# Allocations mémoires,
## Pourquoi et comment profiler ?

gray-scott reloaded / LAPP - Annecy / 07-2024

*Sébastien Valat*

# Plan

# Vi-HPS
## Virtual Institute — High Productivity Supercomputing

**VI-HPS**

https://www.vi-hps.org/

**SCAN ME**

Next : **4 – 6 september 2024 / Ostrava**

× MASTODON × TWITTER/X × IMPRINT × PRIVACY × SITEMAP

**VI-HPS**    VIRTUAL INSTITUTE — HIGH PRODUCTIVITY SUPERCOMPUTING

× ABOUT    × MEMBERS    × ORGANIZATION    × TOOLS    × TRAINING    × SYMPOSIA    × PROJECTS    × NEWS

> Tools

**TOOLS OVERVIEW**

To assist developers with the selection of the appropriate tools provided by VI-HPS, our Tools Guide booklet linked below offers a brief overview of the respective tools. The guide showcases their individual debugging, correctness checking and performance analysis capabilities. Furthermore, it indicates their support for parallel computer systems, programming models and languages.

- VI-HPS Tools Guide

Detailed information about each particular tool can be found on the listed websites in the document above and the respective tool pages below.

**Single Node Performance**

Callgrind
LIKWID
MAQAO

**Parallel Performance**

Dimemas
Extra-P
Linaro MAP (Part of Linaro Forge)
Linaro Performance Reports (Part of Linaro Forge)
mpiP
Open|SpeedShop
Paraver
Scalasca
TAU
Vampir

**Instrumentation**

Opari 2

**Measurement**

Extrae
PAPI
Score-P

**Integration**

Component-based Tool Framework
LaunchMON
P$^n$MPI

**Visualization**

Cube

Sébastien Valat - INRIA / LJK / UGA

# Origin of the tools

➢ **PhD**. On **memory management** for **HPC** (at CEA / UVSQ)

➢ **MALT** : post-doc at Versailles :

Exascale computing research

➢ **NUMAPROF** : side project post-doc work at :

CERN

➢ **URL** :

**Get both on :**
https://memtt.github.io/

Sébastien Valat - INRIA / LJK / UGA

# Motivations

➢ Lot of **issues** today :

- **Huge** memory **space** to **manage** (~TB of memory)
- **Lot more** distinct **allocations** (e.g. 75M in 5 minutes)
- **Multi-threaded** : 256 threads
- **Hidden** into large (**huge**) C/C++/Fortran **codes** (**~1M** lines)

➢ Access:

