



Profiling Tutorial

Soon Yung Jun (Fermilab, SCD/Physics and Detector Simulation Group)
LArSoft Tools and Technology Workshop
20 June 2017, Fermilab

Debugging is finally done! Ready for a test drive?



My program runs, but seems very slow ...

CPU

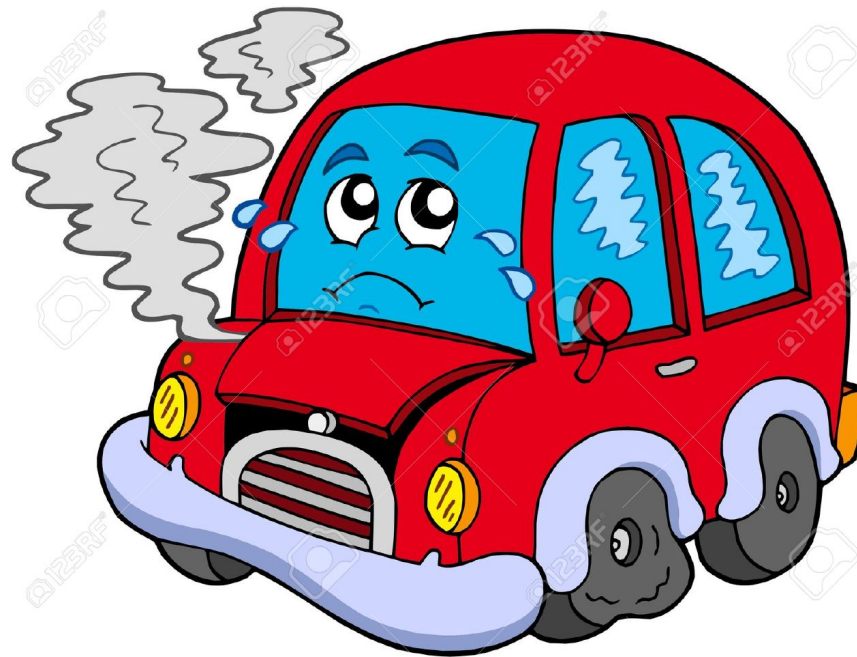
Throughput

Instruction stall

Memory

Latency

Cache misses



I/O

Communication

DB contention

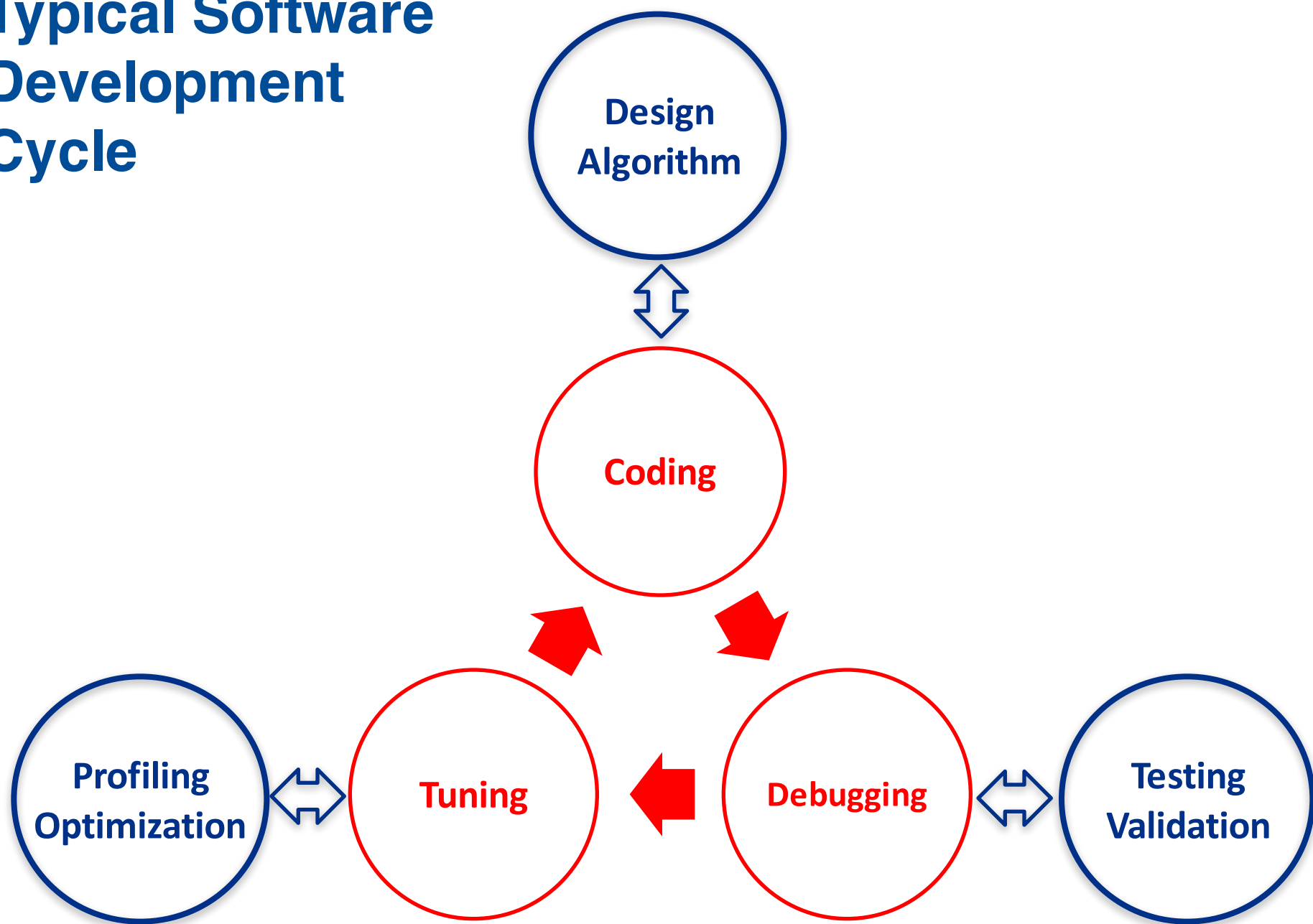
Multithreading

Load balancing

Scalability

Where to start?

Typical Software Development Cycle



This Tutorial

- PART-I
 - A brief introduction to computing performance profiling
 - An overview of selected profiling tools and examples
- PART-II
 - Profiling results of the LArTest application with IgProf and OpenSpeedshop
- Demos and Questions



PART-I

Introduction to Computing Performance Profiling
Overview of Selected Profiling Tools and Examples

Why Profile?

Performance tuning is an essential part of the development cycle

- Free lunch is over as modern hardware architectures are getting more complex and parallel
- HEP applications are usually complicated too
- Every \$/Watt matters (computing with a limited budget)
- Understanding the code performance is responsibility of the software developer
- **Maximize CPU flop rate and minimize memory operations (balancing them is not an easy task)**

Computing Performance Profiling and Analysis

- Performance benchmarking quantifies usage/changes of CPU time and memory (amount required or churn)
- Performance profiling analyzes
 - Hot spots, bottlenecks and efficient utilization of resources
 - Code efficiency (instruction/cycle, latencies, I/O and etc.)
- Identifying opportunities for optimization



Low Hanging Fruit



Tools

After Production

Understanding Computer Performance

- Hardware platform (processors) popularly used in HEP
 - CISC (x86), RISC (ARM), MIC, GPU(SIMT), FPGA
- Speed: cycle vs. frequency
 - cycle time = $1/(\text{clock frequency})$
 - 2.0 GHz = 0.5 ns per cycle
 - CPU Time = $\Sigma(\text{number of clock cycles})/\text{frequency}$
- Memory: latency vs. bandwidth
 - latency: the time interval between the request for information and the access (to the first bit of that information)
 - Bandwidth: the number of bits per second
- Throughput vs. locality: CPI, MIPS, FLOPS, FMO
- Pipelining: instruction throughput, data dependency, ILP, ...

