

A JOURNEY
OVER THE MEMORY MANAGEMENT
STACK
FOR HPC LARGE APPLICATIONS
ON MODERN ARCHITECTURES

Work coming from



www.cern.ch

Sébastien Valat

IXPUG - 25 sept 2019

Plan

Introduction

- I. Analysis of OS paging policy
- II. NUMA allocator for HPC applications
- III. Cost of first touch handler
- IV. MALT & NUMAPROF memory profilers
- V. Conclusion

INTRODUCTION

The “new” memory context

- 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 **per compute** (cores) **power** ?



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



https://de.wikipedia.org/wiki/Datei:PS2_RAM_Module.jpg

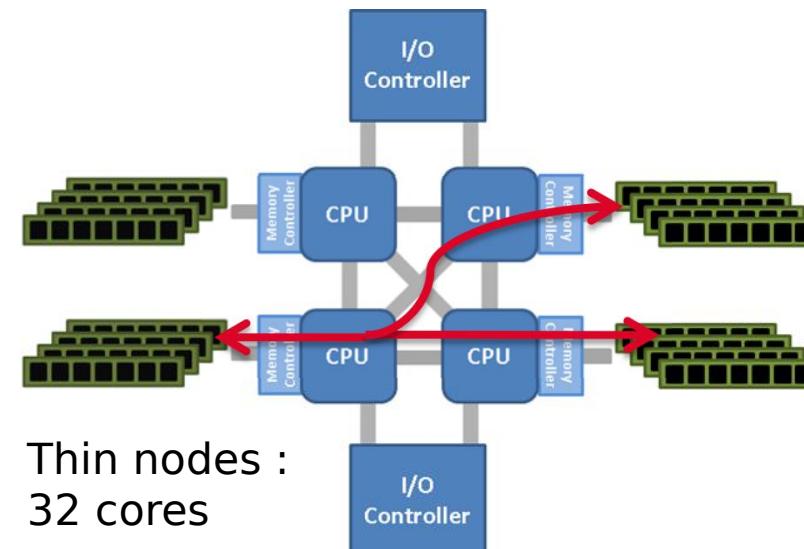
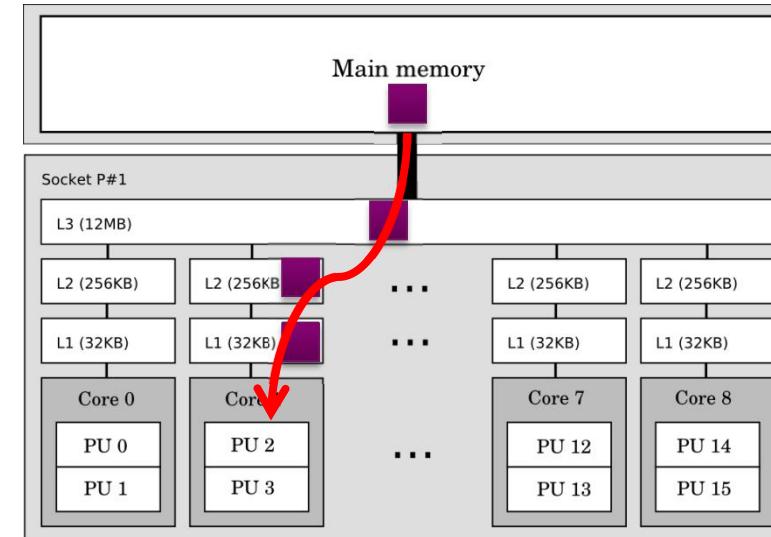
Complex memory hierarchies

- Caches

- NUMA

- MCDRAM ?

- NVMe DIMMs ?



Software memory management layer

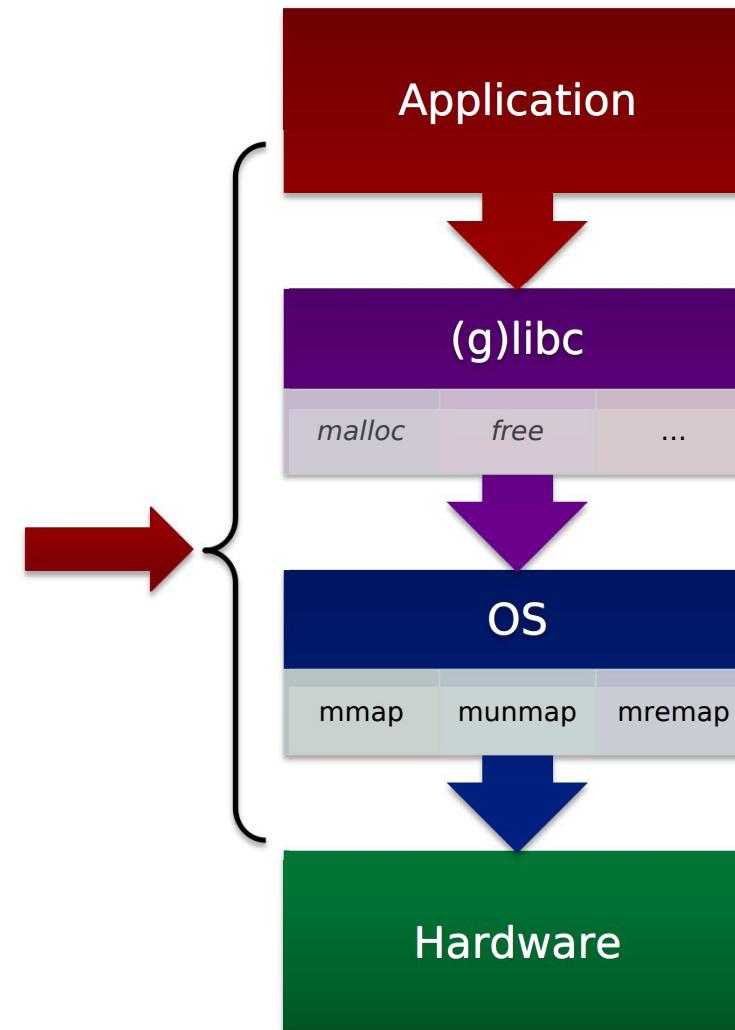
- Impact of memory management mechanisms ?

- Involving two components :

- User space **memory allocator** (malloc)
 - **Operating System** (OS)

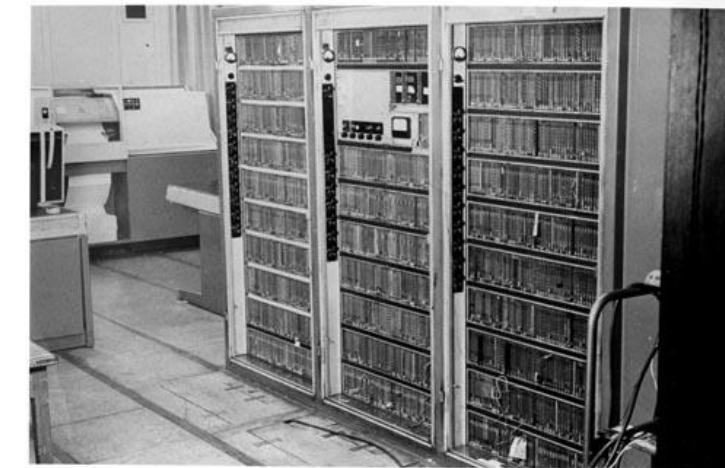
- Focus on :

- Impact on **allocation time**
 - Impact on **access efficiency** (placement)

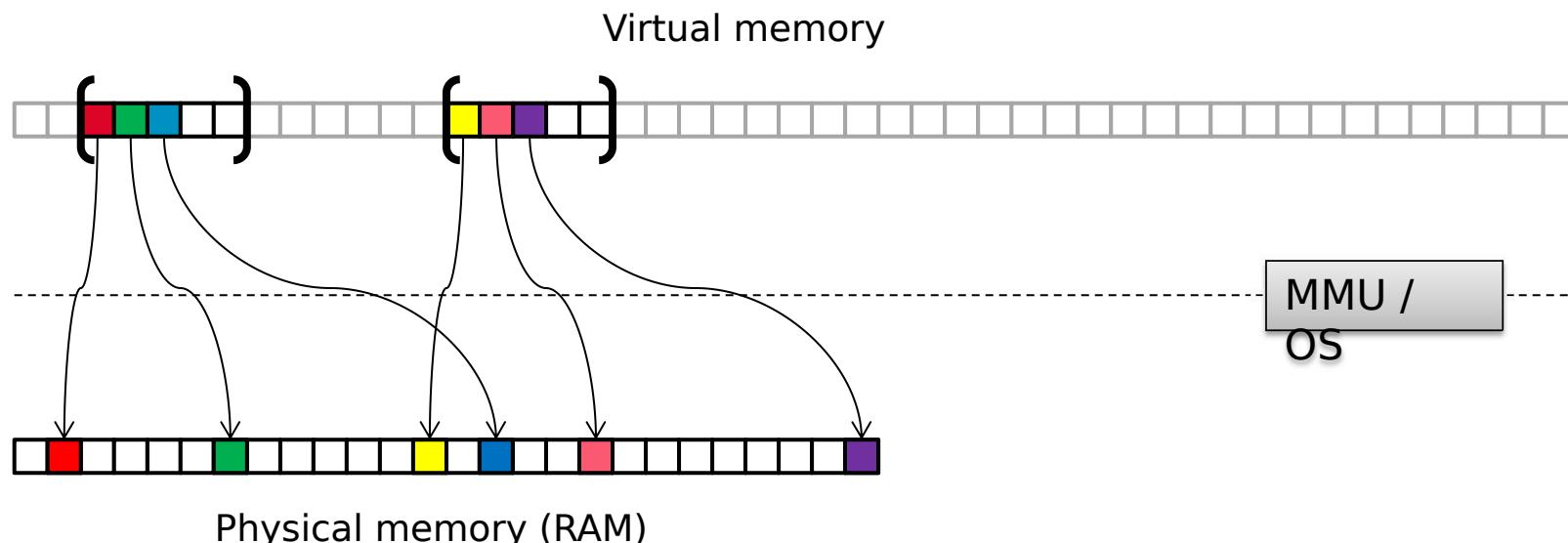


OS virtual / physical address spaces

- Two address spaces : **physical + virtual**
- **Paging** was first used in **1962** on the **ATLAS computer**
- **Area creation with syscalls** : **mmap / munmap / mremap**
- **Malloc** has the responsibility to **hide the pages to developers**