- **NUMA** (Non Uniform Memory Access)
- **Memory wall !**
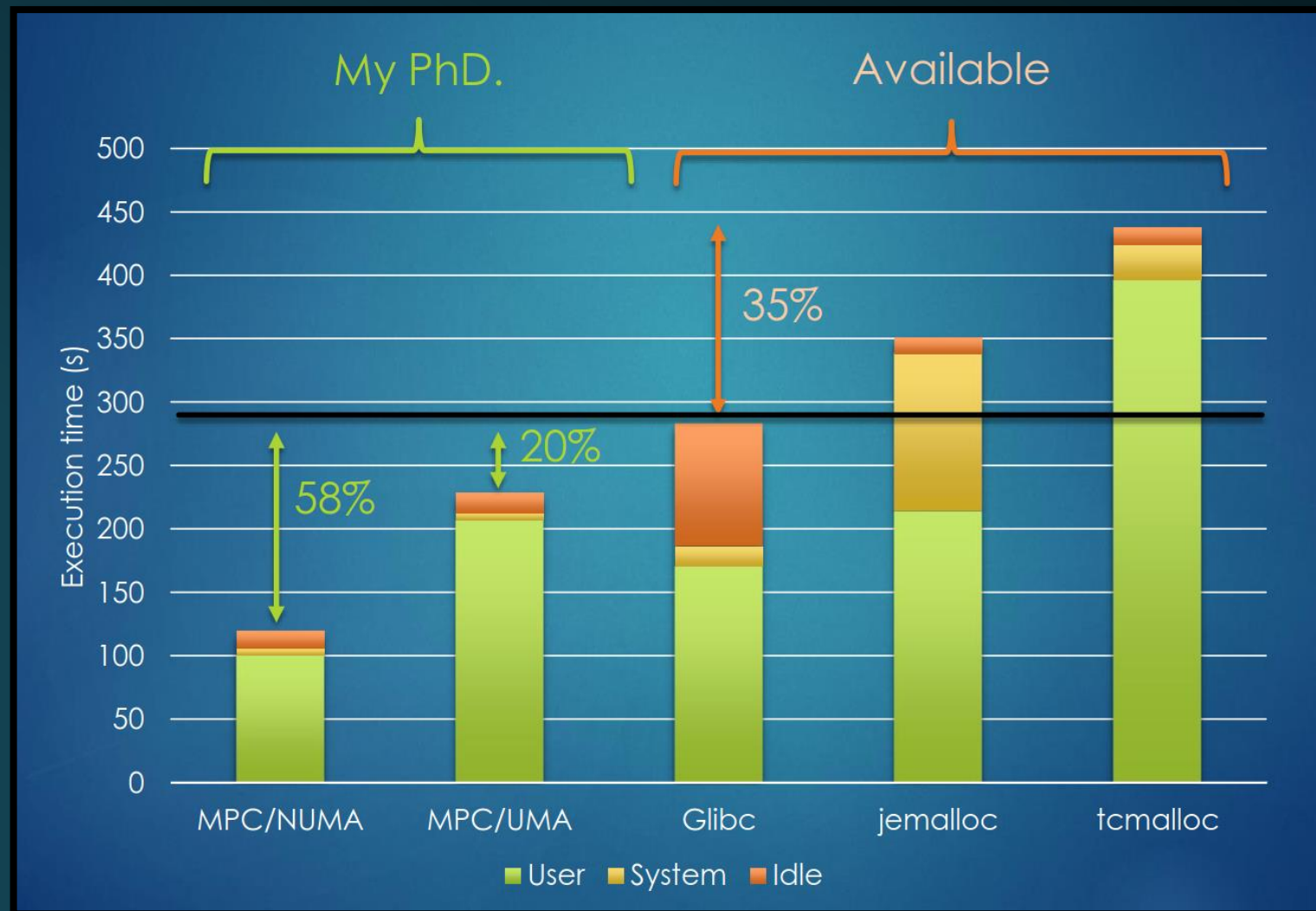
**SuperMicro**, 2024, 5019P-MT
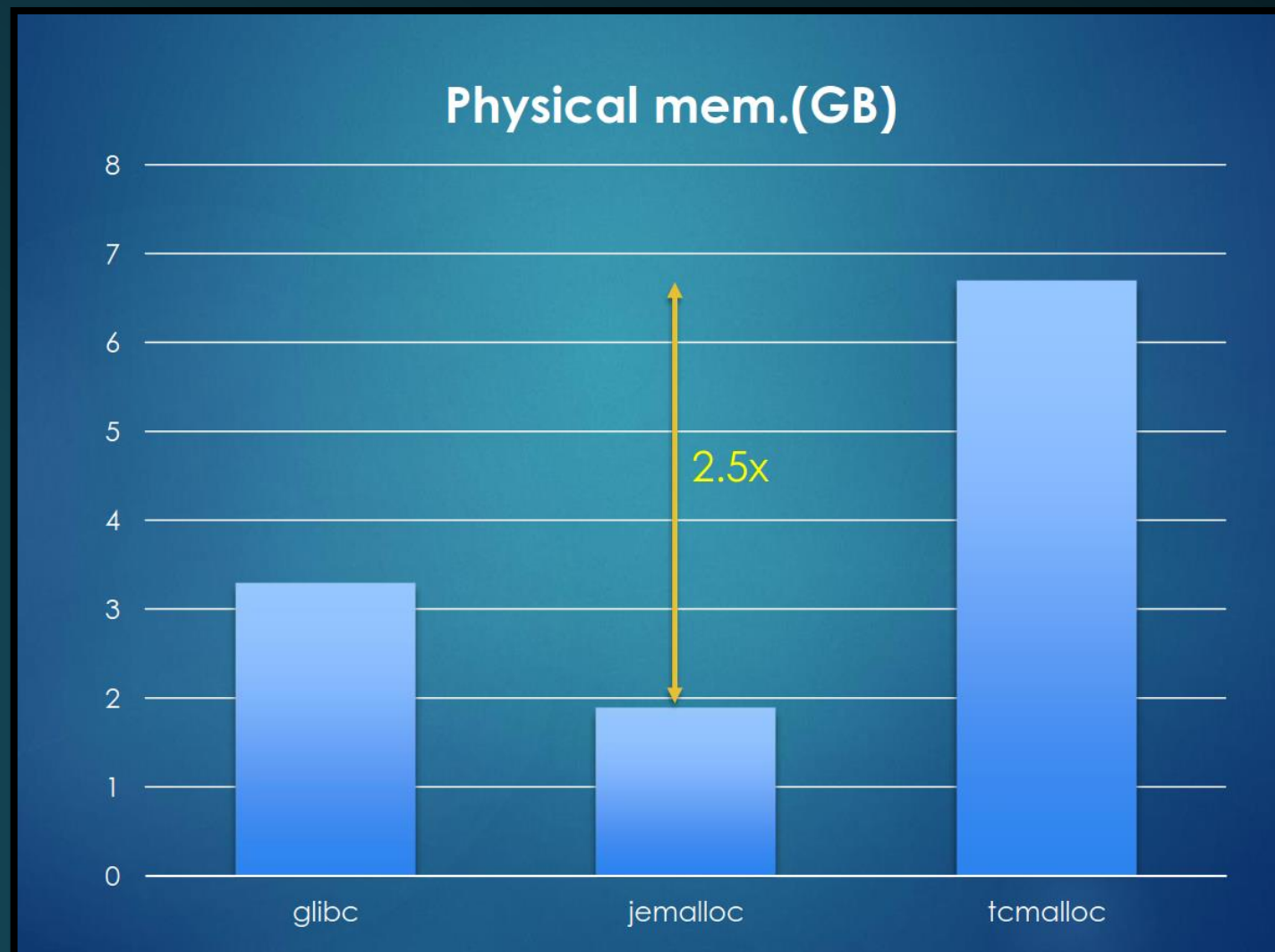**1 TB** ~ 10 000€