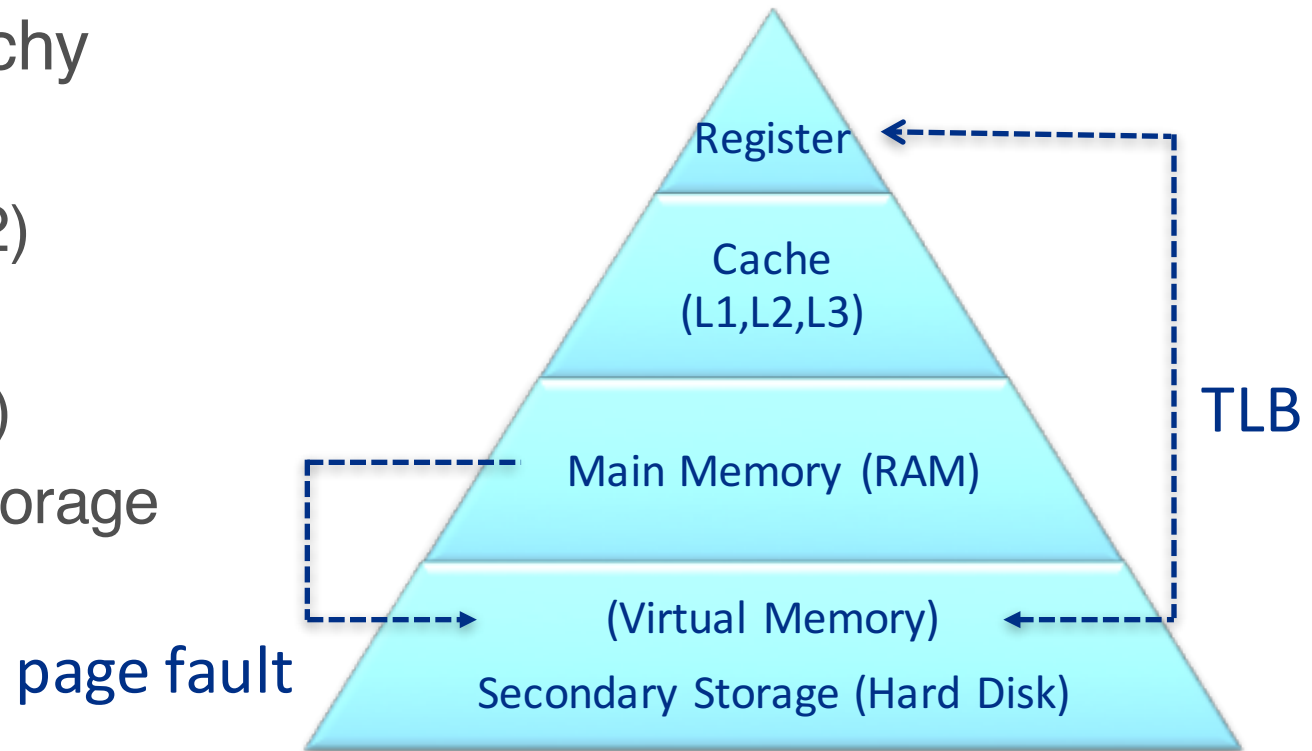
Understanding CPU Performance

- Q1: Which operation takes more cycles?
 1. Integer division
 2. Double division
 3. Function call
 4. `static_cast<int>(double)`
- Strategies
 - Do not mix data type
 - Avoid unnecessary divisions and function calls in the inner most loop

