

Walk over 8 years of research in HPC

- Memory management (tools)
- LHCb 40 Tb/s DAQ R&D

SÉBASTIEN VALAT – ATOS – 04 APRIL 2019

Plan

- ▶ Quick reminder: memory management performance impact
- ▶ MALT : A Malloc Tracer
- ▶ NUMAPROF : A NUMA memory profiler
- ▶ LHCb 40 Tb/s DAQ R&D effort for 2020 upgrade





Quick reminder memory management impact

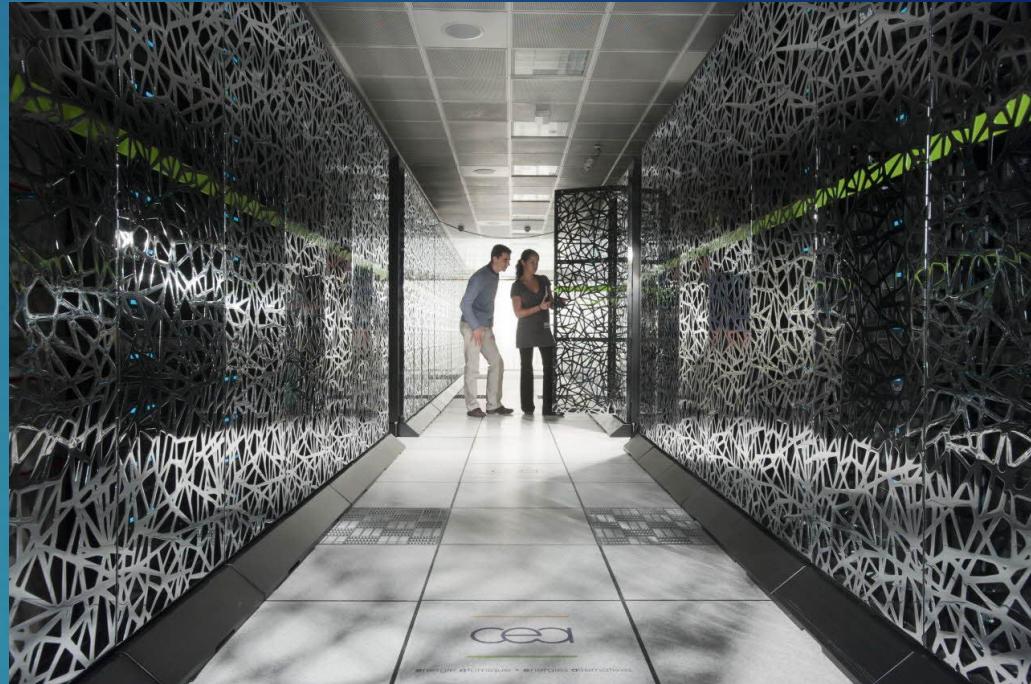
PHD. WORK

A little bit of bibliography

- ▶ **Interesting to read :**
 - ▶ What every programmer should know about memory (Ulrich Drepper)
<https://people.freebsd.org/~lstewart/articles/cpumemory.pdf>
- ▶ For all details from this presentation : look on my PhD. thesis :
 - ▶ PhD: <https://hal.archives-ouvertes.fr/tel-01253537>
 - ▶ PhD slides: <https://svvalat.github.io/docs/2014-07-phd-defense.pdf>
 - ▶ Webcast of similar talk (FR): <https://youtu.be/wp8cCTD7XDI>

Context : HPC

- ▶ **Supercomputers** for numerical simulations
- ▶ At CEA, **Tera 100** :
 - ▶ **6^e** from TOP 500 in 2010
 - ▶ **140 000** cores, **1.05** Pflops.
- ▶ **Massively parallel** machines
 - ▶ Now **million cores**
 - ▶ More and **more cores** !



<http://www.cea.fr/multimedia/Pages/galeries/defense/Tera-100.aspx>

Architecture	Proc.	Cores	Threads
Tera 100 thin nodes	4	32	32
Tera 100 large nodes	16	128	128
Intel KNL	1	64	256
PEZY-SC2	1	2048	16 384

Context : HPC

- ▶ **Memory** becomes a **critical resource**
- ▶ Growing impact on **performance**
- ▶ **Data movements** : speed gap CPU / RAM, **memory wall**.
- ▶ **Management** : now have to handle close to **TB** of memory
- ▶ Decreasing **memory per core**



<http://www.cea.fr/multimedia/Pages/galeries/defense/Tera-100.aspx>

Architecture

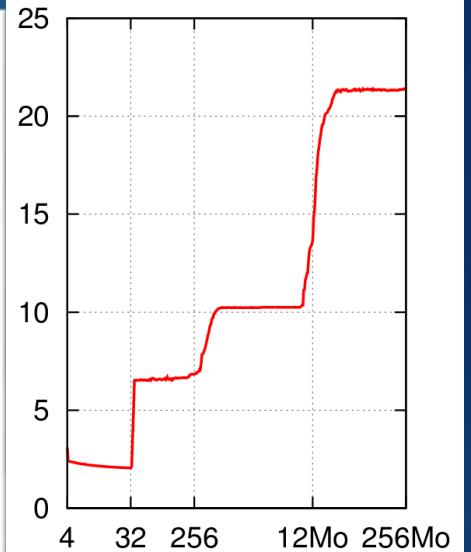
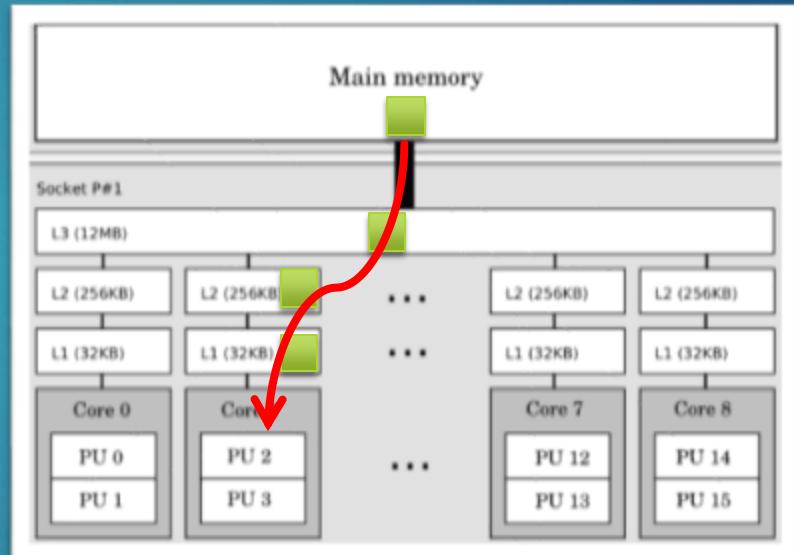
- ▶ Computer science: **operations & data**

- ▶ Multiple **memory levels**

- ▶ Hierarchical **caches**

- ▶ *Pre-fetcher*

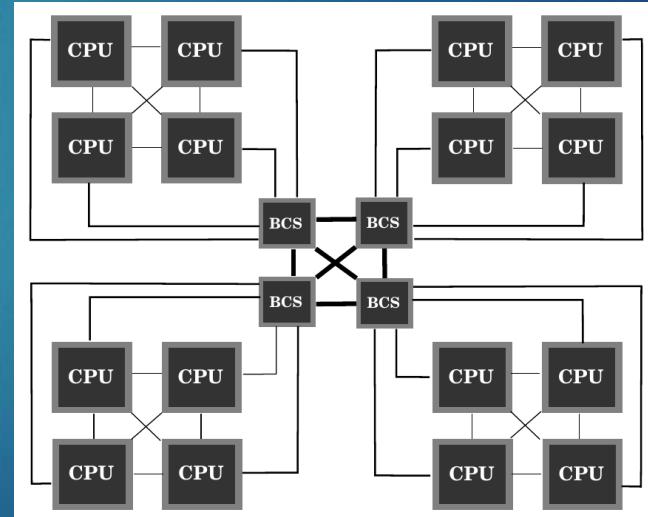
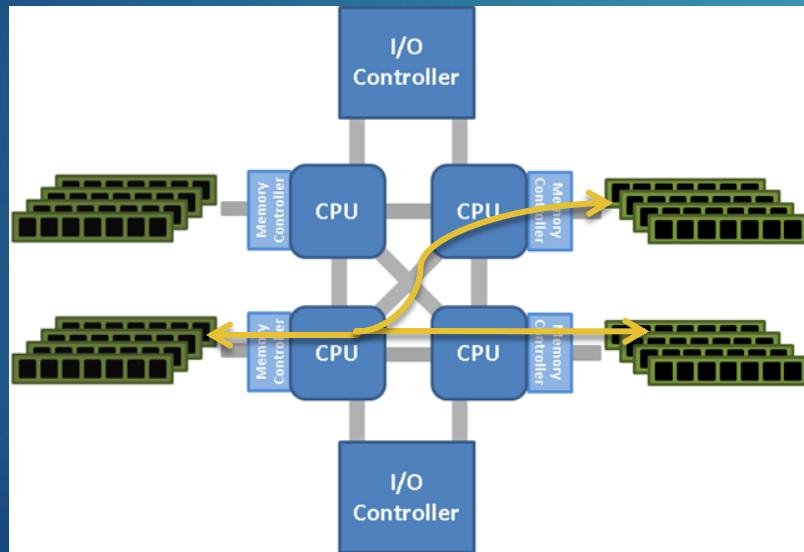
Processor : 8 cores



Architecture

| Slide 8 / 43

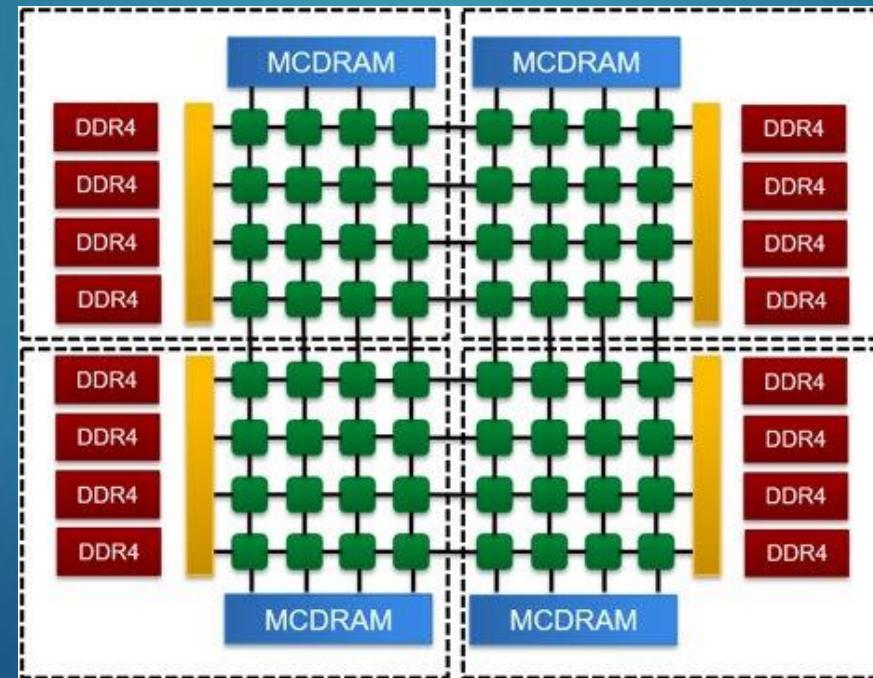
- ▶ Hierarchical memory
- ▶ Remote / local memories (**NUMA : Non Uniform Memory Access**)



Large nodes : 128 cores (BCS)

Now also inside the CPU – Intel KNL

- ▶ Intel KNL (64 cores) can be configured in **2 or 4 NUMA domains**
- ▶ Also add **MCDRAM** (similar idea than GPU GDDR5) **viewed as a N**



- ▶ Or on **AMD CPUs**

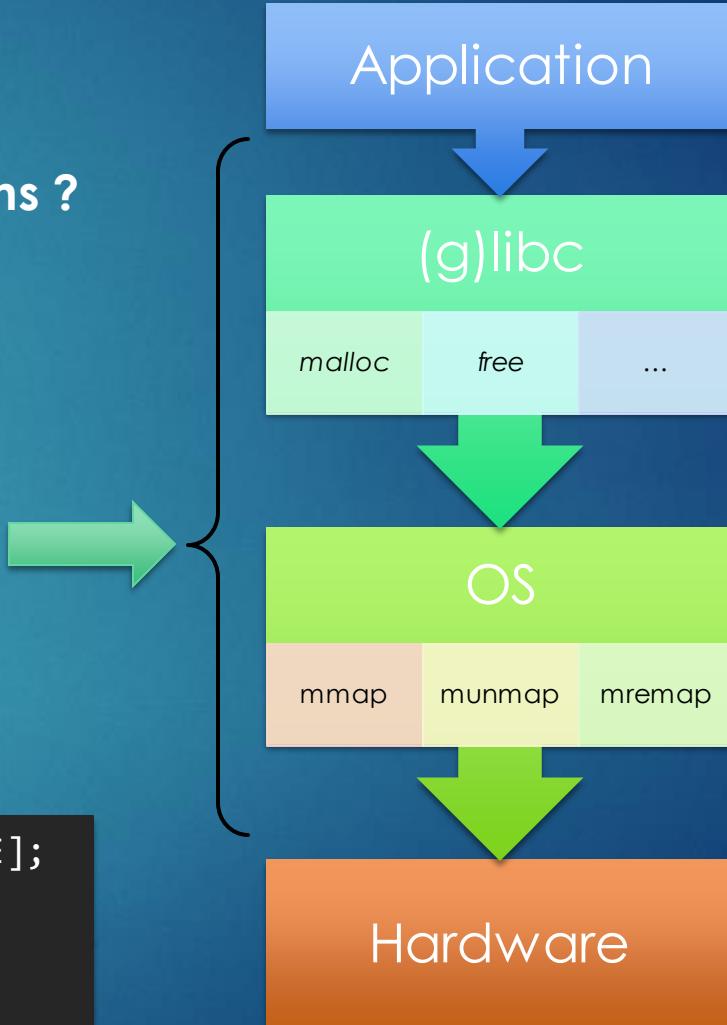
Software memory management layer

| Slide 10 / 43

- ▶ Impact of memory management mechanisms ?
- ▶ Involving **two components** :
 - ▶ User space **memory allocator** (malloc)
 - ▶ **Operating System** (OS)
- ▶ Focus on :
 - ▶ Impact on **allocation time**
 - ▶ Impact on **access efficiency** (placement)
- ▶ Malloc C or C++ interface :

```
float * ptr = malloc(SIZE);
...
ptr = realloc(ptr, NEW_SIZE);
...
free(ptr);
```

```
float * ptr = new float[SIZE];
...
...
delete [] ptr;
```

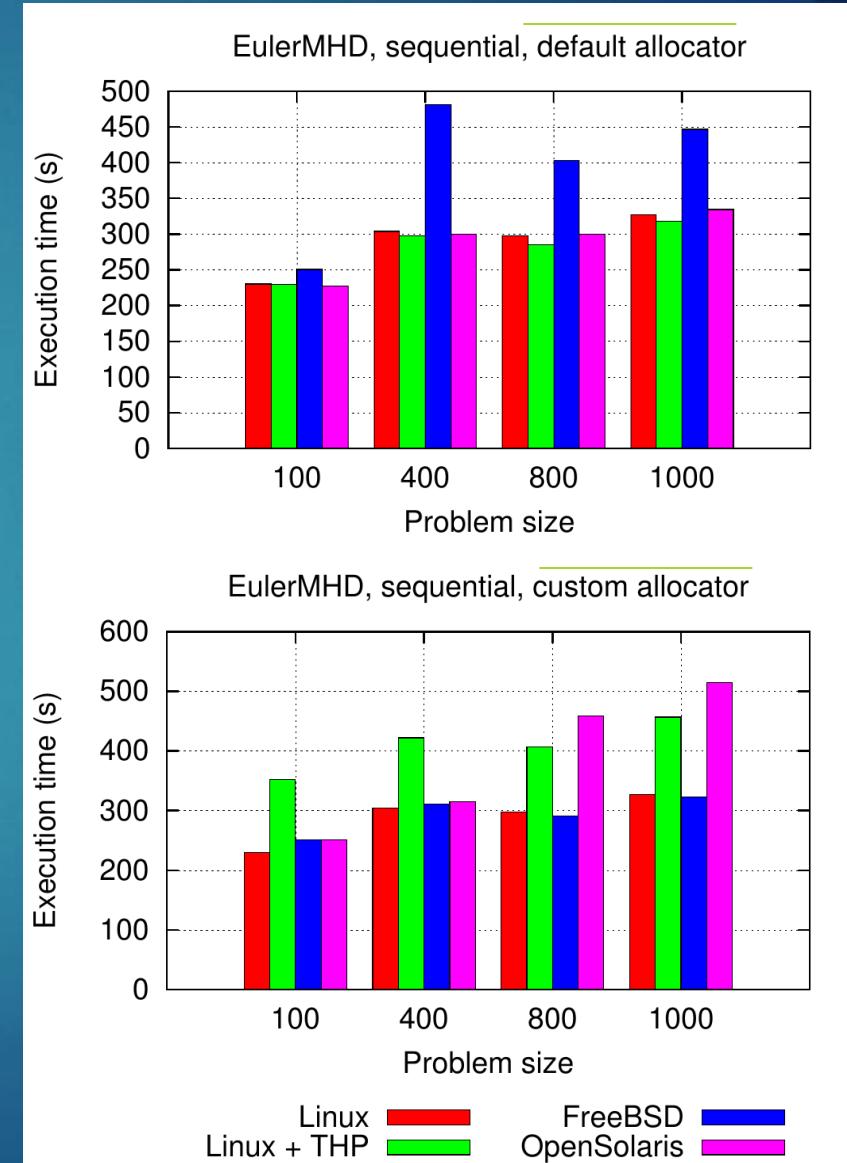


EulerMHD issue

- ▶ EulerMHD (CEA) :
 - ▶ C++ /MPI
 - ▶ Magnéto-hydrodynamic **stencil code**
- ▶ FreeBSD : slowdown of **1.5x**, up to **3x** in **parallel**
- ▶ Impacted function only do compute.
- ▶ Function with **9 arrays pre-allocated** at init. :

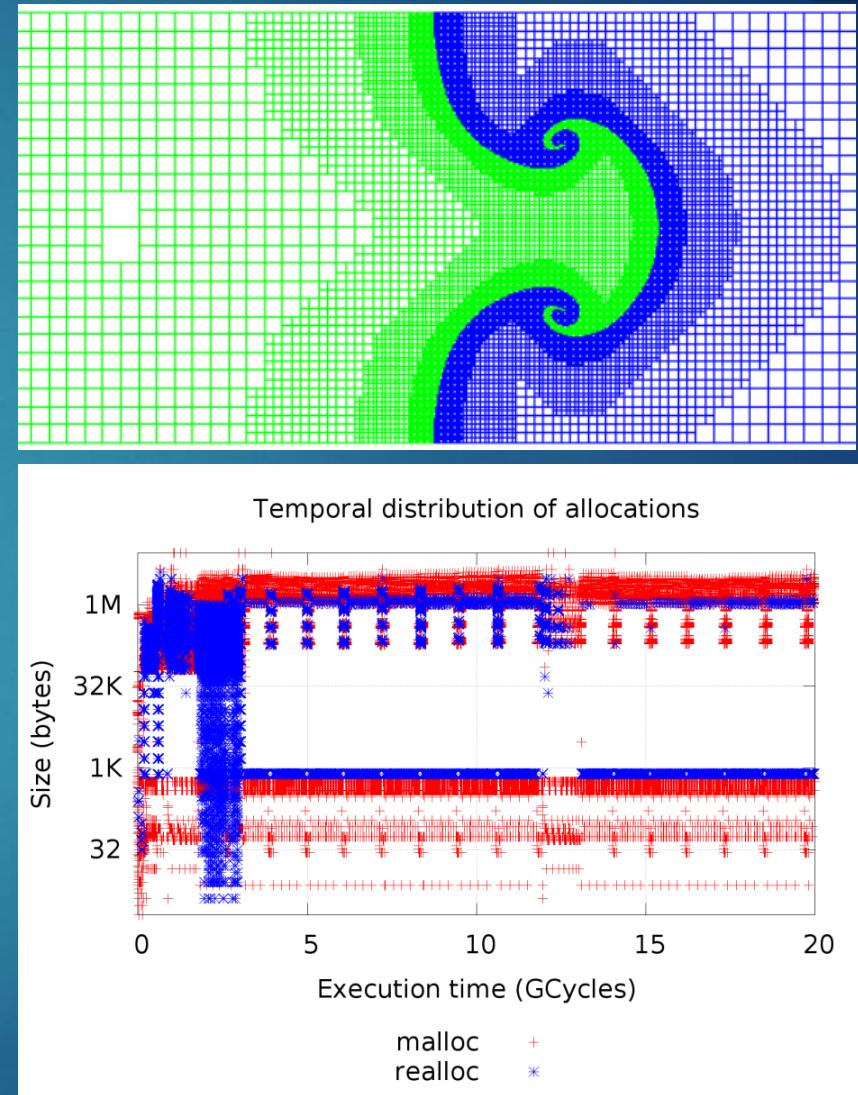
```
for (i = 0 ; i < SIZE ; i++)
    x1[i] = x2[i] + x3[i] ... + x9[i]
```

- ▶ Change between OS's :
 - ▶ User space memory allocator (malloc).
 - ▶ OS paging policy
 - ▶ *(Scheduler)*
- ▶ Effect can be controlled by **changing the allocator**.



