# MALT : MALloc Tracker

## A memory profiling tool

# Possibly huge impact

- **Memory management** can have **huge impact** on performance

- **Extreme case** on a **1.5 million** C++ lines **HPC simulation** app. on a **16 processors** server

## Execution time (s)

4x

Legend: ■ User ■ System ■ Idle

(Bar chart categories: MPC/NUMA, MPC/UMA, Glibc, jemalloc, tcmalloc)

- 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

- I wanted 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;
}
```
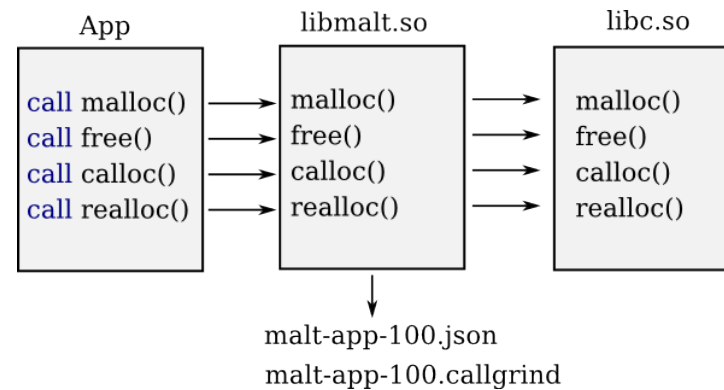
Global variables and TLS

Indirect allocations

Leak

