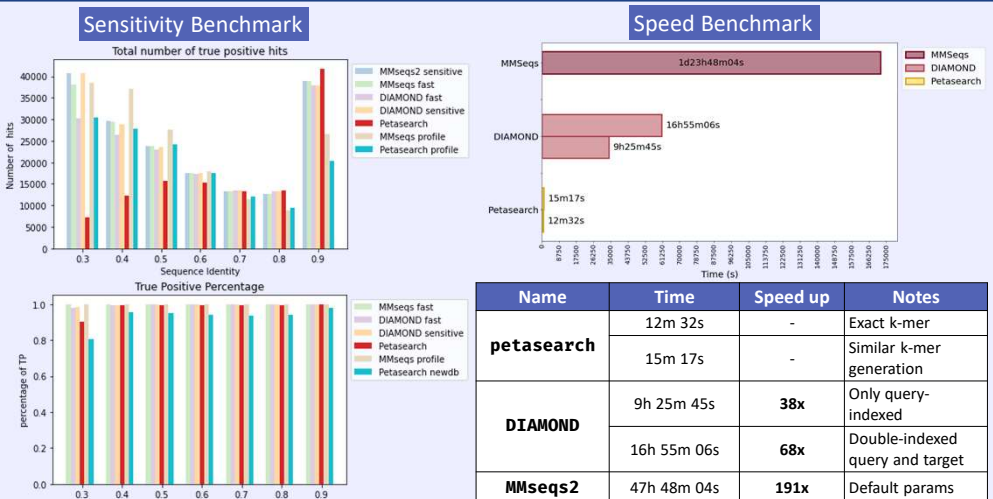
## Sequence Compression



## K-mer Index Compression



## Sensitivity and Speed Benchmarks



## Conclusion

- **petasearch** is a lightning fast searching algorithm with comparable sensitivity at higher sequence identity.
- **petasearch** will enable researchers to exploit the vast amount of evolutionary information available in current and upcoming databases.

## Acknowledgement

This work is supported by XXX research fund program. XXXXXXXX