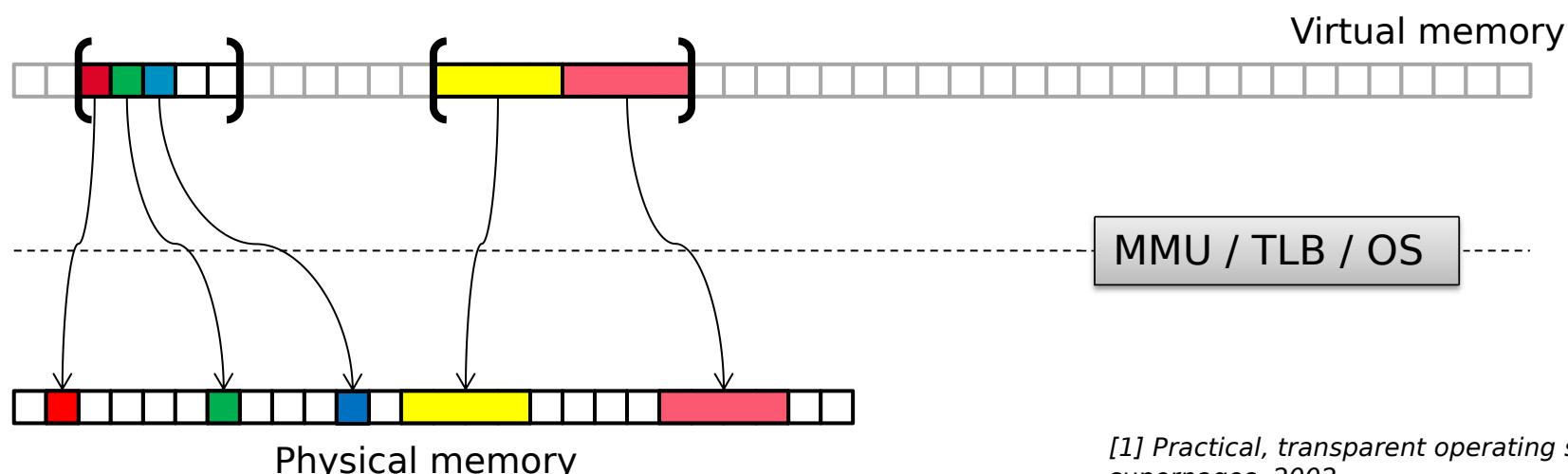


<http://www.computerhistory.org/collections/catalog/102698470>



Huge pages

- Huge pages : 2 MB
- First real support : FreeBSD (superpages, 2002) [1]
- Support Linux : *old HugeTLBfs* then now **Transparent Huge Pages (THP), 2011**

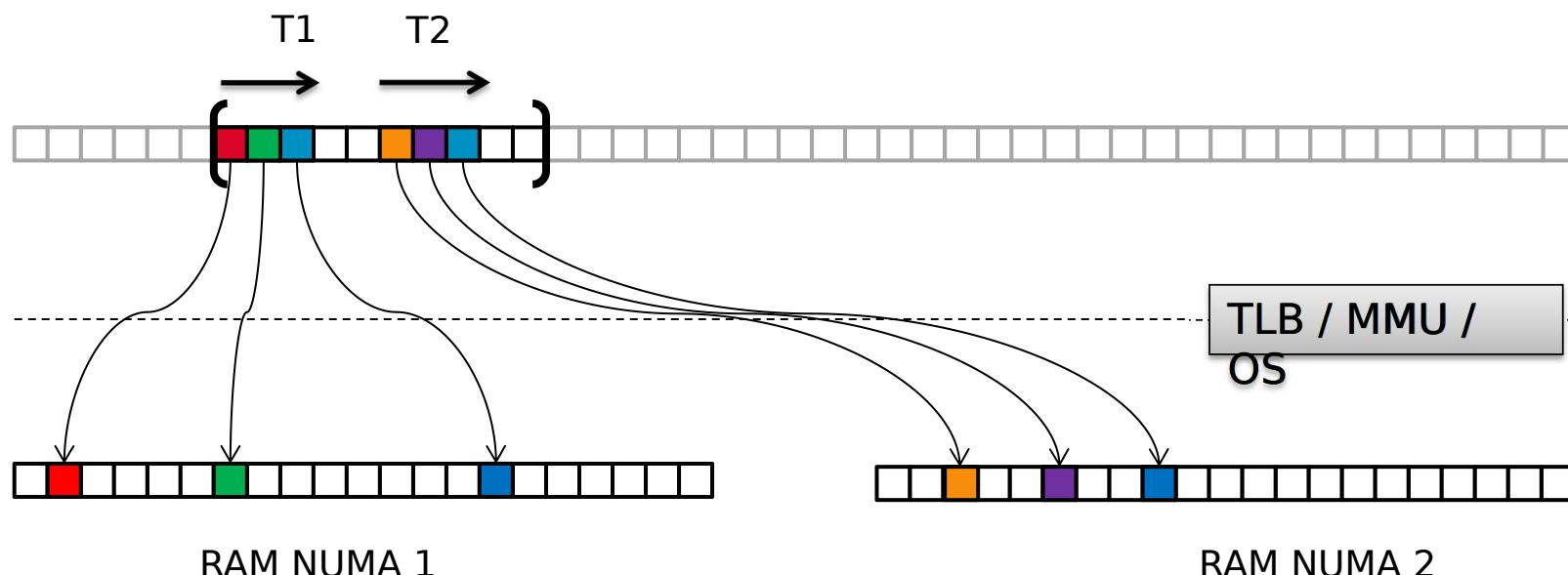


[1] Practical, transparent operating system support for superpages, 2002

Lazy page allocation

- **mmap** creates **pure virtual** area
- First touch creates a **page fault** for each virtual page
- OS provides **physical pages** on **first touch**
- First touch implicitly determines **NUMA placement** of the page

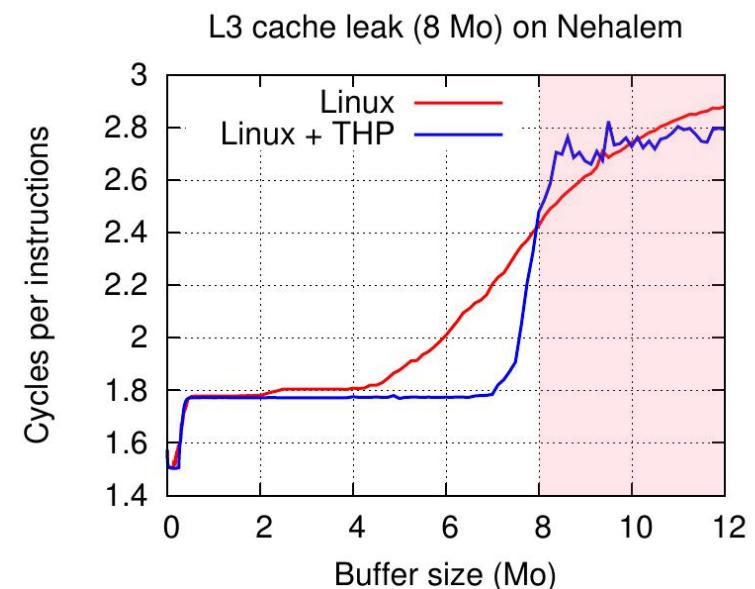
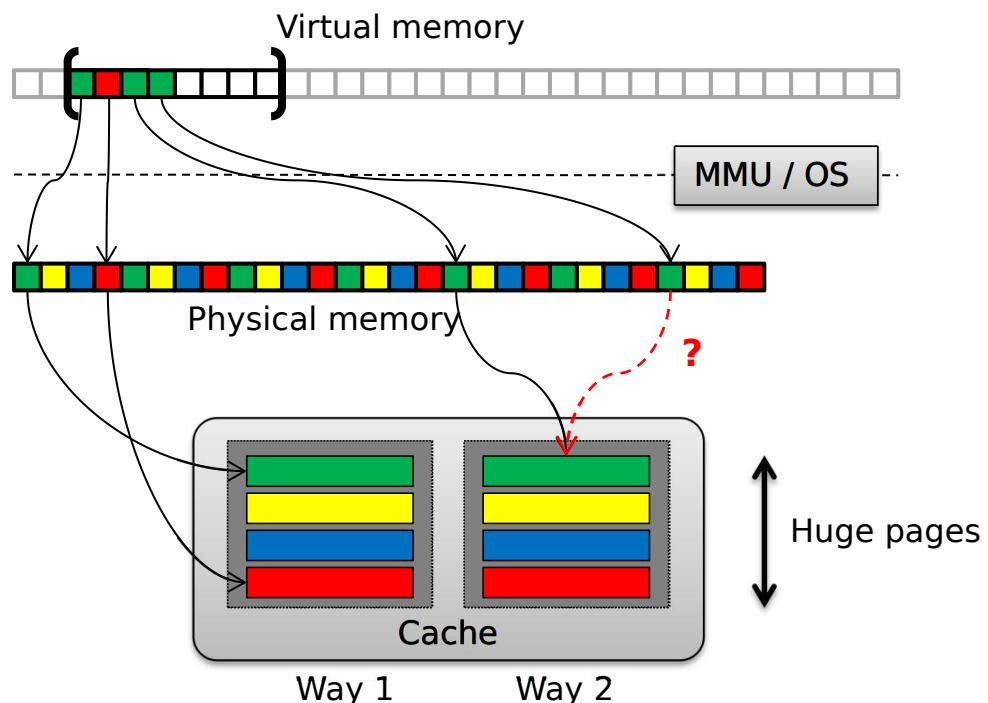
```
ptr = mmap(...,SIZE,...);
#pragma omp parallel for
for (i = 0 ; i < SIZE ; i++)
    ptr[i] = 0;
```



ANALYSIS OF OS PAGING POLICY

Cache associativity

- Data can only be placed in one of the **N lines associated to the address**
- Can create **conflicts** depending on the OS
- Linux “**randomly** chooses” the pages



OS strategies comparison (2010)

- Each **system** has its default paging **strategy**:

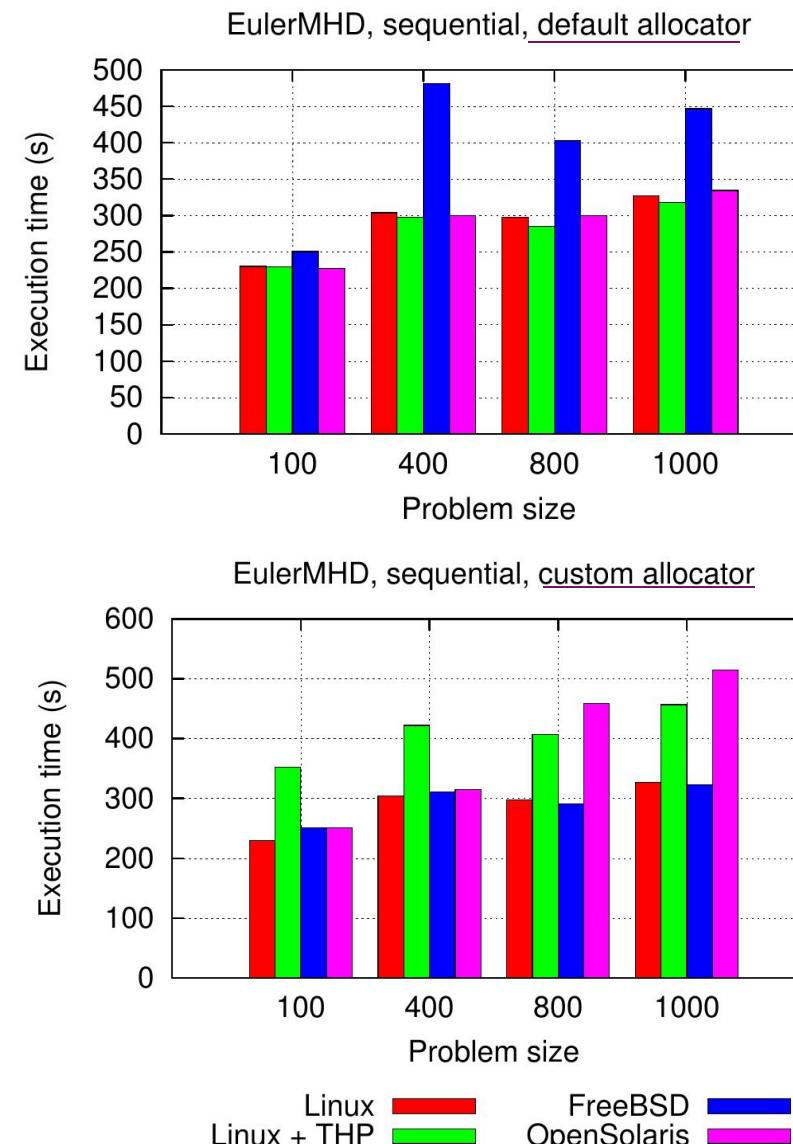
OS	Strategy
Linux	4K random
OpenSolaris	Page coloring
FreeBSD	Huge pages

- Is **Linux** slower due to **random paging** ?
- Tested architecture : Intel **Nehalem bi-socket**
- Use a fixed compile chain : **GCC/Binutils/MPI/BLAS**
- Focus a pathological case

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**.



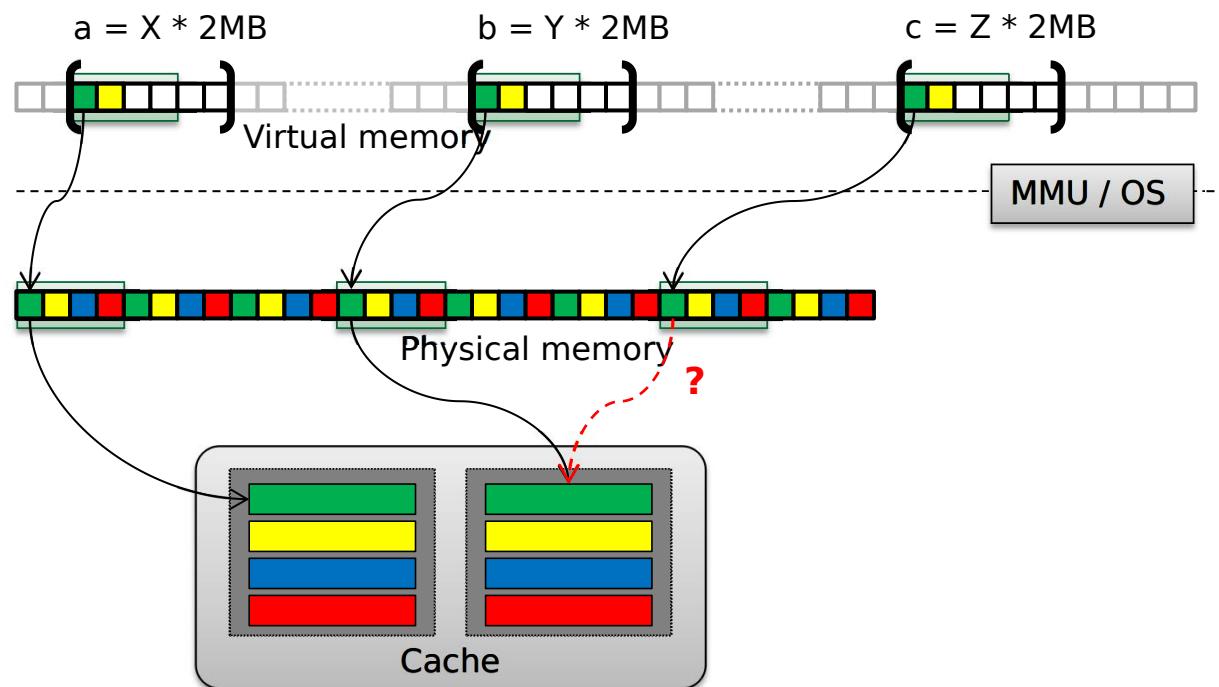
Alignment effect on regular coloring

- Each **malloc** (OS) produces different **alignments**

```
for (i = 0 ; i < SIZE ; i++)  
    a[i] = b[i] + c[i];
```

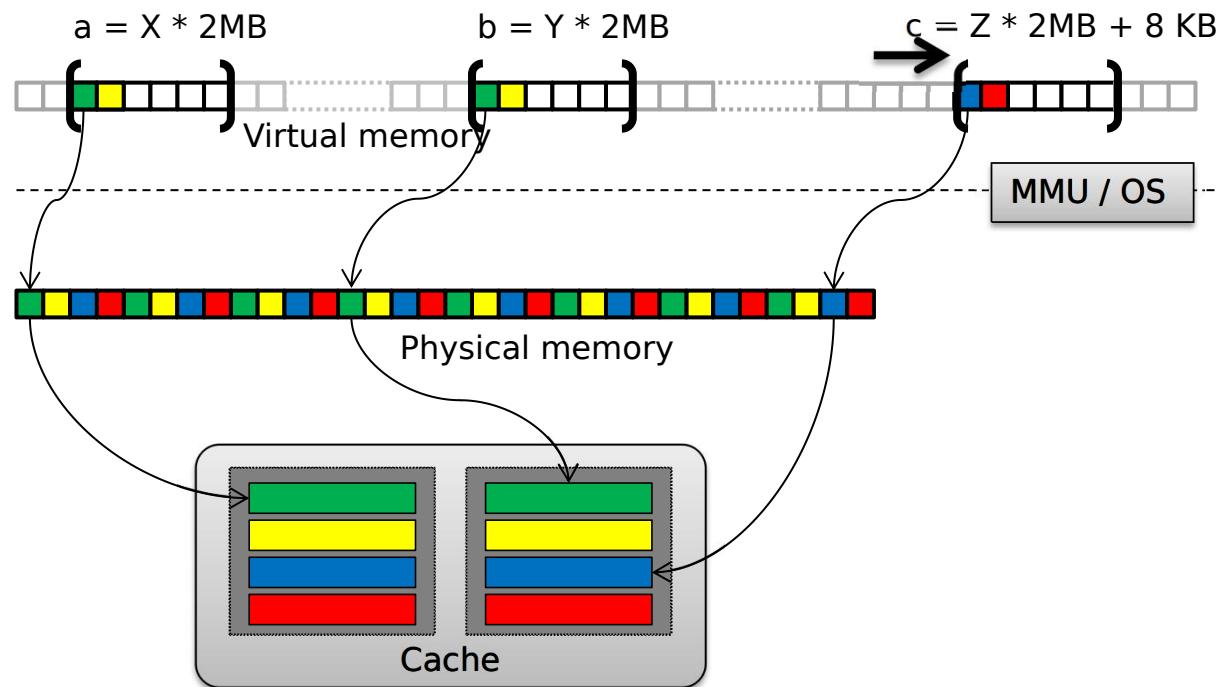
- FreeBSD align **large segments** on 2 MB

- It interferes with **regular patterns** generated by :
 - OpenSolaris coloration method (modulo)
 - Huge pages



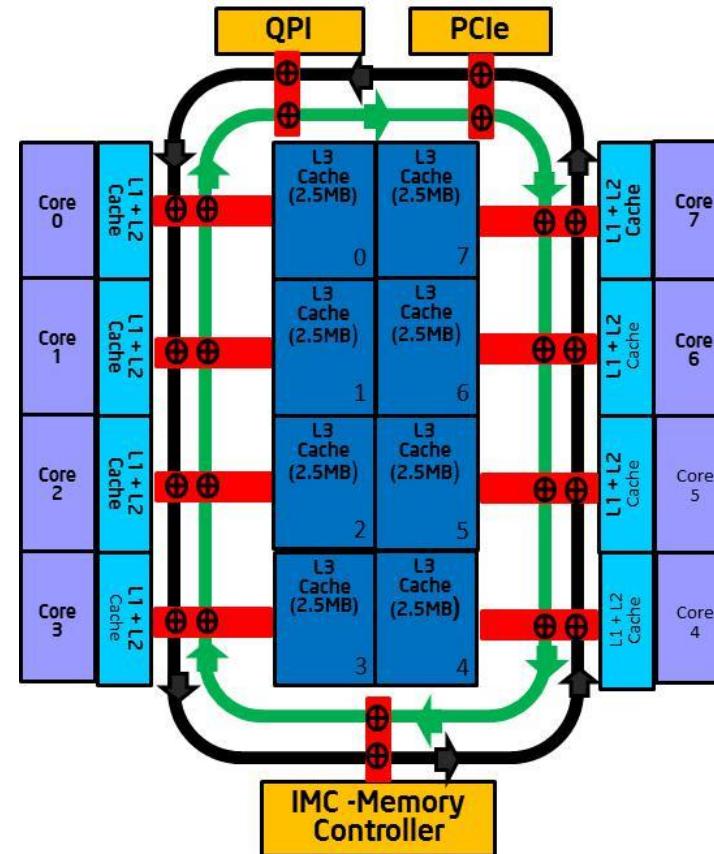
Solution

- Avoid segment **alignments** on **cache way size** (mmap / malloc).
- The **Linux random** approach **prevents pathological cases**