Might lead to swap for large size

C++11 auto induced allocs

Short life allocations

# What I want to provide

- Same **approach** than **valgrind/kcachgind**

- **Mapped** allocations on **sources lines** and **call stacks**

- **Using a web-based GUI**
  - I started with kcachgrind
  - But wanted more flexibility and time charts

- Use **LD_PRELOAD** to intercept **malloc/free/...** as Google heap profiler



- **Map** allocations on **call stacks**

- **Build & consolidate summary** metrics

- Generate **JSON** output file

# Source annotations

Web technology (**NodeJS**, **D3JS**, **Jquery**, **AngularS**)



Inclusive/Exclusive

Metric selector

Per line annotation

Symbols

Details of symbol or line

Call stacks reaching the selected site.

# Call tree view

MALT, Sébastien Valat

- ## Memory consumption over time
  - Physical
  - Virtual
  - Requested (malloced)

## Memory allocated over time

## Size over time



## Lifetime over size

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



Allocation rate
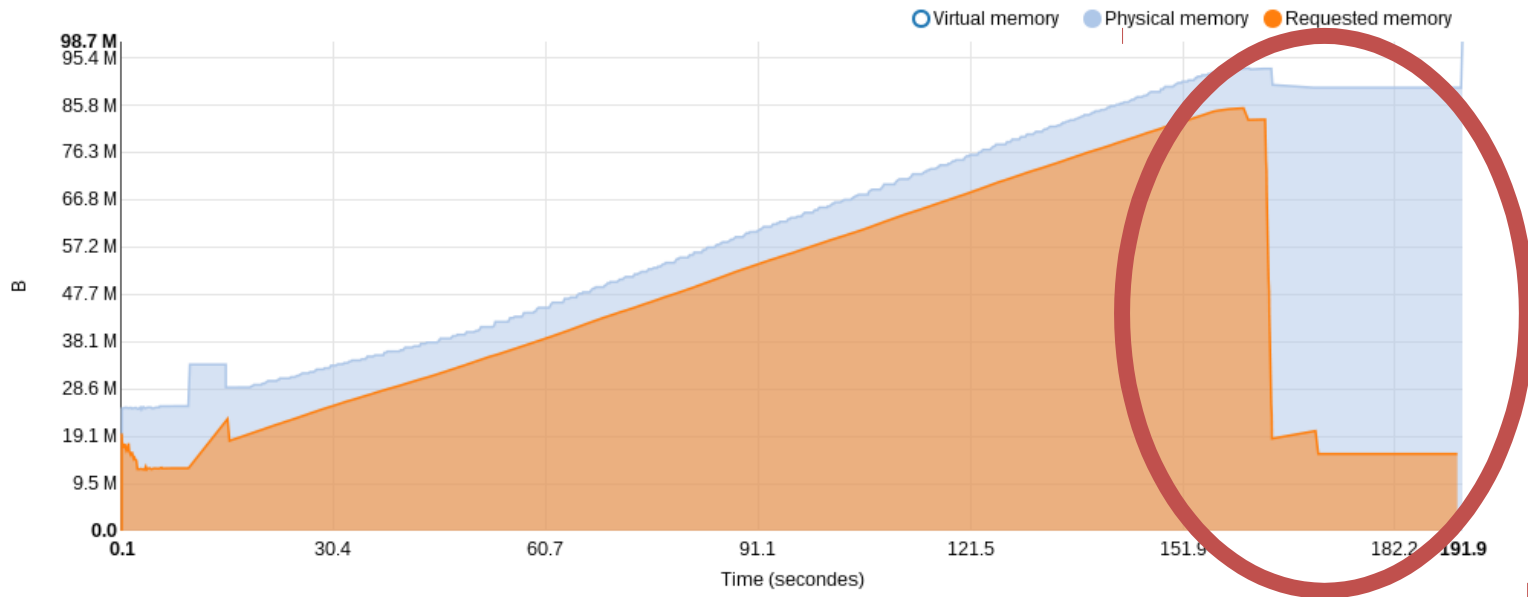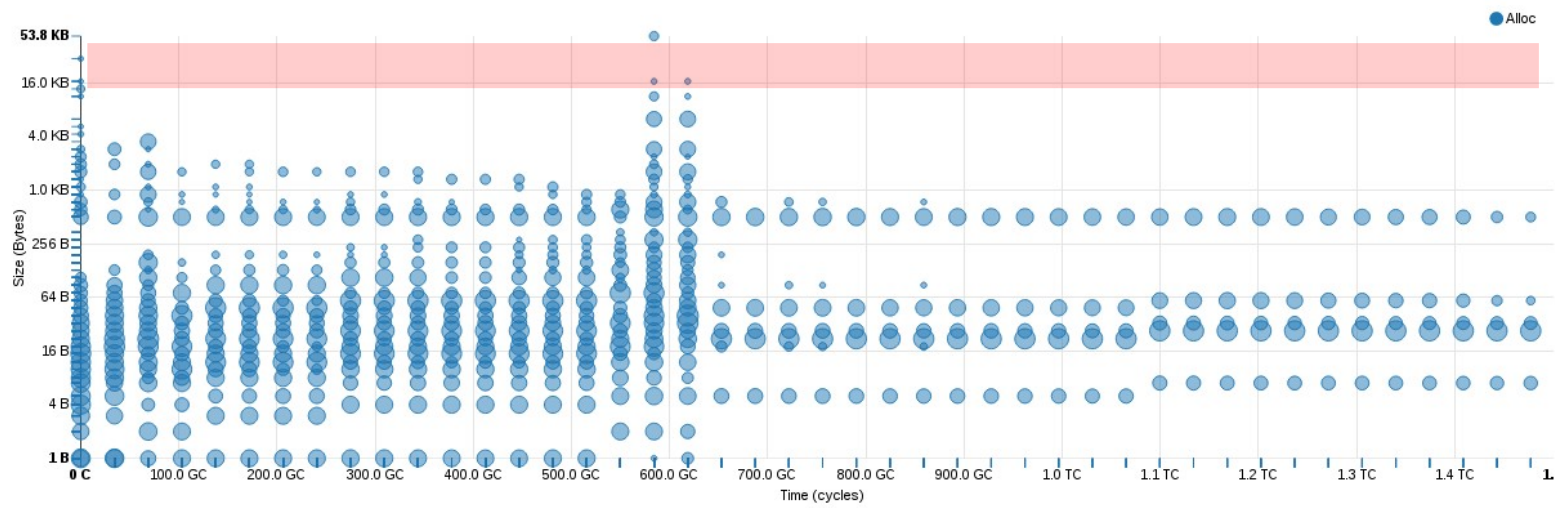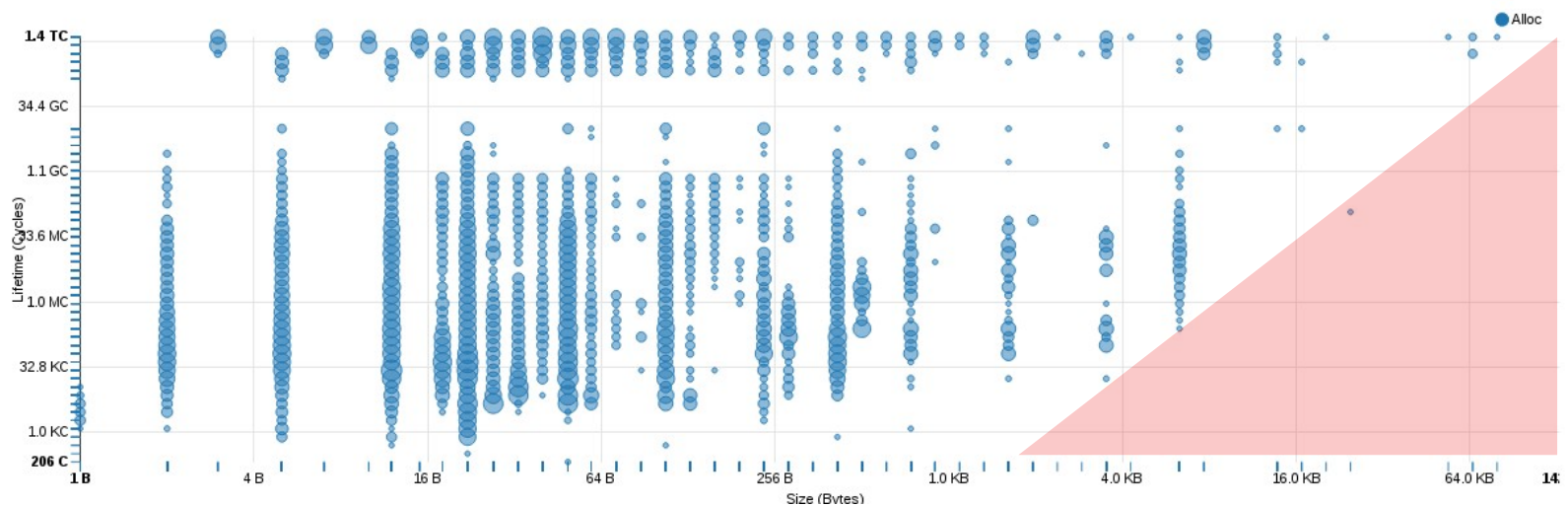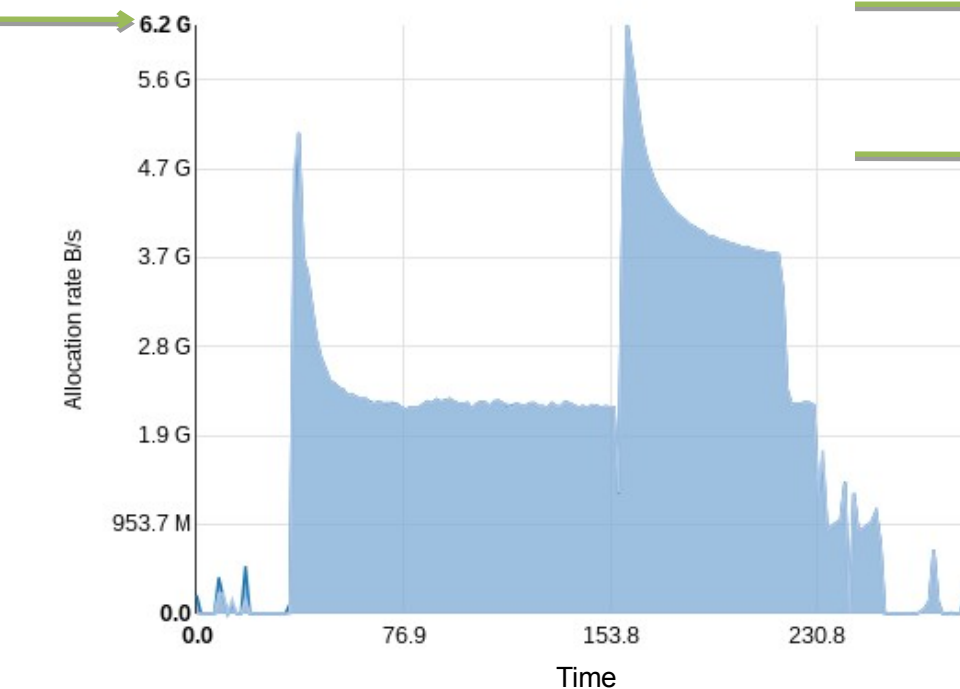
```
 99    CALL assert(capacity==size(array),&
100            'array and capacity variable are not .
101
102    IF (needed_size>capacity) THEN
103        IF (ALLOCATED ( temp) ) DEALLOCATE(temp)
104        ALLOCATE ( temp(capacity))
105
106        DO i=1,capacity
107            temp(i)=array(i)
108        END DO
109
110        DEALLOCATE ( array)
111        ALLOCATE ( array(new_cap))
112
113        DO i=1,capacity
114            array(i)=temp(i)
115        END DO
116
117        capacity=new_cap
118    END IF
119 <
```

57 GB

57 GB

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

Function
▶    _start

- Optionally recompile with debug flags :