# Key today

You need to
**well understand
the memory behavior**
of your (HPC)
application !

# Eg: **>1M lines** C++ simulation.
# On **128** cores / **16** NUMA CPUs

Sébastien Jean Valat. Contribution à l'amélioration des méthodes d'optimisation de la gestion de la mémoire dans le cadre du Calcul Haute Performance. (tel-01253537)

# Same about **memory consumption** on 12 cores



**Physical mem.(GB)**

Bar chart with values: glibc ≈ 3.3, jemalloc ≈ 1.9, tcmalloc ≈ 6.7. Annotation: 2.5x between jemalloc and tcmalloc.

8

# MALT

Malloc Tracker



**MALT**
**A Malloc Tracker**

# Goal

➤ We have **profiling tool** for **timing** (eg. **Valgrind** or **vtune**)

➤ But for **memory usage** ?

➤ Memory can be an issue :
- Failed to run (or swap) due to **lack of memory resource**.
- **Performance impact** of memory management functions.
- Impact due to **memory layout**.

# Some issue examples

We want to help searching :

➢ **Where** memory is allocated.

➢ **Properties** of allocated chunks.

➢ **Bad** allocation **patterns** for performance.

➢ **Leaks**

➢ **Global variables** (TLS)

# Some issue examples

```
Int gblVar[SIZE];
int * func(int size)
{
        child_func_with_allocs();
        void * ptr = new char[size];
        double* ret = new double[size*size*size];
        for (…..)
        {
                double* buffer = new double[size];
                //short and quick do stuff
                delete [] buffer;

        }
        return ret;

}
```

Global **variables** or **TLS**

**Indirect** allocations

**Leak**

Might lead to swap for **large size**

**Short life** allocations

# Existing tools

➢ Valgrind - **memcheck**

➢ Valgrind - **massif**

# Existing tools

➢ **Google heap profiler** (tcmalloc):