New intel L3 cache slices

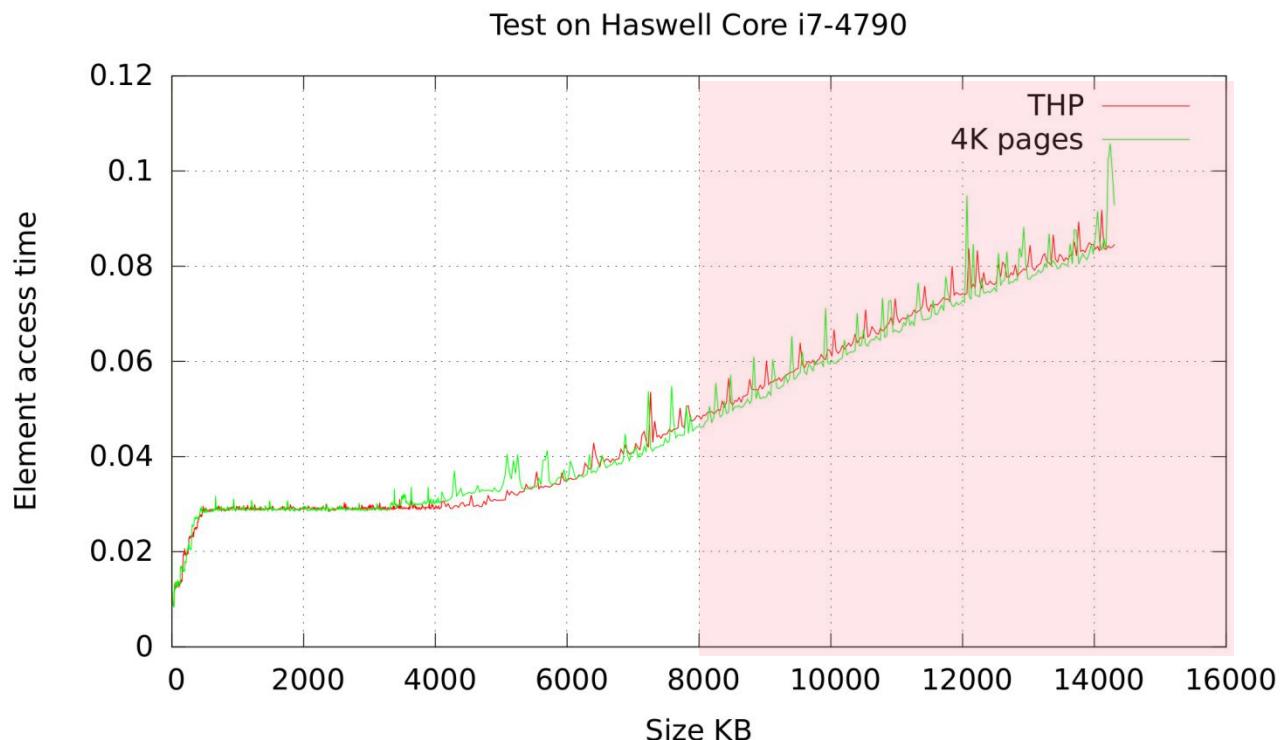
- Since Sandy Bridge
- L3 splits in **slices**
- Slice is selected by **hashing the address**
- Each slice has associativity with **16 ways**
- This fix the **coloring/alignment issue**



<https://software.intel.com/en-us/articles/intel-xeon-processor-e5-26004600-product-family-technical-overview>

On today CPUs

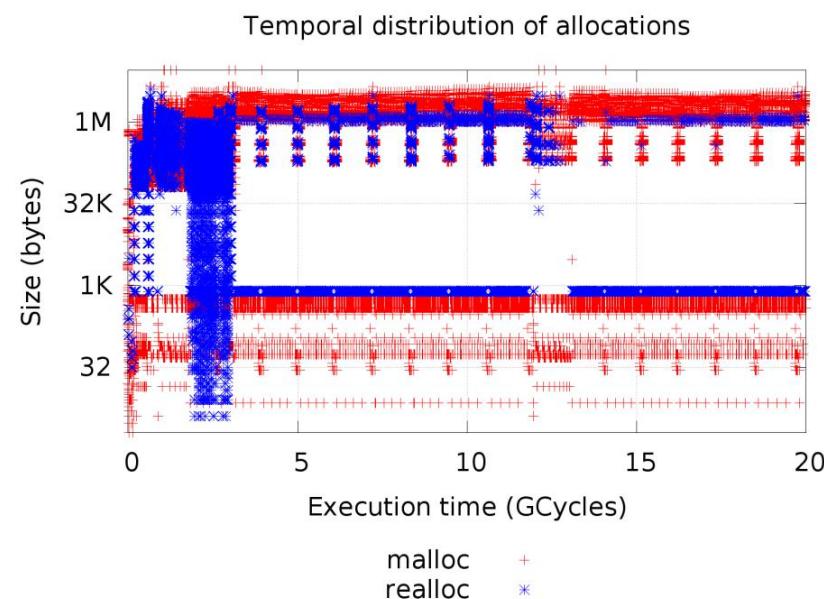
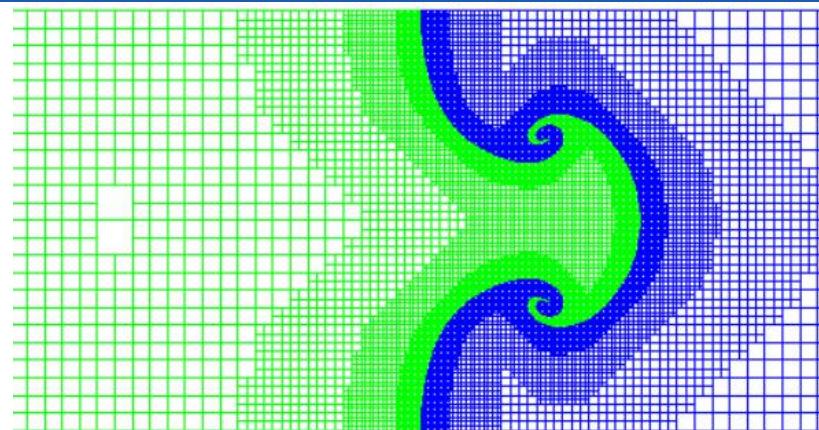
- Not anymore an issue for Intel L3 caches
 - Change of topology : slices
- AMD Zen (Ryzen)
 - Now also use slices
 - Should solve the issue
- Still an issue on IBM power 8
 - L3 cache has 8 ways for 8 MB
 - Issue present
 - Power 9 ? Also “regions” in LLC ?
- For ARM (v7/v8) ?
 - L2 shared associative cache
 - Issue should be present
 - But I never tested
- Issue for L2 of all processors !
 - Think hyperthreading with 8 ways !



NUMA ALLOCATOR FOR HPC APPLICATIONS

Allocator performance on HPC applications

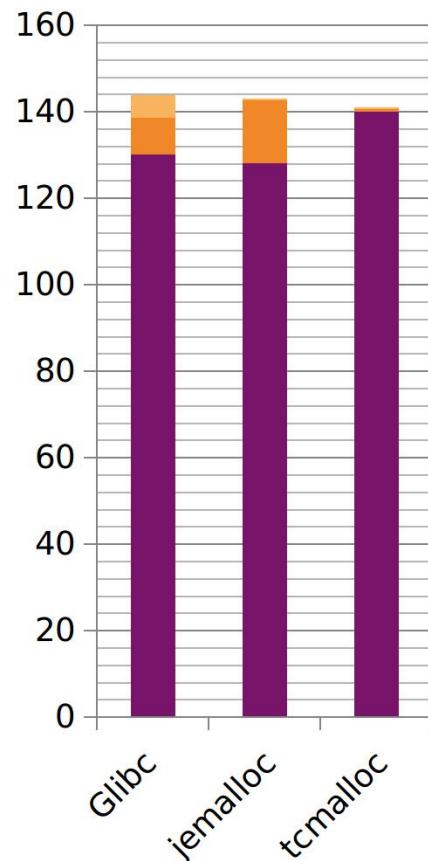
- Main interest : **malloc time cost**
- Test case : **Hera (CEA)**
 - **Adaptive Mesh Refinement (AMR)**
 - **Massive C++/MPI code (~1 million lines).**
- **Large number of memory allocations**
(~75 millions / 5 minutes on 12 cores)
- **Large number of alloc/realloc around ~20 MB**
- **Available allocators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



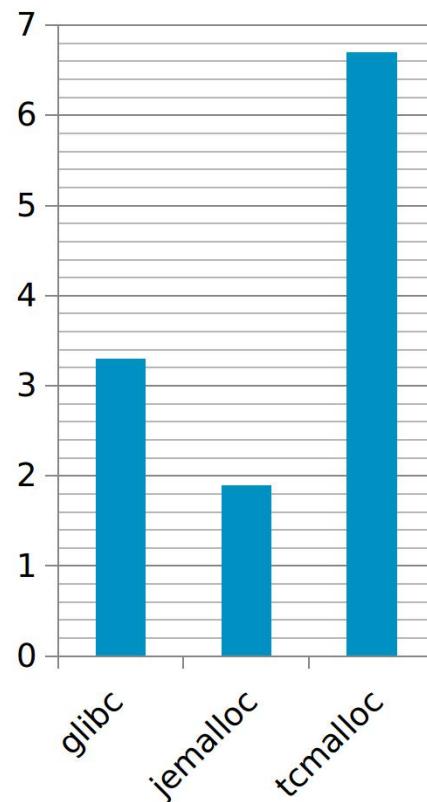
Hera preliminary results

12 cores

Execution time(s)

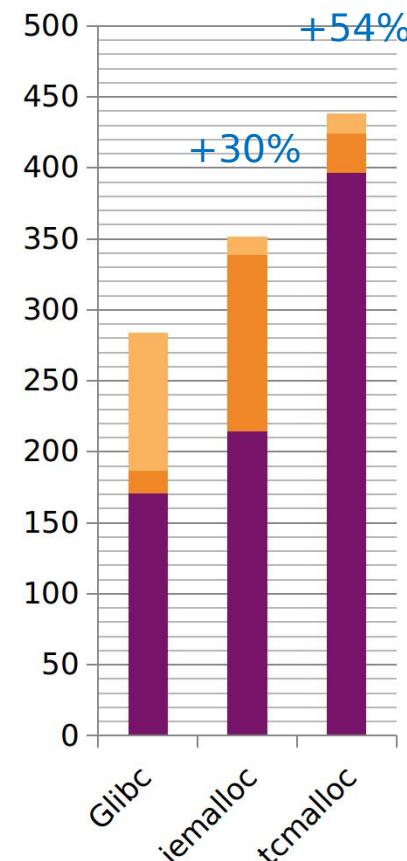


Physical mem.(Go)



128 cores

Execution time(s)



