# SCALABLE DISTRIBUTED COMPUTATION OF CRITICAL PATH

Pierre-Frédérick DENYS
Thursday 8 December 2022

# Agenda

- Introduction

- Proposed algorithm

- Benchmarks

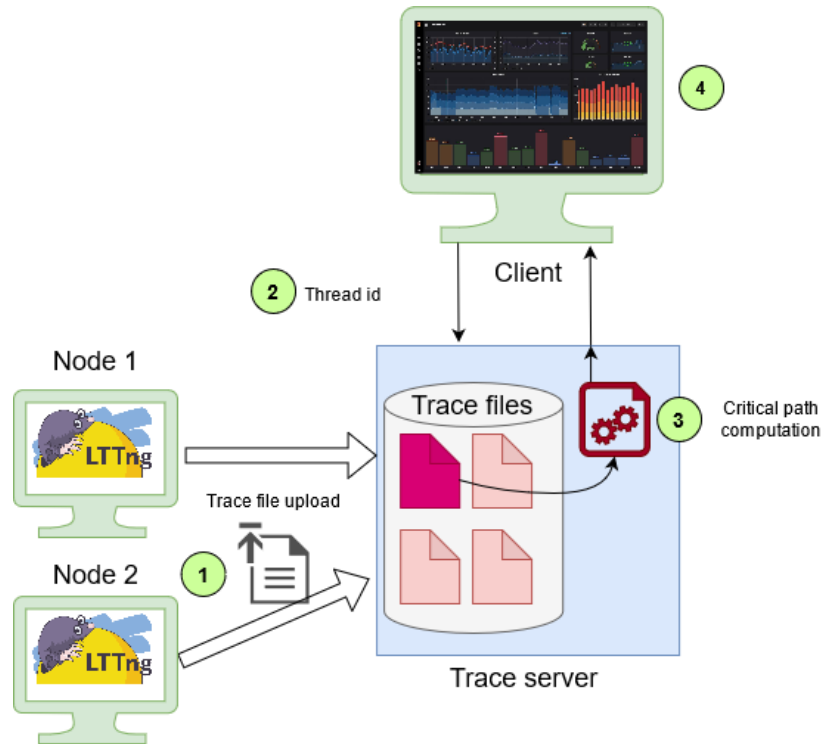- Future work and usecases

- Conclusion

# Critical path usage

- Need for large distributed systems tracing
    - HPC systems
    - MPI clusters
    - Kubernetes and container clusters
- Transfer of trace files on analysis node was mandatory
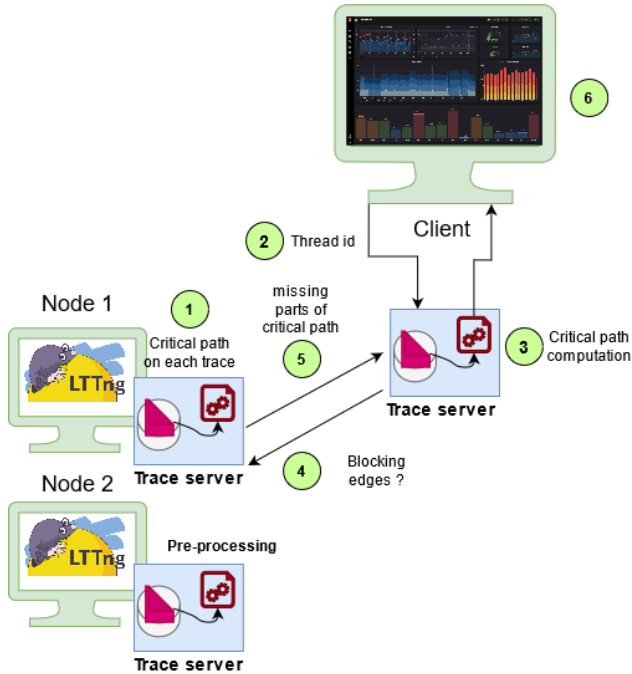- Critical path distributed computation was not optimized