➢ IBM **Purify++** / Parasoft **Insure++**

•1
4

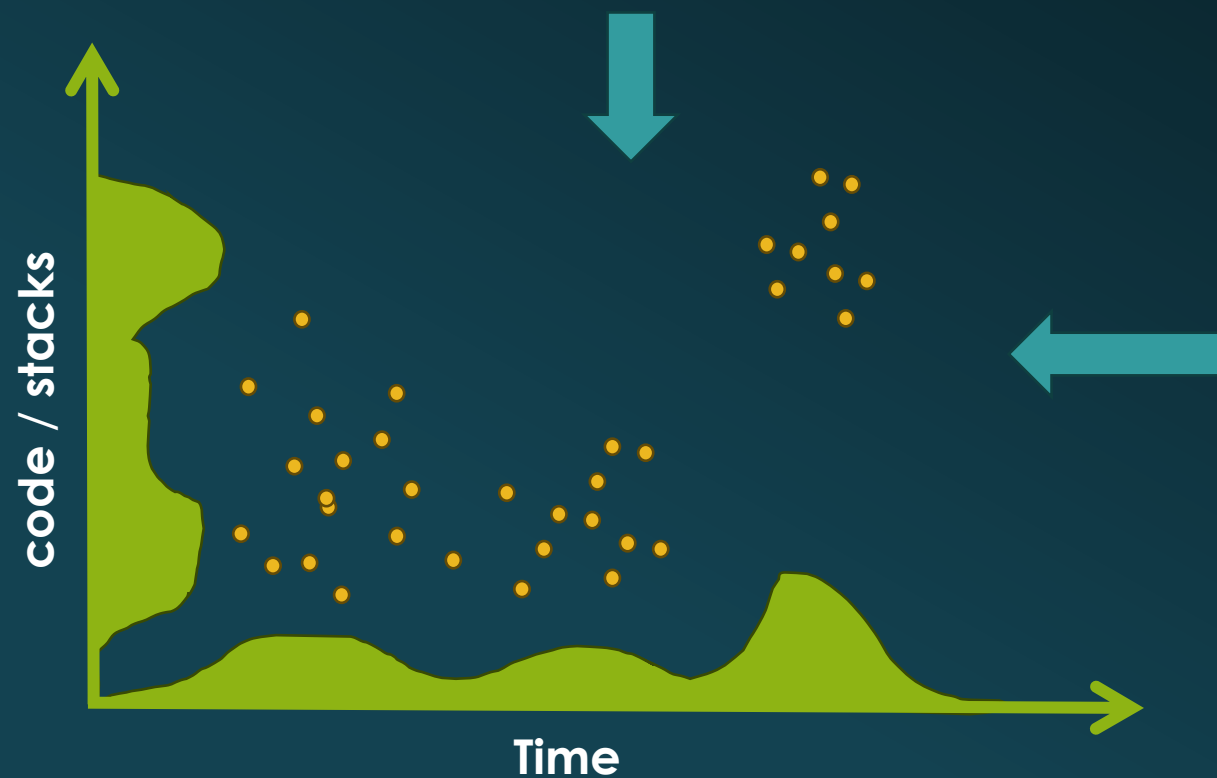# Visual Studio Ultimate memory profiler

# What I want to provide

➢ Same **approach** than **valgrind / kcachegrind**

➢ **Mapped allocations** on **sources lines**

➢ For memory resource usage :
- **Memory leaks** (malloc without free)
- **Peak** and total allocated memory

➢ For performance :
- Allocation **count**
- Allocation **sizes** (min/mean/max)
- Chunk **lifetime** (min/mean/max)

# Profile based

➢ Project on **two axis**:
- Over **code / stacks**
- Over **time**

# Usage

**Profiling**

`malt ./mon_prog`

JSON

**Viewing
via light web server**

`malt-webview …`

`firefox http://localhost:8080`

**Can be forwarded via SSH**

`ssh –L8080:localhost:8080 ………`

# Profiling on a cluster

**Your laptop**

```
ssh –L8080:localhost:8080 .........
```



**Remote cluster / server**

**Profiling**

```
malt ./mon_prog
```

Project sources

JSON

Libraries sources

**Viewing via light web server**

```
malt-webview …
```

# Global summary

| EXECUTION TIME | PHYSICAL MEMORY PEAK | ALLOCATION COUNT | AVAILABLE PHYSICAL MEMORY |
|---|---|---|---|
| 00:00:00.25 | 2.3 MB | 379 | 4.1 Gb |

## Run description

| Executable : | simple-case-finstr-linked |
|---|---|
| Commande : | `./simple-case-finstr-linked` |
| Tool : | matt-0.0.0 |
| Host : | localhost |
| Date : | 2014-11-26 22:40 |
| Execution time : | 00:00:00.25 |
| Ticks frequency : | 1.8 GHz |