```
gcc -g ...
```

- Run

```
malt [--config=file.ini] YOUR_PRGM [OPTIONS]
```

- Use the web view && http://localhost:8080:

```
malt-webview -i malt-{YOUR_PRGM}-{PID}.json
```

- In case there is a QT wrapper embedding NodeJS + Webkit

```
malt-qt -i malt-{YOUR_PRGM}-{PID}.json
```

- **Open sourced** since one year on https://github.com/memtt
- Co-hosted with a **similar tool :**
  **NUMAPROF** for **Non Uniform Memory Access** profiling.



- My **research** on memory management for **HPC** : http://svalat.github.io/

Thank you.

# QUESTIONS ?

# BACKUP

Example from YALES2 with gfortran issue



Many really small allocations

# Output, first idea, kcachegrind
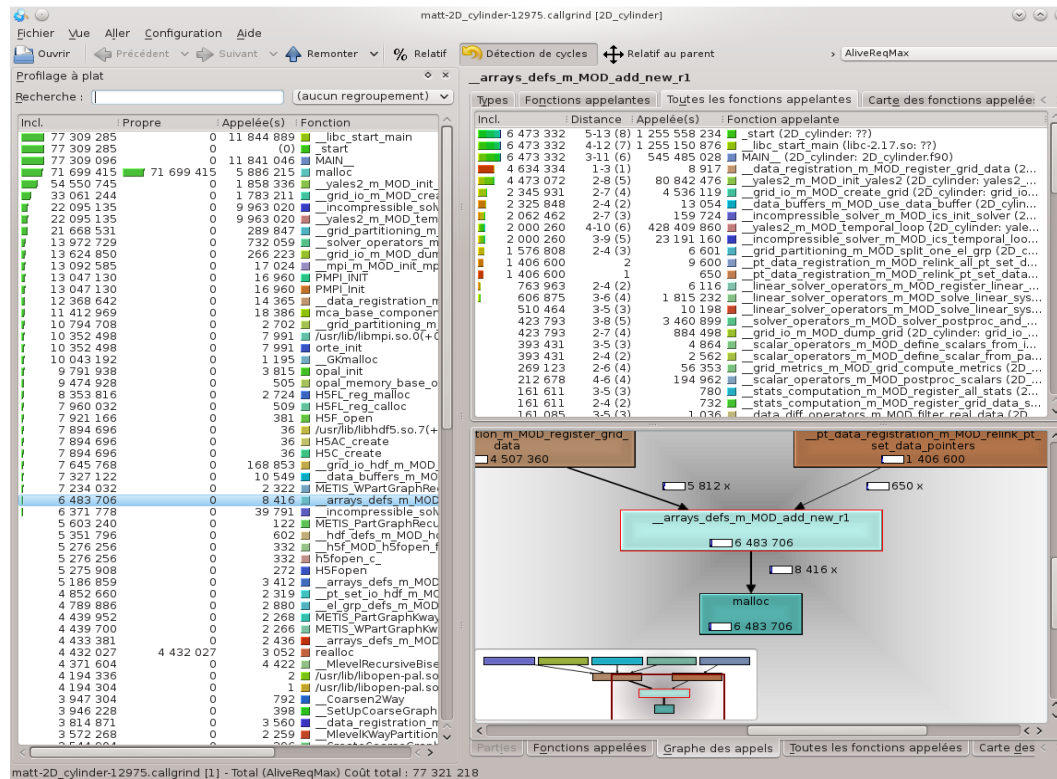
## Callgrind compatibiltiy

- Can use kcachgrind
- Might be usefull for some users, cannot provide all metrics.

- **Started** with **kcacegrind** GUI…. But …

- Display **human readable** units
  – You prefer **15728640** or **15 MB** ?
  – I want to **compare to what I expect**.

- Cannot handle **non sum cumulative metrics**
  – **Inclusive** costs **only rely** on **+ operator**
  – Some mem. metrics **requires max/min** (eg. lifetime)

- No way to express **time charts**

- No way to express **parameter distributions** (eg. sizes).

- Add NUMA statistics

