

Performance Analysis at Scale: The Score-P Tools Infrastructure



Frank Winkler (frank.winkler@tu-dresden.de)



Performance tools will not automatically make your code run faster. They help you understand, what your code does and where to put in work.





Performance engineering workflow



High Performance Computing







SO, YOU HAVE DECIDED TO UNDERSTAND WHAT A PROGRAM EXACTLY DOES?







Score-P: Functionality

- Typical functionality for HPC performance tools
 - Instrumentation (various methods)
 - Sampling (experimental)
- Flexible measurement without re-compilation
 - Basic and advanced profile generation
 - Event trace recording
- Programming paradigms:
 - Multi-process
 - MPI, SHMEM
 - Thread-parallel
 - OpenMP, Pthreads
 - Accelerator-based
 - CUDA, OpenCL, OpenACC



– Hybrid parallelism



Score-P: Architecture











• To see all available options for instrumentation:







Score-P Workflow: Advanced Instrumentation

• For CMake and autotools based build systems it is recommended to use the scorep-wrapper script instances



• Pass instrumentation and compiler flags at make







• Measurements are configured via environment variables

```
$ scorep-info config-vars --full
SCOREP_ENABLE_PROFILING
[...]
SCOREP_ENABLE_TRACING
[...]
SCOREP_TOTAL_MEMORY
Description: Total memory in bytes for the measurement system
[...]
SCOREP_EXPERIMENT_DIRECTORY
Description: Name of the experiment directory
[...]
```

• Example for generating a profile:

```
$ export SCOREP_ENABLE_PROFILING=true
```

```
$ export SCOREP_ENABLE_TRACING=false
```

```
$ export SCOREP_EXPERIMENT_DIRECTORY=profile
```

\$ aprun <instrumented binary>





- Profile analysis tool for displaying performance data of parallel programs
- Originally developed as part of Scalasca toolset
- Available as a separate component of Score-P
- Representation of values (severity matrix) on three hierarchical axes
 - Performance property (metric)
 - Call-tree path (program location)
 - System location (process/thread)
- Three coupled tree browsers









Score-P: Cube Analysis Presentation









Score-P Workflow: Filtering

- Use scorep-score to define a filter
 - Exclude short frequently called functions from measurement
 - For profiling: reduce measurement overhead (if necessary)
 - For tracing: reduce measurement overhead and total trace size



• Filter file:



Reduce measurement overhead and size of event trace to about 40 GB! (Example uses 4 processes)

• Example for generating a trace

```
$ export SCOREP ENABLE PROFILING=false
```

```
$ export SCOREP ENABLE TRACING=true
```

```
$ export SCOREP EXPERIMENT DIRECTORY=trace
```

```
$ export SCOREP TOTAL MEMORY=2G
```

```
$ export SCOREP_FILTERING_FILE=scorep.filt
```

```
$ aprun <instrumented binary>
```





- Show dynamic run-time behavior graphically at a fine level of detail
- Provide summaries (profiles) on performance metrics

Timeline charts

• Show application activities and communication along a time axis

		84.8 s	84.9 s	85.0 s	85.1 s	85,2 s
Process 0	YSU		- 1	W.C.	11 100	-
Process 1		APR WAR	H		1 1	
Process 2		A AME	Wat			
Process 3			species in the second	MPT Wit		
Process 4	950	CUMULUS DRI	VER			
Process 5	the second se	A No America	Wat		11 18	
Process 6		A MPI	Wat			
Process 7		MAG		AHPT Wait		
Process 8	YSU	CUMULUS OR	VER			
Process 9	CUMULUS	DRIVER MAT	Wat			
Process 10	-	AND	Wat /			
Process 11	Contraction of the	77	and the second	HIPT Weat		
Process 12	YSU CUMU	UUS ORIVE		Concession in succession of the local division of the local divisi		N
Process 13	and the second s	/M MAT	Wat			11
Process 14		AND SHE AND	Wait		11 121	
Process 15		and the second s	Support of the local division of the local d	MPT Wat	11 1/11	

Summary charts

• Provide quantitative results for the currently selected time interval







Vampir: Performance Charts



Master Timeline



• Trace visualization of FDS (Fire Dynamics Simulator)



Vampir at Scale: FDS with 8192 cores

DRESDEN

• Fit to chart height feature in Master Timeline



Center for Information Services & High Performance Computing

• Jacobi Example

- Iterative solver for system of equations $U_{old} = U$

$$u_{i,j} = bu_{old,i,j} + a_x(u_{old,i-1,j} + u_{old,i+1,j}) + a_y(u_{old,i,j-1} + u_{old,i,j+1}) - rHs / b$$

 Code uses OpenMP, CUDA and MPI for parallelization







High Performance Computing

• Domain decomposition

- Halo exchange at boundaries:
 - Via MPI between processes
 - Via CUDA between hosts and accelerators



• Documentation at https://www.olcf.ornl.gov/support/software/



Demo: Jacobi Solver / Setup

Connect to Summit-dev and copy sources

- \$ cp /ccs/home/winklerf/scorep_tutorial/jacobi.tar.gz .
- \$ tar xzvf jacobi.tar.gz
- \$ cd jacobi

Change programming environment and load modules

```
$ module load gcc/5.4.0
$ module load cuda
```

```
$ module load scorep
```

• Compile benchmark and submit job

```
$ make
$ cd bin
$ bsub < run.lsf
$ less run.lsf
Jacobi relaxation Calculation: 8192 x 8192 mesh with
2 processes and 6 threads + one Tesla P100-SXM2-16GB for each process.
614 of 4097 local rows are calculated on the CPU to balance the load
between the CPU and the GPU.
0, 0.489197
100, 0.002397
[...]
total: 9.409952 s
```