## Global statistics

Show all details   Show help

| Physical memory peak | 2.3 MB |
|---|---|
| Virtual memory peak | 103.7 MB |

# Global summary

## Global statistics

[Show all details] [Show help]

| | |
|---|---|
| Physical memory peak | 2.3 MB |
| Virtual memory peak | 103.7 MB |
| Requested memory peak | 3.9 KB |
| Cumulated memory allocations | 26.4 KB |
| Allocation count | 379 |
| Recycling ratio | 6.7 |
| Leaked memory | 2.1 KB |
| Largest stack | 6.0 KB |
| Global variables | 41.2 KB |
| TLS variables | 4.0 KB |
| Peak allocation rate | 25.2 MB/s |

# Can give some hints (warnings)

| Ticks frequency : | 2.9 GHz |
|---|---|
| Allocator used : | /usr/lib/x86_64-linux-gnu/libc.so.6 |

## Global statistics

[Show all details] [Show help]

| Physical memory peak | 13.8 MB |
|---|---|
| Virtual memory peak | 189.4 MB |
| Requested memory peak | 12.1 MB |
| **Cumulated memory allocations** | **144.0 GB** ⚠ |
| Allocation count | 49.2 K |
| **Recycling ratio** | **12204.8** ⚠ |
| Leaked memory | 6.1 MB |
| Largest stack | 96 B |
| Global variables | 43.8 KB |
| TLS variables | 40 B |
| **Global variable count** | **789** ⚠ |
| Peak allocation rate | 309.4 GB/s |

# Per thread statistics

# Source annotation

**Inclusive/Exclusive**

**Metric selector**

MALT WebView · Calltree · ... · sizes · Realloc · Global variables · Help

/home/svalat/Projects/malt/src/lib/tests/simple-case.cpp |

Allocated count

Search

```
     79        *(char*)ptr='c';//requir...      e malloc/free
 2   80        free(ptr);
     81        funcC();
     82    }
     83
     84    /****************  FUNCTION  *****************/
     85    void funcA()
     86    {
     87        void * ptr = malloc(16);
110  88        *(char*)ptr='c';//required otherwise new compilers will remove malloc/free
     89        free(ptr);
 99  90        funcB();
     91    }
     92
     93    /****************  FUNCTION  *****************/
     94    void recurseA(int depth)
     95    {
     96        if (depth > 0)
     97        {
 1   98            void * ptr = malloc(64);
 1   99            *(char*)ptr='c';//required otherwise new compilers will remove malloc/free
    100            recurseA(depth-1);
    101            free(ptr);
 1  102        }
    103    }
```

419  _start
417  __libc_start_main_impl
417  __libc_start_call_main
417  main
345  omp_fulfill_event
304  __clone3
304  start_thread
304  MALT::pthreadWrapperStartRou...
230  testRecuseIntervedB(int)
220  testParallelWithRecurse() [c...

« 1 2 3 4 ... »

**Per line annotation**

**Call stacks reaching the selected site.**

**Symbols**

| Inclusive | |
|---|---|
| Allocated memory | 6.2 KB |
| Freed memory | 6.2 KB |
| Local peak | 1.1 KB |
| Leaks | 0 |
| 99 alloc | [ 64 B , 64 B , 64 B ] |
| 99 free | [ 64 B , 64 B , 64 B ] |
| Lifetime | [ 18.5 K , 1.0 M , 12.3 M ] (cycles) |

| Function | Metric |
|---|---|
| ▼ _start | 99 |
| ▼ __libc_start_main_impl | 99 |
| ▼ __libc_start_call_main | 99 |
| ▼ main | 99 |
| ▼ recurseA(int) | 99 |
| ▼ recurseA(int) | 99 |
| malloc | 11 |
| ▼ recurseA(int) | 88 |
| malloc | 11 |
| ▼ recurseA(int) | 77 |
| malloc | 11 |

**Details of symbol or line**

24

# Metrics

**Bloc sizes**

**Frees**

**Peaks**

**Short life** allocations

**Recycling**

**Realloc**



↓≡  %  |  Allocated count ▾

Allocated mem.
Allocated count
Min. alloc size
Mean alloc size          rt_main_impl
Max. alloc size
Freed mem.               rt_call_main
Free count
Memory ops.
Local peak               l_event
Global peak
Leaks
Max lifetime             ead
Min lifetime             readWrapperStartRouti…
Recycling ratio
Realloc count            eIntervedB(int)
Realloc sum
                         220  testParallelWithRecurse() [clo…