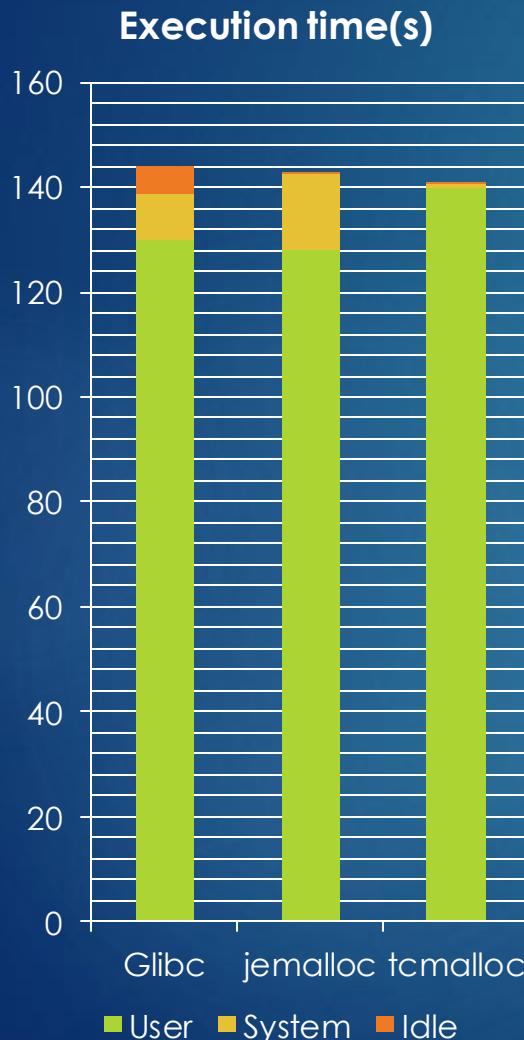
Allocator performance on HPC applications

- ▶ Main interest : **malloc time cost**
- ▶ Test case : **Hera (CEA)**
 - ▶ Adaptive Mesh Refinement (AMR)
 - ▶ Massive C++/MPI code (~1 million lines).
- ▶ **Large number of memory allocations** (~75 millions / 5 minutes on 12 cores)
- ▶ **Large number of alloc/realloc** around ~20 MB
- ▶ **Available allocators :**
 - ▶ Doug Lea / PTMalloc : libc Linux
 - ▶ Jemalloc : FreeBSD / Firefox / Facebook
 - ▶ TCMalloc : Google
 - ▶ Hoard

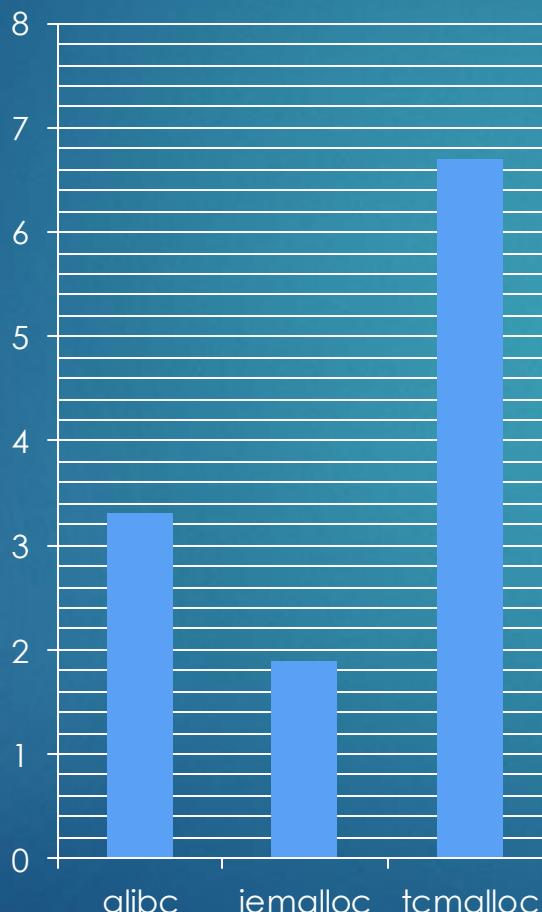


Hera preliminary results

12 cores

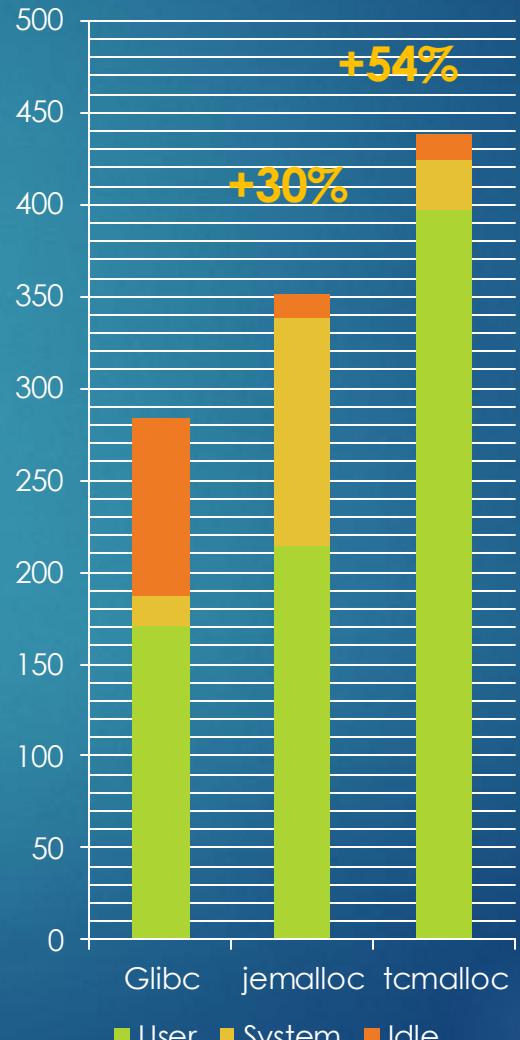


Physical mem.(Go)

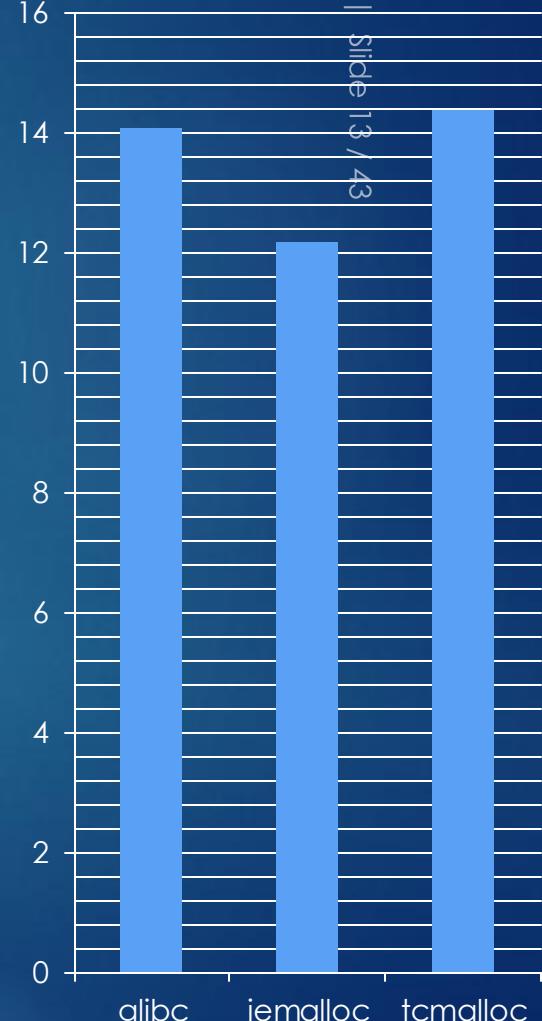


128 cores

Execution time(s)



Physical mem.(Go)



How to measure malloc time

| Slide 14 / 43

- ▶ Measurement method :

```
T0 = clock_start();
ptr = malloc(SIZE);
T1 = clock_end();
```

- ▶ Ok for **small blocks**, but not for **large** one :

```
T0 = clock_start();
ptr = malloc(SIZE);
for ( i = 0 ; i < SIZE ; i += PAGE_SIZE)
    ptr[i] = 0;
T1 = clock_end();
```

- ▶ **Lazy page allocation.**
- ▶ **Page faults** on first access.

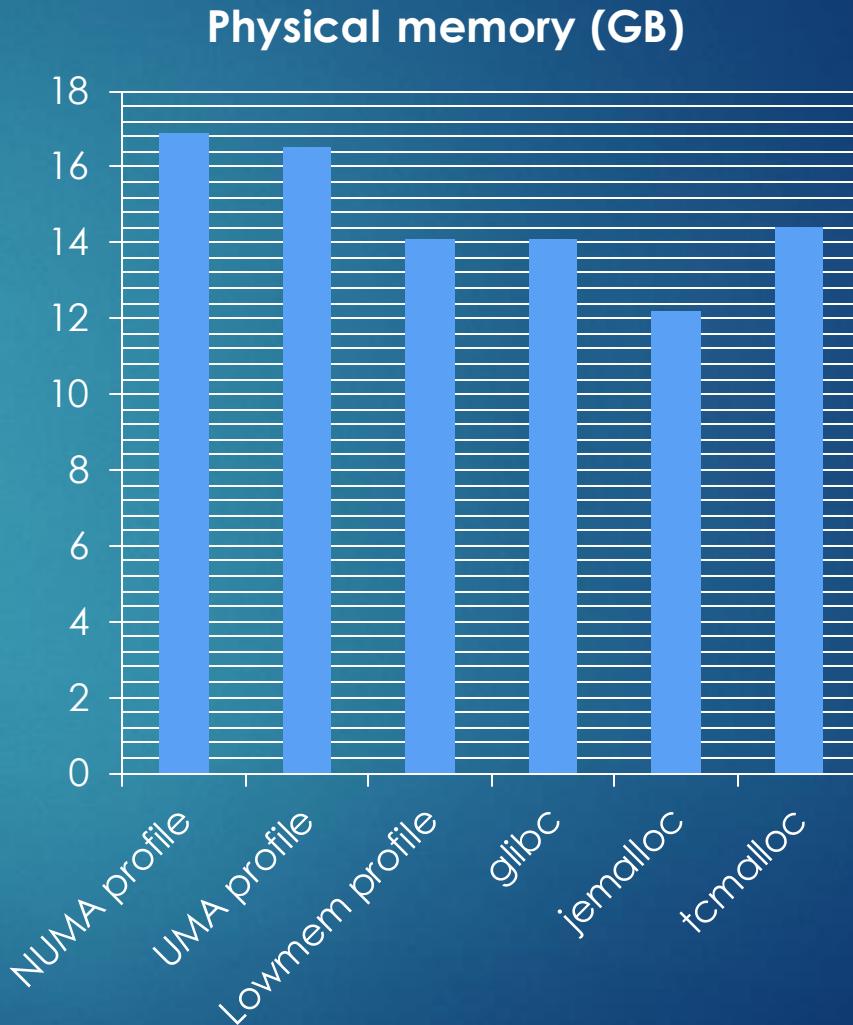
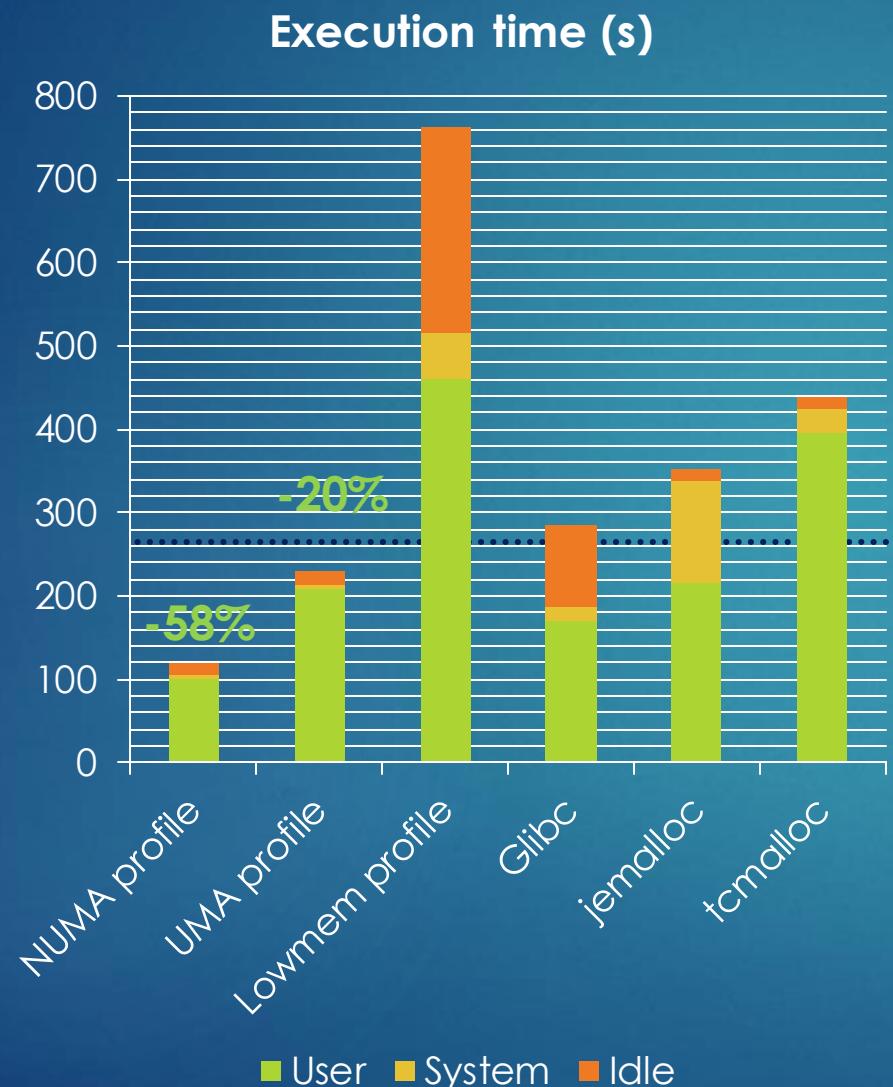
For 4GB	Malloc	First access
Time (M cycles)	0,008	1 217

Large allocations

| Slide 15 / 43

- ▶ Small allocation **well handled** by most allocators, **best is jemalloc / tcmalloc.**
- ▶ Cost for **large allocation : page faults.**
- ▶ **Commonly neglected**, literature mainly discuss small allocations
- ▶ Direct call to **mmap/munmap**
- ▶ **HPC applications** (expected to) use **large arrays**

Hera on Nehalem-EP (128 : 4*4*8 cores)



MALT

A MALLOC TRACKER



UNIVERSITÉ DE
VERSAILLES
ST-QUENTIN-EN-YVELINES



Questions

- ▶ We have **good profiling tool** for **timings**(eg. Valgrind or vtune)
- ▶ But for what **memory profiling**?
- ▶ Memory can be an issue :
 - ▶ **Availability** of the resource
 - ▶ **Performance**
- ▶ Three main questions :
 - ▶ How to reduce **memory footprint** ?
 - ▶ How to improve overhead of **memory management** ?
 - ▶ How to improve **memory usage** ?

Some issue examples

- ▶ We want to point :
 - ▶ **Where** memory is allocated.
 - ▶ **Properties** of allocated chunks.
 - ▶ **Bad** allocation **patterns** for performance.

```
_thread int gblVar[SIZE];  
int * func(int size)  
{  
    child_func_with_allocs();  
    void * ptr = new char[size];  
    double * ret = new double[size*size*size];  
    for (auto it : iter_items)  
    {  
        double * buffer = new double[size];  
        //short and quick do stuff  
        delete [] buffer;  
    }  
    return ret;  
}
```

Global variables and TLS

Indirect allocations

Leak

Might lead to swap for large size

Auto missing reference

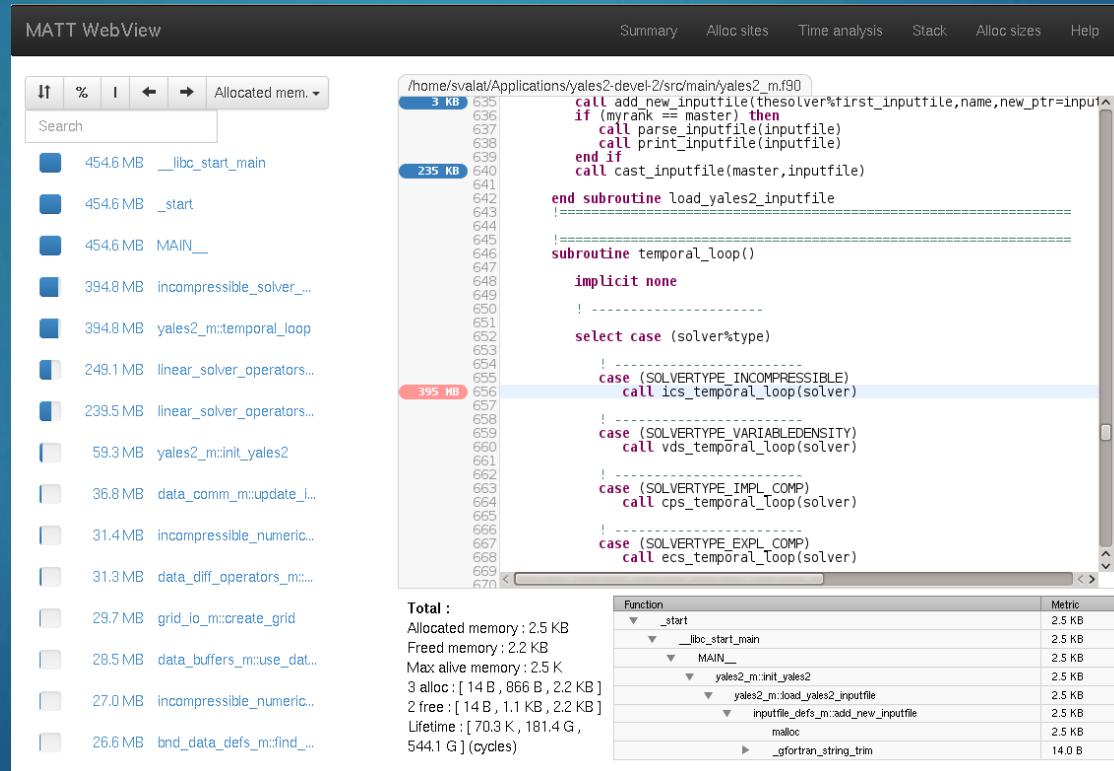
Short life allocations

What I want to provide

- ▶ Same **approach** than **valgrind/kcachgind**
- ▶ **Mapped** allocations on **sources lines** and **call stacks**
 - ▶ **Using profile** approach, **not snapshots**
- ▶ For **memory resource usage** :
 - ▶ Memory **leaks**
 - ▶ **Memory on peak**
- ▶ For **performance** :
 - ▶ Allocation **count** and **cumulated size**
 - ▶ Allocation **sizes** (min/average/max)
 - ▶ Chunk **lifetime** (min/average/max)

Our GUI : based on web tech.

- ▶ Web technology (NodeJS, D3JS, Jquery, AngularJS)
- ▶ Easier for remote usage



Global summary

- ▶ Provide a small summary
- ▶ Provide some warnings

Show all details	Show help
Physical memory peak	66.7 MB
Virtual memory peak	158.1 MB
Requested memory peak	6.1 MB
Cumulated memory allocations	11.5 MB
Allocation count	172.2 K
Recycling ratio	1.9
Leaked memory	743.7 KB
Largest stack	0 B
Global variables	10.0 MB 
TLS variables	48 B
Global variable count	421.0 K 
Peak allocation rate	37.8 MB/s

Global summary : top 5 functions

- ▶ Summarize **top functions** for some metrics
- ▶ Points to check
- ▶ Examples on YALES2

Alloc count

Ratio	Allocs	Function
	911.9 K	data_comm_m::copy_int_comm_to_data
	896.4 K	data_comm_m::copy_data_to_int_comm
	853.2 K	data_comm_m::update_int_comm
	484.9 K	sponge_layer_m::calc.sponge_layer_mask
	296.0 K	incompressible_numerics_m::ics_diffuse_velocity_rk_4th

Allocated memory

Ratio	Allocs	Function
	202.4 MB	linear_solver_operators_m::solve_linear_system_deflated_pcg
	26.6 MB	bnd_data_defs_m::find_bnd_data
	21.8 MB	linear_solver_operators_m::solve_el_grp_pcg
	19.0 MB	data_comm_m::copy_int_comm_to_data
	18.1 MB	data_comm_m::update_int_comm

Peak memory

Source annotations

Inclusive/Exclusive

Metric selector

Per line annotation

Symbols

Details of symbol or line

Call stacks reaching the selected site.

The screenshot shows the MATT WebView interface. On the left, a sidebar lists symbols with their sizes: __libc_start_main (28.4 KB), _start (28.4 KB), main (28.2 KB), testMaxAlive() (12.5 KB), recurseA(int) (6.9 KB), testThreads() (6.3 KB), funcB() (1.0 KB), testRecuseInterv... (1.0 KB), testRecuseInterv... (1.0 KB), funcC() (704.0 B), testParallelWithRecur... (704.0 B), OutOfMainAlloc (128.0 B), __cxx_global_var_init1 (128.0 B), global constructors ke... (128.0 B), and __libc_csu_init (128.0 B). The main area displays a portion of the source code with memory usage annotations above each line. A call stack table at the bottom shows the memory usage for each function and its children. A tooltip 'Inclusive/Exclusive' points to the metric selector in the top right. A tooltip 'Metric selector' points to the metric selector dropdown. A tooltip 'Per line annotation' points to the annotated source code. A tooltip 'Symbols' points to the sidebar. A tooltip 'Details of symbol or line' points to the call stack table. A tooltip 'Call stacks reaching the selected site.' points to the call stack table.