• Build instrumented executable

```
$ make clean
$ make scorep
scorep --cuda cc ... -o bin/jacobi mpi+openmp+cuda
```

Submit job for profiling run

```
$ less run_profile.lsf
[...]
export SCOREP_ENABLE_PROFILING=true
export SCOREP_ENABLE_TRACING=false
export SCOREP_EXPERIMENT_DIRECTORY=jacobi_profile
export SCOREP_CUDA_ENABLE=yes
[...]
mpirun -n 2 ./jacobi_mpi+openmp+cuda 8192 8192 0.15
$ bsub < run_profile.lsf
$ less jacobi.o[JOB_ID]
Jacobi relaxation Calculation: 8192 x 8192 mesh with
2 processes and 6 threads + one Tesla P100-SXM2-16GB for each process.
[...]
total: 10.678350 s
```

Demo: Jacobi Solver / Profile Analysis

• Perform flat profile analysis with cube_stat

\$ cd bin								
<pre>\$ cube_stat -t 10 -p jacobi_mpi+openmp+cuda_profile/profile.cubex</pre>								
cube::Region	NumberOfCalls	ExclusiveTime	InclusiveTime					
!\$omp for @jacobi_cuda.c:188	32000.000000	131.797289	131.797289					
!\$omp implicit barrier	32000.000000	104.298683	104.298683					
!\$omp for @jacobi_cuda.c:258	32000.000000	42.999056	50.568642					
[]								

Perform call-path profile analysis with Cube

cube jacobi profile/profile.cubex \$ Cube-4.3.4: acobi mai+esemmp-icuda profile/profile cubes Ets Duplay Engine Help Restore Setting . Sale Settings Absolute . Absolute Add (Colors . Mattic tree Califier Fatures System bas | Boulftat 3 Haft Visits mett 0.00 main Description over 1997 🖬 0.06 init_mpi - Caloriat: Iow 0.00 Minimum Inclusive Time Iseci-0.01 Nandia command line arguments 🗈 🗖 – catinat lä 16.69 Maximum Inclusive Time Issuit - cept 0 1.19 init, heat 0 0 types_pit dyteic 0.37 init_codu - Di - Blasie II 🗋 0 bytes, get dytes) i di bió mateix. i - D - attic 8448 5 31v8 ALLOCATION_SIZE (bytes) IS DO MP1_IRamai - nic 10208 5.3748 DEALLOCATION_SIDE INVALID 0.00 start times - node rid03734 C bytes jeaked bytest E 45 29 MPL Rate: 2 3 48 (81)(8) E 19 wild taurich, jacobic, kemel, any klocolet Roat, Nutr. Vol. Int. Nutr. 2 thet maximum_heap_memory_attacated bytes: 0-0 - Mare 7 2 60x8 Process inemory usage distant to he thomp parallel @secator_could in TB4 - C) - 2010 8447 8.10547 Sylas_asti Sylas) CI - AL 10201 8.10x7 bytes_received drytesi 8.11-Bong-critical @jacobi_stada C213 C - node md037p3 C 415 20 MP1 Mare 1 C 02 Barry colical struck @pacets_cuda c 213 3 102 00 Borry orgholt karner @pacelor, custa c 217 C 14 Road walk, pacebo, transactioned Road*1 11 11.11 D 14 MP1 Adresius . 14 [4] 4 4 All (36 eleitmetts) 305 (190 10%) 100 ID 00 130 59 (44.10%) 296.000 10.00 130.59 **Salected** "Time"

Demo: Jacobi Solver / Scoring

• Do we need a filter? (Overhead and memory footprint)

							No filtoring	
\$ scor Estima	required.							
Estimated requirements for largest trace buffer (max buf): 5MB								
Estimated memory requirements (SCOREP_TOTAL_MEMORY): 41MB								
(hint: When tracing set SCOREP_TOTAL_MEMORY=41MB to avoid intermediate								
flushes or reduce requirements using USR regions filters.)								
flt	type	max_buf[B]	visits	time[s]	time[%]	time/visit[us]	region	
	ALL	4,924,060	310,504	308.53	100.0	993.63	ALL	
	OMP	4,135,850	256,417	287.31	93.1	1120.46	OMP	
	CUDA	494,338	38,025	10.40	3.4	273.53	CUDA	
	COM	156,260	12,020	10.46	3.4	870.58	COM	
	MPI	137,222	4,012	0.30	0.1	73.96	MPI	
	MEMORY	260	20	0.06	0.0	2972.15	MEMORY	
	USR	130	10	0.00	0.0	10.26	USR	





• Submit job for tracing run

```
Ścd..
$ less run trace.lsf
[...]
export SCOREP ENABLE PROFILING=false
export SCOREP ENABLE TRACING=true
export SCOREP EXPERIMENT DIRECTORY=jacobi trace
export SCOREP CUDA ENABLE=yes
export SCOREP TOTAL MEMORY=50MB
[...]
mpirun -n 2 ./jacobi mpi+openmp+cuda 8192 8192 0.15
$ bsub < run trace.lsf</pre>
$ less jacobi.o[JOB ID]
Jacobi relaxation Calculation: 8192 x 8192 mesh with
 2 processes and 6 threads + one Tesla P100-SXM2-16GB for each process.
 614 of 4097 local rows are calculated on the CPU to balance the load
 between the CPU and the GPU.
    0, 0.489197
 100, 0.002397
  [...]
  900, 0.000269
 total: 9.895828 s
```

Demo: Jacobi Solver / Trace Analysis

• Perform analysis on the trace data with Vampir

- \$ cd bin
- \$ module load vampir
- \$ vampir jacobi_trace/traces.otf2