Understanding Memory Transaction

- Memory hierarchy

- Registers
- Cache (L1/L2)
- DRAM (rss)
- Virtual (vsize)
- Secondary storage

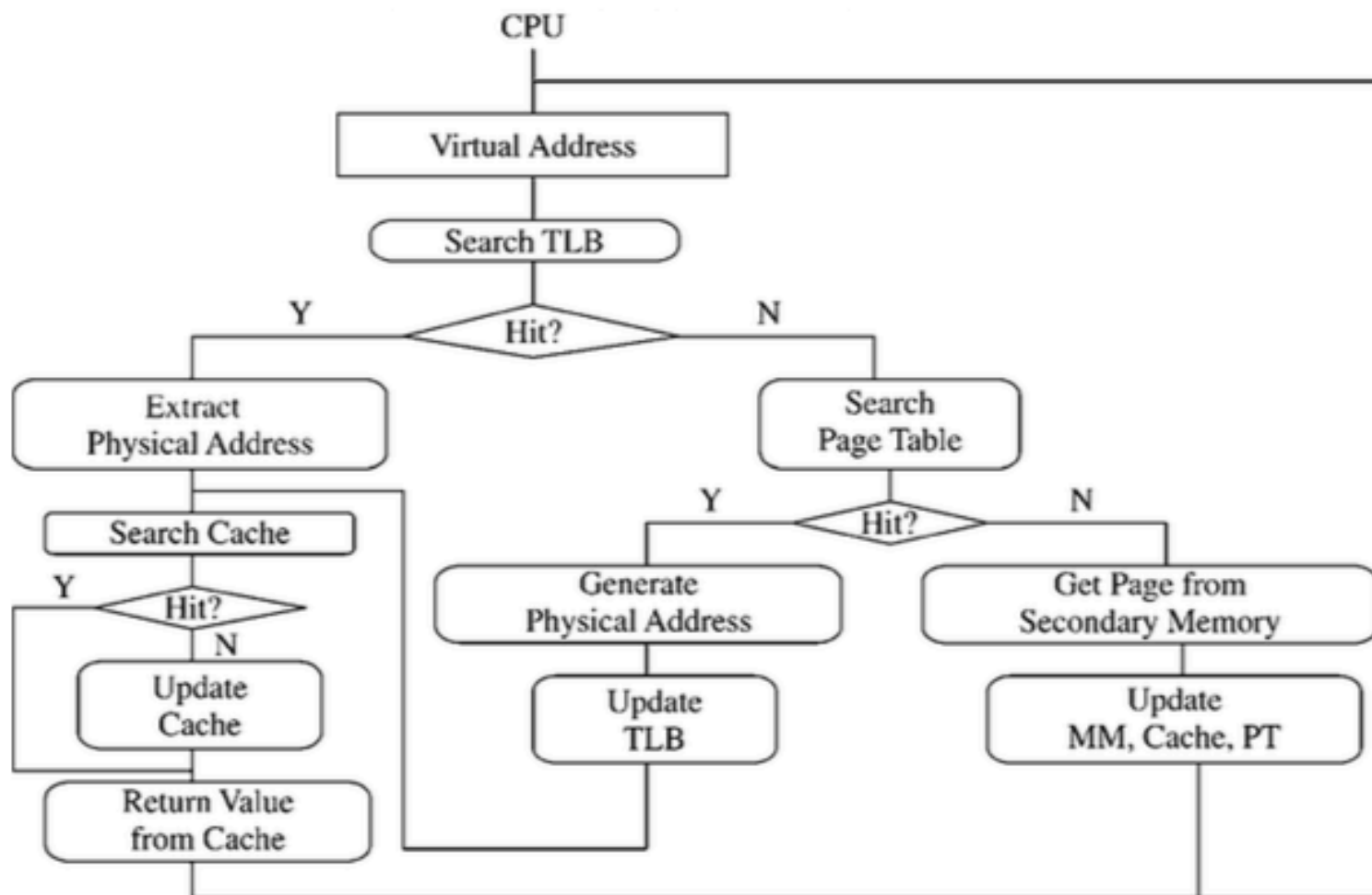


- Caching

- Spatial locality (data storage, coalescence)
- Temporal locality (data reusability in near future)
- Replacement policies
- TLB (translation look-aside buffer, the most recent page access)