Summary Alloc sites Time analysis Stack Alloc sizes Help

```

/home/svalat/Projects/matt/src/lib/tests/simple-case.cpp
53 int * ptr = new int[16];
54 *(char*)ptr = 'c';//required otherwise new compilers will remove malloc/free
55 delete [] ptr;
56 }
57 **** FUNCTION ****
58 void funcB()
59 {
60     void * ptr = malloc(32);
61     *(char*)ptr='c';//required otherwise new compilers will remove malloc/free
62     free(ptr);
63     funcC();
64 }
65 **** FUNCTION ****
66 void funcA()
67 {
68     void * ptr = malloc(16);
69     *(char*)ptr='c';//required otherwise new compilers will remove malloc/free
70     free(ptr);
71     funcB();
72 }
73 **** FUNCTION ****
74 void recurseA(int depth)
75 {
76     if (depth > 0)
77     {
78         void * ptr = malloc(64);
79         *(char*)ptr='c';//required otherwise new compilers will remove malloc/free
80         free(ptr);
81         recurseA(depth-1);
82     }
83 }
84 **** FUNCTION ****
85
86
87
88 **** FUNCTION ****

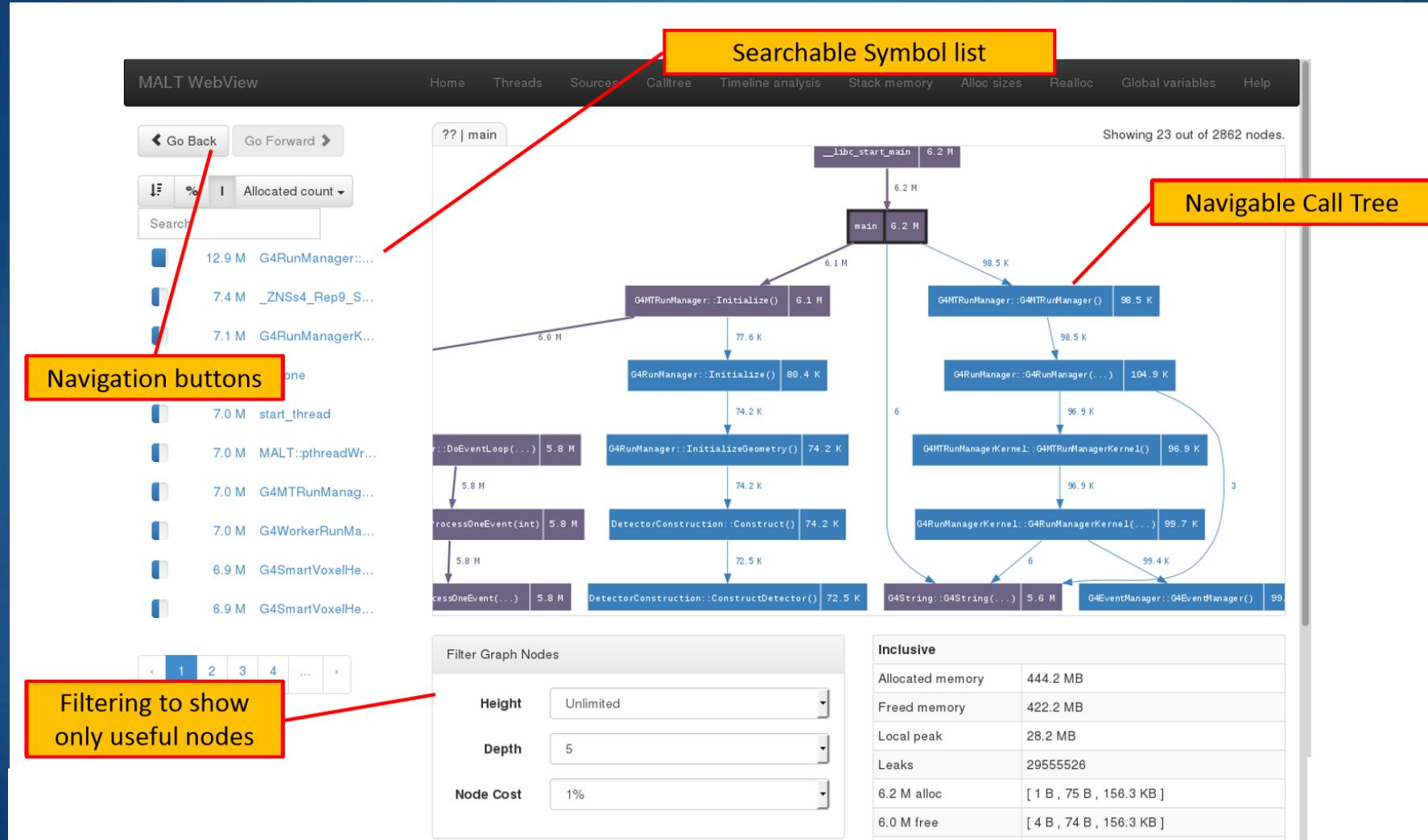
```

Total :
Allocated memory : 96 B
Freed memory : 96 B
Max alive memory : 96
2 alloc : [32 B , 48 B , 64 B]
2 free : [32 B , 48 B , 64 B]
Lifetime : [41.3 K , 42.1 K , 42.9 K] (cycles)

Function	Metric
_start	96.0 B
__libc_start_main	96.0 B
main	96.0 B
funcA()	96.0 B
funcB()	96.0 B
malloc	96.0 B
funcC()	96.0 B