- Provide virtual/physical ratio

- Estimate page fault costs

- Exploit traces in GUI for deeper analysis
  - Alive allocations at a certain time
  - Fragmentation analysis
  - Time charts from call sites
  - Usage over threads for call sites

**Exascale** ∞
computing research

| 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 |
| Requested memory peak | 2.9 KB |

- **Profile over time** :
  - Allocation **rate**
  - **Physical / Virtual / Requested** memory
  - **Stack size** for each **thread** (require function instrumentation)

- **Example on YALES2 with gfortran** :

# EXISTING TOOLS

# Existing tools

- ## Valgrind (massif)
  - Memory **over time** (snapshots) & **fu**
  - Memory per function **at peak**
  - Has a simple GUI

- ## Valgrind (memchek)
  - **Leaks**
  - No real GUI

- ## Google heap profiler (tcmalloc)
  - Memory **over time**  (snapshots)
  - Faster then valgrind
  - No GUI

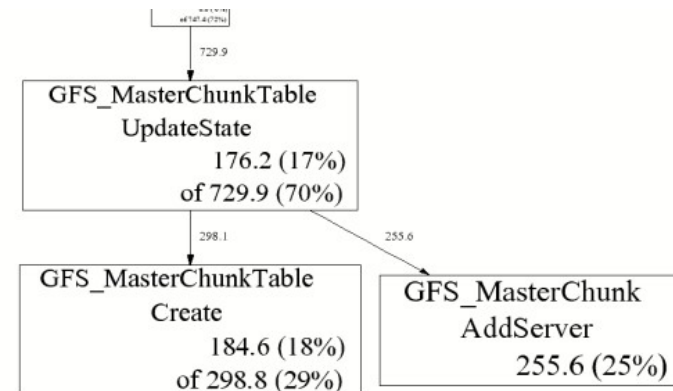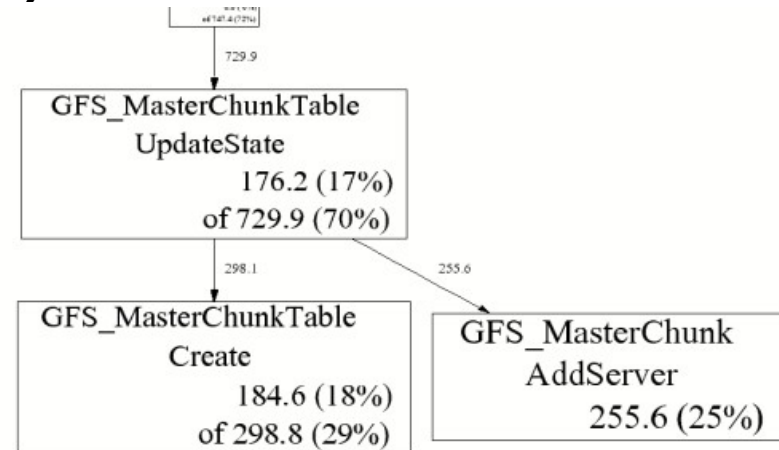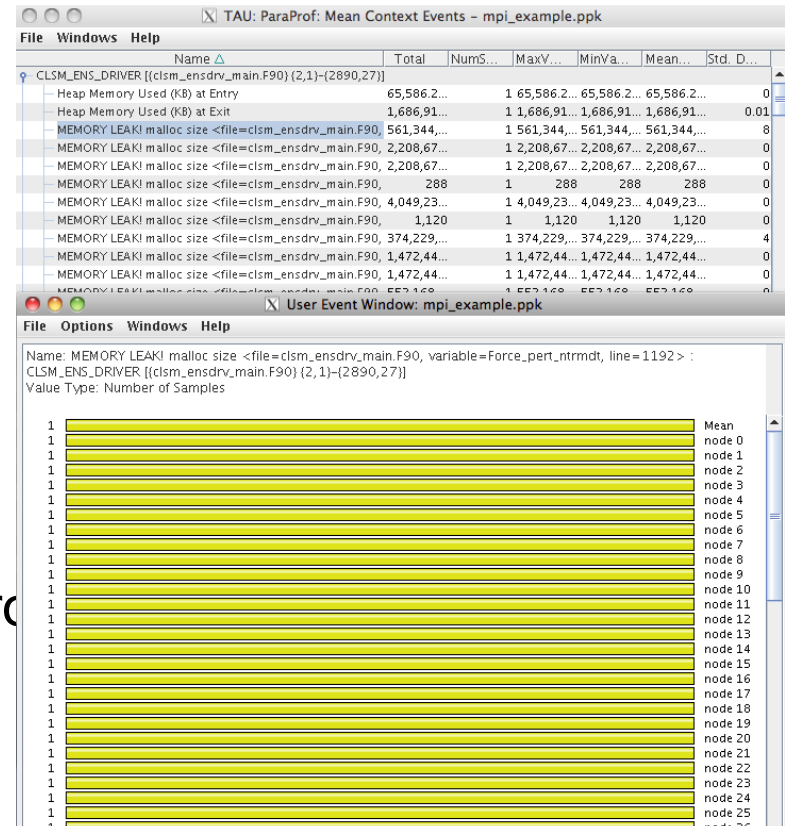- **Google heap profiler (tcmalloc):**
  - **Small overhead.**
  - Similar metric than massif
  - Only provide snapshots of **alloca** **memory per stacks**.
  - Peak might not be captured.
  - Lack of a real GUI to use it.