Understanding Memory Transaction

- Example of memory accesses scenarios



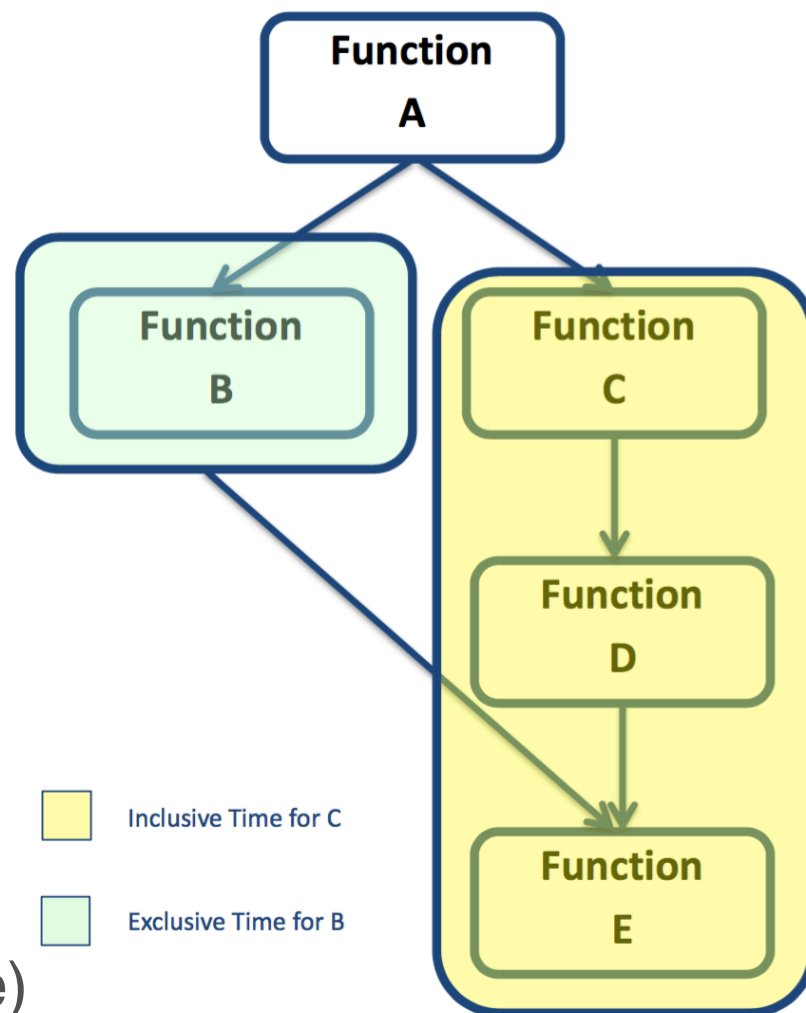
- Do not over-optimize by yourself, but rely on profiling first

Understanding Memory Performance

- Q2: Which ratio is the biggest in memory access?
 1. L1 Cache/Register
 2. L2 Cache/L1 Cache
 3. RAM/L2 Cache
 4. Virtual Memory/RAM
- Strategies
 - Try to fit everything in RAM
 - Try to fit essential calculations in cache

Basic Concepts of Performance Profilers

- Program segments:
 - Code
 - Stack (program)
 - Heap
- Collecting program events
 - Hardware interrupts
 - Code instrumentation
 - Instruction set simulation
 - Tracing (when)
- Periodic sampling
 - Top of the stack (exclusive)
 - Anywhere in the stack (inclusive)



Classification of Profilers by Techniques used

- **Instrumentation:** inserts extra code at each function call to count how many times the function is called and how much time it takes.
- **Sampling:** The profiler tells the operating system to generate an interrupt and counts how many times an interrupt occurs in each part of the program
 - no modification of the program
 - time-based
 - event-based
- **Debugging tools:** The profiler inserts temporary debug breakpoints at every function or every code line (valgrind)