■ User ■ System ■ Idle

How to measure malloc time

- 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();
```

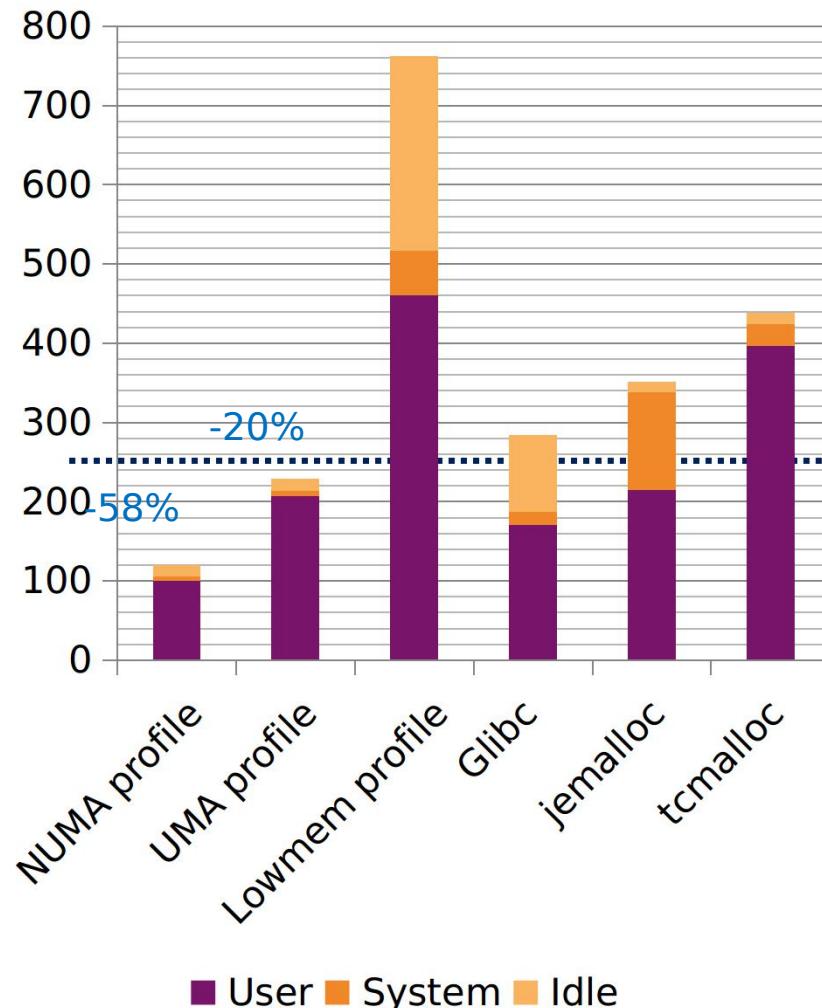
- **Lazy page allocation.**

- **Page faults** on first access.

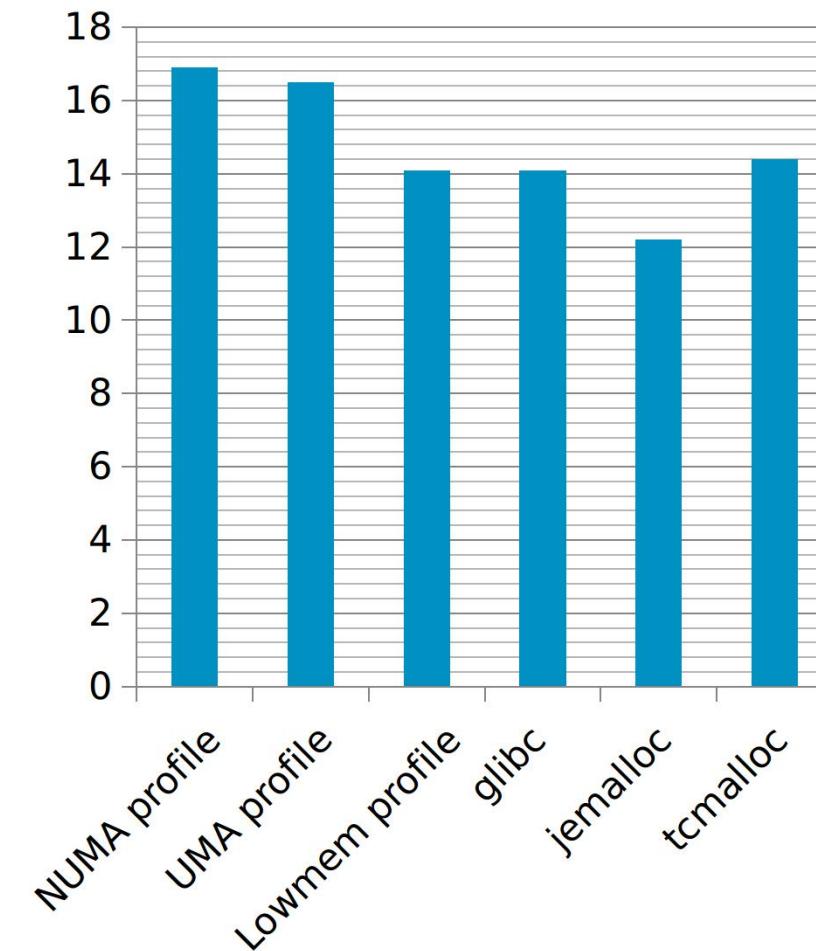
For 4GB	Malloc	First access
Time (M cycles)	0,008	1 217

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

Execution time (s)



Physical memory (GB)



COST OF FIRST TOUCH HANDLER

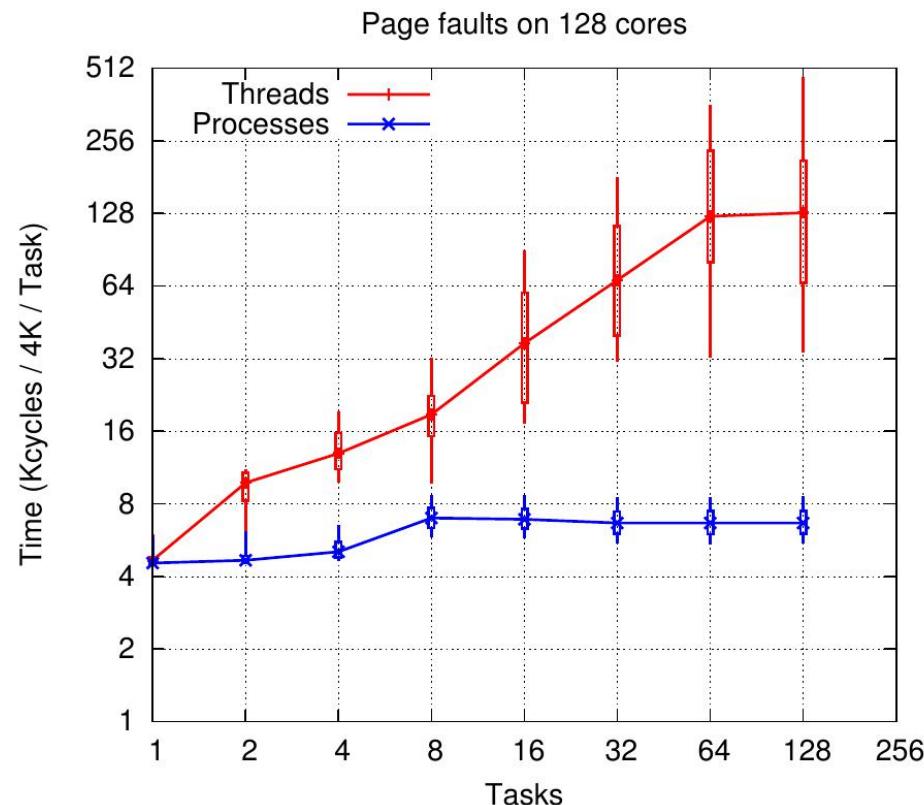
Benchmarking page faults

- Page faults are an issue for allocation performance
- We previously limit them with large segment recycling
- Can we improve fault performance of large allocations?
- Micro-benchmark :