```
% pprof gfs_master profile.0100.heap
  255.6  24.7%  24.7%    255.6  24.7% GFS_MasterChunk::AddServer
  184.6  17.8%  42.5%    298.8  28.8% GFS_MasterChunkTable::Create
  176.2  17.0%  59.5%    729.9  70.5%
GFS_MasterChunkTable::UpdateState
  169.8  16.4%  75.9%    169.8  16.4% PendingClone::PendingClone
   76.3   7.4%  83.3%     76.3   7.4%
__default_alloc_template::_S_chunk_alloc
   49.5   4.8%  88.0%     49.5   4.8% hashtable::resize
```

**Exascale** *computing research*



- **TAU memory profiler**
  - Provide profiles
  - Follow stacks
  - Track leaks
  - Parallel, done for HPC/MPI
  - Lack easy matching with sourc

- **FOM**

- **IBM Purify++ / Parasoft Insure++**
  - Commercial
  - Leak detection, access checking, memory debugging tools.
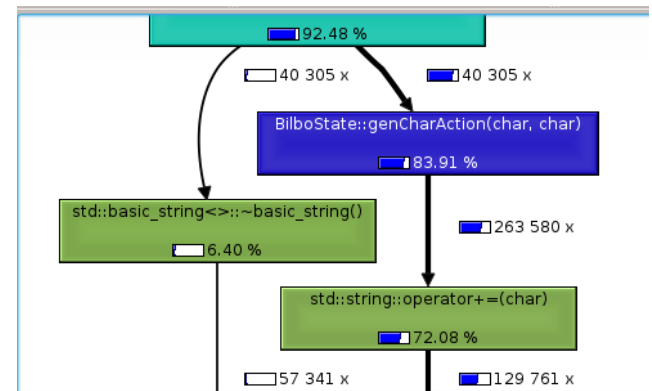  - Use binary or source instrumentation.
  - Windows / Redhat
- **Visual Studio Ultimate Edition Memory profiler**
  - Nice but windows only and commercial

- Two approach implemented : **backtrace** and **instrumentation**

- **Backtrace** (default) :
  - Work out of the box
  - Manage all dynamic libraries
  - **Slow** for **large number of calls** (~>10M)

- **Instrumentation :**
  - Need source **recompilation** (available) : *-finstrument-function*
  - Or tools for **binary instrumentation** : MAQAO / Pintool (experimental)
  - Faster for really large number of calls to malloc
  - **Only** provide stacks for the **instrumented** binaries

# What is good in kcachgrind

Exascale ∞ computing research

- List of **functions** with **exclusive/inclusive** costs

- Nice **call tree**

- **Annotated** sources

# SOME VIEWS

- 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 |

- 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

# Tracking stack memory



**Exascale** computing research

MATT WebView     Summary   Alloc sites   Time analysis   Stack   Alloc sizes   Help

Display largest stack for thread ID

Stack space used by functions on peak

Thread ID : 0

2.1 KB | 54.0 KB grid_io_hdf_m::dump_grid_data_to_hdf | 29.8 KB grid_io_hdf_m::dump_grid_to_hdf | 2.1 KB | 3.6 KB | 7.9 KB | solver_incompressible

Thread ID

0

Stack size over time

77.4
■ 0  26.6 K

**Exascale** computing research

Example from YALES2



Many really small allocations

## Distribution over binaries

○ Grouped ● Stacked                                                                      ● Global variables   ● TLS variables

| Binary | |
|---|---|
| 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

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

# REAL CASES

# Allocatable arrays on YALES2

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

MATT WebView

Search intensive alloc functions

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

Search

911.9 K  data_comm_m::copy_i...

896.4 K  data_comm_m::copy_...

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)
Owns :
Allocated memory : 9.5 MB

And mostly really small allocations !

- 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.

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

```
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
```

1  B  535

- 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**

### Memory allocated over time



○ Virtual memory  ● Physical memory  ● Requested memory