Call tree view

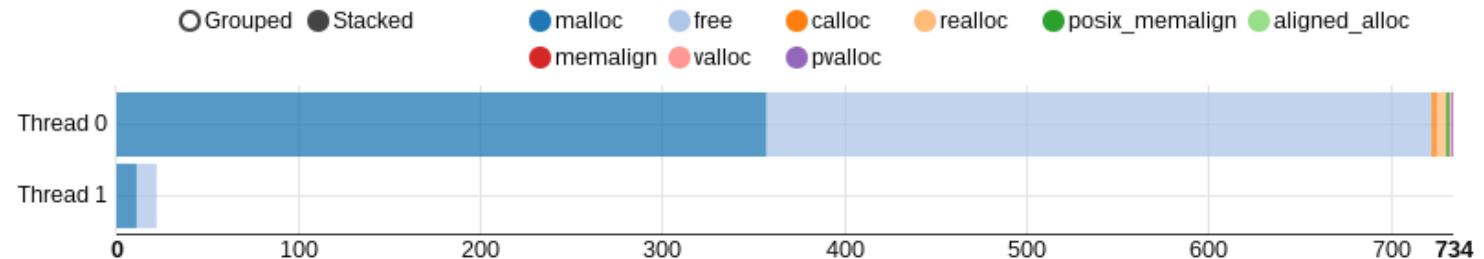
Work by CERN/Openlab summer student



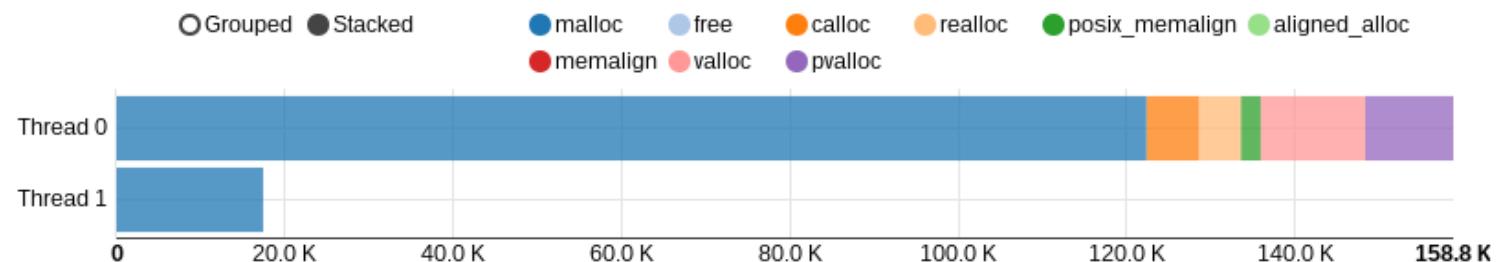
Per thread statistics

CERN-IT-MALT, Sébastien Valat
13/03/2018

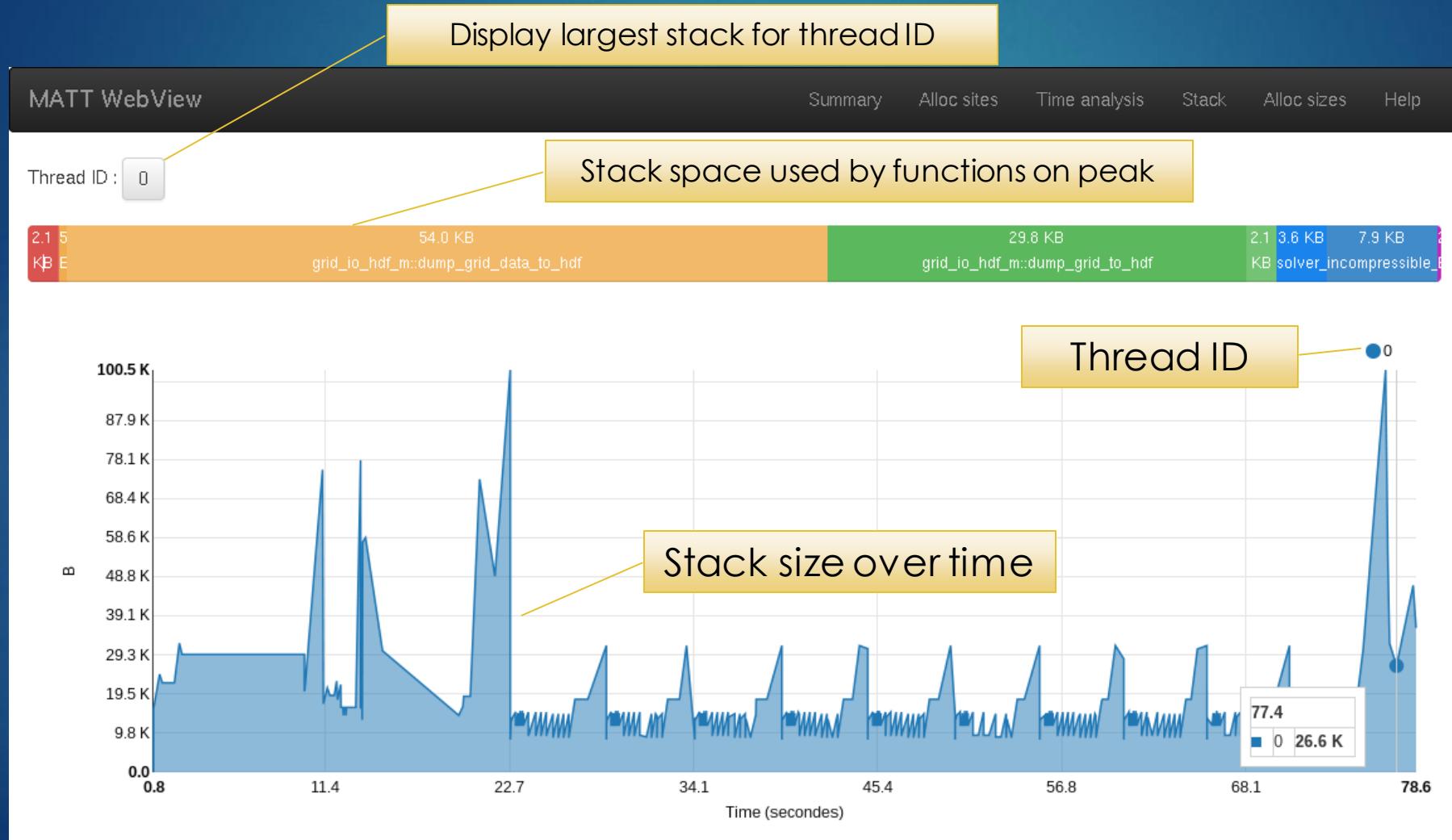
Call per thread



Time per thread

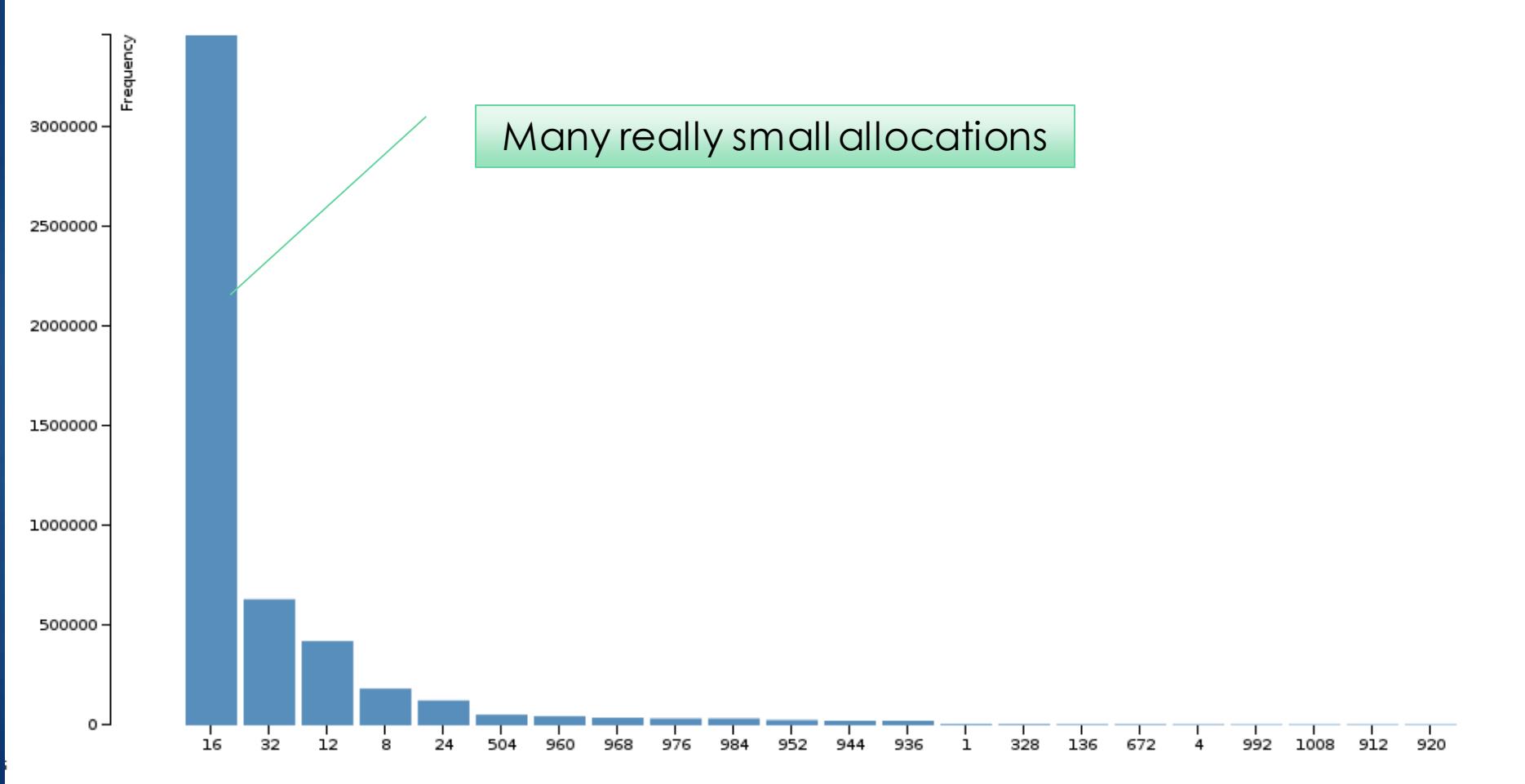


Tracking stack memory

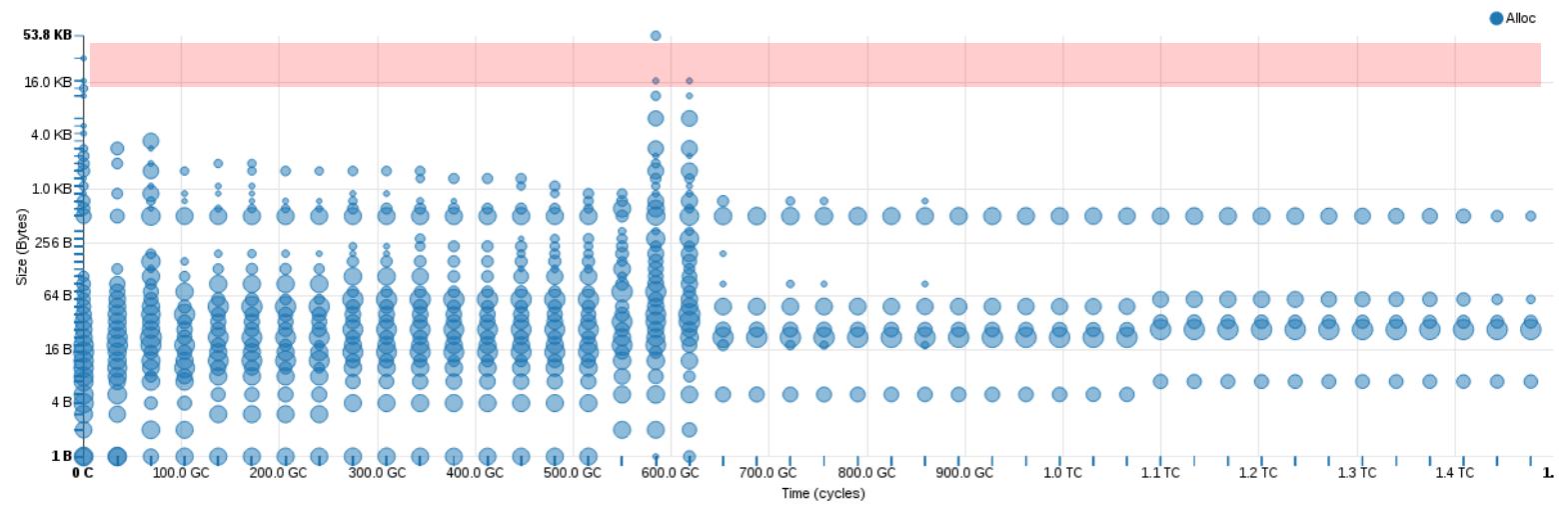


Chunk size distribution

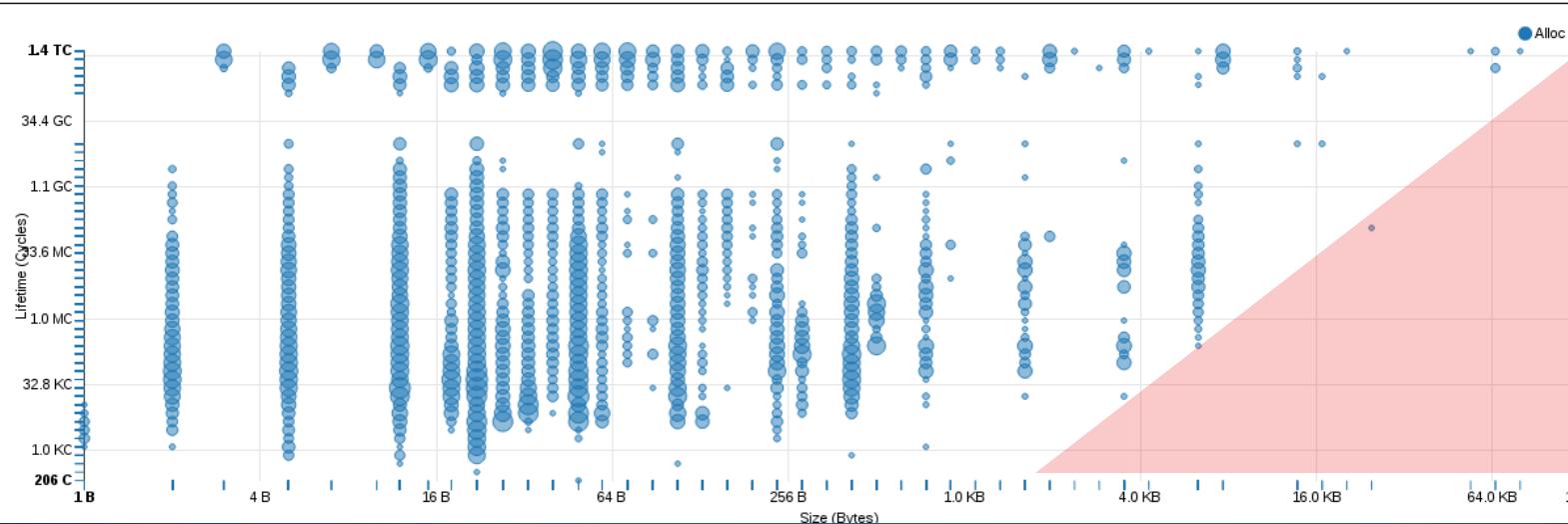
Example from YALES2



Size over time

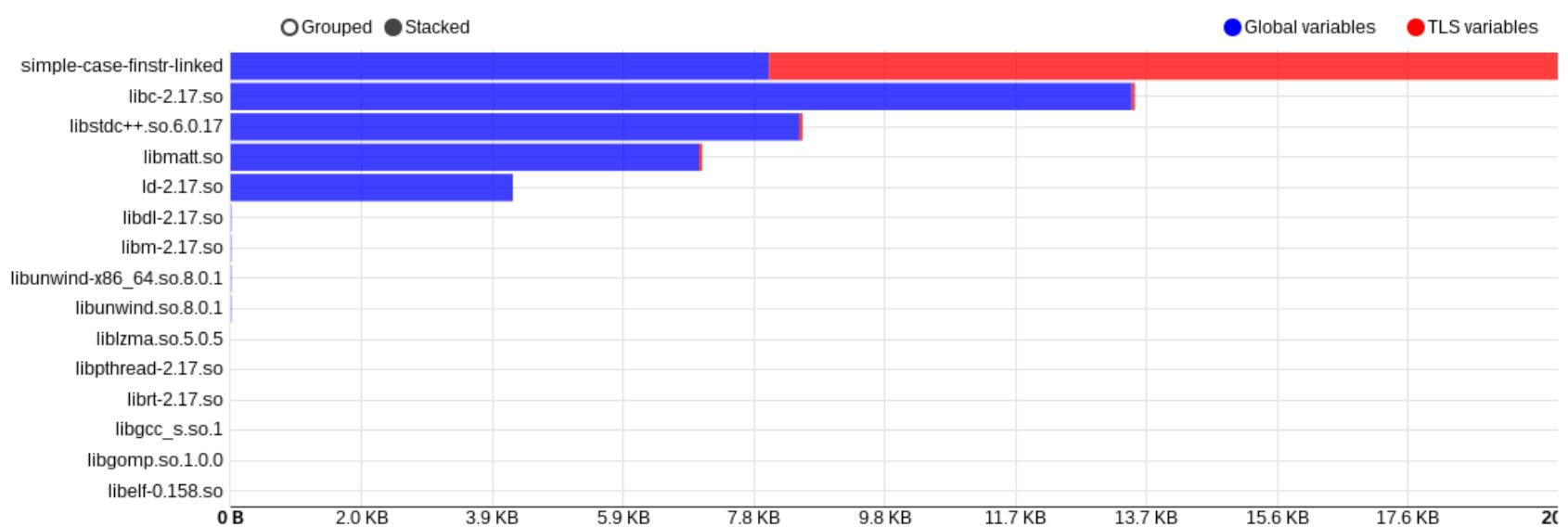


Lifetime over size

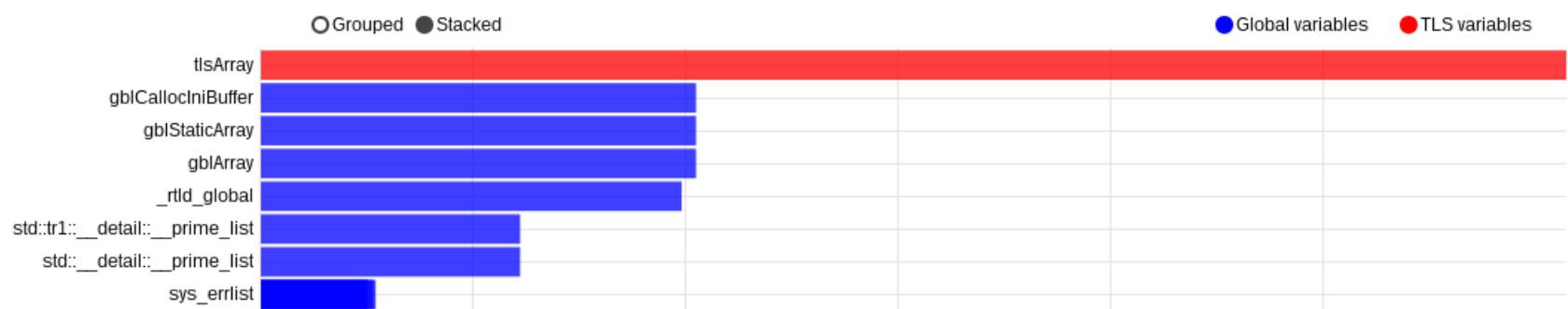


Global variables

Distribution over binaries

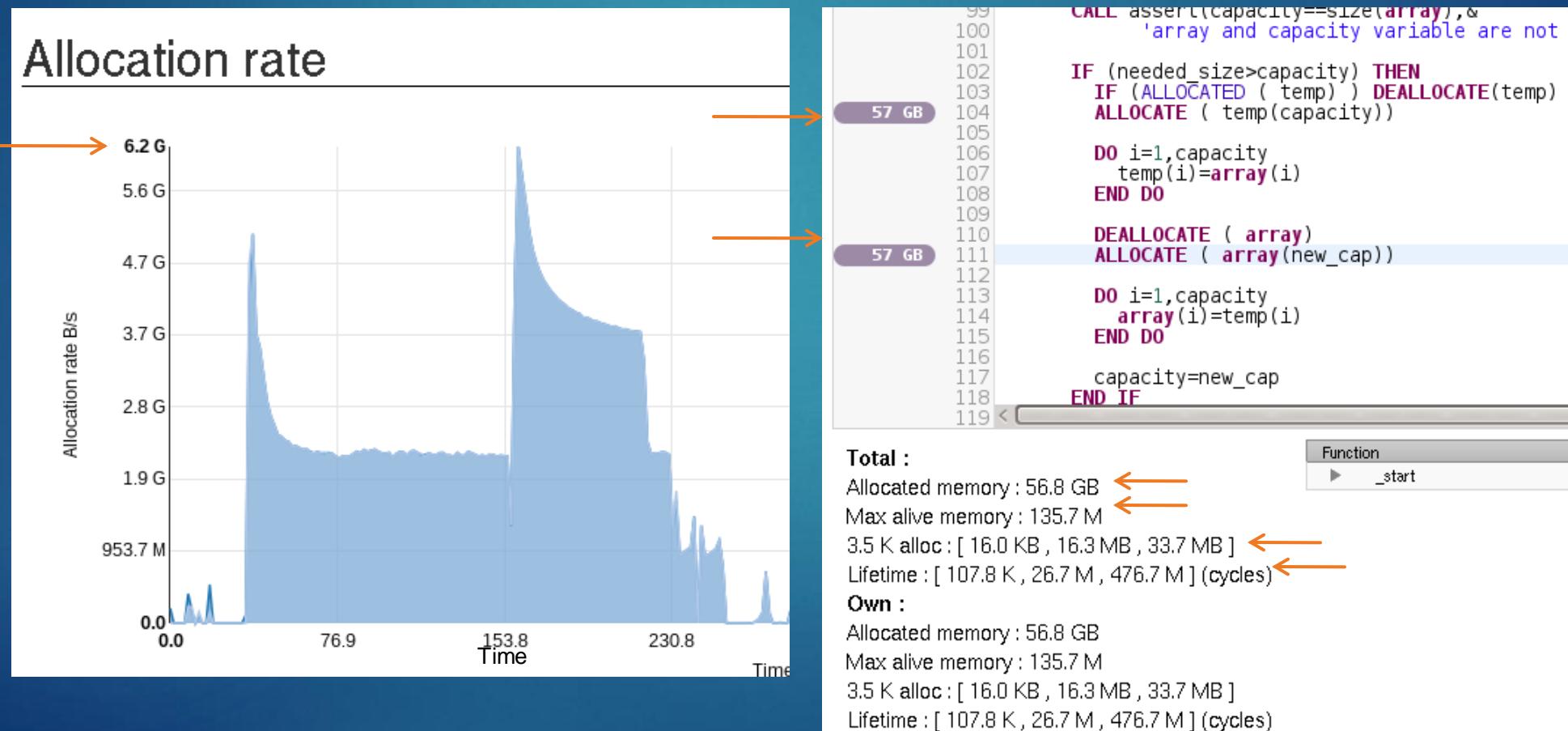


Distribution over variables



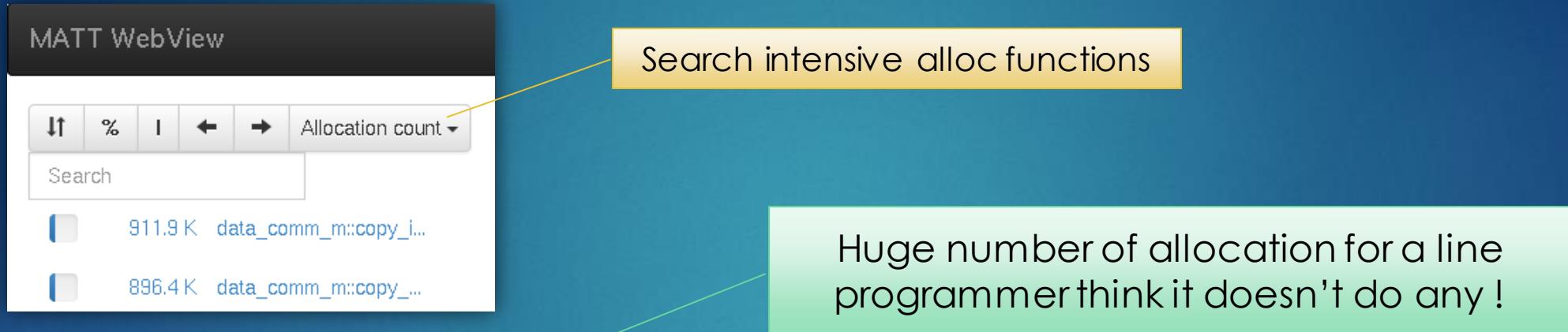
Example on AVBP init phase

- Issue with **reallocation** on init
- Detected with **allocation rate & cumulated allocated mem.**



Allocatable arrays on YALES2

- Issue only occur with **gfortran**, ifort uses stack arrays.



608 K

Total :
Allocated memory : 9.5 MB
Freed memory : 9.5 MB
Max alive memory : 432
608.0 K alloc:[16 B , 16 B , 16 B]
608.0 K free:[16 B , 16 B , 16 B]
Lifetime : [24.5 K , 39.9 K , 37.8 M] (cycles)
Own :
Allocated memory : 9.5 MB

And mostly really small allocations!

We can found allocs of 1B !

- Examples on YALES 2, small allocations :

MATT WebView

↑ % | ← → Min. size ▾

Search

1.0 B	/usr/lib/gcc/x86_64-p...
1.0 B	_strdup
1.0 B	data_defs_m::resize_...

Search for the minimal chunk size.

1.0 B /usr/lib/gcc/x86_64-p...
1.0 B __strup
1.0 B data_defs_m::resize_...

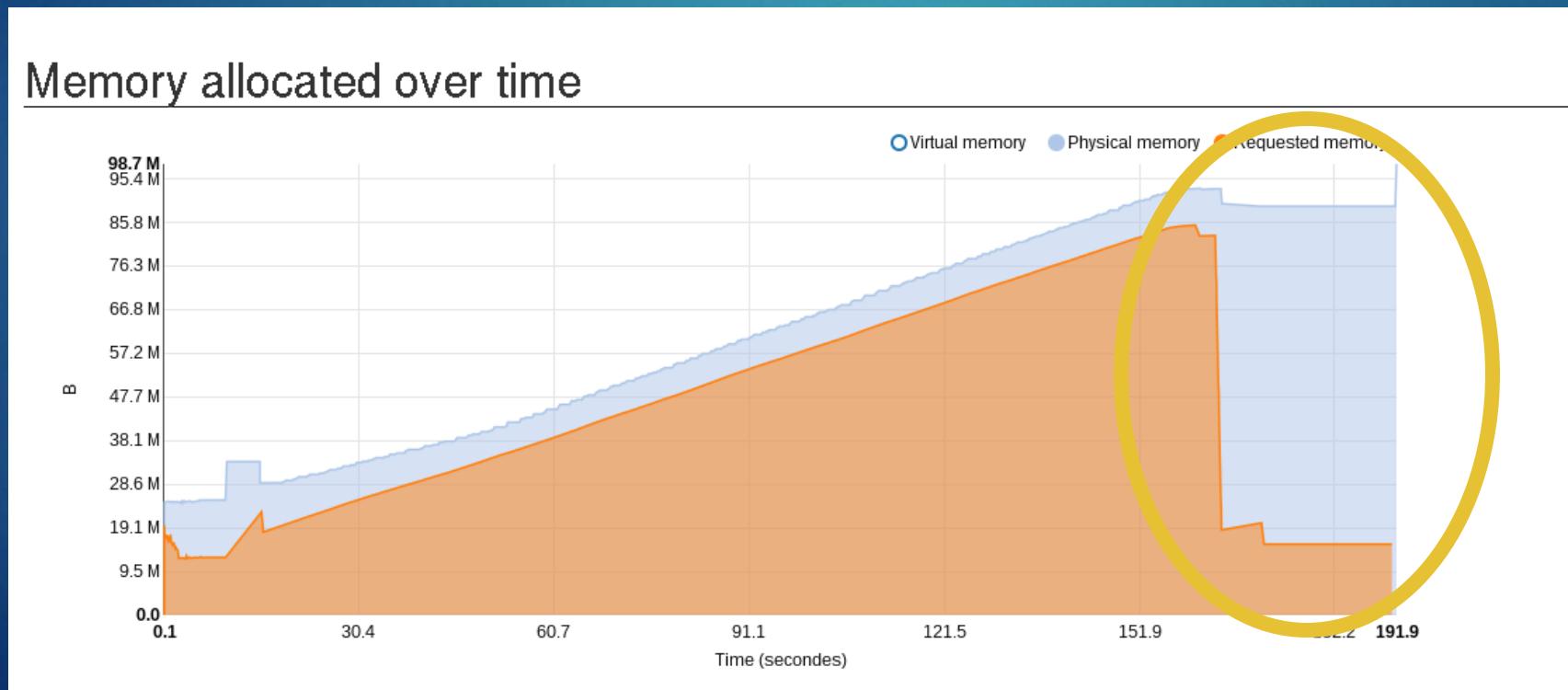
Many codes produce allocations of 1B.
OK with moderation.

1 B

```
530 case (DATATYPE_REAL_NODE_VECTOR,DATATYPE_REAL_ELEM_VECTOR, &
531 DATATYPE_REAL_FACE_VECTOR,DATATYPE_REAL_PAIR_VECTOR)
532 if (associated(data_ptr%r2_ptrs)) then
533   deallocate(data_ptr%r2_ptrs)
534 end if
535 allocate(data_ptr%r2_ptrs(nel_grps))
536 do n=1,nel_grps
537   NULLIFY(data_ptr%r2_ptrs(n)%ptr)
538 end do
539
```

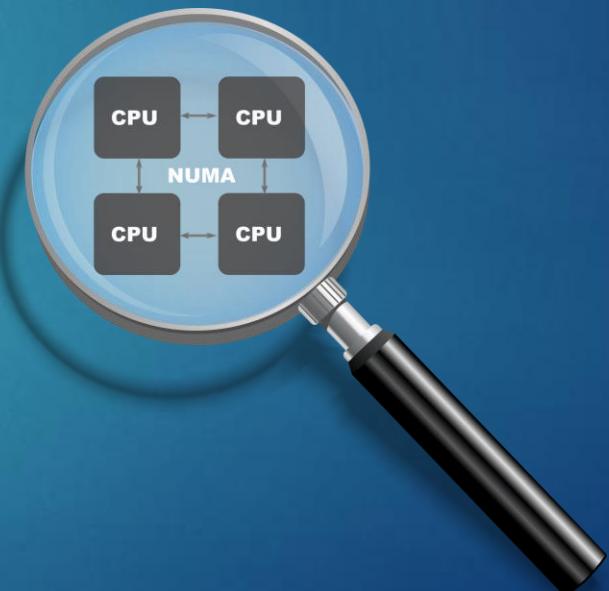
Fragmentation issue

- ▶ Example of **fragmentation** detection
- ▶ Using the time chart with **physical**, **virtual** and **requested memory**
- ▶ **Solution : avoid interleaved allocation of chunks with different lifetime.**
- ▶ Looking on **source annotation** : most of them **can be avoided**.



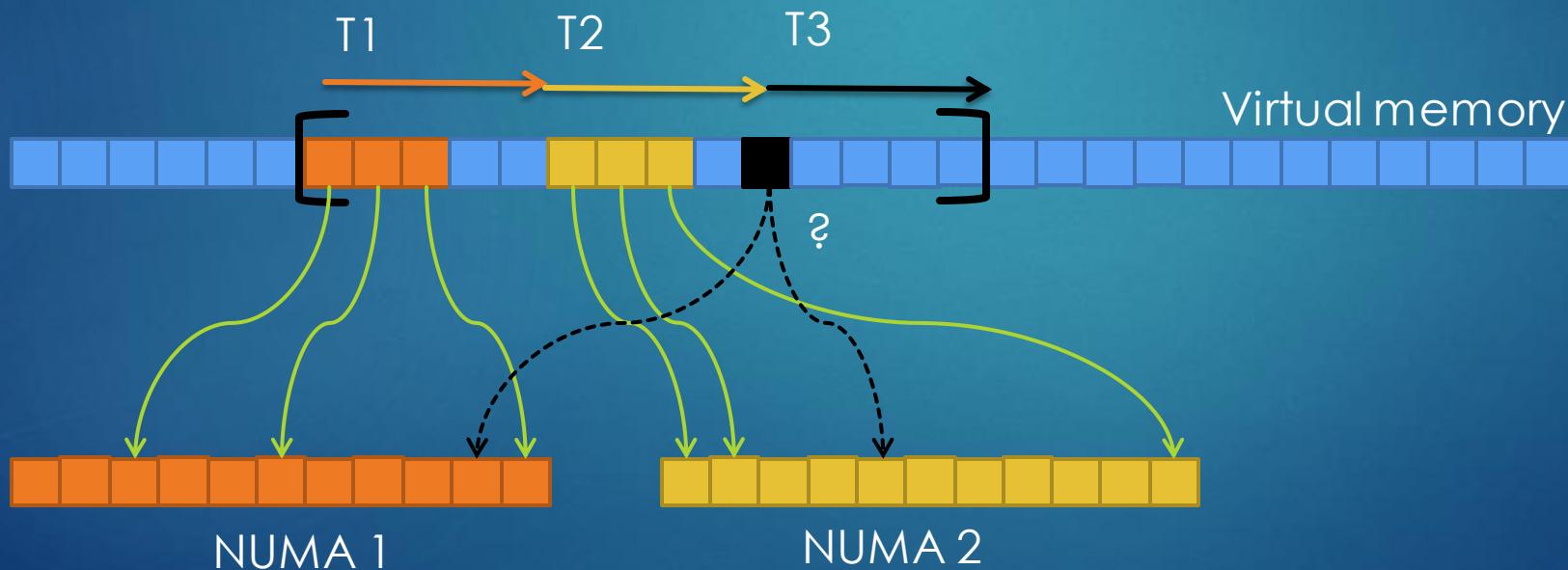
NUMAPROF

A NUMA MEMORY PROFILING TOOL



Implicit binding : first touch

- ▶ New allocated segments are **physically empty**
- ▶ They are filled on **first touch**
- ▶ Page selection **depend** of the **thread position**



Typical OpenMP mistake

- ▶ Make first **init outside of OpenMP** (in thread 1)
- ▶ So **each pages** will be first touched **on NUMA 1**

```
#pragma omp parallel for
for (int i = 0 ; i < SIZE ; i++)
    array[i] = 0;
```

- ▶ Then access

```
#pragma omp parallel for
for (int i = 0 ; i < SIZE ; i++)
    array[i]++;
```

- ▶ **Bad performance** due to remote accesses !

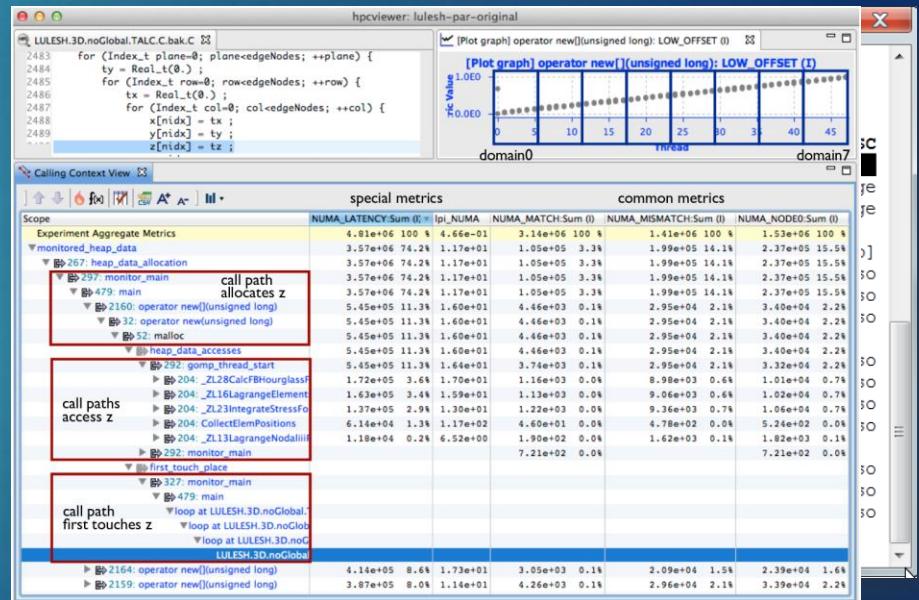
Wish list for a profiling tool...

- ▶ We want to know if we make **remote accesses**
- ▶ Ideally we need to know **where...**
- ▶ We can dream, we want to know **which allocation contain issues**
- ▶ We want to know **where** the **first touch** has been done
- ▶ On KNL we want to check **MCRAM accesses**



Existing tools

- ▶ **MemProf** [1], a research paper from Grenoble
 - ▶ Hardware feature provided by AMD & kernel module, No GUI
- ▶ **SNPERF** [2], again a research project
 - ▶ Link utilization on time chart
- ▶ **NUMAgrind** [3]
 - ▶ Looks nice but not available
- ▶ **Numatop**
 - ▶ Similar to top, global profiling
- ▶ **HPCToolkit** [4]
 - ▶ Also have a nice interface, but hw counters & sampling



[1] Renaud Lachaize and Al. MemProf: A Memory Profiler for NUMA Multicore Systems.

[2] U. Prestor and Al., "An application-centric ccNUMA memory profiler,"

[3] Profiling Directed NUMA Optimization on Linux Systems: A Case Study of the Gaussian Computational Chemistry Code

[4] A Tool to Analyze the Performance of Multithreaded Programs on NUMA architectures

NUMAPROF

- ▶ Take back the idea from **MALT**
 - ▶ **Web interface**
 - ▶ **Source annotation**
 - ▶ **Global metrics**
- ▶ Use intel **Pin**
 - ▶ Permit to **instrument** all **memory accesses**
 - ▶ **Parallel** opposite to valgrind
 - ▶ **Difficulty:** we cannot easily use libs inside the tool
 - ▶ I would have used hwloc and libnuma.....

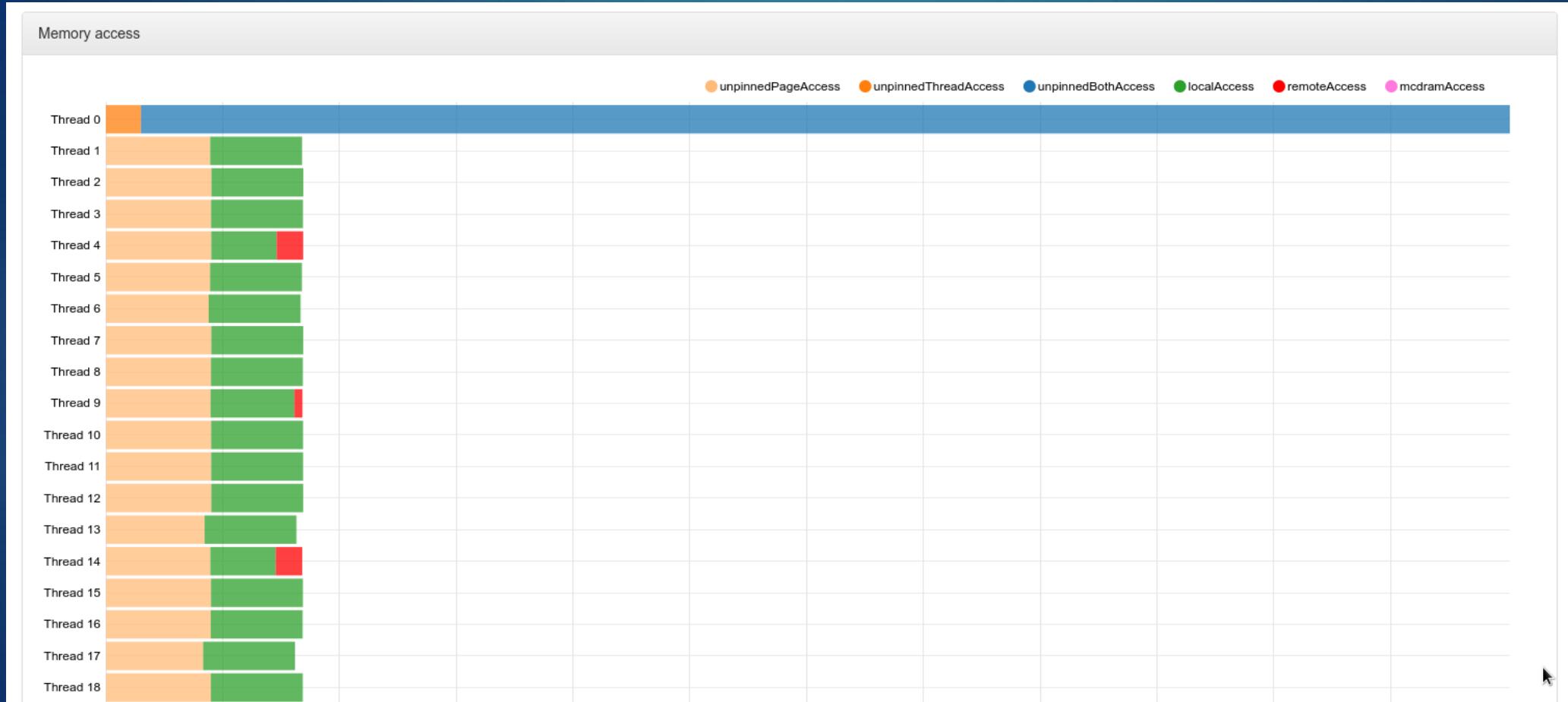
On access we need...