```
ptr = mmap(SIZE);
#pragma omp parallel for
for ( i = 0 ; i < SIZE ; i += PAGE_SIZE)
{
    TIME_DISTRIBUTION(ptr[i] = 0);
}
```

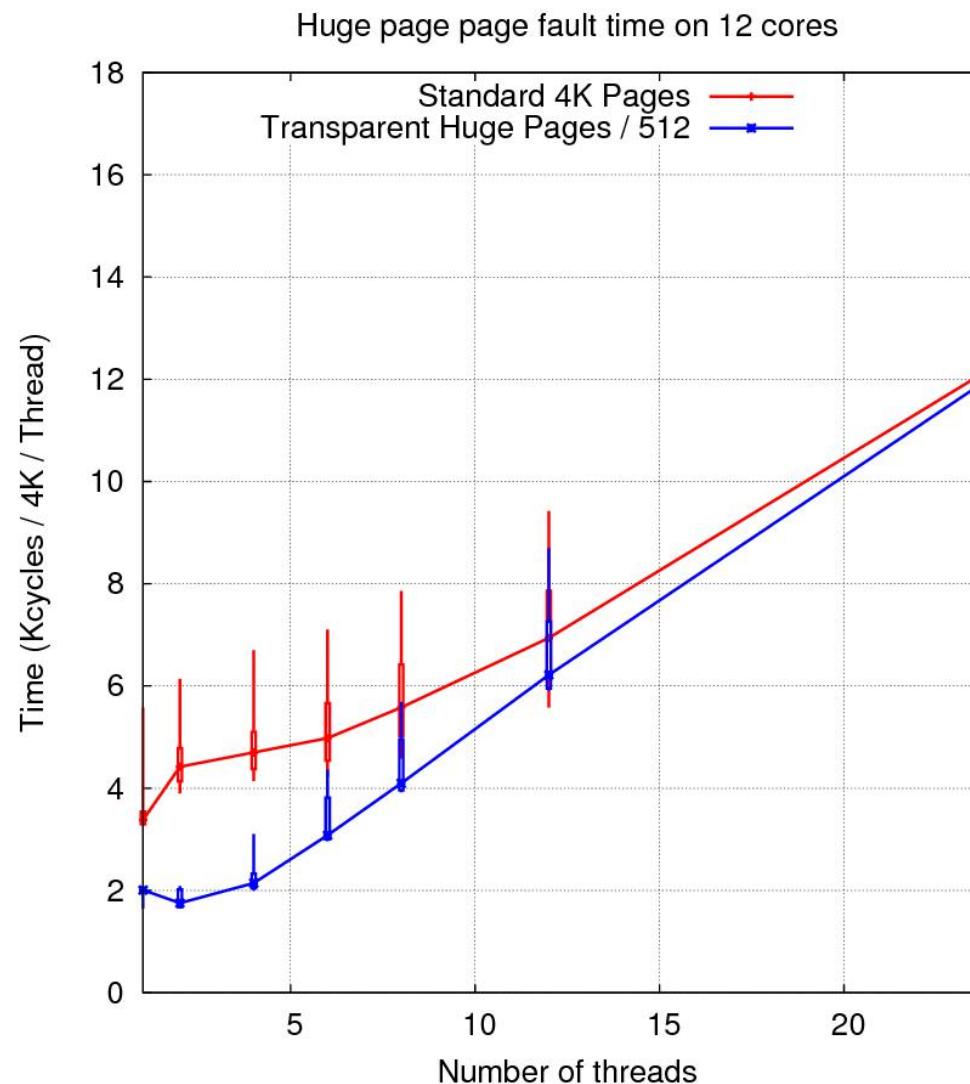
Page fault scalability

- Are page faults scalable ? Over threads or processes.
- Mesurement on **4*4 Nehalem-EP** (128 cores) and on **Xeon Phi** (60 cores)
- Get scalability issue !



Can huge pages solve this issue ?

- Standard pages: **4K**
- Huge pages (x86_64): **2M**
- **Divide number of faults by 512**
- Impact on performance ?
 - Sequential : **only 40%**
 - Parallel : **No**
- Why ?

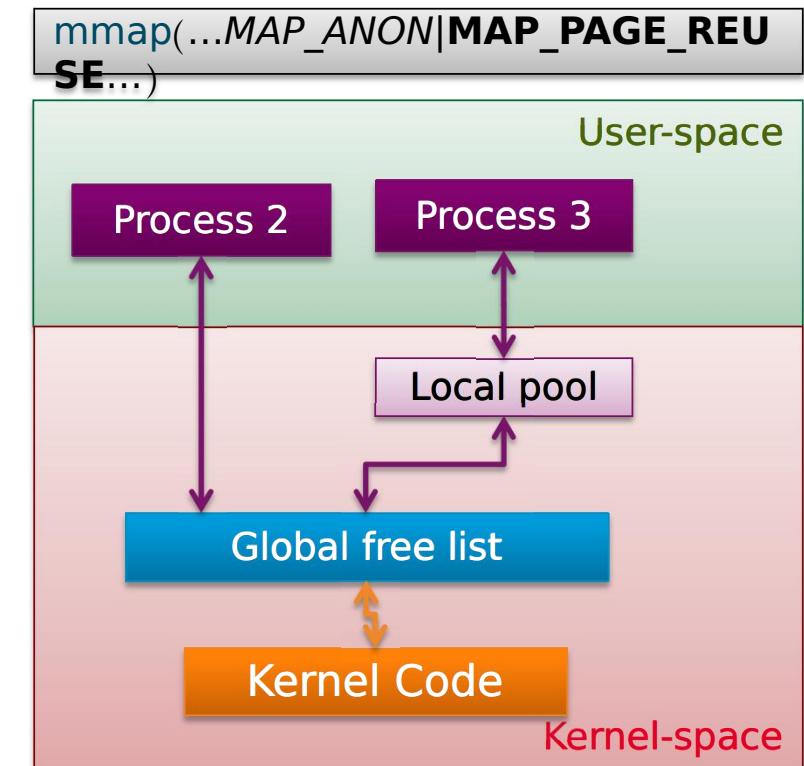


What happens on first touch page fault ?

- Hardware generates an interruption to the OS
 - Take **locks** on **page table**
 - Check reason of the fault
 - Is **first touch** from **lazy allocation**
 - Request a free page to NUMA **free lists**
 - Clear the page content
 - Map the page, update the **page table**
 - Release locks
-
- Possible issue on large NUMA domains
- $\sim 1400 / 3400$ cycles
40%
99% for THP !
- Locks, but hard to fix
(some work from
A.T. Clement ASPLOS12)

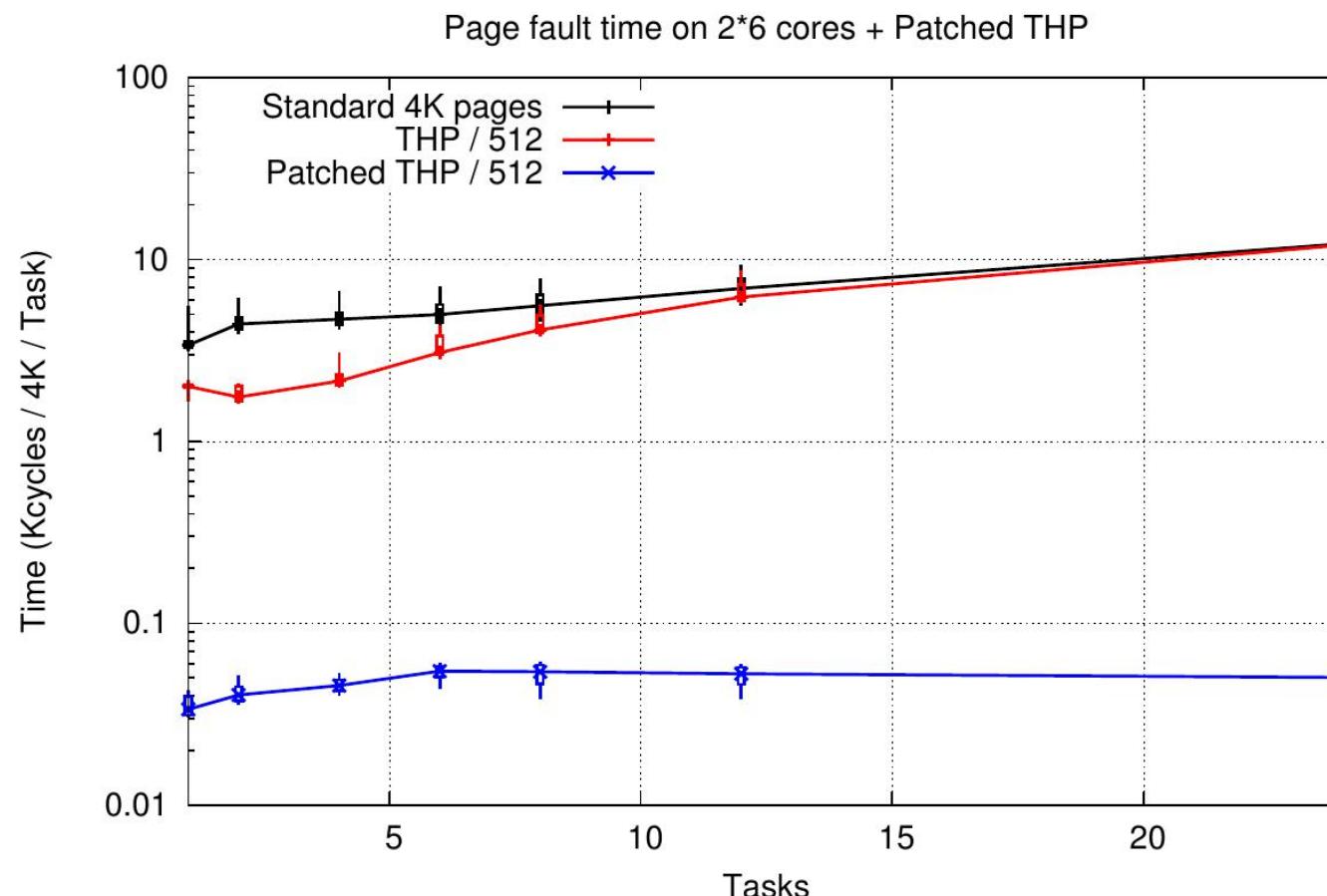
How to avoid page zeroing cost ?

- Microsoft approach :
 - Windows uses a **system thread** to clear the memory
 - So its done **out of critical path**
- But **zeroing**:
 - Implies **useless work**
 - Consumes CPU **cycles** so **energy**
 - Consumes **memory bandwidth**
- Why not **avoid them** ?
 - **Skip them (security ?)**
 - Use a **per process memory in kernel (published)**
 - Do in DIMM hardware



Performance impact on huge pages

- Huge pages (2 MB) faults become **47** times faster, **60** in parallel.
- New interest for huge pages.





MALT - A MALLOC TRACKER

The question

- 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 : list)
    {
        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

“compiler added allocations”

Short life allocations

Source annotations

MATT WebView

Inclusive/Exclusive

Metric selector

Summary Alloc sites Time analysis Stack Alloc sizes Help

Allocated mem. ▾

Search

28.4 KB __libc_start_main

28.4 KB _start

28.2 KB main

12.5 KB testMaxAlive()

6.9 KB recurseA(int)

6.3 KB testThreads()

1.0 KB funcB()

1.0 KB testRecuseInterv...
edA(i...)

1.0 KB testRecuseInterv...
edB(i...)

704.0 B funcC()

704.0 B testParallelWithRecur...
sion()

128.0 B OutOfMainAlloc

128.0 B __cxx_global_var_init1

128.0 B global constructors ke...

128.0 B __libc_csu_init

704 B 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(16);
61 *(char*)ptr='c';
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 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)

Per line annotation

Call stacks reaching the selected site.

S. Valat | IXPUG | 25 sept 2019

32

Symbols

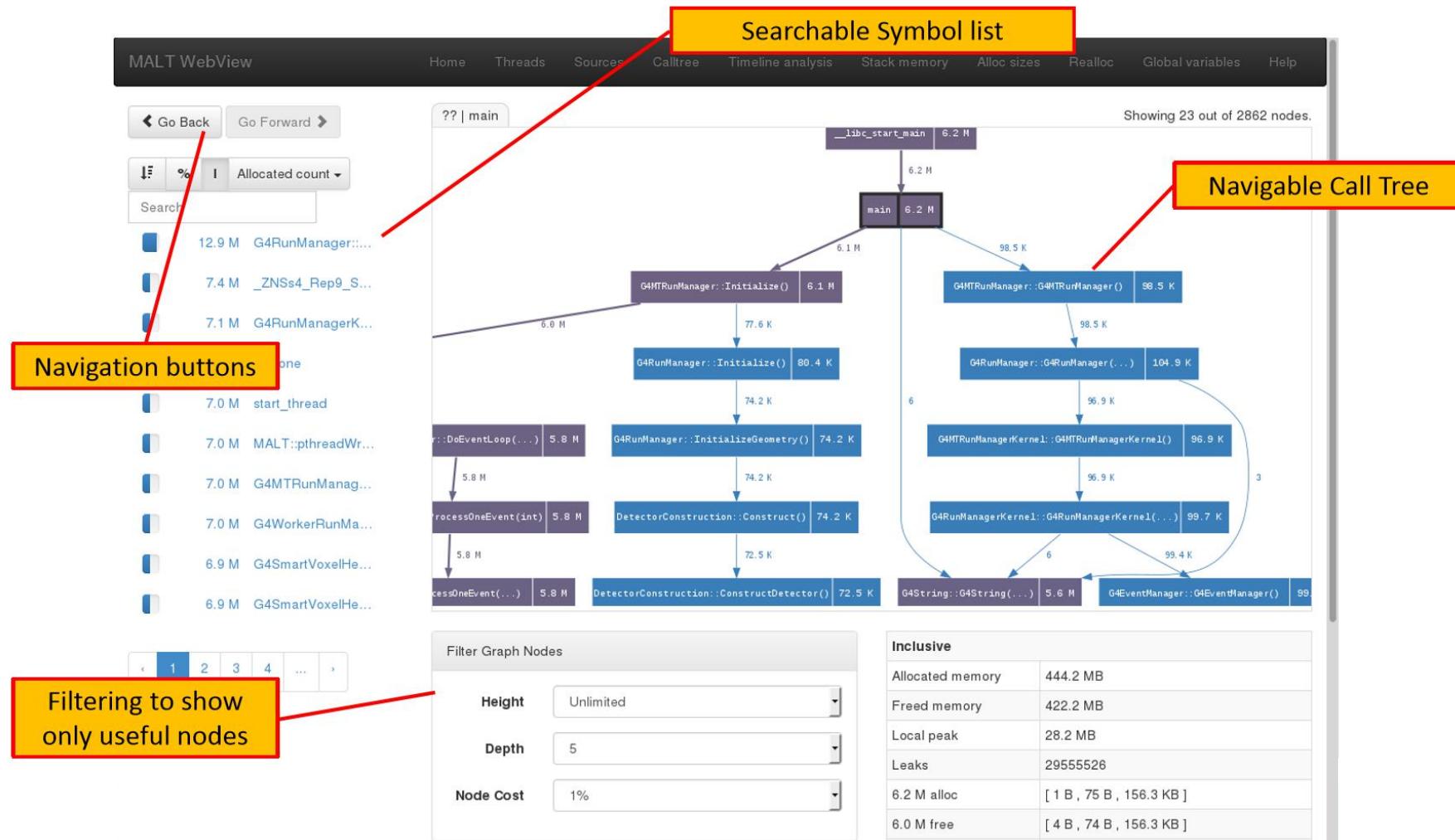
Details of symbol or line

Function

- _start
- __libc_start_main
- main
- funcA()
- funcB()
- malloc
- funcC()

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	32.0 B
funcC()	0 B

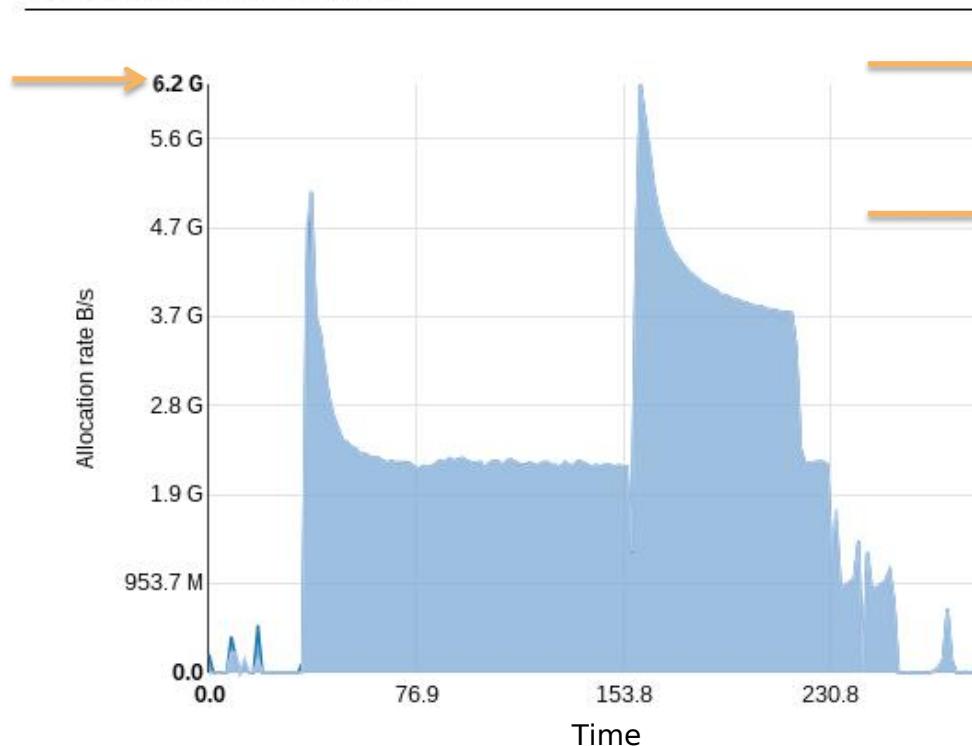
Call tree view



The question

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

Allocation rate



