

### BENCHMARKING OF I/O ACTIVITIES IN HPC APPLICATIONS WITH STRACE INSPECTOR

34TH POP WEBINAR | JUNE 2, 2025 | ARAVIND SANKARAN



Mitglied der Helmholtz-Gemeinschaft



- About me
- The Story Application benchmarks.
- Application Benchmarking with Directly-Follows-Graph (DFG).
- Synthesis of large amounts of I/O related system call traces.
- Final words.



### **The Linear Algebra Mapping Problem**

The mapping of high level expressions to an optimized sequence of basic linear algebra operations.





#### **The Linear Algebra Mapping Problem**





### **The Linear Algebra Mapping Problem**

The mapping of high level expressions to an optimized sequence of basic linear algebra operations.



Kernels are highly optimized.



### The Linear Algebra Mapping Problem

The mapping of high level expressions to an optimized sequence of basic linear algebra operations.





# THE LAAB STORY

### The Linear Algebra Aware Benchmarks

Sample CB report

- Developed benchmarks to evaluate the linear algebra awareness of the build of popular ML framework (v2022).
- Designed to run in a continuous benchmarking setup to automatically generate reports (v2025 – yet to release).
- Source code: <u>https://github.com/HPAC/LAAB-</u>
  <u>Python</u>

#### LAAB-Python | LA Awareness-CPU | PyTorch/2.1.2-foss-2023a | HPC2N\_x86\_64

This report evaluates whether the software build performs operations equivalent to those of optimized math libraries (e.g., OpenBLAS, MKL), and whether it leverages linear algebra techniques to accelerate CPU computations. Unless stated otherwise, all benchmarks use matrices of size  $3000 \times 3000$  and are executed on a single CPU core of AMD EPYC 7413 24-Core Processor.

#### 1) PyTorch vs. BLAS for matrix multiplication:

	Call	time (s)
$A^T B$	t(A)@B	0.5094 🔽
"	<pre>linalg.matmul(t(A),B)</pre>	0.5072 🔽
Reference	sgemm	0.4942

#### 2) Elimination of common sub-expression:

 $E = A^T B + A^T B.$ 

Expr	Call	time (s)
E	t(A)@B + t(A)@B	0.5251 🔽
Reference	2*(t(A)@B)	0.5271

#### $E_1 = (A^T B)^T (A^T B)$ and $E_2 = (A^T B)^T A^T B$ .

Expr	Call	time (s)		
$E_1$	t(t(A)@B)@(t(A)@B)	1.0191 🔽		
$E_2$	t(t(A)@B)@t(A)@B	1.5224 🗙		
Reference	S=t(A)@B; t(S)@S	1.0188		



JÜLICH SUPERCOMPUTING CENTRE

8

## THE LAAB STORY

### The Linear Algebra Aware Benchmarks

#### Sample CB report

#### 3) Choosing the optimal order to evaluate a matrix chain:

Given *m* matrices of suitable sizes, the product  $M = A_1A_2...A_n$  is known as a matrix chain. Because of associativity of matrix product, matrix chain can be computed in many different ways, each identified by a specific paranthesization. The paranthesization that evaluates *M* with the least number of floating point operations (FLOPs) is considered optimal. PyTorch provides the method torch.linalg.multi\_dot specifically for evaluation of matrix chains.

#### a) Right to left:

The input matrix chain is  $H^T H x$ , where x is a vector of lenth 3000. The reference implementation, evaluating from right-to-left - i.e.,  $H^T (H x)$ , avoids the expensive  $O(n^3)$  matrix product, and has a complexity of  $O(n^2)$ .

Expr	Call	time (s)
$H^T H x$	t(H)@H@x	0.5100 🗙
	<pre>linalg.multi_dot([t(H), H, x])</pre>	0.0064 🔽
Reference	t(H)@(H@x)	0.0054

For more info, refer: A. Sankaran, N. A. Alashti, C. Psarras and P. Bientinesi, "Benchmarking the Linear Algebra Awareness of TensorFlow and PyTorch," 2022 IEEE International Parallel and Distributed Processing Symposium Workshops (IPDPSW), Lyon, France, 2022, pp. 924-933, doi: 10.1109/IPDPSW55747.2022.00150.



- Developed benchmarks to evaluate the linear algebra awareness of the build of popular ML framework (v2022).
- Designed to run in a continuous benchmarking setup to automatically generate reports (v2025 – yet to release).
- Source code: <u>https://github.com/HPAC/LAAB-</u> Python

Mitglied der Helmholtz-Gemeinschaft

9

# **BENCHMARKING WITH DFG**

### The algorithms use case

Identification of root causes of performance differences among algorithm variants.



