



Extreme-scale Protein Similarity Search PASTIS (ExaBiome/ExaGraph) GPUs for Science, NERSC

Oct 25, 2022 Oguz Selvitopi (AMCR, LBNL)



Protein similarity search & clustering

▷ Comparative genomics

Functional or taxonomic contents of collected samples

Similarity search between two sets of protein sequences

- Functional annotation
- Gene localization
- Studying protein evolution

Common use case

- Protein similarity networks
- Detection of protein families



Motivation and Goals

- ▷ Distributed many-against-many protein similarity search
 - Parallel protein similarity search
 - Clustering
 - Metagenomics
 - Days, weeks, or even months

Existing software: MMseqs2, LAST, DIAMOND, BLASTP

• Optimized for shared-memory

PASTIS: Protein Alignment via Sparse Matrices for SEARCH
HipMCL: High Performance Markov Clustering for CLUSTERING 3



PASTIS Workflow



- \triangleright Fixed size k-mers to discover overlapping sequences
 - Exact, substitute variants
- Large-scale parallelism via CombBLAS (distributed sparse matrix library)
- Alignment via external libraries (on-node)
- > **Optimizations** for load balancing, memory management, overlapping, etc₄



Sparse matrices in PASTIS





PASTIS: Parallelism

- Scalable distributed-memory implementation of the components
 - Hybrid MPI-OpenMP
- ▷ Libraries used in PASTIS
 - CombBLAS
 - SeqAn/ADEPT
- PASTIS itself also makes use of MPI-OpenMP





Challenges of similarity search on huge datasets

Computational patterns

- Search operations on sequences
- Pairwise alignments
 - Edit distance computations

Memory requirements

- Many-against-many search
- Example: 100 million sequences
 - Candidate pairs ~ trillions (10¹²)
 - Alignments ~ hundreds of billions (10¹¹)
 - Similar pairs ~ billions (10⁹)



Techniques

> **Blocked (incremental) formation** of the similarity graph

Blocked 2D Sparse SUMMA

▷ Load balancing

- Index-based
- Triangularity-based
- Interleaved with blocked formation

▷ Pre-blocking

- GPU accelerators
- Utilize all resources on nodes simultaneously



Blocked 2D Sparse SUMMA



Advantages

- Bounds the maximum memory utilization
- Enables in-memory search for huge datasets
- Opens up the path for several further optimizations

Disadvantages

- Increases time compared to doing all at once
- Increased communication



Load balancing

- > **Symmetricity** of the overlap matrix
 - Can avoid alignments + perhaps sparse matrix computations
 - $\circ \text{ Important } \rightarrow \text{ Alignments are expensive}$
- How to achieve load balancing in blocked overlap detection?
- Approach #1: Symmetricity based on triangularity
- Approach #2: Symmetricity based on indices



half of the processes would stay idle if used in blocked multiplication



Load balancing comparison



Index-based scheme has better load balance

Triangularity-based method can save from computations



Pre-blocking

- Incremental formation of the similarity graph
 - Alignments \rightarrow on GPUs
 - Sparse computations & other \rightarrow on CPUs
- \triangleright CPUs stay idle
 - Can go ahead and prepare the next batch(es)
- ➢ Hide the overhead of sparse computations
 - Distributed sparse computations
 - Avoid collective communication + memorybound low-intensity computations

Runtime (sec)

#blocks	Plain	Pre-blocking	Total%
10	1555	1090	70
20	1606	1123	70
30	1659	1163	70
40	1724	1203	70
50	1774	1245	70

Overlap efficiency: >95%



Strong scaling & weak scaling





Perlmutter porting

- PASTIS was already running on Summit
 - Smooth porting

▷ Remarks

- PASTIS needs and uses both CPU and GPU resources
- PASTIS benefited from a faster CPU as well as faster GPUs on a node on Perlmutter

Small-scale performance 64 nodes





Protein similarity search at scale

⊳ <u>Summit</u>

- o 4600+ nodes
- o 42 CPU cores (512 GB)
- 6 GPUs (V100) (16 GB)

Dataset size

- o 313 million (April)
- o 405 million (July)

▷ Test scale

- o 2025 nodes
- o 3364 nodes

- ⊳ <u>Perlmutter</u>
 - 1500+ nodes
 - 0 64 CPU cores (256 GB)
 - 4 GPUs (A100) (40 GB)

Dataset size

- o 157 million
- o 200 million (July)
- ▷ Test scale
 - o 1024 nodes



Summit runs

Alignments per second: 320 million \rightarrow 691 million Cell updates per second: 143.9 TCUPs → 176.3 TCUPs

Problem size: ~1.67x Run size: ~1.66x

- ▷ 2025 nodes ▷ 313 million sequences
- **Parameters**
 - Blocking factor: 14 x 14 0
 - Load balancing: triangularity 0
 - Pre-blocking 0

▷ Results

- Discovered candidates: 53T 0
- Performed alignments: **4.5T** 0
- Similar pairs: **215B** 0
- Output size: **5.4 TB** 0
- **3.89** hours 0

- ▷ 405 million sequ^r
- **Parameters**

▷ 3364 nodes

- Blocking ' 0
- , of a , alarity Load ¹ 0 P۲ 0

cir

- $\triangleright \mathbf{R}$ 0 ⊿ candidates: 96T
 - Perfo. ied alignments: 8.6T 0
 - Similar pairs: **1.1T** 0
 - Output size: 27 TB 0
 - 3.44 hours 0



Perlmutter runs

▷ 1024 nodes

▷ 200 million sequences

▷ Parameters

- Blocking factor: 20 x 20
- Load balancing: triangularity
- Pre-blocking

▷ Results

- Discovered candidates: 18T
- Performed alignments: 1.6T
- Similar pairs: 76B
- Output size: 1.5 2TB
- 2.48 hours
- Excluding IO: 2.01 hours

	Alignments per sec (per node)		
	Summit	Perlmutter	
Overall	205294	173325	
Excluding IO	218809	213844	

	Summit	Perlmutter
Sparse (candidate disc.)	7427	5539
Alignment	.37×10168	1.16x ₆₄₀₂
IO%	5%	19%



Problem size comparison

- \triangleright Tools with distributed memory search option
 - MMSeqs2, DIAMOND
- ▷ DIAMOND (reported in 2021)
 - 281 million sequences vs. 39 million sequences
 - 520 nodes (Cobra at Max Planck Society)
 - ~23 billion alignments
 - 5.42 hours (very sensitive) 17.77 hours (ultra sensitive)

▷ Alignment rate

• **1.2 million** alignments per sec vs. **690.6 million** alignments per sec

15.0x

• Per node: 2.3k vs. 205k

▷ Search space

• 281 x 39 x 10¹² vs. 405 x 405 x 10¹²

575.5x



Conclusions

Compute-intensive and huge memory footprint

- Accelerators for alignments
- Algorithmic techniques for staying in-memory
 - No intermediate IO
- Optimizations to take advantage of heterogeneous node architecture
 - Use all resources on the node
- **> HPC for bioinformatics**
 - Cutting time-to-solution from days/weeks to hours



Attend our talk at SC22!

Thank you QUESTIONS