- ▶ On each access we want to know if it is
 - ▶ **Remote** access
 - ▶ **Local** access
 - ▶ **MCDRAM** access
 - ▶ **Page is pinned**
 - ▶ **Thread is pinned**
- ▶ So, we need to know
 - ▶ Where is **the page**
 - ▶ Where is **the current thread**
- ▶ We can **skip** accesses to **local stack** (overhead 80x -> 40x)

Global summary



Statistics per thread

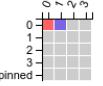


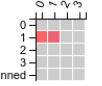
Details per thread

Numaprof Home Threads **Details** Sources Assembler Help

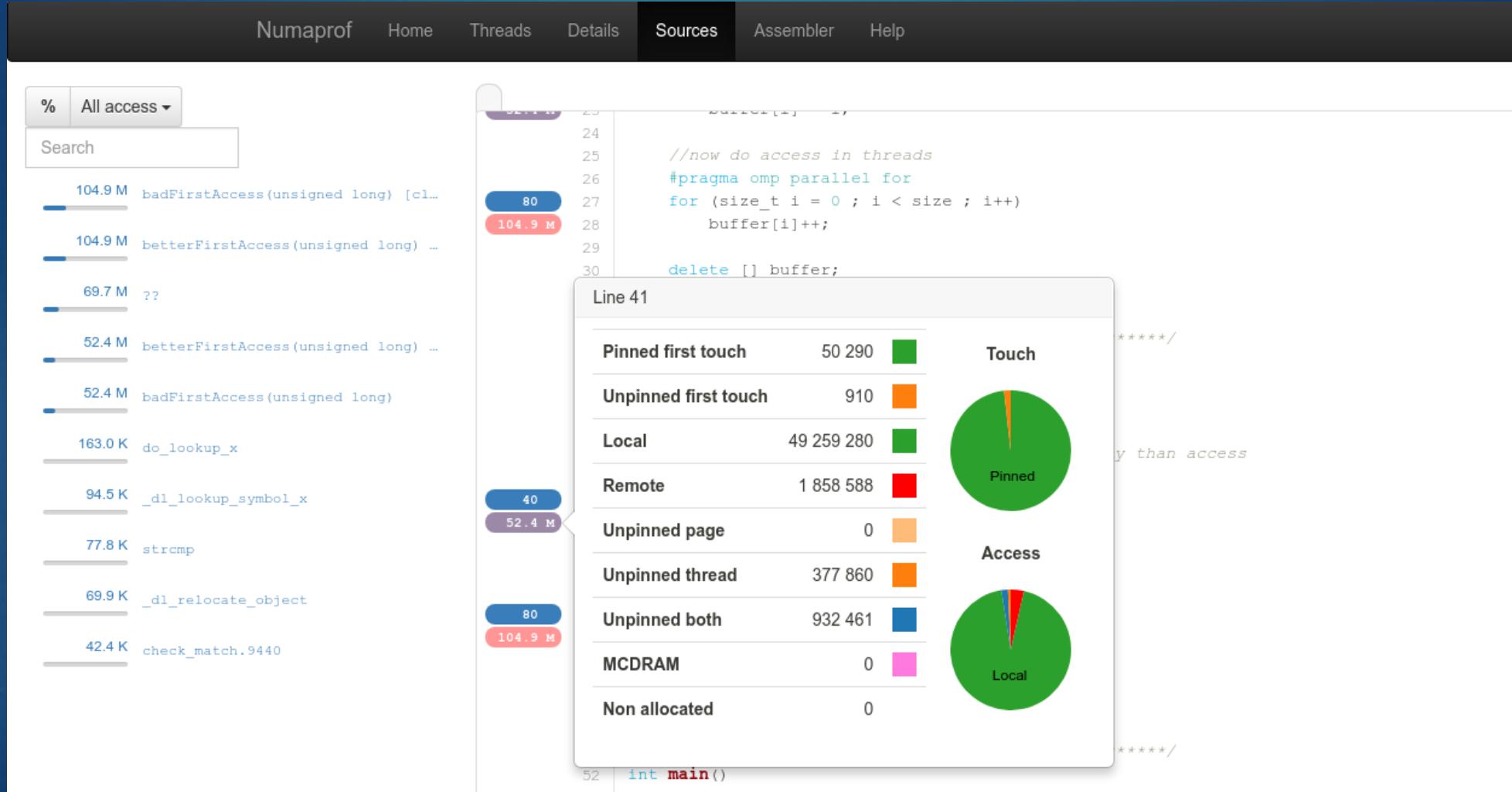
← « 1 2 3 4 5 6 7 8 9 10 » →

5

Thread 4	
Lifetime	64.31% → 90.65%
CPU thread binding	4
Numa thread binding	0
Numa mem. policy	MPOL_DEFAULT on -1 considered as NO_BIND
First touch	 Pinned 1,024.00
Accesses	
Accesses	
Pinning log	At 64.31%, pin thread on node 0 At 64.31%, do memory binding MPOL_DEFAULT on -1 considered as NO_BIND

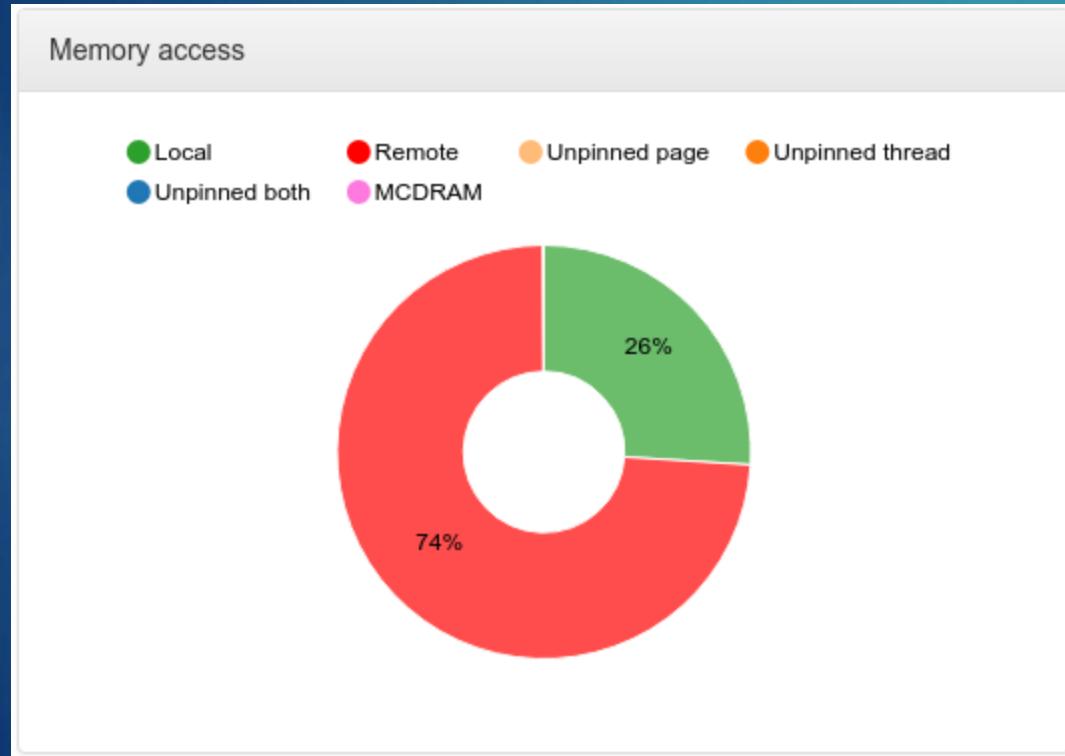
Thread 5	
Lifetime	64.32% → 90.91%
CPU thread binding	5
Numa thread binding	1
Numa mem. policy	MPOL_DEFAULT on -1 considered as NO_BIND
First touch	
Accesses	
Accesses	
Pinning log	At 64.32%, pin thread on node 1 At 64.32%, do memory binding MPOL_DEFAULT on -1 considered as NO_BIND

Source & asm annotations

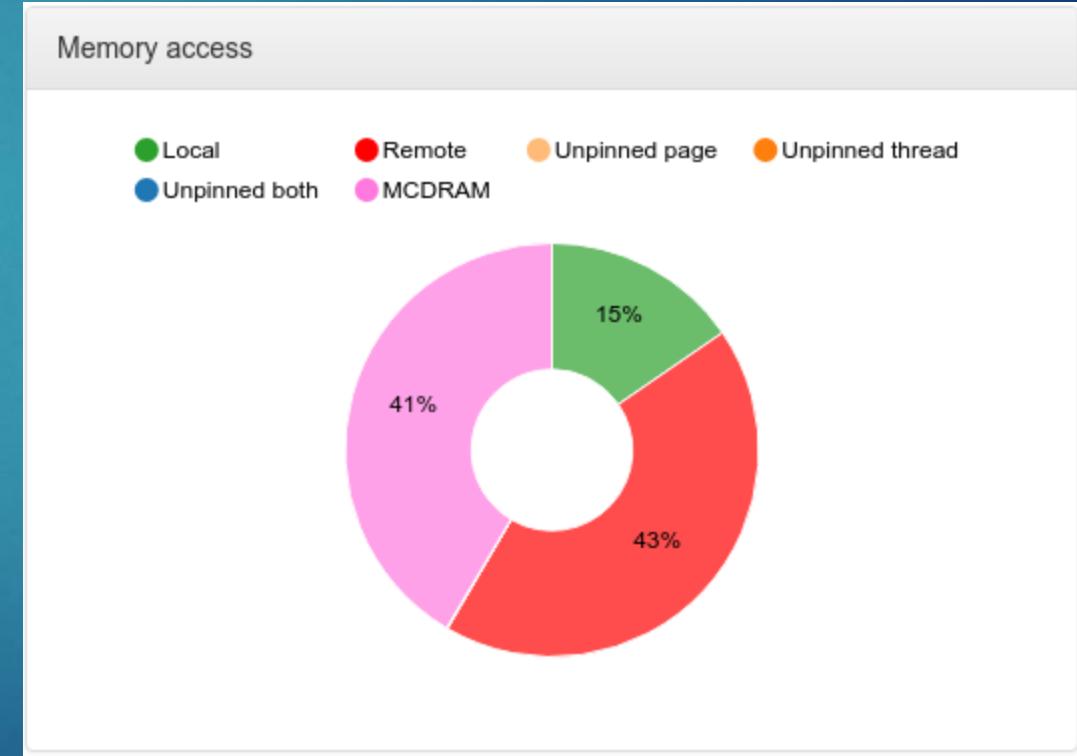


Code Hydro

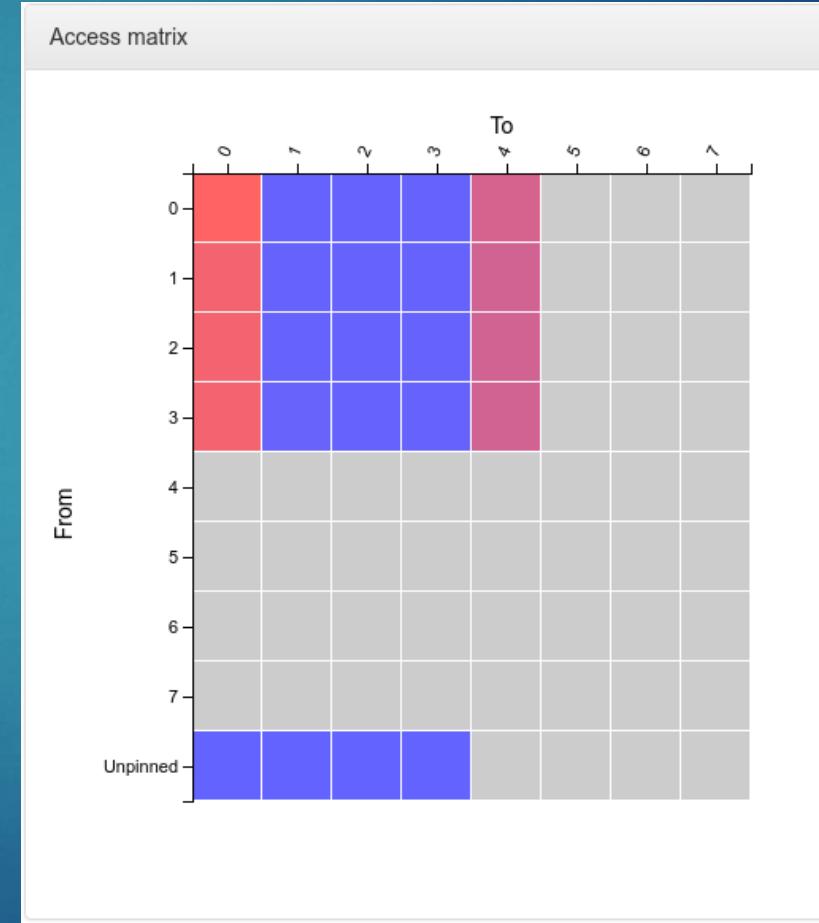
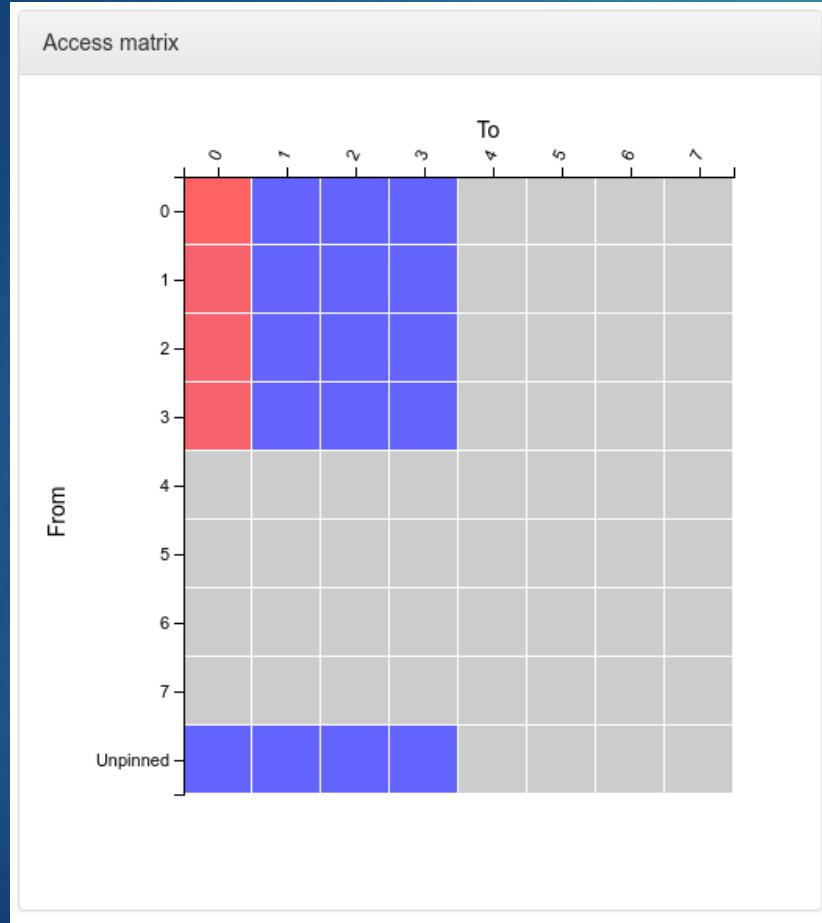
▶ KNL Without HBM



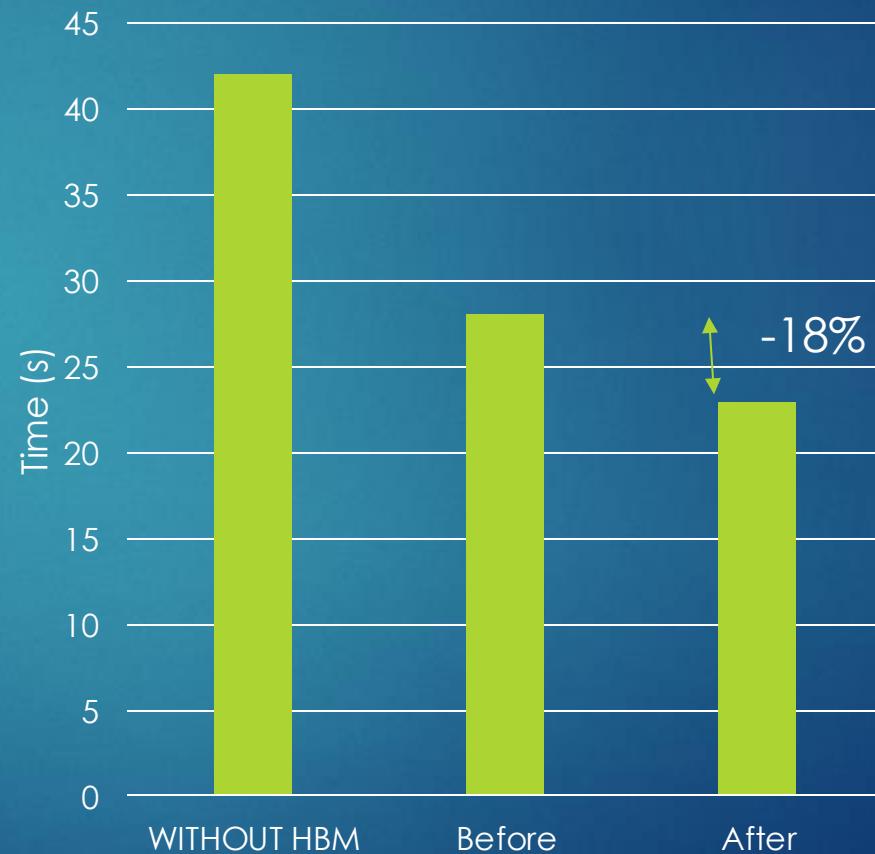
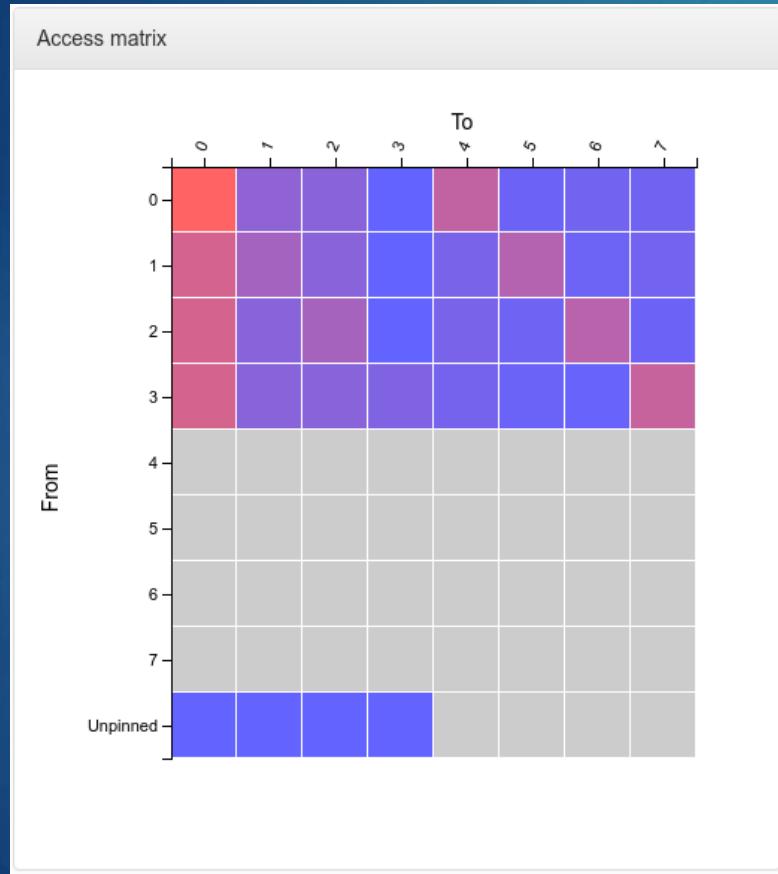
WITH HBM



Original Hydro access matrix



Speed up obtained on Hydro



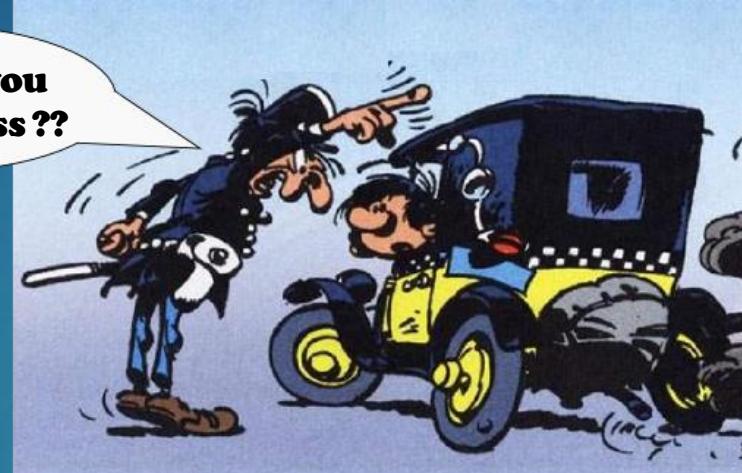
Conclusion on memory

- ▶ Need to understand memory behavior of apps
- ▶ Tools are need
- ▶ I provide two on <http://memtt.github.io/> un CeCILL-C (LPLC-Like)
 - ▶ MALT
 - ▶ NUMAPROF
- ▶ Still work needed in OS layers to progress

Questions ?

THANK YOU.

**Who give you
this address ??**



M'ENFIN !

What you want ??
Some users expect
0x1A8E3 as answer to
malloc(0) !! ???





Studying a 40 Tb/s
data acquisition system for
the LHCb experiment

LHCb, an upgrade for 2018-2020

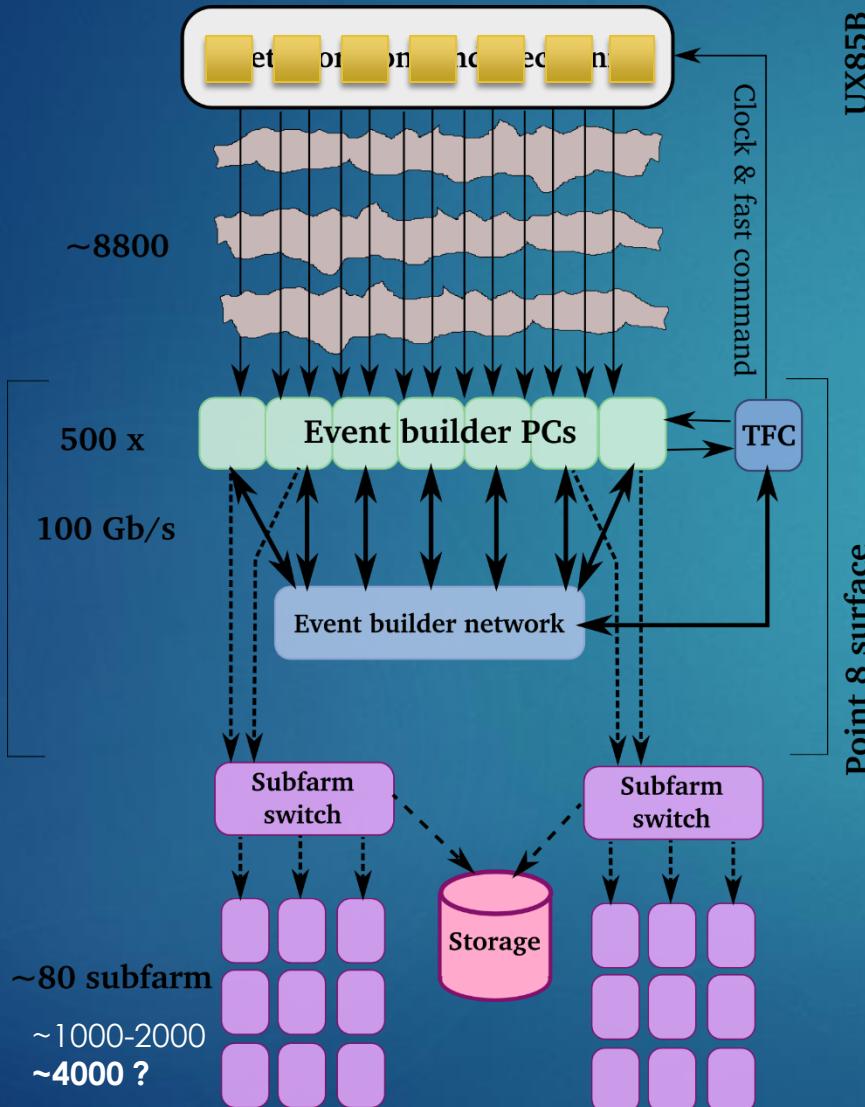
- ▶ Update of sub-detectors
- ▶ Removal of hardware trigger
 - ▶ Currently in **FPGA** and **analogic** components
 - ▶ Have to answer in **real time**
 - ▶ **Hard to maintain** and update
- ▶ New status :
 - ▶ Larger **event rate** (1 Mhz to **40 Mhz**)
 - ▶ Larger **event size** (50 KB to \sim 100 KB)
- ▶ Much more data for DAQ & Trigger ($x80 \Rightarrow 40 \text{ Tb/s}$)