Variant	Kernel sequence
$alg_0$	potrf, trsm, trsv, syrk, gemv, potrf, trsv, trsv
$alg_1$	potrf, trsv, trsm, syrk, potrf, gemv, trsv, trsv
$alg_2$	potrf, trsm, trsv, syrk, gemv, qr, gemv, trsv
$alg_3$	potrf, trsv, trsm, gemm, potrf, gemv, trsv, trsv
$alg_4$	potrf, trsm, trsv, gemm, gemv, potrf, trsv, trsv
$alg_5$	potrf, trsm, trsv, gemm, potrf, gemv, trsv, trsv
$alg_6$	transpose, potrf, trsm, trsv, syrk, potrf, trsv, trsm, gemv, trsv
$alg_7$	transpose, potrf, trsm, syrk, potrf, trsv, trsv, trsm, gemv, trsv
$alg_8$	transpose, potrf, trsm, syrk, potrf, trsm, trsm, trsv, trsv, gemv
alg <sub>9</sub>	transpose, potrf, trsm, syrk, potrf, trsv, trsv, trsm, trsm, gemv

Algorithmic variants of the generalized least square problem.



## **BENCHMARKING WITH DFG**

### The algorithms use case





Variant	Kernel sequence
alg <sub>0</sub>	potrf, trsm, trsv, syrk, gemv, potrf, trsv, trsv
alg <sub>1</sub>	potrf, trsv, trsm, syrk, potrf, gemv, trsv, trsv
alg <sub>2</sub>	potrf, trsm, trsv, syrk, gemv, qr, gemv, trsv
alg <sub>3</sub>	potrf, trsv, trsm, gemm, potrf, gemv, trsv, trsv
alg <sub>4</sub>	potrf, trsm, trsv, gemm, gemv, potrf, trsv, trsv
alg <sub>5</sub>	potrf, trsm, trsv, gemm, potrf, gemv, trsv, trsv
alg <sub>6</sub>	transpose, potrf, trsm, trsv, syrk, potrf, trsv, trsm, gemv, trsv
alg <sub>7</sub>	transpose, potrf, trsm, syrk, potrf, trsv, trsv, trsm, gemv, trsv
alg <sub>8</sub>	transpose, potrf, trsm, syrk, potrf, trsm, trsm, trsv, trsv, gemv
alg <sub>9</sub>	transpose, potrf, trsm, syrk, potrf, trsv, trsv, trsm, trsm, gemv

**Ref:** Sankaran, A., Karlsson, L. & Bientinesi, P. Ranking with ties based on noisy performance data. *Int J Data Sci Anal* (2025). https://doi.org/10.1007/s41060-025-00722-1

![](_page_9_Picture_6.jpeg)

# AND NOW,

The picture is incomplete without I/O information.

### The Problem:

- How to include I/O in LAAB reports?
- Investigate the use of DFG method to *synthesize* the large amounts of I/O performance data.

![](_page_10_Picture_5.jpeg)

### WHICH I/O DATA TO USE?

![](_page_11_Figure_1.jpeg)

![](_page_11_Picture_2.jpeg)

## WHICH I/O DATA TO USE?

![](_page_12_Figure_1.jpeg)

![](_page_12_Picture_2.jpeg)

## **TRACING WITH STRACE**

### Basic Usage:

strace [COMMAND]

#### **Example:**

### \$ strace ls

![](_page_13_Picture_6.jpeg)

## **TRACING WITH STRACE**

Example:

#### strace -f -tt -T -y -e read, write ls

![](_page_14_Figure_3.jpeg)

![](_page_14_Picture_4.jpeg)

## **TRACING WITH STRACE**

![](_page_15_Figure_1.jpeg)

Trace file: host\_9042.st

![](_page_15_Picture_3.jpeg)

## **THE CHALLENGE**

	pid	call	start	duration	bytes	fs	case	end
10	22085	read	1900-01-01 18:54:16.207116	0.000015	832	/p/software/fs/jusuf/stages/2024/software/HDF5	st.jsfc134.22051.log	1900-01-01 18:54:16.207131
33	22085	read	1900-01-01 18:54:16.211619	0.000017	832	/p/software/fs/jusuf/stages/2024/software/psmp	st.jsfc134.22051.log	1900-01-01 18:54:16.211636
221	22085	read	1900-01-01 18:54:16.246555	0.000008	832	/usr/lib64/libc.so.6	st.jsfc134.22051.log	1900-01-01 18:54:16.246563
231	22085	read	1900-01-01 18:54:16.247709	0.000015	832	/p/software/fs/jusuf/stages/2024/software/IME/	st.jsfc134.22051.log	1900-01-01 18:54:16.247724
235	22085	read	1900-01-01 18:54:16.248466	0.000016	832	/p/software/fs/jusuf/stages/2024/software/Szip	st.jsfc134.22051.log	1900-01-01 18:54:16.248482