```
CALL assert(capacity==size(array),&
           'array and capacity variable are not :
IF (needed_size>capacity) THEN
  IF (ALLOCATED ( temp) ) DEALLOCATE(temp)
  ALLOCATE ( temp(capacity))

  DO i=1,capacity
    temp(i)=array(i)
  END DO

  DEALLOCATE ( array)
  ALLOCATE ( array(new_cap))

  DO i=1,capacity
    array(i)=temp(i)
  END DO

  capacity=new_cap
END IF
```

Total :

Allocated memory : 56.8 GB

Max alive memory : 135.7 M

3.5 K alloc : [16.0 KB , 16.3 MB , 33.7 MB]

Lifetime : [107.8 K , 26.7 M , 476.7 M] (cycles)

Own :

Allocated memory : 56.8 GB

Max alive memory : 135.7 M

3.5 K alloc : [16.0 KB , 16.3 MB , 33.7 MB]

GERN-IT - MALT, Sébastien Valat

Lifetime : [107.8 K , 26.7 M , 476.7 M] (cycles)

Function

▶ _start

NUMAOROF - A NUMA PROFILER



Typical NUMA example

- 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...

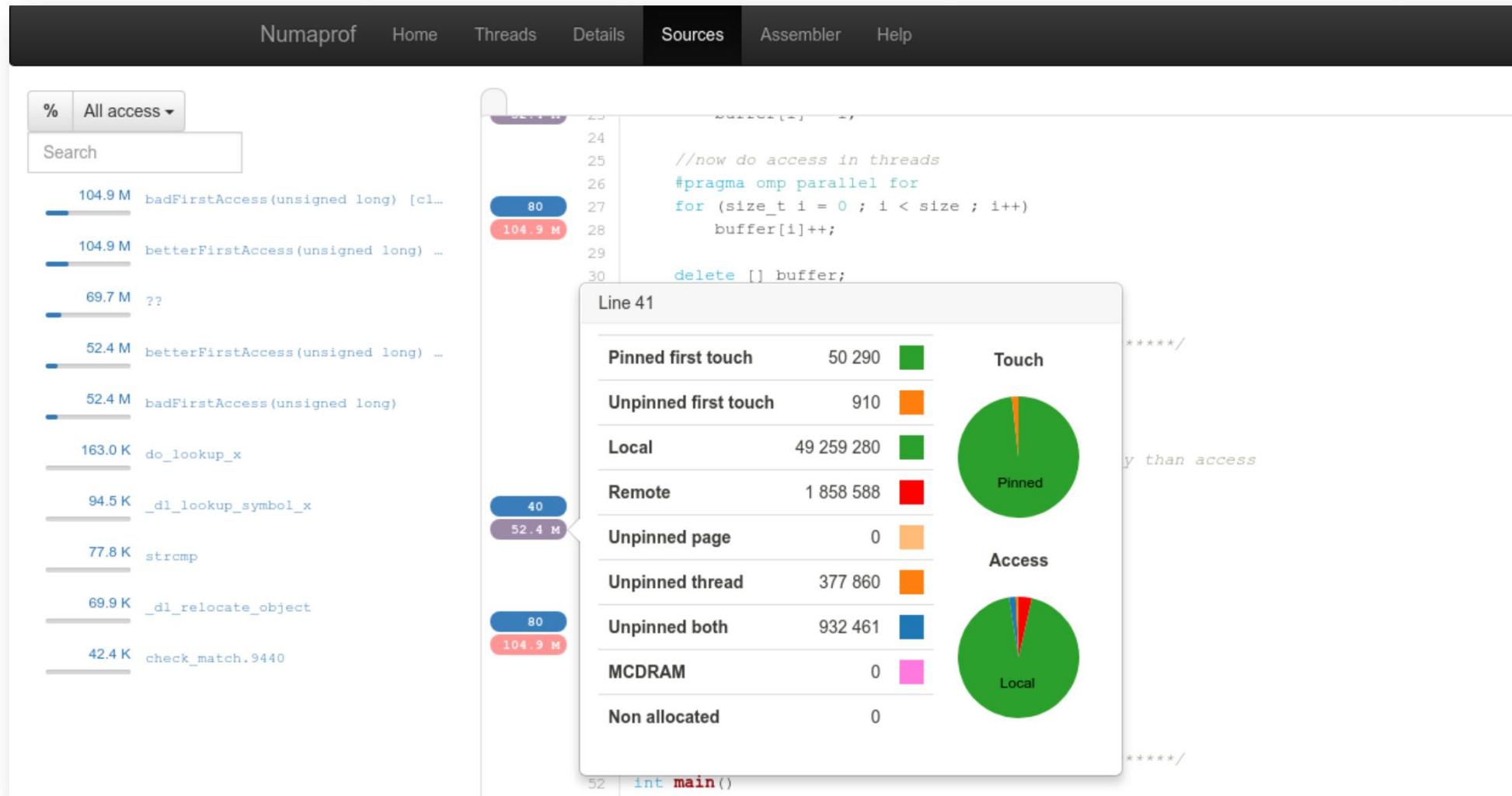
- We want to know if we make **remote accesses**
- Ideally we need to know **where...**
- We can dream, we want to know **which allocation contain issues**
- We want to know **where** the **first touch** has been done
- On KNL we want to check **MCRAM accesses**



Global summary



Source & asm annotations

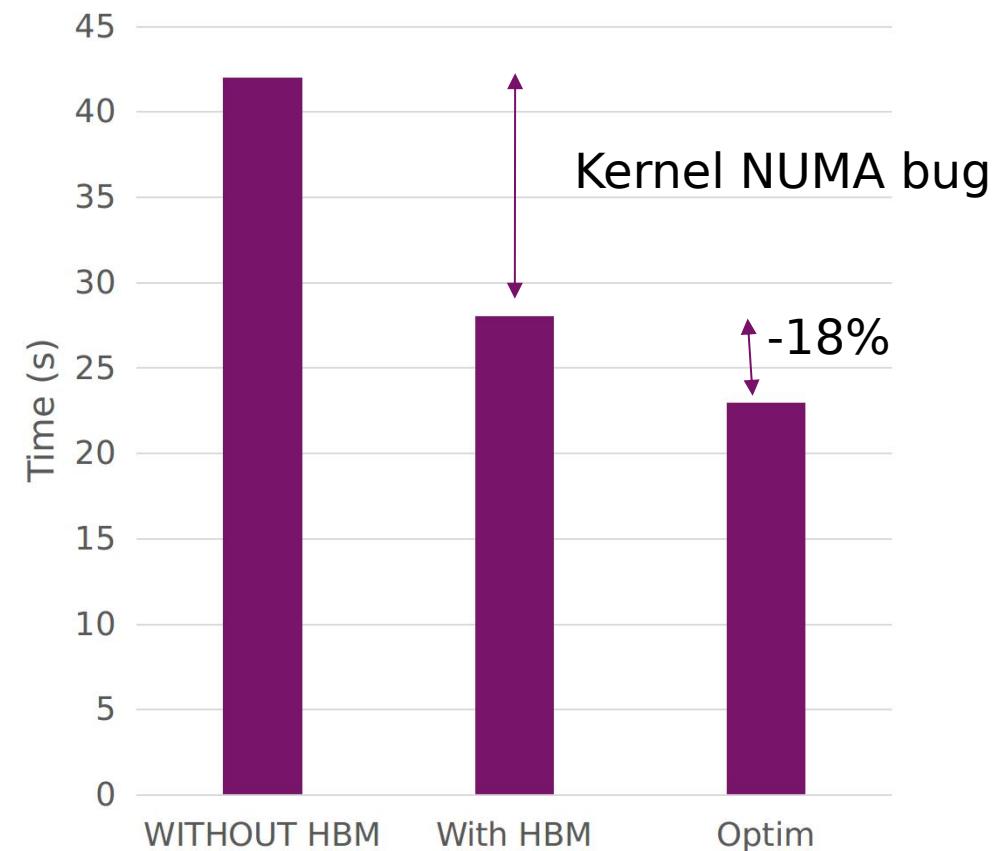
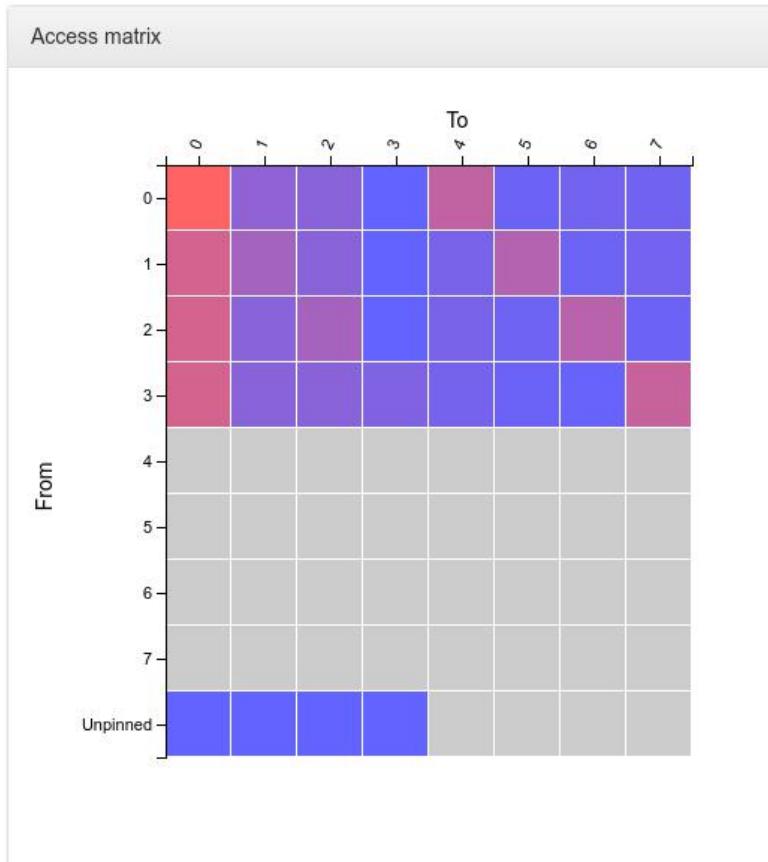


Non parallel allocations