40 Tb/s for LHCb in 2020



Data flow



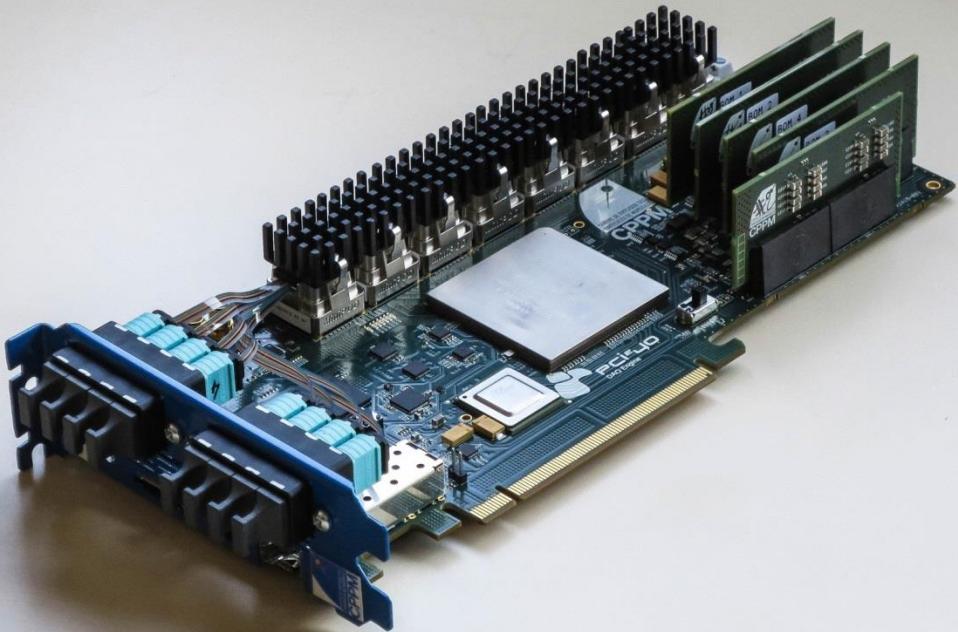
► Numbers

- ~8800 optical links going out from detector to the surface (~100 m) and up to ~4.8 Gb/s each.
- ~500 readout nodes (up to 24 input links each)
- Filter farm of O(1000) nodes
- 40 Tb/s to transport
 - Can run over 512 nodes
 - Need a 100 Gb/s network
 - Considering net. load at 80%
 - So 80 Gb/s
 - Current estimation is 70 Gb/s

Starting point : the acquisition board

55

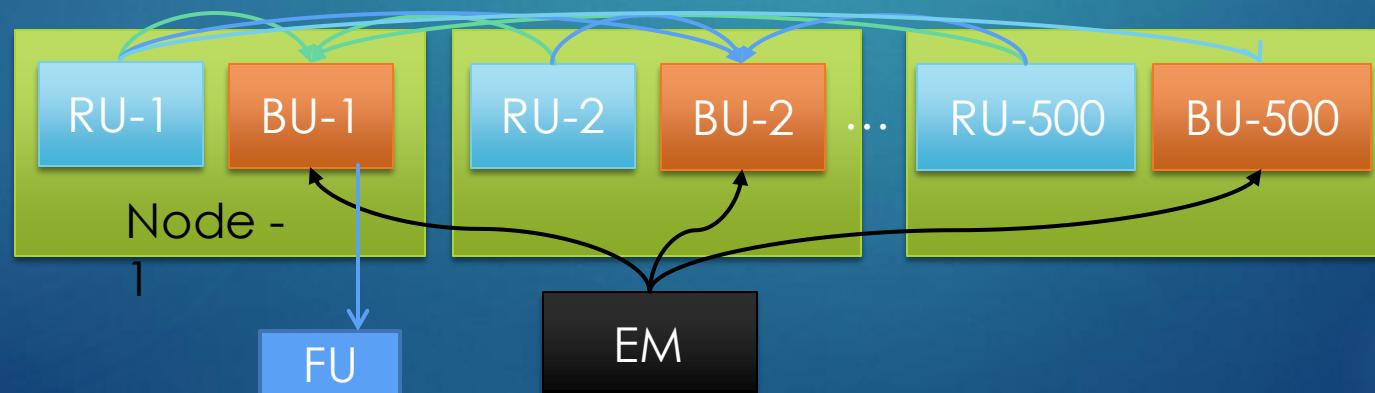
Studying a 40 Tb/s DAQ - Sébastien Valat
29/01/2018



Up to 48 fibers and 100 Gb/s

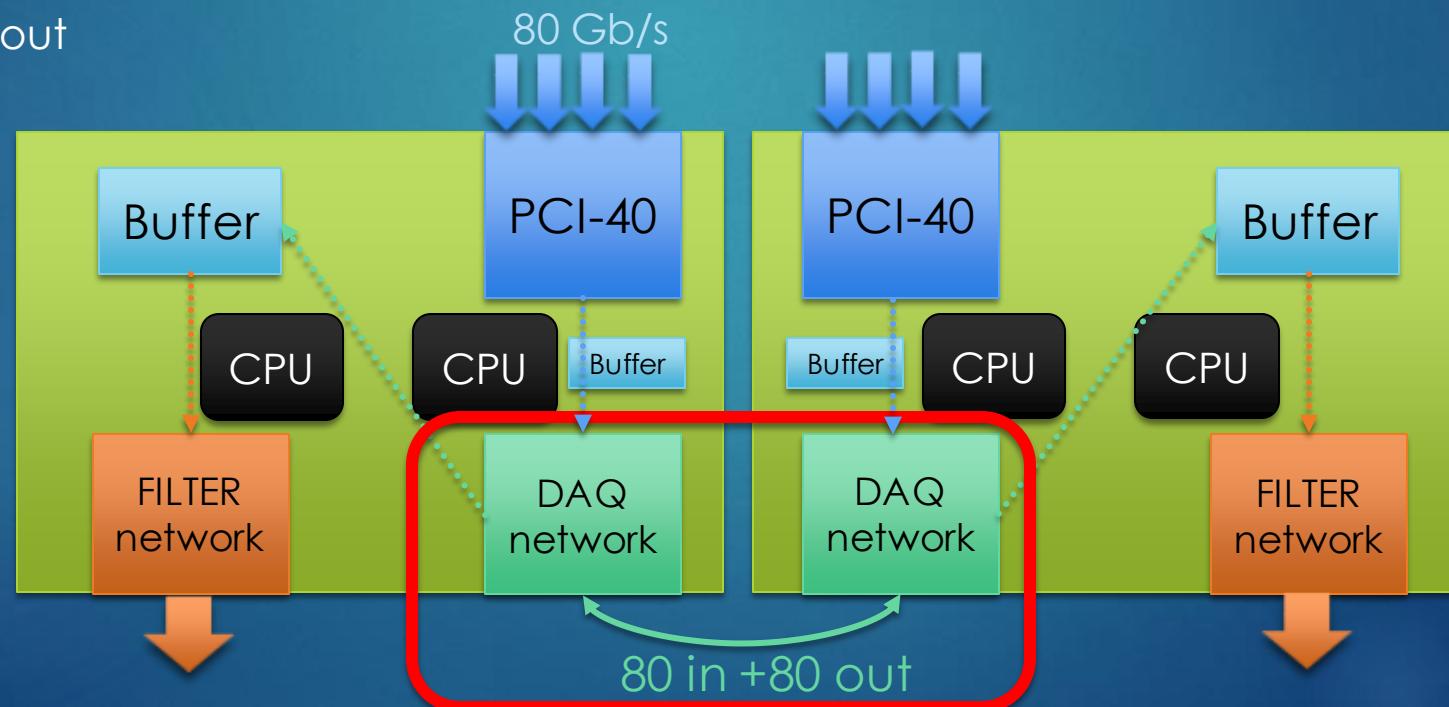
Event building network topology

- ▶ Work with 4 units:
 - ▶ **Readout unit (RU)** to **read** data from PCI-E readout board
 - ▶ **Builder unit (BU)** to **merge** data from **Readout unit** and send to **Filter unit**.
 - ▶ **Filter unit (FU)** to select the interesting data (**not considered here**)
 - ▶ **Event manager (EM)** to dispatch the work over **Builder unit** (credits)
- ▶ RU/BU mostly does a “all-to-all”
 - ▶ To **aggregate** the data chunks from each collision



IO nodes hardware

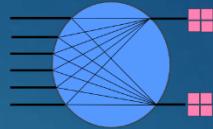
- ▶ Three IO boards at 100 Gb/s per node:
 - ▶ PCI-40 for fiber input
 - ▶ Event building network
 - ▶ Output to filter farm
- ▶ Also stress the memory (320 Gb/s of total traffic)
 - ▶ 80 in + 80 out



The DAQ network technologies

- ▶ We expect to use HPC networks
 - ▶ 100 Gb/s (expect running at 80 Gb/s)
- ▶ Technologies we looked on:
 - ▶ **Infiniband** EDR (HDR 200 Gb/s ?)
 - ▶ Intel **Omni-path** (~~version 2, 200 Gb/s ?~~)
 - ▶ 100 Gb/s **Ethernet**
- ▶ Not as HPC apps
 - ▶ We need to **fully fill the network continuously**
 - ▶ In **full duplex**
 - ▶ Bad pattern : **many gathers** (all-to-all) !
- ▶ Think using a fat-tree topology



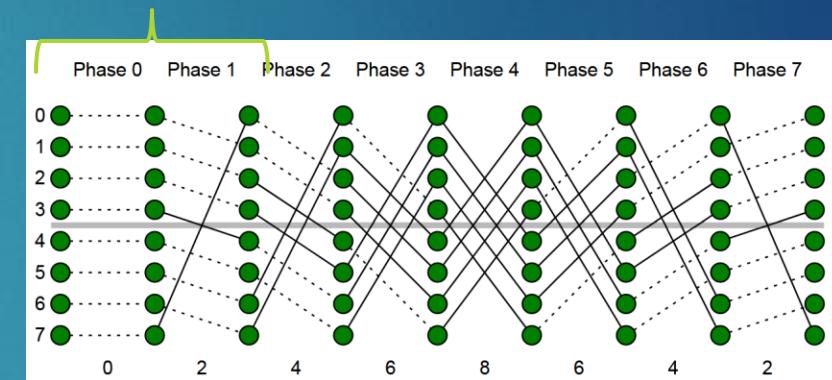


DAQPIPE

- ▶ A benchmark to evaluate event building solutions
 - ▶ DAQ Protocol Independent Performance Evaluator
 - ▶ Provides EM/RU/BU units
- ▶ Support various APIs
 - ▶ **MPI**
 - ▶ Libfabric
 - ▶ PSM2
 - ▶ Verbs
 - ▶ TCP / UDP
 - ▶ RapidIO
- ▶ Multiple protocols
 - ▶ PUSH
 - ▶ **PULL**

Communication control parameters

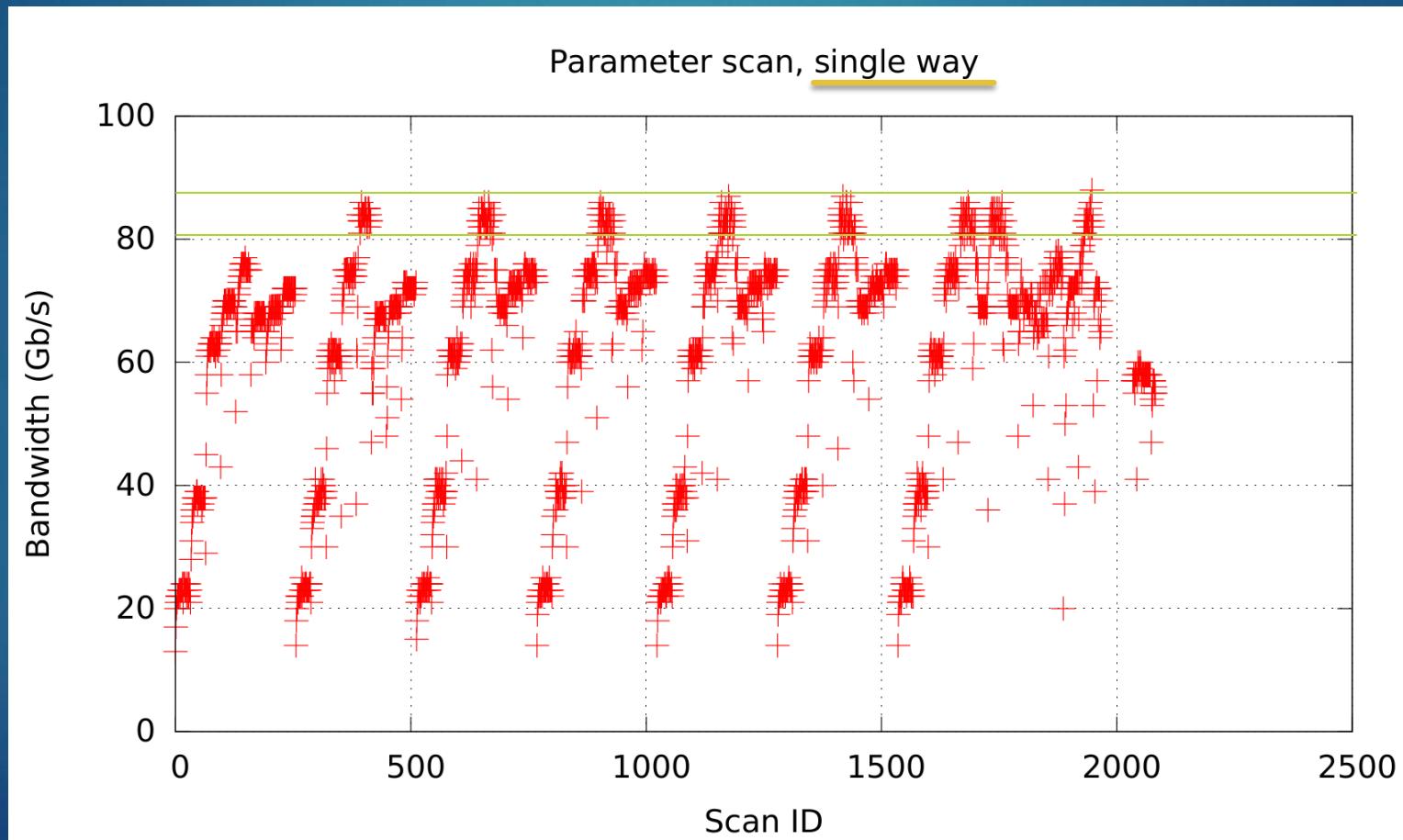
- ▶ Manage communication scheduling models
 - ▶ **Barrel shift** ordering (with N on-fly)
 - ▶ **Random** ordering (with N on-fly)
 - ▶ Send **all in one go**
- ▶ Messages size
 - ▶ best ~**512 KB**
- ▶ Parallel gathers (credits)
 - ▶ best ~**2**
- ▶ Parallel communication per gather
 - ▶ With how many nodes we communicate
 - ▶ best ~**8**
 - ▶ So in total **2*8** communication in flight
- ▶ Process per nodes, 1 or 2
 - ▶ best **2**



* From Bandwidth-optimal all-to-all exchanges in fat tree networks

What we have seen here on our cluster with 16 nodes OPA

87 Gb/s = **2%** of the points

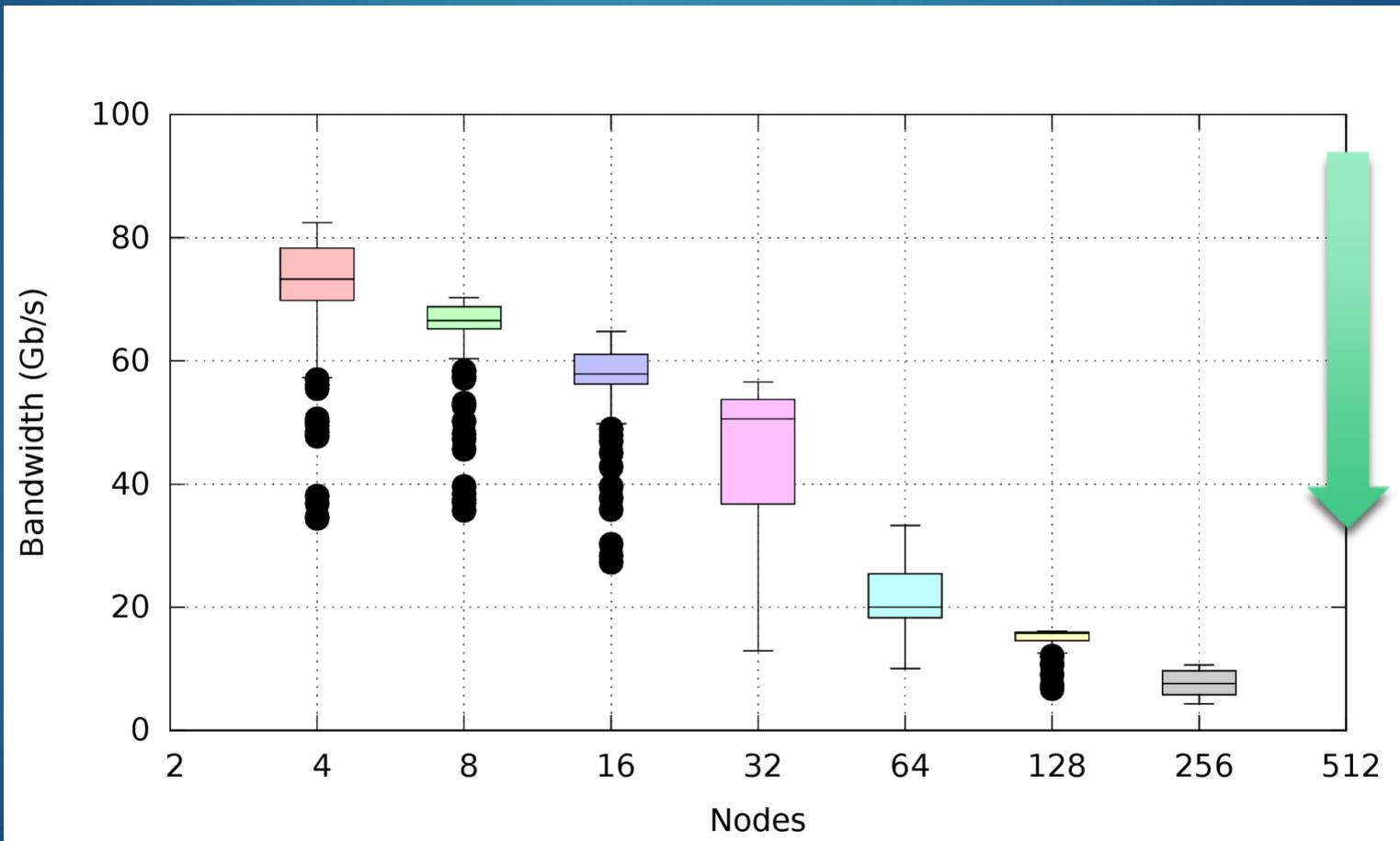




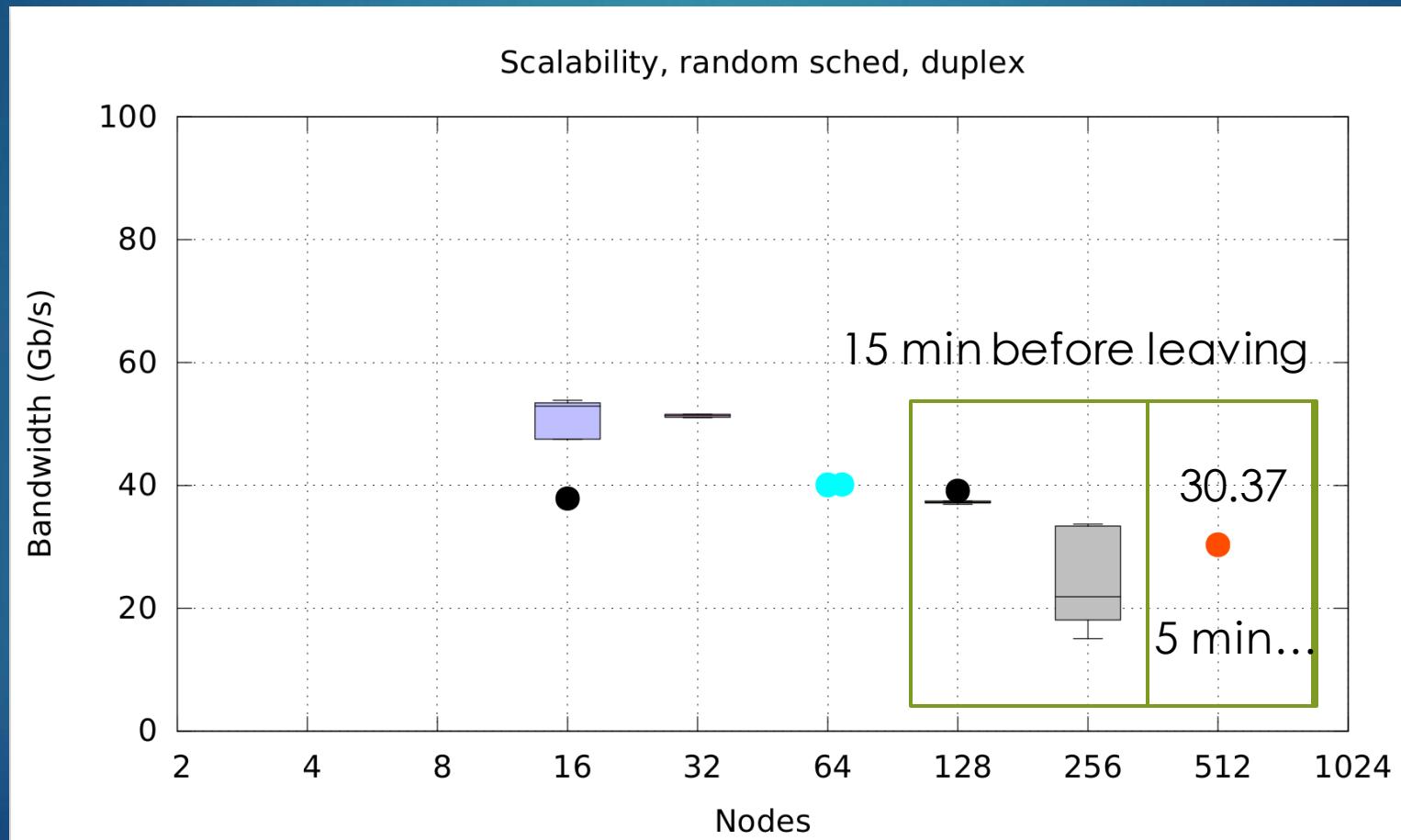
From <http://www.flickr.com/photos/tags/supercomputer/interesting/>