**Challenge:** How do you extract useful information from large amounts of information in the system call traces?

![](_page_16_Picture_3.jpeg)

Idea:

- Map each row to a string that helps answer your question. We call this string "Activity".
- Apply grouping based on activities and compute statistics.
- Identify the directly-follows relation between the activities to build a Directly-Follows-Graph (DFG).

pid	call	start	duration	bytes	fs	Activity
22085	read	1900-01-01 18:54:16.207116	0.000015	832	/p/software/fs/jusuf/stages/2024/software/HDF5 :	→ read+/p/software
22085	read	1900-01-01 18:54:16.211619	0.000017	832	/p/software/fs/jusuf/stages/2024/software/psmp	→ read+/p/software
22085	read	1900-01-01 18:54:16.246555	0.000008	832	/usr/lib64/libc.so.6 :	→ read+/usr/lib64

![](_page_17_Picture_6.jpeg)

pid	call	start	duration	bytes	fs	<u>Activity</u>
22085	read	1900-01-01 18:54:16.207116	0.000015	832	/p/software/fs/jusuf/stages/2024/software/HDF5 :	read+/p/software
22085	read	1900-01-01 18:54:16.211619	0.000017	832	/p/software/fs/jusuf/stages/2024/software/psmp	read+/p/software
22085	read	1900-01-01 18:54:16.246555	0.000008	832	/usr/lib64/libc.so.6 :	→ read+/usr/lib64
				Ť		

- This data is can be formalized as an event-log in Process Mining.
- Process mining introduces scalable techniques to translate event logs into different types of dependency graphs, including Directly-Follows Graph.
- Ref: W. M. P. Van Der Aalst, "Foundations of process discovery," in Process Mining Handbook, DOI: https://doi.org/10.1007/978-3-031-08848-3\_2

![](_page_18_Picture_5.jpeg)

![](_page_19_Figure_1.jpeg)

![](_page_19_Picture_2.jpeg)

![](_page_20_Figure_1.jpeg)

![](_page_20_Figure_2.jpeg)

![](_page_20_Picture_3.jpeg)

![](_page_21_Figure_1.jpeg)

![](_page_21_Figure_2.jpeg)

![](_page_21_Picture_3.jpeg)

![](_page_22_Figure_1.jpeg)

![](_page_23_Picture_0.jpeg)

![](_page_23_Picture_1.jpeg)

Mitglied der Helmholtz-Gemeinschaft

### **SNAPSHOT OF LAAB REPORT WITH I/O**

![](_page_24_Figure_1.jpeg)

![](_page_24_Picture_2.jpeg)

## THE IMPLEMENTATION: STRACE INSPECTOR

- One log file is generated per node used.
- The parsing of log files and mapping to activities are done in parallel (*in v1.1.0*).
- The construction of DFG requires one pass through the activity log O(n), where n is the number of lines in the filtered event log after mapping.
- The complexity of rendering the DFG is O(m2), where m is the number activities. For all intents and purposes, m should be small.

#### Source code: <u>https://gitlab.jsc.fz-juelich.de/st-inspector/st\_inspector</u>

**Ref:** A. Sankaran, I. Zhukov, W. Frings and P. Bientinesi, "Inspection of I/O Operations from System Call Traces using Directly-Follows-Graph," *SC24-W: Workshops of the International Conference for High Performance Computing, Networking, Storage and Analysis*, Atlanta, GA, USA, 2024, pp. 1562-1575, doi: 10.1109/SCW63240.2024.00196

![](_page_25_Picture_7.jpeg)

### **Acknowledgements**

### Thank you for you attention.

- This research was financially supported by the Juelich Supercomputing Center at Forschungszentrum Juelich and the BMBF project 01-1H1-6013 AP6 NRW Anwenderunterstutzung SiVeGCS.
- Additionally, supervision support from RWTH Aachen University through the DFG project IRTG-2379 is gratefully acknowledged.
- Compute time on JUWELS Booster at JSC and Kebnekaise at HPC2N (Sweden) is gratefully acknowledged.

![](_page_26_Picture_6.jpeg)

Modern

Inverse

Problems

M!P

![](_page_26_Picture_7.jpeg)

![](_page_26_Picture_8.jpeg)