

Allocations mémoires, pourquoi et comment profiler sébastien valat - Jeudis de Gray scott - Février 2024





Origin of the tools

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

MALT : post-doc at Versailles :



NUMAPROF : side project post-doc work at :





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

Motivation

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 !

Key today

4

You need to well understand memory Behavior of you (HPC) Application !

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



Same about **memory consumption** on 12 cores





MALT MALLOC TRACKER

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

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 help searching :

- ▶ Where memory is allocated.
- Properties of allocated chunks.
- **Bad** allocation **patterns** for performance.
- Leaks
- Global variables (TLS)

Some issue examples

Global variables or TLS

Indirect allocations

Leak

Might lead to swap for large size

Short life allocations

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;
```

Existing tools

Valgrind - massif

Valgrind - memcheck









Existing tools

Google heap profiler (tcmalloc):



IBM Purify++ / Parasoft Insure++





What I want to provide

Same approach than valgrind/kcachgind

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)

Global summary

EXECUTION TIME 00:00:00.25	PHYSICAL MEMORY PEAK	ALLOCATION COUNT	AVAILABLE PHYSICAL MEMORY 4.1 Gb		
Run description					
Executable :	simple-case-finstr-lin	ked			
Commande :	./simple-case-f	./simple-case-finstr-linked			
Tool :	matt-0.0.0				
Host :	localhost				
Date :	2014-11-26 22:40	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					
Participated memory peak	20.42					

Per thread statistics



Source annotations



•16

Time charts

17

Sébastien Valat - INRIA / LJK 15/02/2024



System memory



JSON data | CSV data | Save as SVG

Chunk size distribution Example from YALES2 with gfortran issue



•18

Global variables

Distribution over binaries



Distribution over variables





Sébastien Valat - INRIA / LJK 15/02/2024

Real cases

Cerfacs - AVBP- CFD simulator



Example on AVBP init phase

Issue with reallocation on init



Sébastien Valat - INRIA / LJK 15/02/2024 100 101 IF (needed_size>capacity) THEN
 IF (ALLOCATED (temp)) DEALLOCATE(temp) 102 103 104 105 ALLOCATE (temp(capacity)) 57 GB 106 D0 i=1,capacity 107 temp(i)=array(i) END DO 108 110 DEALLOCATE (array) 111 ALLOCATE (array(new cap)) 57 GB 112 113 DO i=1.capacity array(i)=temp(i) 114 115 FND DO 116 117 capacity=new cap 118 END IF 119 Total : Function _start 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)

•22

Coria – YALES2 - Combustion



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



•24

We can found allocs of 1B!







Sébastien Valat - INRIA / LJK 15/02/2024

Usage & conclusion

Usage



Backtrace mode :

Optionally recompile with debug flag to get source lines :
cc -g ...
Run your program
\${PREFIX}/bin/malt [--config=file.ini] YOUR_PRGM [OPTIONS]

Use the web view :

#Launch the server
malt-webserver -i malt-{YOUR_PRGM}-{PID}.json
Connect with your browser on http://localhost:8080



Sébastien Valat - INRIA / LJK 15/02/2024

Reminder on NUMA

Today topology

Example current Intel Knight Landing, mode SNC2 or SNC4

Also add fast memory MCDRAM presented as NUMA or LLC cache



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;</pre>

Then access

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

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...



32

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

NUMPROF HOW TO KNOW IF WE ARE RIGHT IN A REAL APPLICATION ?

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.....

Overhead and scalability

Of course overhead is large: ~30x

► But is scale

Example code hydro on KNL:



36

Sébastien Valat - INRIA / LJK 15/02/2024

GUI and example

Global summary



Statistics per thread



Source & asm annotations



Code Hydro



Sébastien Valat 15/02/2024

► KNL Without HBM



WITH HBM



Original Hydro access matrix





Ordering issue

42

%	Alloc remote access -		/data/svalat/Projects/Hydro/HydroC/HydroCplusMPI/Domain.cpp Domain::setTiles() [clone		
			699		
Search		700	<pre>m_localDt = AlignedAllocReal(m_nbtiles);</pre>		
11.9 G gomp_malloc			756.6 K 701	<pre>m_tiles = new Tile *[m_nbtiles];</pre>	
			702	<pre>#pragma omp parallel for private(i) if (m_numa) SCHEDULE</pre>	
			703	<pre>for (int32_t i = 0; i < m_nbtiles; i++) {</pre>	
	8.4 G Domain::setTiles() [cloneomp_fn.0]		8.4 G 704	<pre>m_tiles[i] = new Tile;</pre>	
				}	
	1.3 G _{Soa::Soa(in}	Soa::Soa(int, int, int)		// Create the Morton holder to wander around the tiles	
•				<pre>m_morton = new Matrix2 < int32_t > (mortonW, mortonH);</pre>	
2.6 M	2.6 M milessipitm	Tile::initTile(Soa*)	708	// cerr << mortonW << " " << mortonH << endl;	
	file::initi		709	<pre>m_mortonIdx = m_morton->listMortonIdx();</pre>	
	1.9 M			<pre>assert(m_mortonIdx != 0);</pre>	
_	no main		711		

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

Non parallel allocations

	%	Alloc re	note access 、		/data/svala	t/Proje	cts/Hydro/HydroC/Hy
	_					22	asing namespac
	Sea	rch				23	ThreadBuffers:
		5.6.0				24	{
	5.6 G	gomp_malloc		25	int32 t lg		
	959.4 M		Domain::setTiles() [cloneomp		26	lgx = (xma	
		959.4 M			27	lgy = (yma	
					28	lgmax = lg	
	_	206.1 M	<pre>Soa::Soa(int, int, int)</pre>	1		29	if (lgmax ·
						30	lgmax
		19.5 M	ThreadBuffers::ThreadBu	iffers(int		31	
					232.2 K	32	$m_q = new$
80.0		1.4 M	Domain::setTiles()		121.4 K	33	$m_q xm = ner$
	_				121.4 K	34	$m_qxp = ner$
-2.8		773.2 K	Tile::initTile(Soa*)		230.7 K	35	m_dq = new
	_				121.4 K	36	m_qleft =
		769.0 K	main		121.4 K	37	m_qright =
	_				121.4 K	38	m_qgdnv =
		128.3 K	Domain	toTootCooo()		39	
			DomainCreaterestCase()	11.5 M	40	m_c = new	
		72 9 14			7.0 M	41	m_e = new
		12.0 N	<pre>matrix2<int>::listMorto</int></pre>	utax()		42	

/droCplusMPI/ThreadBuffers.cpp | ThreadBuffers::ThreadBuffers(int, in ເ ງເພ, :ThreadBuffers(int32 t xmin, int32 t xmax, int32 t ix, lgy, lgmax; ix - xmin); x - ymin); IX; < lgy) = lgy;Soa(NB VAR, lgx, lgy); w Soa(NB VAR, lgx, lgy); w Soa(NB VAR, lgx, lgy); Soa(NB VAR, lgx, lgy); new Soa(NB VAR, lgx, lgy); new Soa(NB VAR, lgx, lgy); new Soa(NB VAR, lgx, lgy); Matrix2 < real t > (lgx, lgy); Matrix2 < real t > (lgx, lgy);

Parallel allocations

Original

for (int32_t i = 0; i < m_numThreads; i++) {
 m_buffers[i] = new ThreadBuffers(...);
 assert(m_buffers[i] != 0);</pre>

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