First test at scale (512):
hitting the wall ☹

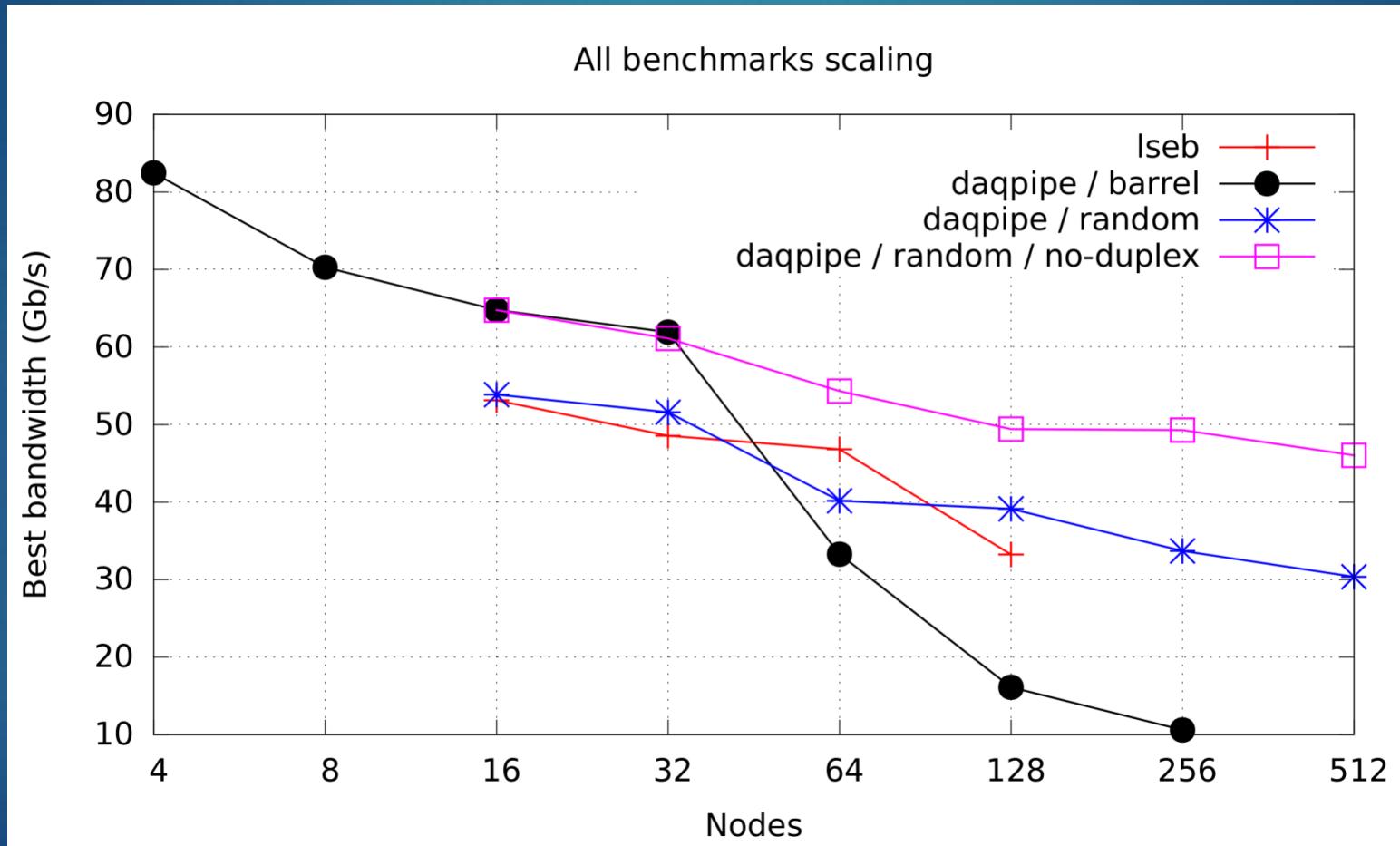
Summary scalability on Day-1 😞



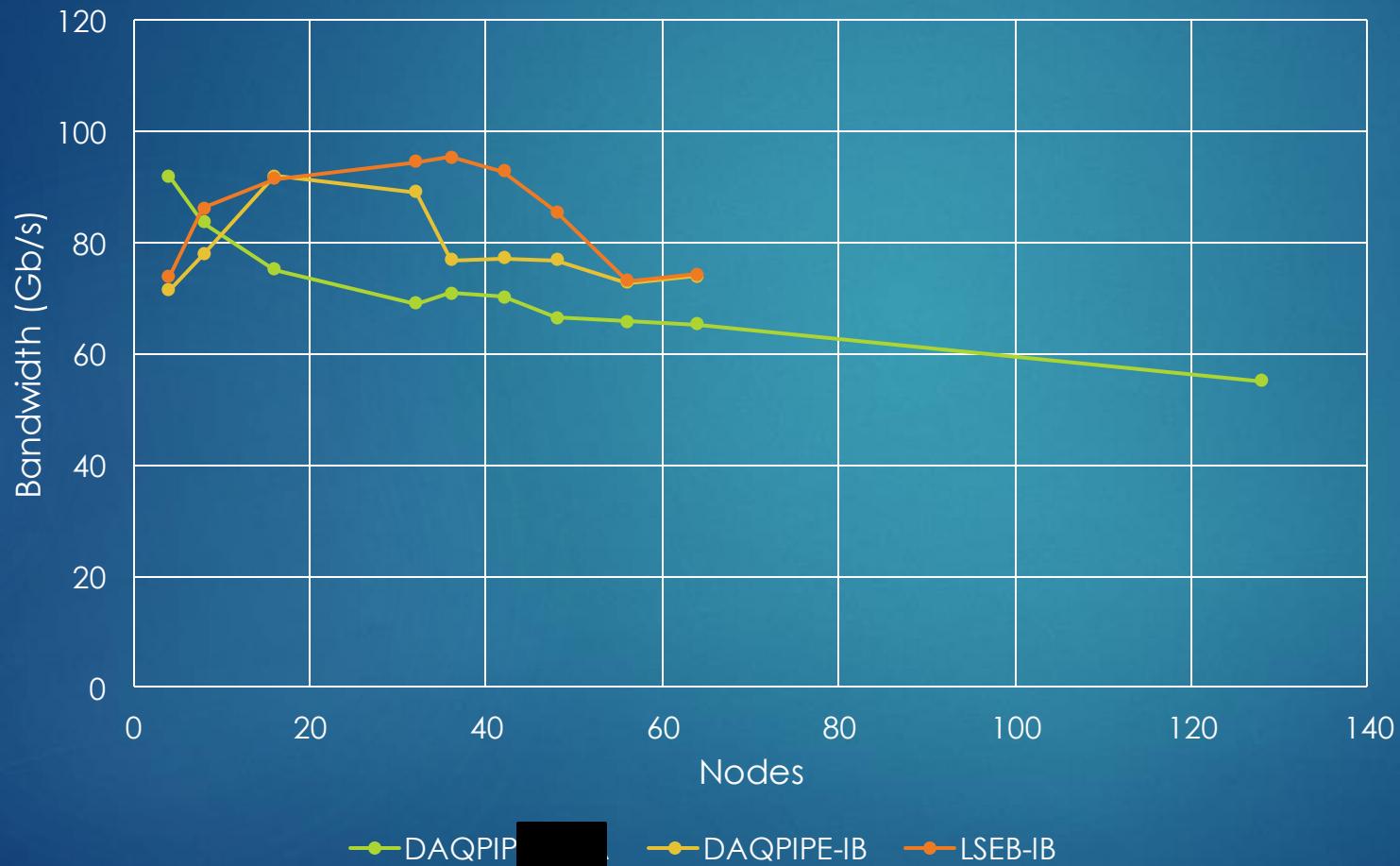
Day-2, 2 hours before leaving try random scheduling



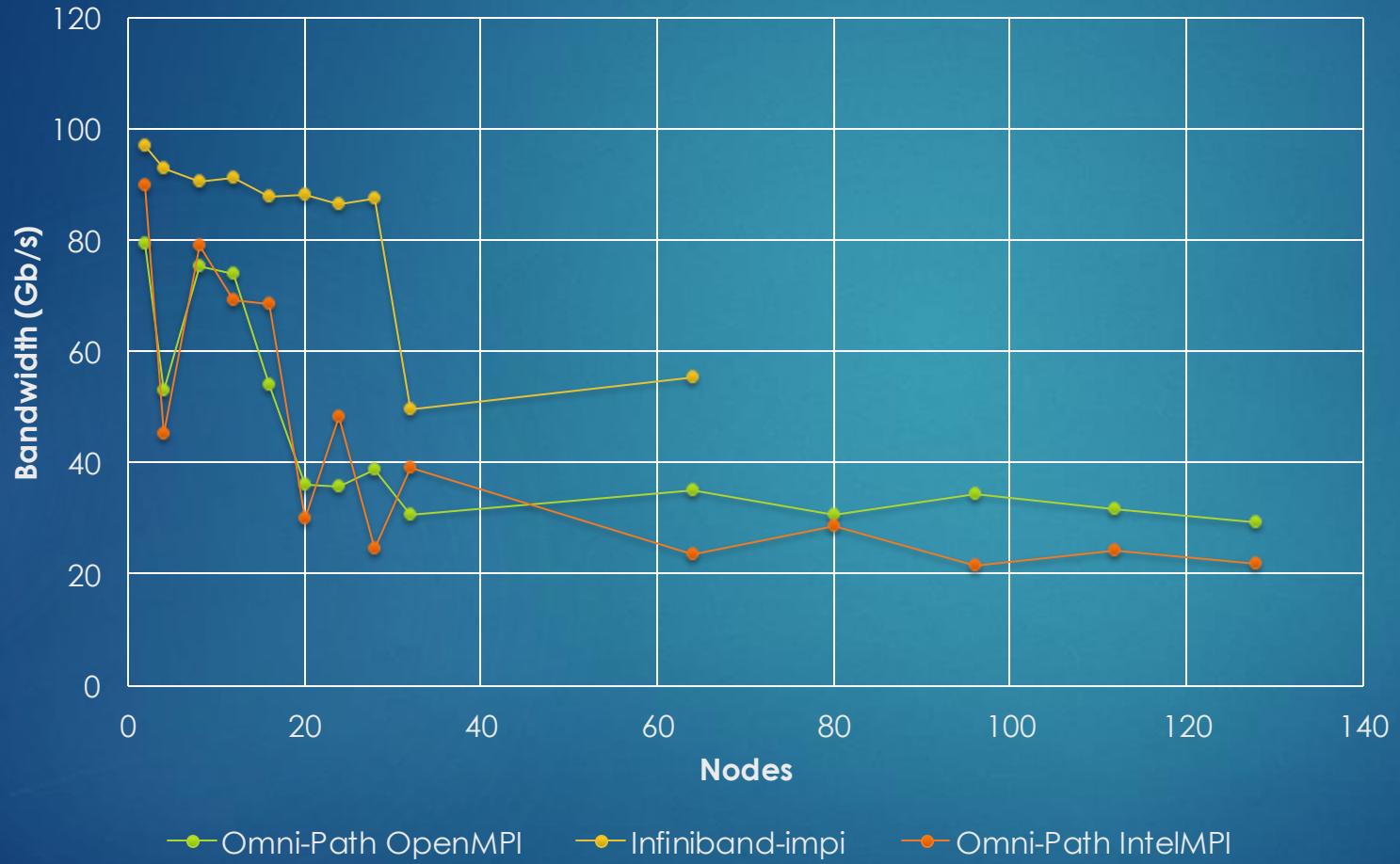
Final summary



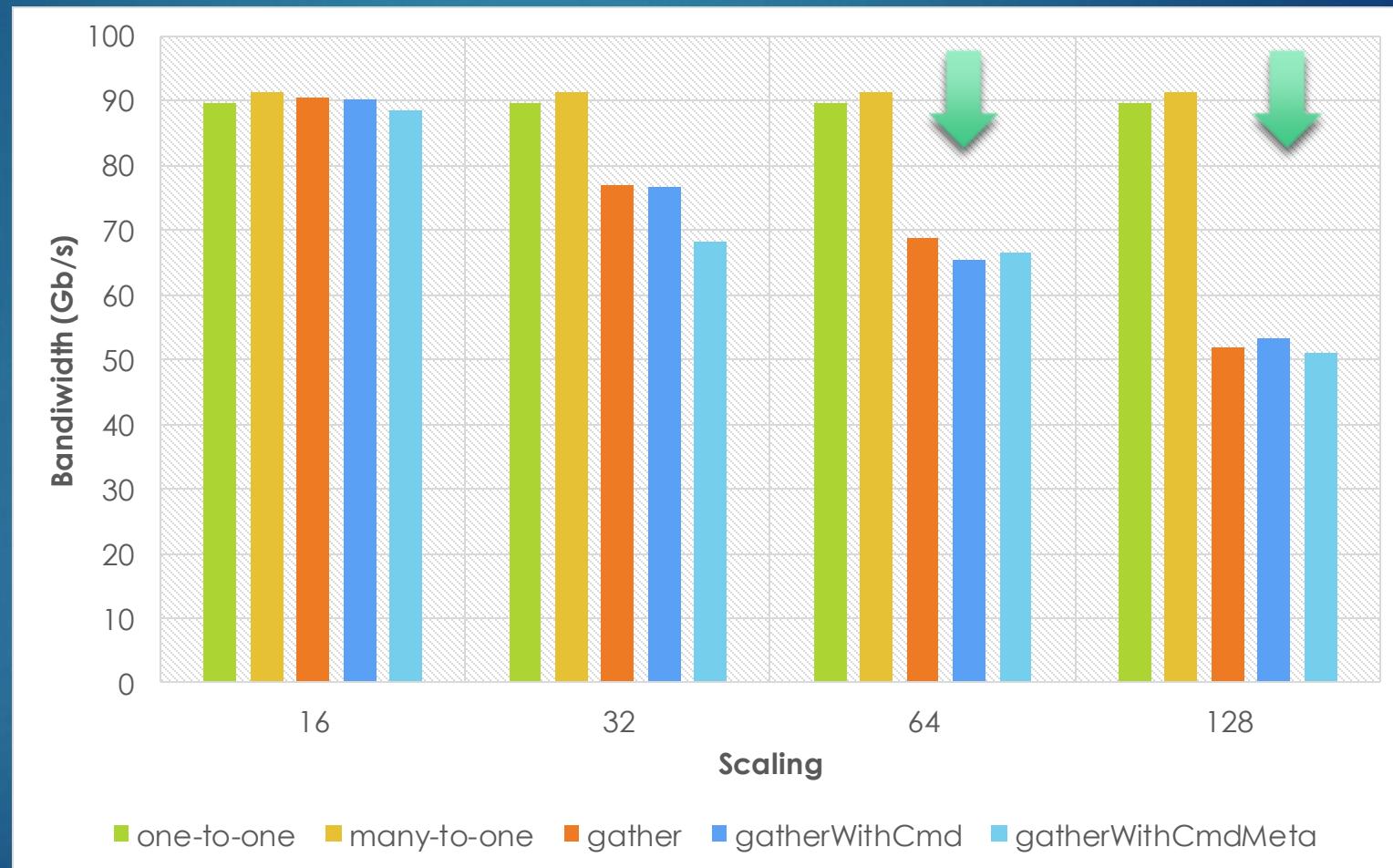
Another test on two clusters



MPI_Alltoall

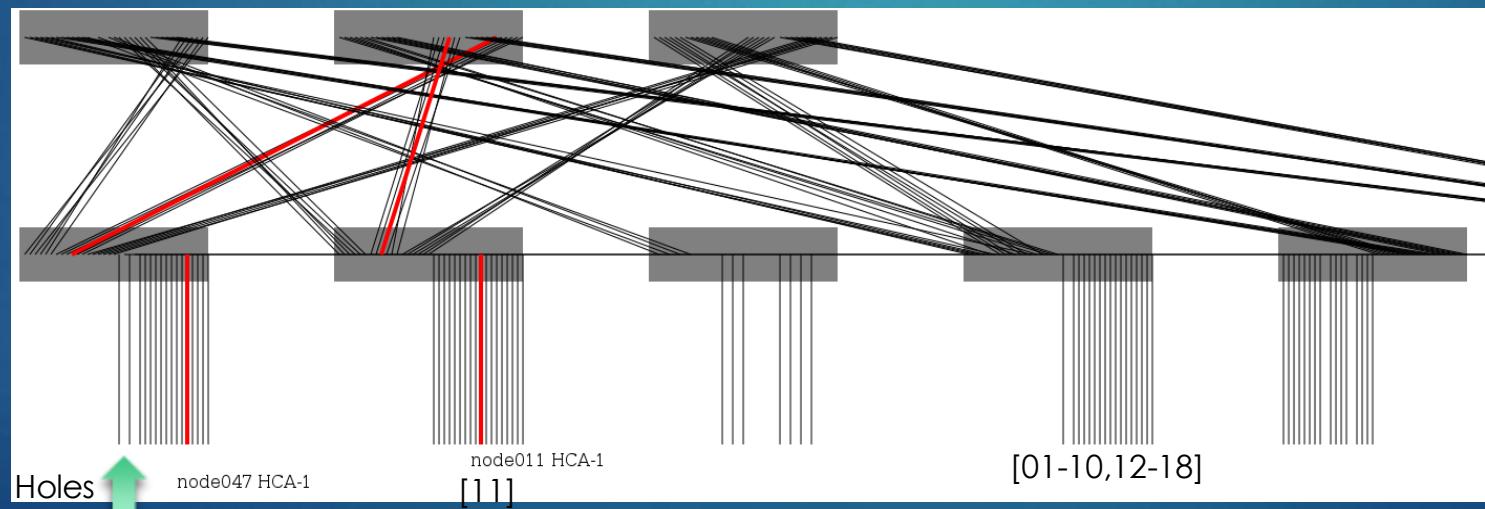
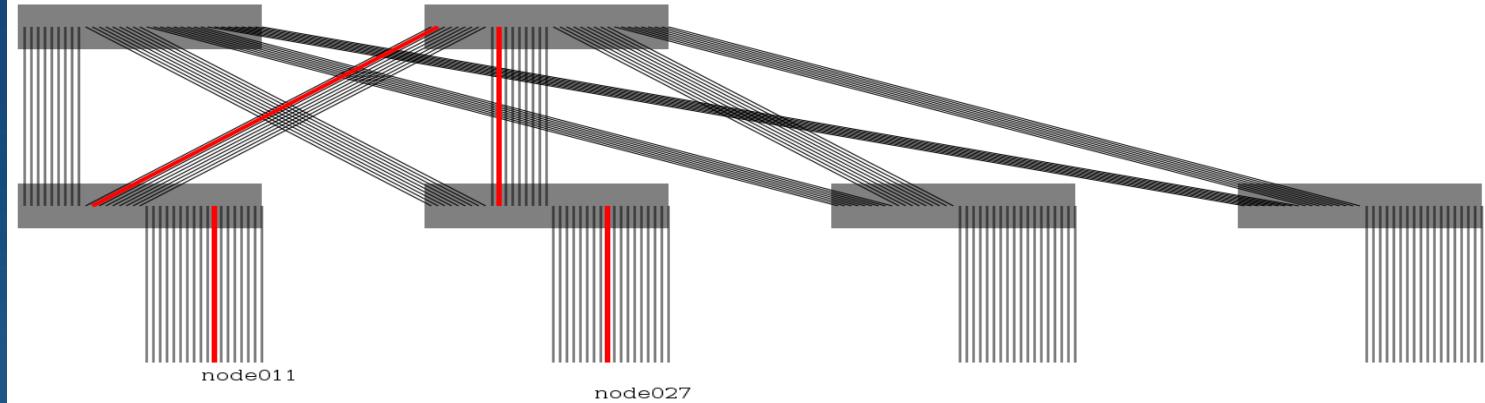


Micro-benchmarks on OPA



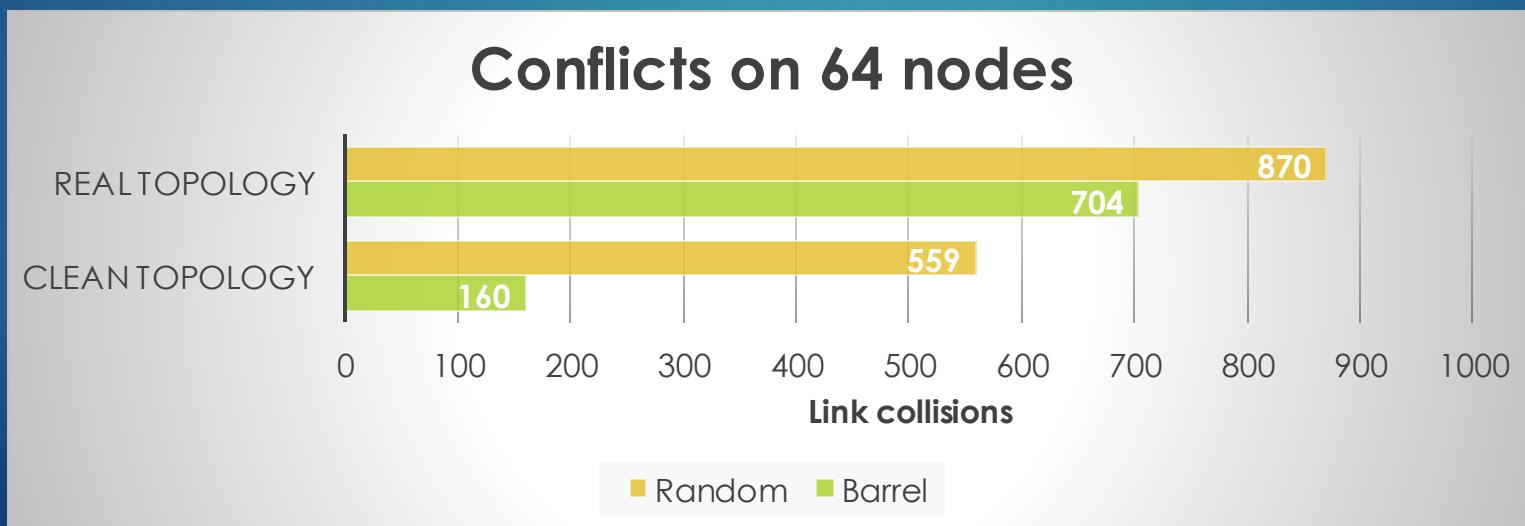
Ideal and real topologies

Studying a 40 Tb/s DAQ - Sébastien Valat
29/01/2018



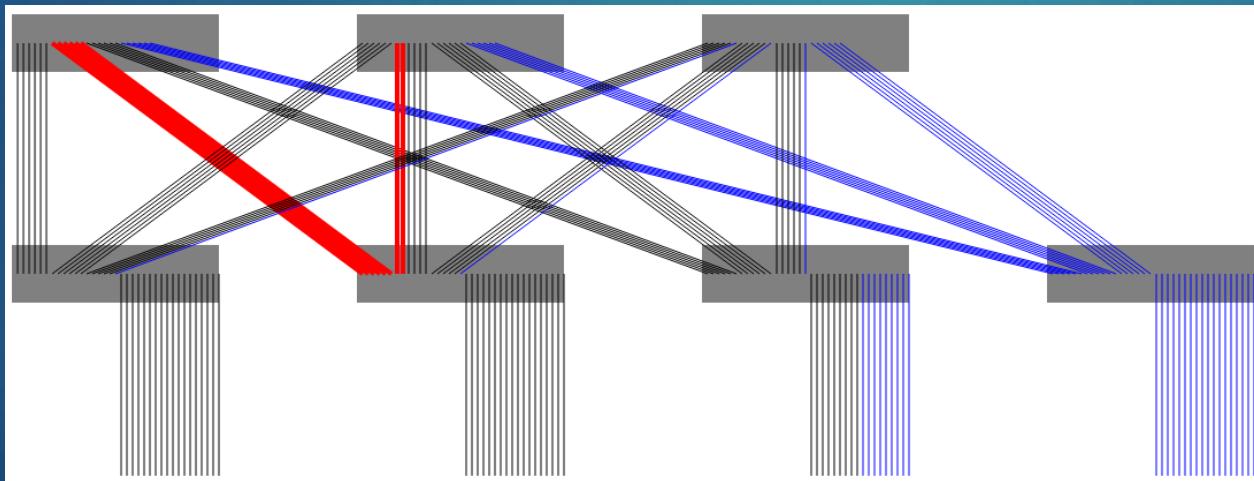
Conflicts

- ▶ The real topology implies many conflicts
- ▶ Due to :
 - ▶ routing table
 - ▶ missing nodes in leaf switches
- ▶ 64 nodes imply 4096 communications
 - ▶ 870 conflicts => **21%**, 704 conflicts => **17%**, 160 conflicts => **3%**