# Call stack tree

*Jaffery, Syed Mehdi Raza / CERN*

# Time charts

# Chunk size distribution
## Example from YALES2 with gfortran issue



Many really small allocations

# Global variables

## Distribution over binaries

○ Grouped ● Stacked

● Global variables ● TLS variables

- simple-case-finstr-linked
- libc-2.17.so
- libstdc++.so.6.0.17
- libmatt.so
- ld-2.17.so
- libdl-2.17.so
- libm-2.17.so
- libunwind-x86_64.so.8.0.1
- libunwind.so.8.0.1
- liblzma.so.5.0.5
- libpthread-2.17.so
- librt-2.17.so
- libgcc_s.so.1
- libgomp.so.1.0.0
- libelf-0.158.so

0 B    2.0 KB    3.9 KB    5.9 KB    7.8 KB    9.8 KB    11.7 KB    13.7 KB    15.6 KB    17.6 KB    20

## Distribution over variables

○ Grouped ● Stacked

● Global variables ● TLS variables

- tlsArray
- gblCallocIniBuffer
- gblStaticArray
- gblArray
- _rtld_global
- std::tr1::__detail::__prime_list
- std::__detail::__prime_list
- sys_errlist

# Memory used by stacks

Need to recompile
the app with
-finstrument-functions

| LARGEST STACK | THREADS |
|---|---|
| 3.1 MB | 1 |

## Stack peaks



● Stack

0

0 B                                                                                    3.1 MB

Maximum stack size

## Stack of thread 0



● Thread 0

3.1 MB
2.9 MB

2.4 MB

1.9 MB

1.4 MB

976.6 KB

488.3 KB

32 B
0.7 ms      0.9 ms      1.0 ms      1.2 ms      1.4 ms      1.5 ms      1.7 ms      1.9 ms      2.1 ms      2.2 ms      2.4 ms

Time (secondes)

## Stack memory used by functions for thread 0



● Frame size

depth_3()
depth_2(int)
depth_1()
doing_things_with_memory()

# Real cases

# Cerfacs - AVBP- CFD simulator

ugal

# Example on AVBP init phase

➢ Issue with reallocation on init

# Coria – YALES2 - Combustion

# Allocatable arrays on YALES2
## Issue only occur with **gfortran**, **ifort** uses stack arrays.

MATT WebView

| ↕ | % | | | ← | → | Allocation count ▾ |

Search

▮ 911.9 K data_comm_m::copy_i...

▮ 896.4 K data_comm_m::copy_...

Search intensive alloc functions

Huge number of allocation for a line programmer think it doesn't do any !

```
      892    do i=1,nitem_el_grp
      893        el_grp_ind = el_grp_index2int_comm_index%val(1,i)
      894        int_comm_ind = el_grp_index2int_comm_index%val(2,i)
608 K 895        el_grp_r2%val(1:dim1,el_grp_ind) = int_comm_r2%val(1:dim1,int_comm_ind)
      896    end do
```

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 !

35

# We can found allocs of 1B !

MATT WebView



Search for the minimal chunk size.

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

# Overhead

Sébastien Valat - INRIA / LJK / UGA

# Reminder on NUMA

Non Uniform Memory Access

# What is NUMA ?

➢ Each CPU has its **own memory**

➢ Access to **remote memory** we need to **go through the owner** CPU

# Architecture

➢ **A bit more complex**







Large nodes : 16 proc - 128 cores
Bull BCS ~ 2010

# Today topology

- ➢ **Intel Knight Landing,**

  mode **SNC2** or **SNC4**

- ➢ Also add fast memory

  **MCDRAM (HBM)** presented

  as **NUMA** or **LLC cache**

# Today topology

Safe usage :

Each **process** bound over
a **single NUMA node**

# 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**

```
memset(array, 0, SIZE)
```

➢ Then access

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

➢ **Bad performance** due to remote accesses !

# Performance impact

➢On a real machine (2 NUMA) :

| | Value |
|---|---|
| Memset | 35 |
| OMP Loop | 15 |

Chart axis: 0, 5, 10, 15, 20, 25, 30, 35, 40

# 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**

46

*Je voudrais…*

*Othman Bouizi*

Sébastien Valat - INRIA / LJK / UGA

# NumaProf

How to know if we are right In a real application ?