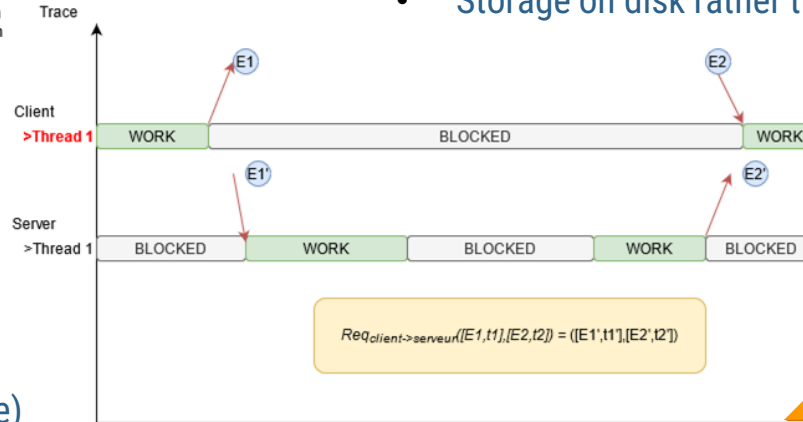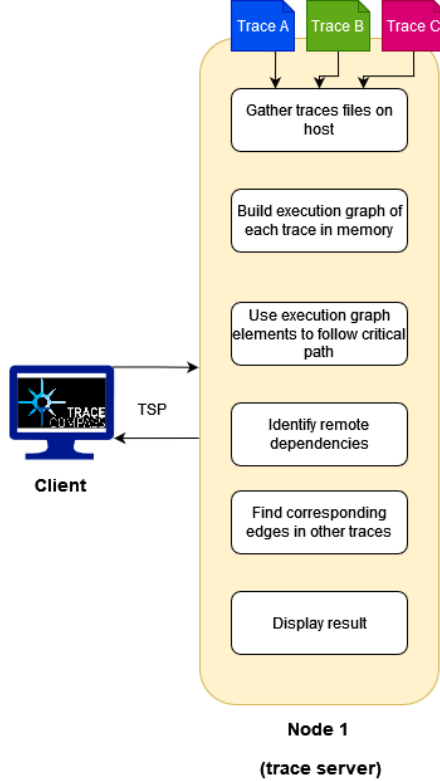# " Critical path computation evolution

-

- Pre-processing of critical path on each node

- On client request, process the critical path of the trace, and ask only the missing parts of the path to other nodes

- Distributed processing, suitable for large number of nodes, less network load

- Storage on disk rather than in memory*

(*related work done by Arnaud and Geneviève)



Node 1

Critical path on each trace

Trace server

Node 2

Pre-processing

Trace server

Thread id

missing parts of critical path

Client

Trace server

Critical path computation

Blocking edges ?

Trace

Client
>Thread 1   WORK   BLOCKED   WORK

E1   E2

Server
>Thread 1   BLOCKED   WORK   BLOCKED   WORK   BLOCKED

E1'   E2'

