

LProf

Version 2.2 May 2025

www.maqao.org

AQAO TUtorial series

1 Introduction

MAQAO Lightweight Profiler (LProf) is the MAQAO module which allows to easily profile an application to detect hot functions and loops in two steps:

1) Data collection using sampling

LProf uses hardware counters to profile large-scale parallel applications (2000+ cores) with a very low overhead.

It is also possible to provide a custom list of hardware counters to sample.

2) Data display

LProf output allows to quickly identify time-consuming functions and loops, observe the amount of time spent by the application between different categories (I/O, Runtime, etc...) and detect load balancing issues.

2 Running MAQAO LProf

2.1 Sequential Run Command



2.2 Parallel Run Command (version 2.5+)

Interactive runs: Interactive runs:



Runs with launch script (typically to submit a job) since version 2.20:



Older versions (2.5 to 2.19): use --batch-script and --batch-command

In jobscript, application executable and its arguments have to be replaced by <run_command>.

```
$ cat jobscript.sh
...
mpirun -n 4 <run_command> # instead of mpirun -n 4
<application> [args]
# <mpi command> <run command> # if mpi-command used
```

Since 2.12.0, you can (and must) inform Lprof about the maximum number of processes per node (if greater than 1), allowing it to set correct internal settings: --maximum-processes-per-node

Starting from 2.14.5, it is autodetected when missing but it is still recommended to set --maximum-processes-per-node if known and > 1.

2.3 Kernel samples exclusion

Since 2.12.0, kernel samples are not collected by default (recent Linux distributions do not allow this by default). To collect them:

 If sysctl kernel.perf_event_paranoid returns 2 or more, this step must be performed first:

```
$ sudo sysctl -w kernel.perf_event_paranoid=1
# lost after reboot
$ sudo sh -c `echo kernel.perf_event_paranoid=1 >>
/etc/sysctl.d/local.conf
# persists after reboot
```

• If sysctl kernel.perf_event_paranoid returns 1 or less:

```
$ maqao lprof -include-kernel ...
```

2.4 Options (collect step)

To list all options along with their descriptions:

maqao lprof --help

Options in gold color can be used to mitigate sampling overhead. Options in light red can be used to override default behavior and workaround profiling issues.

Main options (collect step)						
Name	Short Description	Values				
include-kernel	No effect with the 'no- perf' engine. Count kernel samples (requires perf-event- paranoid level 1 or less)	(no value)				
-mc/mpi- command=	Specify command for interactive MPI run or replacement value for <mpi_command> in job script</mpi_command>	Ex:"mpirun -n 4"				
-ls/launch- script=	Script used to launch application. If it is executable (user-exec perm. and shebang), LProf directly executes it, otherwise Lprof executes <launch- command> <launch- script></launch- </launch- 	Path to script (string)				
-lc/launch- command=	Interpreter (e.g bash, python) or batch submission command to launch <launch- script> (as prefix). If omitted, guessed from</launch- 	Ex: "sbatch"				

	<launch-script> extension</launch-script>	
stdin-path	Defines a file for redirection to stdin.	path to a stdin-redirection file
sampling-rate=	Number of collected samples per second	 highest (2000 Hz, btm=off recommended) high (1000 Hz, avoid btm=stack) medium (200 Hz, default) low (50 Hz) lowest (10 Hz)
-ldi=	Scan debug information into all or specified (provided list) library(ies) to get loops details	on (all) off (default) or r list of libraries ('lib1, lib2,')
start-after- seconds=	Setup sampling start delay	 0 (default) -1: never starts n positive integer: delay
pause-resume- at-SIGTSTP	Each time LProf receives SIGTSTP (from the target application or CTRL+Z), sampling is paused (if running) or resumed (if paused).	(no value)
stdout-start- keywords=	Each time LProf detects a keyword from stdout, sampling is started (no effet if already running)	Comma-separated list of words
stdout-stop- keywords=	Each time LProf detects a keyword from stdout, sampling is stopped (no effet if	Comma-separated list of words

MAQAO Tutorial series: LProf 6

	already paused)	
-btm=	Select backtraces (callchains) collection method	 fp (default on AArch64, recompile application with - fno-omit-frame-pointer) stack (default on x86-64, higher overhead but no need to recompile application) branch (not really callchains but branch history, HW-dependent) off (no callchains, lowest overhead)

Advanced/other Options (collect step)						
Name	Short Description	Values				
use-OS- timers	Use OS timers instead of hardware events. Needed in case of unavailable HW counters or undetected processor. With autotuning features	(no value)				
cpu-clock- MHz	[perf-* engines] Override the "cpu-clock" perf-event rate (in MHz) measured by a calibration loop.	integer value				
ref- cycles-MHz	[perf-* engines] Override the "ref-cycles" perf-event rate (in MHz) measured by a calibration loop.	integer value				
replace	Overwrites an already existing output directory (reuse it). Remark: no effect on a not yet existing directory.	(no value)				
-tpp/	[perf-high-ppn only] Maximum	integer value				