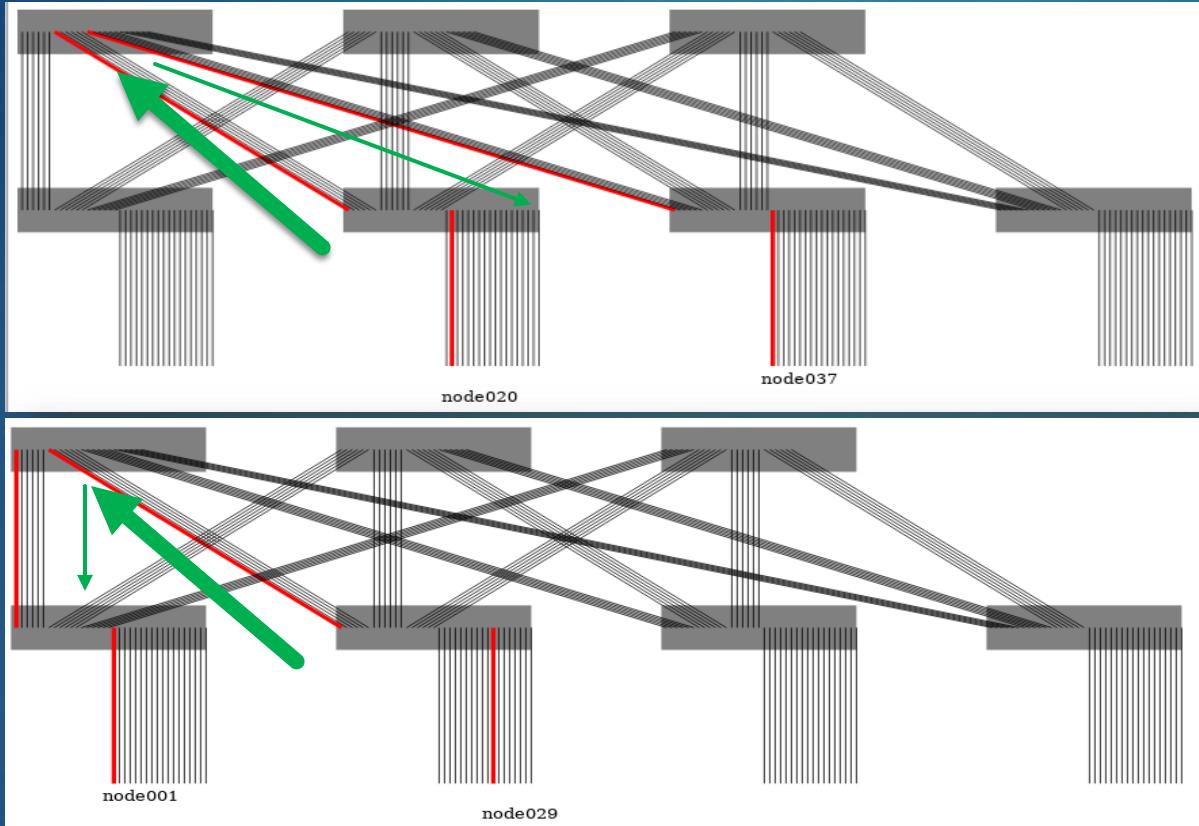


Example on clean topology

- ▶ Switch **radix : 36**
- ▶ Run on **45 nodes** ($2 * 18 + 9$)
- ▶ Look on **step 28 => 9 conflicts**



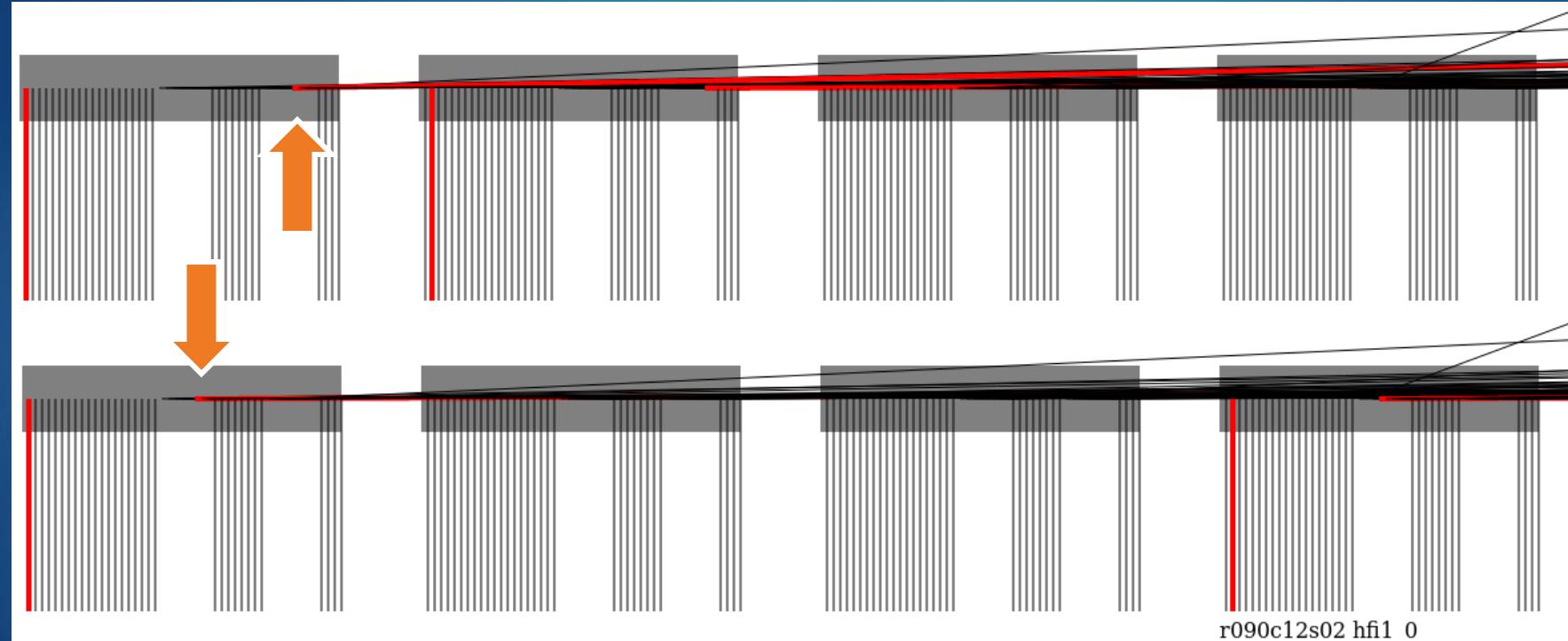
Extract two conflicts to understand



- ▶ Example, **2 conflicts of step 28**
 - ▶ **20 -> 37**
 - ▶ **29 -> 1**
- ▶ They use the same **UP** link

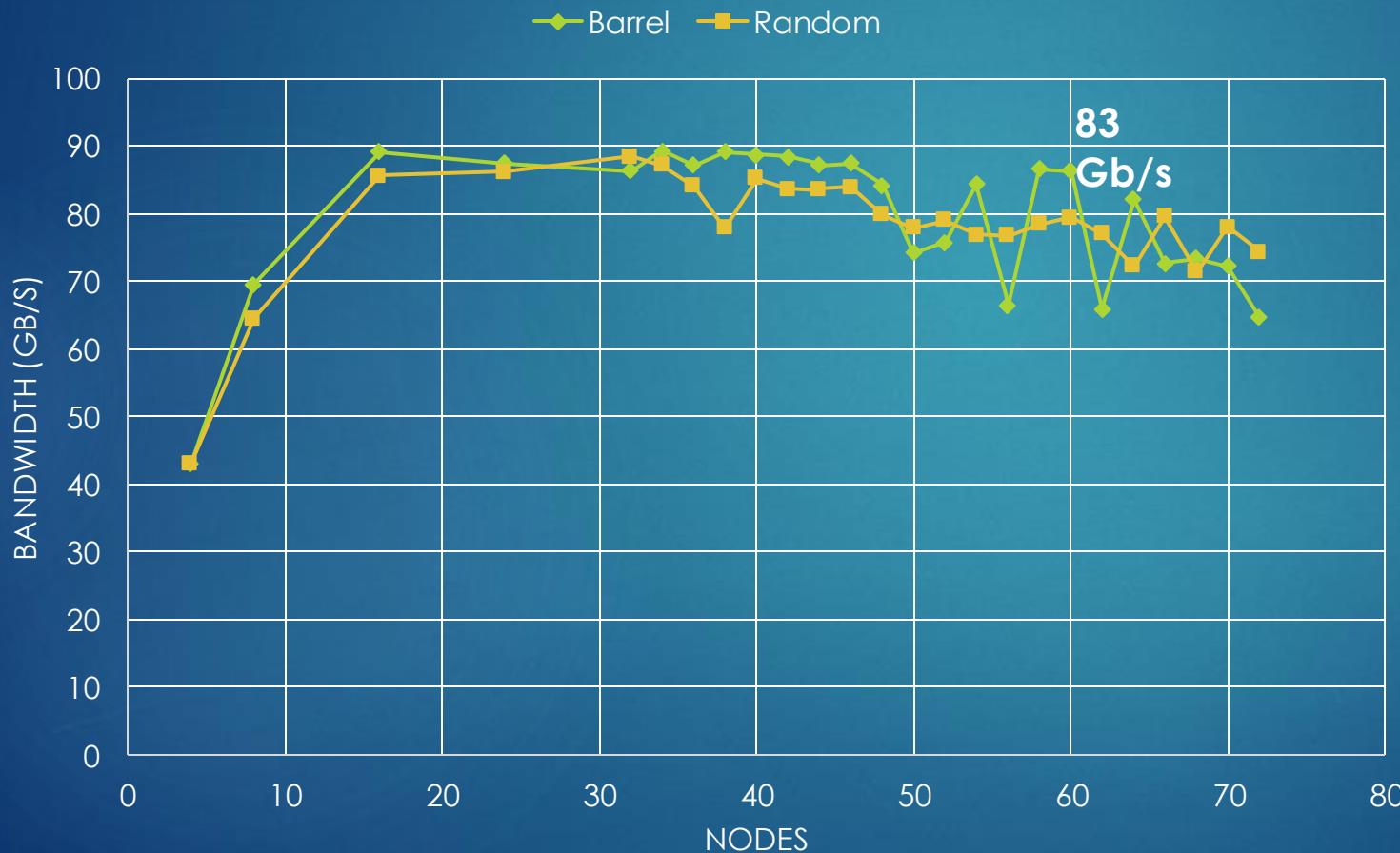
Default routing tables OPA

Not same up-link used to reach N-eme node on leaf switches

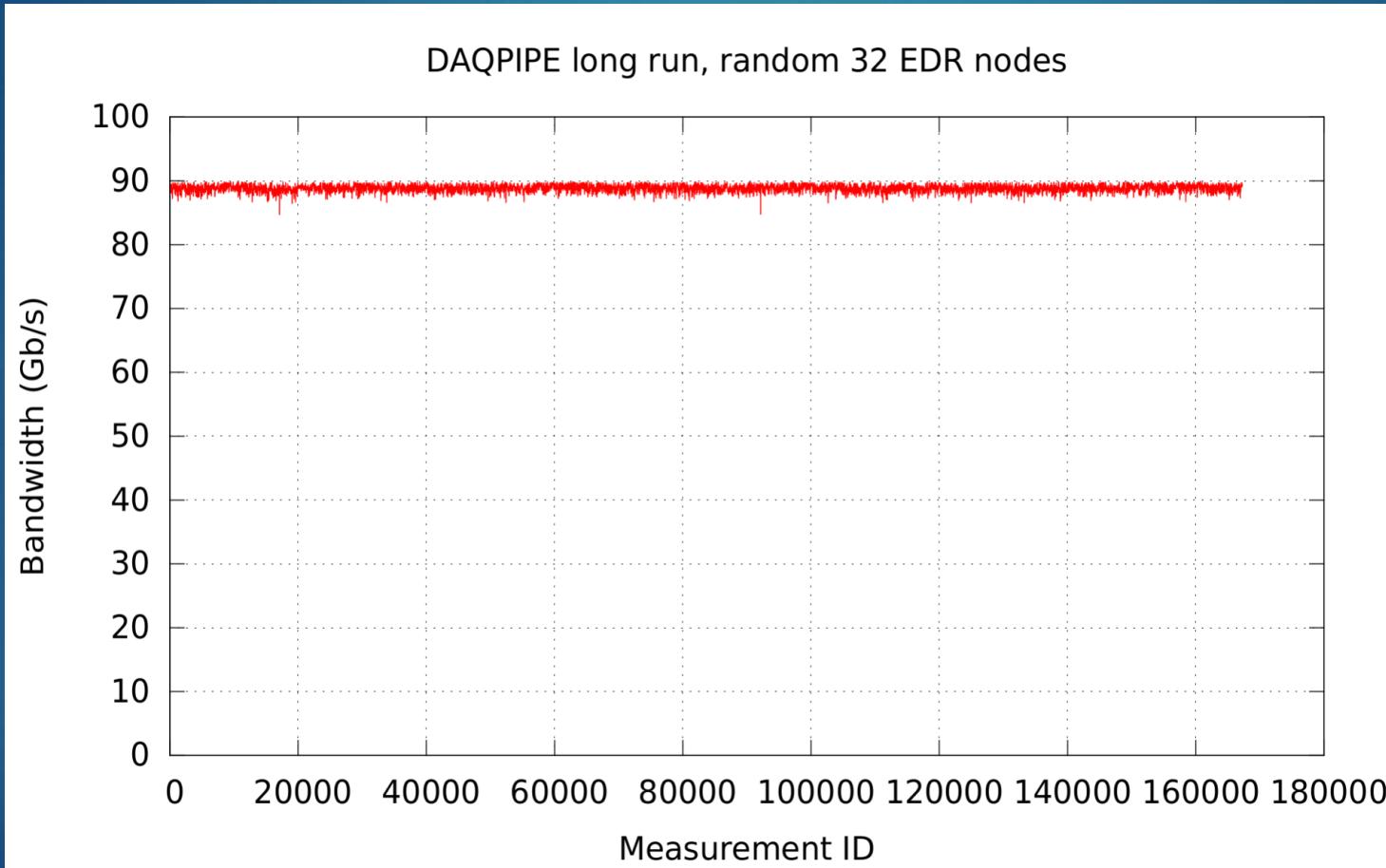


Last results on Tech B (LENOVO)

This time with close to perfect cabling

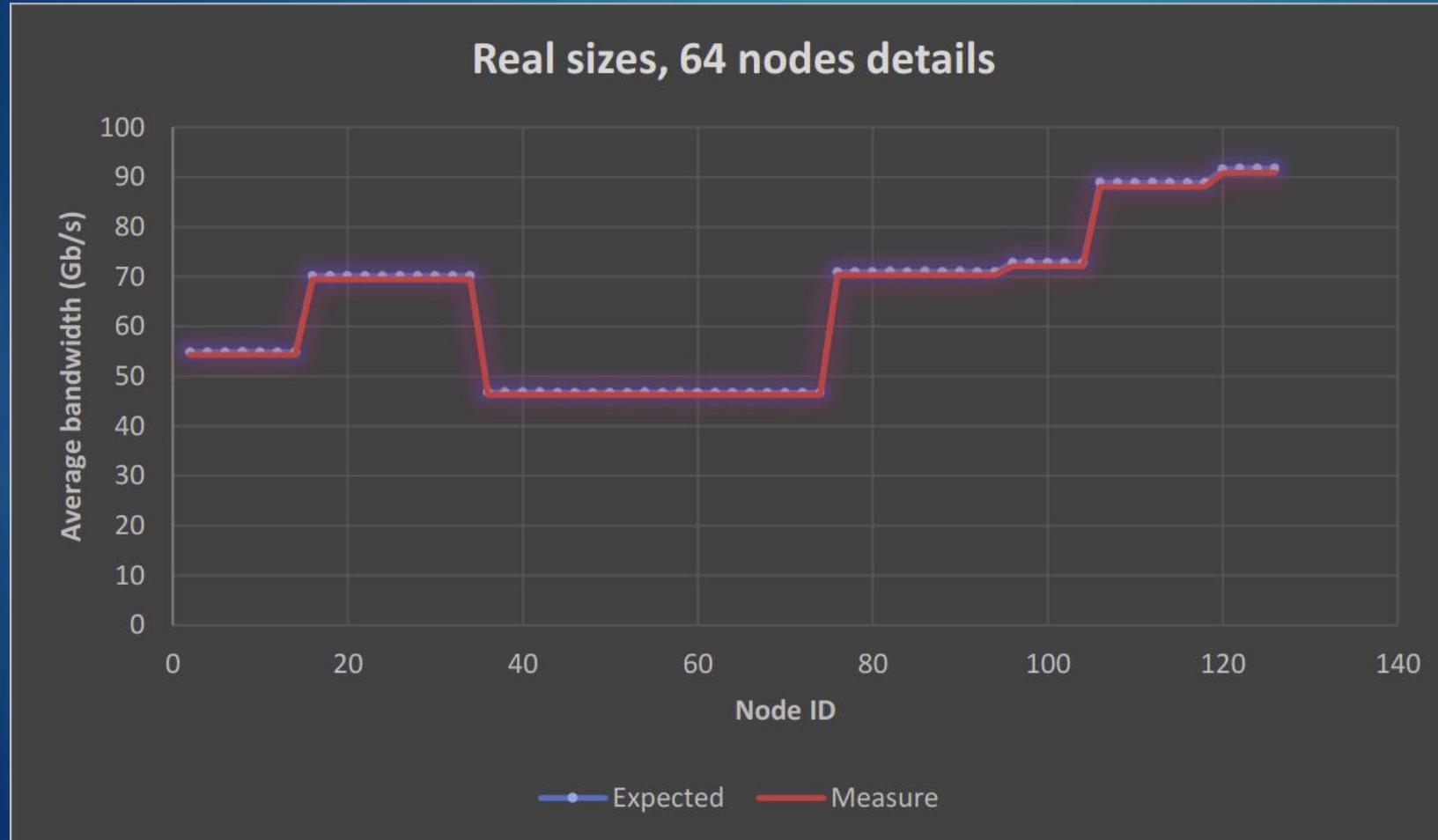


Stability over time (3 hours)



Real unbalanced **sub-detectors** bandwidth

Finally that's should be OK with IB/EDR.... but not so much margins



Conclusion

- ▶ That was a long road
- ▶ Finally IB/EDR should be OK
- ▶ Would sleep better if IB/HDR (200 Gb/s) will be available in time
- ▶ OPA need work on routing (intel did it, need scale tests)
- ▶ We also ave failure recovery
- ▶ Is needed ? MPI restart ~30s ?
- ▶ Current prospective about 100Gb/s ethernet (via iwap/RoCE)

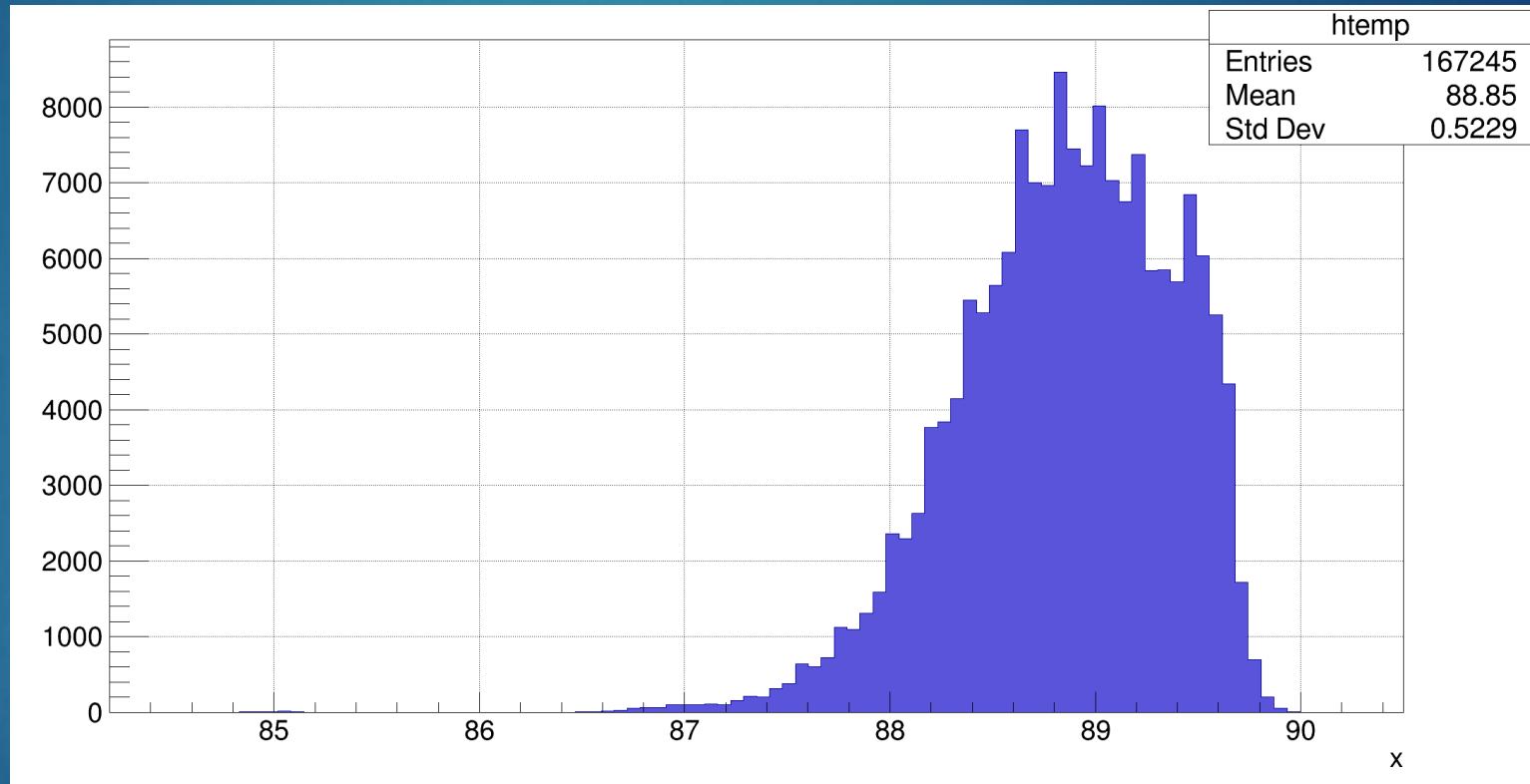
Backup

Htopml integration

Communication scheduling



Stability over time (3 hours)



RDMA

Remote Direct Memory Access

- › **Communication with Ethernet**
 - Traverse all **IP stack** in OS
 - Internal **memory copies** (use mem. **bandwidth & CPU**)
 - Higher latency
- › **OS bypass in InfiniBand or Omni-Path**
 - **RDMA** : Remote Direct Memory Access
 - No memory copies
 - **Lower latency**
 - **CPU not involved**