$Req_{client->serveur}([E1,t1],[E2,t2]) = ([E1',t1'],[E2',t2'])$
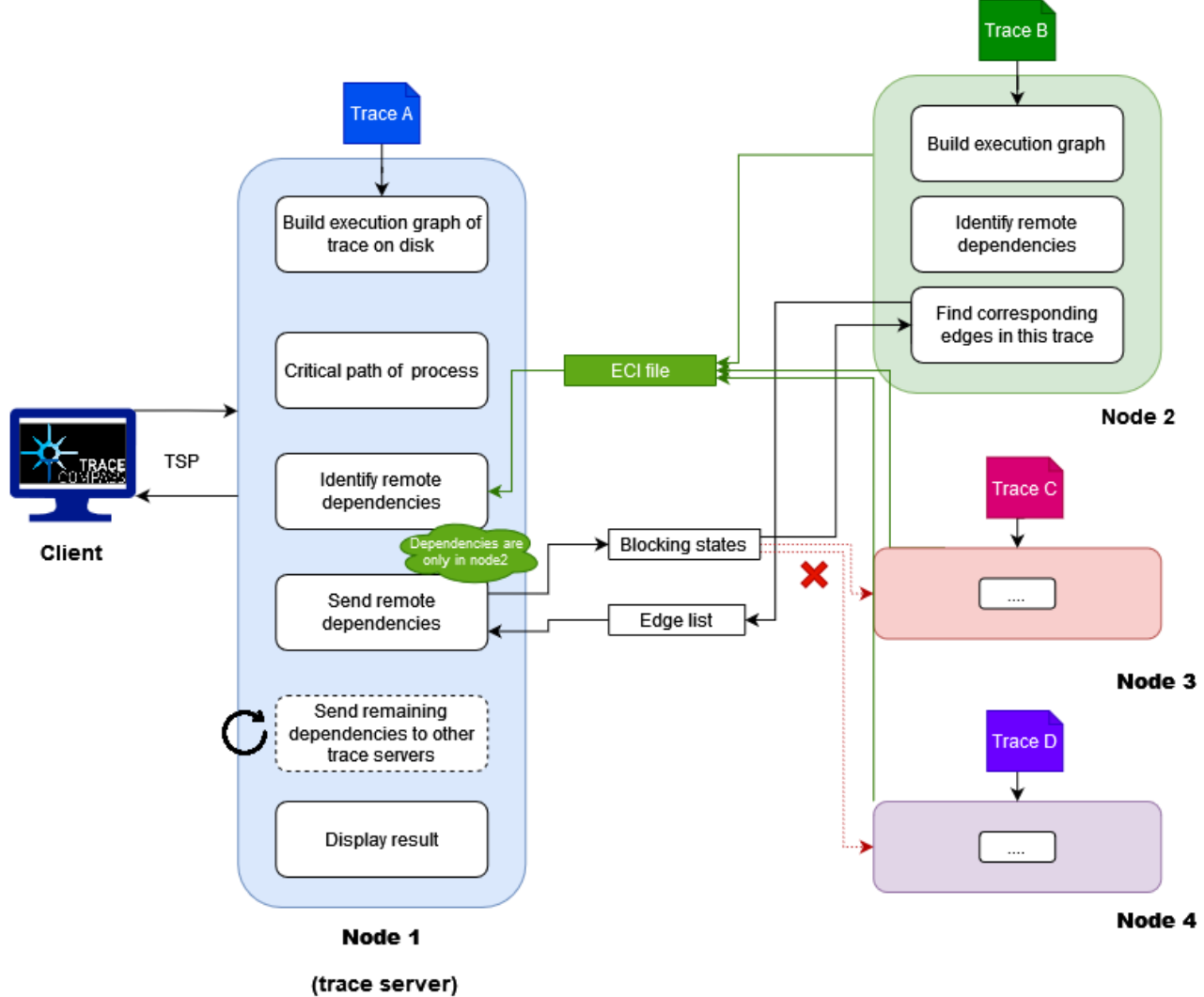
Time

AL1 vs AL2

Actual algorithm (AL1)          Distributed version (AL3)

# My work (AL 3) : External communication Index

- Improve algorithm to :

Introduce External communication Index (ECI) usage

- Index file exchanged after execution graph processing

Identify remote dependencies and location of remote trace

- Remove the need of broadcast communication on computing nodes to find remote dependencies

Pre-identification of remote dependencies for faster critical path processing

```
1  {
2    ...
3    {nodeId : "0xc0a81e02", traceFile : "8949844
       vdfd98", nodeLocation : "0xc0a81e19"},
4    {nodeId : "0xc0a81e04", traceFile : "8949844
       vdfd98", nodeLocation : "0xc0a81e19"}
5    ...
6  }
```

"