With the help of **Othman Bouizi**

# 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)

# Overhead and scalability

➢ Of course overhead is large: **~30x**

➢ But is **scale**

➢ Example code hydro on **KNL**:

**Overhead**

# GUI and example

# Global summary

# Statistics per thread

# OMP and huge pages

➤ **Huge pages** & thread splitting

➤ Most of the time do not match exactly

➤ Not a big issue if limited

# Details per thread

# Source & asm annotations

# Code Hydro

➤ KNL Without HBM                           WITH HBM

# Original Hydro access matrix

# Ordering issue



```
#pragma omp parallel for private(i) if (m_numa) SCHEDULE
        for (int32_t i = 0; i < m_nbtiles; i++) {
                int t = m_mortonIdx[i];
                m_tiles[t] = new Tile;
        }
```

# Non parallel allocations

# Parallel allocations
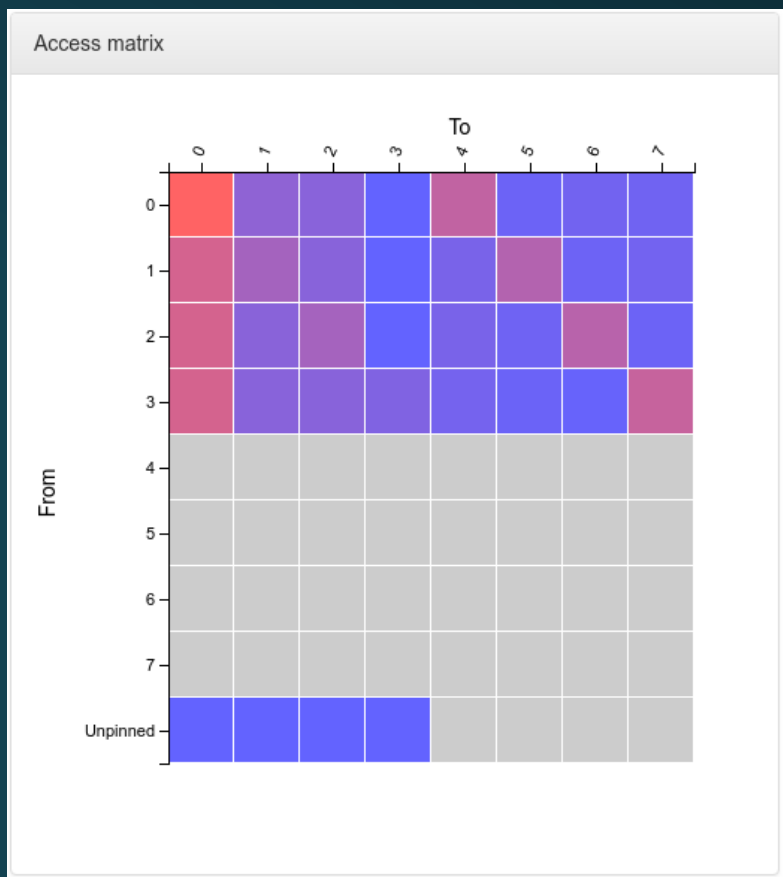
➢ Original

```
for (int32_t i = 0; i < m_numThreads; i++) {
    m_buffers[i] = new ThreadBuffers(…);
    assert(m_buffers[i] != 0);
}
```

➢ Modified

```
#pragma omp parallel
{
    int i = omp_get_thread_num();
            #pragma omp critical
    m_buffers[i] = new ThreadBuffers(..);
    assert(m_buffers[i] != 0);
}
```
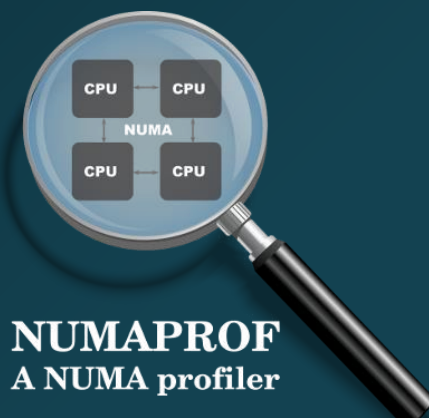
# Speed up obtained on Hydro

# Conclusion

➢ **Memory** is **not trivial** to handle **in large programs**

➢ Need to be taken in account

➢ Some **tiny mistakes** sometimes **cost a lot**
  ▪ **Possibly everywhere** in the program (**global impact**)

➢ Be able to get a view is a first help

MALT
A Malloc Tracker

NUMAPROF
A NUMA profiler

https://memtt.github.io

# Questions ?

SCAN ME