```
/data/svalat/Projects/Hydro/HydroC/HydroCplusMPI/ThreadBuffers.cpp | ThreadBuffers::ThreadBuffers(int, in  
22  
23     ThreadBuffers::ThreadBuffers(int32_t xmin, int32_t xmax, int32_t  
24     {  
25         int32_t lgx, lgy, lgmax;  
26         lgx = (xmax - xmin);  
27         lgy = (ymax - ymin);  
28         lgmax = lgx;  
29         if (lgmax < lgy)  
30             lgmax = lgy;  
31  
32         m_q = new Soa(NB_VAR, lgx, lgy);  
33         m_qxm = new Soa(NB_VAR, lgx, lgy);  
34         m_qxp = new Soa(NB_VAR, lgx, lgy);  
35         m_dq = new Soa(NB_VAR, lgx, lgy);  
36         m_qleft = new Soa(NB_VAR, lgx, lgy);  
37         m_qright = new Soa(NB_VAR, lgx, lgy);  
38         m_qgdnv = new Soa(NB_VAR, lgx, lgy);  
39  
40         m_c = new Matrix2 < real_t > (lgx, lgy);  
41         m_e = new Matrix2 < real_t > (lgx, lgy);  
42
```

40 minutes optimization on HydroC



CONCLUSION

Conclusion

- Memory is one of the **key for performance**
- **Old management** need to be carefully **looked again**
- Performance **gaps** can be **integer factors**
- Dedicated **tools** can **help a lot**
- It requires **flexible software** to reach **global optimization**
 - Unit tests ?

QUESTIONS ?

On access we need...

- Intercept the memory accessse (Intel PIN)
- Track thread location
- Intercepts malloc
- We can skip accesses to local stack
 - overhead 80x -> 40x
- Overhead on 256 KNL threads : 60x

BACKUPS

Ideal view of HPC memory management stack

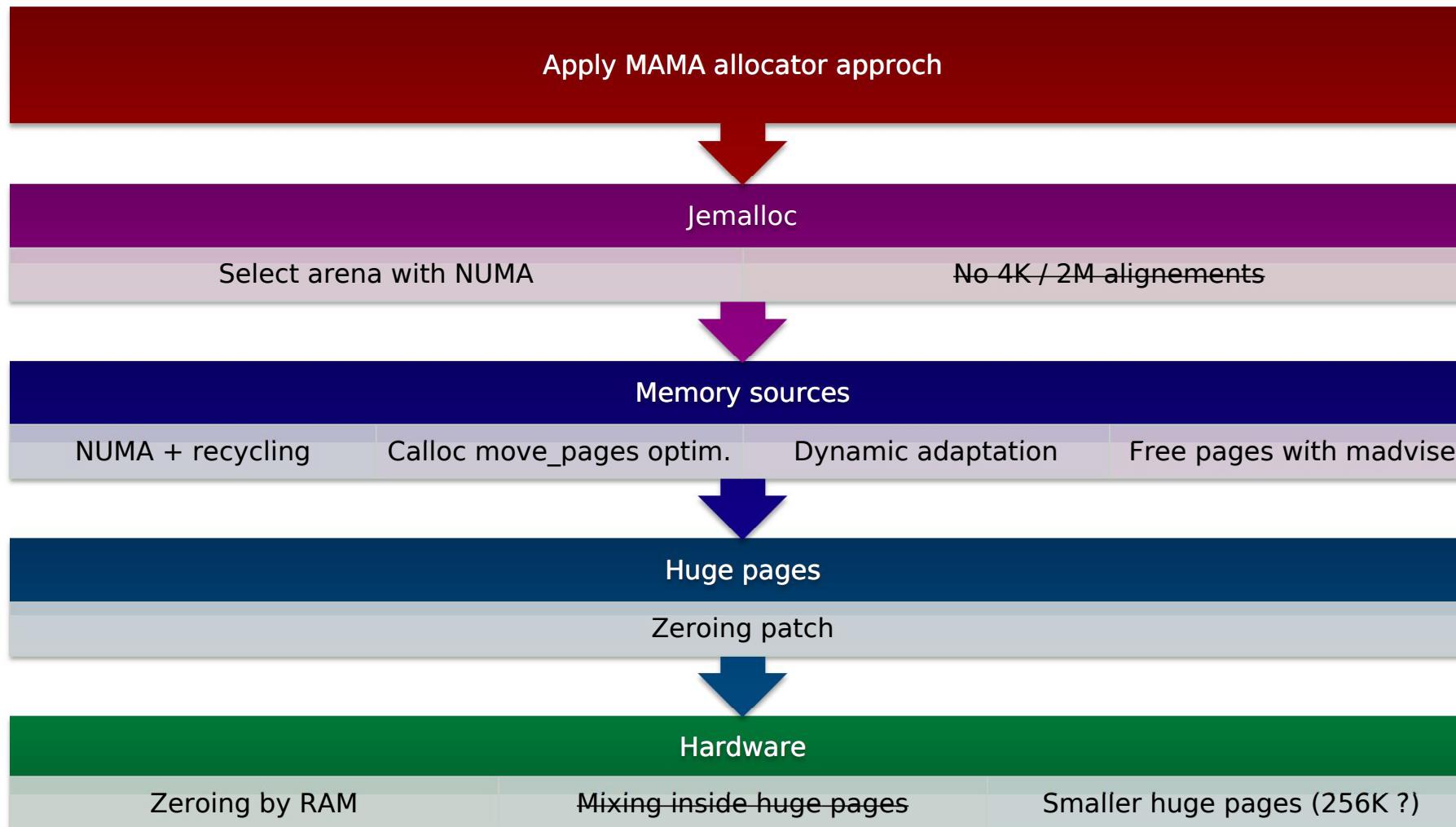


Table 1. Summary of the results of the ANOVA analysis of the effect of culture on the perceived predictability of market changes

Note: The dependent variable is the perceived predictability of market changes. The independent variables are culture (C), gender (G) and age (A).

Source: Own calculations based on the data from the survey of 1,000 respondents.

Table 2. Summary of the results of the ANOVA analysis of the effect of culture on the perceived predictability of market changes

Note: The dependent variable is the perceived predictability of market changes. The independent variables are culture (C), gender (G) and age (A).

Source: Own calculations based on the data from the survey of 1,000 respondents.

Table 3. Summary of the results of the ANOVA analysis of the effect of culture on the perceived predictability of market changes

Note: The dependent variable is the perceived predictability of market changes. The independent variables are culture (C), gender (G) and age (A).

Source: Own calculations based on the data from the survey of 1,000 respondents.

Table 4. Summary of the results of the ANOVA analysis of the effect of culture on the perceived predictability of market changes

Note: The dependent variable is the perceived predictability of market changes. The independent variables are culture (C), gender (G) and age (A).

Source: Own calculations based on the data from the survey of 1,000 respondents.

Table 5. Summary of the results of the ANOVA analysis of the effect of culture on the perceived predictability of market changes

Note: The dependent variable is the perceived predictability of market changes. The independent variables are culture (C), gender (G) and age (A).

Source: Own calculations based on the data from the survey of 1,000 respondents.

Table 6. Summary of the results of the ANOVA analysis of the effect of culture on the perceived predictability of market changes

Note: The dependent variable is the perceived predictability of market changes. The independent variables are culture (C), gender (G) and age (A).

Source: Own calculations based on the data from the survey of 1,000 respondents.

Table 7. Summary of the results of the ANOVA analysis of the effect of culture on the perceived predictability of market changes

Note: The dependent variable is the perceived predictability of market changes. The independent variables are culture (C), gender (G) and age (A).

Source: Own calculations based on the data from the survey of 1,000 respondents.