# Benchmarks

# SCP benchmark

| | File size | Trace file on (SD) each Computing node | Trace file on (SD) main node | Events Main | Events Dest | Events IRQ | Events Network |
|---|---|---|---|---|---|---|---|
| | 100 Mo | 28 Mo (1.86) | 61 Mo (1.68) | 1,026,089 | 2,205,019 | 9.10 % | 1.00 % |
| 1 | 1 Go | 181 Mo (1.92) | 277 Mo (1.83) | 6,757,000 | 9,981,329 | 8.90 % | 1.60 % |
| | 10 Go | 1.8 Go (0.33) | 3.8 Go (0.07) | 69,414.305 | 124,240,618 | 9.50 % | 1.50 % |
| | 100 Mo | 68 Mo (1.93) | 105 Mo (2.12) | 261,926 | 388,210 | 7.50 % | 1.00 % |
| 2 | 1 Go | 346 Mo (1.18) | 252 Mo (1.56) | 12,494,815 | 9,113,923 | 8.20 % | 1.20 % |
| | 10 Go | 3.2 Go (0.12) | 2.4 Go (0.56) | 118,495,717 | 86,933,604 | 9.20 % | 1.40 % |
| | 100 Mo | 71 Mo (1.58) | 75 Mo (1.96) | 2,644,924 | 2,793,396 | 8.40 % | 1.10 % |
| 3 | 1 Go | 443 Mo (1.12) | 266 Mo (1.36) | 15,984,498 | 9,604,204 | 9.20 % | 1.30 % |
| | 10 Go | 4.4 Go (0.18) | 2.4 Go (0.45) | 158,781,607 | 88,181,239 | 8.60 % | 1.70 % |
| | 100 Mo | 83 Mo (1.69) | 78 Mo (1.85) | 2,993,919 | 285,664 | 8.80 % | 1.80 % |
| 10 | 1 Go | 4.1 Go (0.09) | 262 Mo (0.66) | 43,285,500 | 9,565,111 | 7.20 % | 1.90 % |
| | 10 Go | 8.8 Go (0.06) | 6.8 Go (0.05) | 317,427,586 | 248,481,236 | 9.10 % | 1.10 % |
| | 100 Mo | 120 Mo (1.20) | 71 Mo (1.81) | 3,205,774 | 269,877 | 7.80 % | 1.40 % |
| 20 | 1 Go | 9.1 Go (0.08) | 271 Mo (1.76) | 334,953,487 | 9,854,156 | 9.40 % | 1.20 % |
| | 10 Go | 19.5 Go (0.05) | 5.2 Go (0.06) | 698,340,689 | 222,199,310.2 | 8.70 % | 1.90 % |
| | | | | | Average | 8.64 % | 1.41 % |

# SCP benchmark

| | File size | Mean Processing Time*(s) | | | Time overhead (SD) |
|---|---|---|---|---|---|
| | | AL1 (SD) | AL2 (SD) | AL3 (SD) | |
| 1 | 100Mo | 7 | 3 | 3.06 | 2.11% (0.10) |
| | 1Go | 43 | 15 | 15.30 | 2.02% (0.07) |
| | 10Go | 570 | 210 | 214,43 | 2.11% (0.05) |
| 2 | 100Mo | 23 | 5 | 5,11 | 2.25% (0.23) |
| | 1Go | 69 | 19 | 19.39 | 2.03% (0.12) |
| | 10Go | 632 | 235 | 240.66 | 2.41% (0.06) |
| 3 | 100Mo | 212 | 68 | 69.37 | 2.01% (0.18) |
| | 1Go | 687 | 210 | 214.56 | 2.17% (0.09) |
| | 10Go | 958 | 297 | 303.30 | 2.12% (0.06) |
| 10 | 100Mo | 685 | 80.60 | 79.02 | 2.03% (0.19) |
| | 1Go | 9890 | 356 | 364.58 | 2.41% (0.12) |
| | 10Go | NP | 1,916 | 1,957.19 | 2.15% (0.08) |
| 20 | 100Mo | 1360 | 75 | 76.61 | 2.14% (0.21) |
| | 1Go | NP | 38 | 38.78 | 2.06% (0.14) |
| | 10Go | NP | 3,780 | 3,867.70 | 2.32% (0.07) |

*Total CPU Time

# SCP benchmark

# MPI benchmark

- 20 nodes with i5-9400 at 2.9Ghz with six cores and 16Go of memory

# MPI benchmark



- **SOMA :** Offers Monte-Carlo Acceleration is a High-Performance Computing Monte-Carlo simulation for soft coarse-grained polymers. The variable load is the *number of simulated polymers*.
- **Tealeaf :** solves the linear heat conduction equation on a spatially decomposed regular grid. The variable load is the *number of cells*.
- **Minisweep:** is a Nuclear Engineering and radiation transport simulation. The variable load is the *grid cell size*.
- **SPH-EXA:** performs hydrodynamical and computational fluid dynamics simulations. The variable load is the *number of particles to the cube*.