maximum- threads-per- process	number of concurrent threads per process. Default is OMP_NUM_THREADS. Used to set buffers and files size.	
 ppn/maximum- processes- per-node	Since 2.12.0, mandatory when usingmpi-command Optional but recommended starting from 2.14.5 if ppn > 1	Ex on single node: lprof mpi- command="mpirun -n 32" ppn=32
maximum- buffer- megabytes	Allow to override Lprof memory footprint (default is 50 MB per CPU)	Maximum amount per node (Megabytes)
maximum- tmpfiles- megabytes	Limit total temporary files size to X Megabytes per node. Default is 100 MB per CPU (HW thread).	Integer value
-e/evts	Provide custom list of events to sample (CF maqaolist-events)	evt1_name@sample_period, or evt1_code@sample_period,
cnt-evts	[EXPERIMENTAL] Provide a custom list of events (CF maqao list-events) to profile (counting). Use only dynamic PMU events (not counted by the CPU cores PMUs), which requires 0 or negative paranoid level (sudo sysctl -w kernel. perf_event_paranoid=0).	Ex:cnt- evts=RAPL_ENERGY_CORES,UNC_M_CA S_COUNT_IMC0.RD
-p/evts- profiles	Use ready-to-use lists of events. Not yet supporting more than one profile.	string
cnt-evts- profiles	[EXPERIMENTAL] Use ready- to-use lists of events (counting). Presently supported: ENERGY, DRAM_READS and	string

	DRAM_WRITES	
cnt- metrics	[EXPERIMENTAL] Counting metrics. Presently supported: - ENERGY_{PKG,DRAM} (add ENERGY into cnt-evts-profiles) - DRAM_{READS,WRITES} (add DRAM_{READS,WRITES} into cnt-evts-profiles)	string
max-	Maximum callchain length	Positive integer
callchain- length	(default: 20), useful to reduce btm=stack overhead.	
stack-size	Size (in bytes) of stack to dump on samples (default: 8192). Using a smaller size (typically 4096) reduces profiling overhead but may cut (or loose) callchains. Using a bigger size (typically 16384) increases profiling overhead but should guarantee minimal callchains loss.	Positive integer
mmap-pages	Overrides autotuned number of mmap pages for ring buffer payload.	Positive integer
collect- calls-info	Collects source file/line information for callchain nodes (calls). To display them, add use-calls-info=on at display step.	on (default)/off
engine	Use another perf-events based sampling engine	 perf-low-ppn (selected by default when perf-events are available with max 4 processes per node) perf-high-ppn (selected by default when perf-events are available with more than 4

		 processes per node) no-perf (selected by default when perf-events are not available)
include- sleep-time	[no-perf only] Include sleep time (walltime).	(no value)
keep- external- threads	[perf-high-ppn engine only] Profile threads with a different command line than the monitored application.	on/off (default)
keep- indirect- threads	[perf-high-ppn engine only] Profile threads that are not direct children of the monitored application.	on (default)/off
-cpu/cpu- list	Set CPU affinity for the target process. Ex: 0,2 to use CPU0 and CPU2.	comma-separated list of integers
ignore- signals	[no-perf and perf-high-ppn engines] Prevents signals from being interpreted as termination signals. Allows to adapt no-perf and perf-high-ppn to various runtimes. Remark: for ignored signals also specified in set- exit-signals or set-abort-signals, evaluation order is set-abort- signals, set-exit-signals and then ignore-signals.	comma-separated list of integers
set-exit- signals	[no-perf and perf-high-ppn engines] Interpret signals as normal application exit. Allows to adapt no-perf and perf-high- ppn engines to various runtimes. Remark: for exit signals also specified in ignore- signals or set-abort-signals,	comma-separated list of integers

	evaluation order is set-abort- signals, set-exit-signals and then ignore-signals.	
set-abort- signals	[no-perf and perf-high-ppn engines] Interpret signals as abnormal application exit. Allows to adapt no-perf and perf-high-ppn engines to various runtimes. Remark: for abort signals also specified in ignore-signals or set-exit- signals, evaluation order is set- abort-signals, set-exit-signals and then ignore-signals	
legacy- maps	[ADVANCED] Use only if unknown functions coverage is high for executable or libraries. Collect maps via legacy method (out of perf-events) after <legacy-maps> milliseconds and fallback to them in case of unresolved addresses.</legacy-maps>	Positive integer (number of milliseconds)
collect- CPU-time- intervals	[ADVANCED] [perf-low-ppn and perf-high-ppn engines] Collect per-thread CPU-time intervals. Allows to trace when and where (CPU) threads was running, and display them by adding - verbose at display step.	(no value)

2.5 Collect step hints

In case of multiple application processes (typically MPI ranks), use collectcalls-info=off to limit LProf memory footprint when dumping to disk source file/line for each call listed in callchains.

3 Display

The two common display modes are text (default) and HTML.

3.1 Concepts

LProf relates *code regions* contributions to *system-levels*. User must then specify which code regions he is interested in and at which system level/granularity.