| Benchmark | Parameter | Exec. Time (s) | Trace size | | Event Irq (%) | Event Network (%) |
|---|---|---|---|---|---|---|
| | | | Trace avg. per node | Total | | |
| Soma | 1400 | 6 | 7.5 Mo (3.52) | 158 Mo (1.21) | 12.3 | 1.08 |
| | 14000 | 21 | 29 Mo (1.82) | 613 Mo (1.38) | 12.9 | 1.14 |
| | 140000 | 261 | 256 Mo (1.42) | 5.4 Go (1.24) | 13.1 | 1.85 |
| | 1400000 | 2355 | 2.83 Go (1.02) | 63.8 Go (0.88) | 13.6 | 1.90 |
| Tealeaf | 512 | 61 | 73 Mo (2.12) | 1.5 Go (1.28) | 13.6 | 1.91 |
| | 1024 | 78 | 93 Mo (1.21) | 1.9 Go (1.02) | 12.6 | 1.92 |
| | 4096 | 411 | 494 Mo (1.52) | 10.45 Go (0.56) | 13.2 | 1.34 |
| | 32768 | 3120 | 3.94 Go (0.83) | 83.2 Go (0.08) | 12.9 | 1.41 |
| Minisweep | 32 | 75 | 90 Mo (2.04) | 1.9 Go (1.21) | 13.6 | 1.75 |
| | 64 | 150 | 182 Mo (1.84) | 3.85 Go (0.81) | 13.6 | 1.28 |
| | 128 | 315 | 379 Mo (1.62) | 20.3 Go (0.93) | 12.9 | 1.71 |
| | 768 | 1890 | 2.27 Go (1.13) | 48.2 Go (0.38) | 12.3 | 1.05 |
| SPH-EXA | 20 | 246 | 296 Mo (1.57) | 6.26 Go (1.09) | 13.6 | 1.00 |
| | 40 | 470 | 565 Mo (1.84) | 11.9 Go (1.28) | 12.8 | 1.02 |
| | 60 | 705 | 846 Mo (1.92) | 17.9 Go (0.52) | 13.7 | 1.13 |
| | 120 | 1410 | 1.69 Go (1.07) | 36 Go (0.31) | 13.0 | 1.54 |

| Benchmark | Parameter | Processing time* | | | Time overhead (%) (SD) |
|---|---|---|---|---|---|
| | | AL1 | AL2 | AL3 | |
| Soma | 1400 | 143 | 68 | 69 | 1.8 (0.12) |
| | 14000 | 501 | 203 | 206 | 1.6 (0.15) |
| | 140000 | 6223 | 1244 | 1267 | 1.9 (0.11) |
| | 1400000 | NA | 2807 | 2857 | 1.8 (0.09) |
| Tealeaf | 512 | 1455 | 715 | 729 | 1.9 (0.16) |
| | 1024 | 1860 | 695 | 711 | 2.3 (0.17) |
| | 4096 | 9800 | 1960 | 1999 | 2.0 (0.14) |
| | 32768 | NA | 3719 | 3787 | 1.8 (0.08) |
| Minisweep | 32 | 1788 | 894 | 912 | 2.0 (0.15) |
| | 64 | 3577 | 1192 | 1213 | 1.7 (0.16) |
| | 128 | 7511 | 2504 | 2555 | 2.1 (0.19) |
| | 768 | NA | 2253 | 2295 | 1.9 (0.12) |
| SPH-EXA | 20 | 5866 | 2704 | 2757 | 1.9 (0.06) |
| | 40 | 11207 | 3735 | 3825 | 2.4 (0.10) |
| | 60 | 16810 | 3362 | 3443 | 2.4 (0.14) |
| | 120 | NA | 4689 | 4779 | 1.9 (0.13) |

*Total CPU Time

# MPI benchmark

"

# Future work and usecases

# What remains to be done ?

- Remote time synchronisation of traces

- Better protocol for graph elements exchanges

- Automatic coordination between nodes

Target usecases :

- **MPI cluster :** follow a MPI task between computing nodes

- **Kubernetes cluster :** follow a request in a distributed web application

- **ZeroMQ communication :** follow a message exchange between several containers

"

# Conclusion

# Conclusion

- Improvement and benchmark of critical path Distributed computation

- Next step : Integration of Critical path in Trace Server Protocol (for Theia and Grafana viewers)

- Extend parallelisation to other kind of analysis

Thank you for listening !