3.1.1 Code regions (hotspots)

From bigger to smaller:

- Application: set of modules
- Module: set of functions
- Function: set of loops
- Loop: set of blocks
- Block: basic block (compilation concept)

3.1.2 System levels

From bigger to smaller:

- Cluster: set of nodes (machines)
- Node: set of (system) processes
- Process: set of (system) threads
- Thread

3.2 Text Output

3.2.1 Functions Hotspots

To display summary view (at cluster level):

<pre>maqao lprof -df xp=<experiment_directory></experiment_directory></pre>						
hotspots' lev	vel: functions	experin	nent dire	ectory (name or data to display	path) contai /	ning
######################################	Module	Source Info	######################################	6) Time Min(s) [TID]	######################################	Time w.r.t Walltime (s)
<pre>####################################</pre>	bt-mz.A.4 bt-mz.A.4 bt-mz.A.4 bt-mz.A.4 bt-mz.A.4 bt-mz.A.4 bt-mz.A.4 bt-mz.A.4 bt-mz.A.4 bt-mz.A.4	solve_subs.f:206 z_solve.f:45 solve_subs.f:56 y_solve.f:45 x_solve.f:48 rhs.f:33	23.90 13.89 13.39 13.14 12.39 12.12	1.98 [12705] 1.14 [12693] 1.12 [12699] 1.02 [12693] 0.84 [12706] 0.98 [12698]	2.58 [12693] 1.48 [12692] 1.48 [12692] 1.56 [12698] 1.56 [12698] 1.42 [12705] 1.28 [12707]	2.25 1.31 1.26 1.24 1.16 1.14

Figure 1 - LProf Output: Summary View (Functions)

To display view for a lower system level, use -dn (resp. dp, dt) for node (resp. process, thread). For instance, to display thread view:

maqao lprof -df xp=<EXPERIMENT_DIRECTORY> -dt

	Thread	ID	Hostna	ame		Pro	ocess ID		V	Valltime			
			~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~	BACRI - P	ROCESS #1	2693							
			T	nread #12693	- 9.40 s	econd(s							
##	*****	******	*****	,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,	******		*****		******	,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,	*****	****	*****
#	Function Name	۲ I	odule	Source	Info	Covera	age (%)	Time v	I.F.t Wa	illtime (s)		CP1 r	atio #
##	************************	******	*****	*****	*******	*######	*****	*****	########	**********	*####	#####	#######
#	binvcrhs	bt-mz.A	4	solve_subs	.f:206	27.45		2.58				0.68	#
#	matmul_sub	bt-mz.A	4	solve_subs	.f:56	15.74		1.48				0.69	#
#	compute_rhsomp_fn.0	bt-mz.A	.4	rhs.f:33		12.77		1.20				0.99	#
#	z_solveomp_fn.0	bt-mz.A	4	z_solve.f:	45	12.13		1.14				0.71	#
#	x solve . omp fn.0	bt-mz.A	.4	x solve.f:	48	11.91		1.12				0.74	#
#	y solve . omp fn.0	bt-mz.A	.4	y solve.f:	45	10.85		1.02				0.67	#
#	matvec_sub	bt-mz.A	4	solve_subs	.f:27	2.98	i.	0.28			i.	0.62	#

Figure 2 - LProf Output: Thread View (Functions)

#### **3.2.2 Loops Hotspots**

To display summary view (at cluster level):



Figure 3 – LProf Output: Summary View (Loops)

The above figure is truncated. In the actual output, four more columns are available on the right (same as functions mode):

# Coverage (%), Time Min (s), Time Max (s) and Time w.r.t Walltime (s).

As for functions, use -dn/dp/dt to select a lower system level. For instance, to display thread view:



Figure 4 - LProf Output: Thread View (Loops)

## 3.3 Display Options

Basic Options (display step)					
Name	Short Description	Values			
-df/-dl	Display functions/loops	(no value)			
-db	Display basic blocks (for finer granularity than loops)	(no value)			
-dn	Display per-node profiles (instead of cluster by default)	(no value)			
-dp	Display per-process profiles (instead of cluster by default)	(no value)			
-dt	Display per-thread profiles (instead of cluster by default)	(no value)			
-lec/ libraries- extra- categories	Consider specified libraries as extra categories	libraries names as given by 'ldd <application>'</application>			
-of/ output- format	Output results in a file of the given format (default if omitted: console output)	html or csv			
-cc/ callchain	<ul> <li>Specify objects for callchains analysis:</li> <li>exe: display the callchain (if available) for each function with a scope limited to the application.</li> <li>lib: extend the callchain scope to external libraries function calls.</li> <li>all: display the callchain with no limited scope (application + libraries + system calls).</li> <li>off: disable callchains</li> </ul>	exe, lib, all or off			

	analysis. Some OpenMP/MPI functions/loops will no more be correctly categorized. Use this only when display takes too much time/memory.	
-ct/ cumulative- threshold	Display the top loops/functions up to a given cumulated coverage (e.g: ct=50).	integer between 0 and 100

### 3.4 HTML Output

### 3.4.1 Generation of HTML results



This command generates an 'index.html' file into the *<EXPERIMENT_PATH>/html/* directory. Open this file into a web browser to see the results.

#### **3.4.2 Interpretation of the Results**

Refer to the Oneview tutorial: <u>https://maqao.org/documentation/MAQAO.Tutorial.ONEVIEW.pdf</u>