Examples of Profilers

- Basic OS tools:
 - gprop/perf
 - cachegrind/callgrind
- Hardware counter
 - PAPI and tools set
- Vendor tools
 - Intel VTune Amplifier XE, Inspector, Advisor, ITAC
 - AMD CodeAnalyst
 - Allinea (map and DDD)
- ASCR tools (Open source)
 - HPCToolkit (Rice Univ.)
 - TAU (Oregon Univ.)
 - OpenSpeedshop (Krell)
- HEP
 - FAST (FNAL)
 - IgProf
 - Gooda

gprof: demo

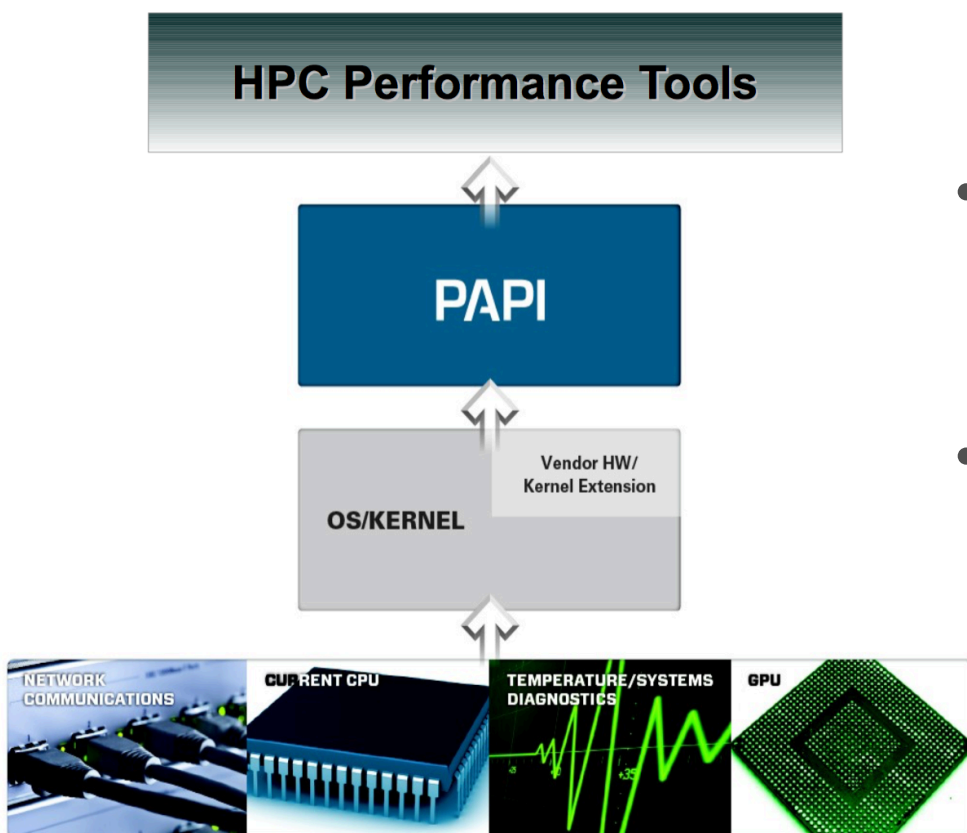
- Compile your program with gcc using `-pg` flag
- Run your program (as usual) – will produce `gmon.out`
- Run `gprof`

```
➤ wget https://g4cpt.fnal.gov/g4p/demos/demo.cc
➤ g++ -pg demo.cc -o demo
➤ time ./demo
➤ gprof ./demo
```

Each sample counts as 0.01 seconds.

%	cumulative	self		self	total	
time	seconds	seconds	calls	ms/call	ms/call	name
36.49	0.51	0.51	1	514.51	514.51	Function_C()
34.34	1.00	0.48	1	484.25	998.76	Function_B()
30.05	1.42	0.42	1	423.72	423.72	Function_A()
0.00	1.42	0.00	1	0.00	0.00	global constructors keyed to _Z10Function_Cv
0.00	1.42	0.00	1	0.00	0.00	__static_initialization_and_destruction_0(int, int)

PAPI (Performance API)



- A standard API to access hardware performance counters
- Relation between software performance and processor events
- Event metrics
 - FLOPS, Load/Store
 - cache hit/miss, TLB miss
 - power consumption (MuMMI)
 - platform specific metrics

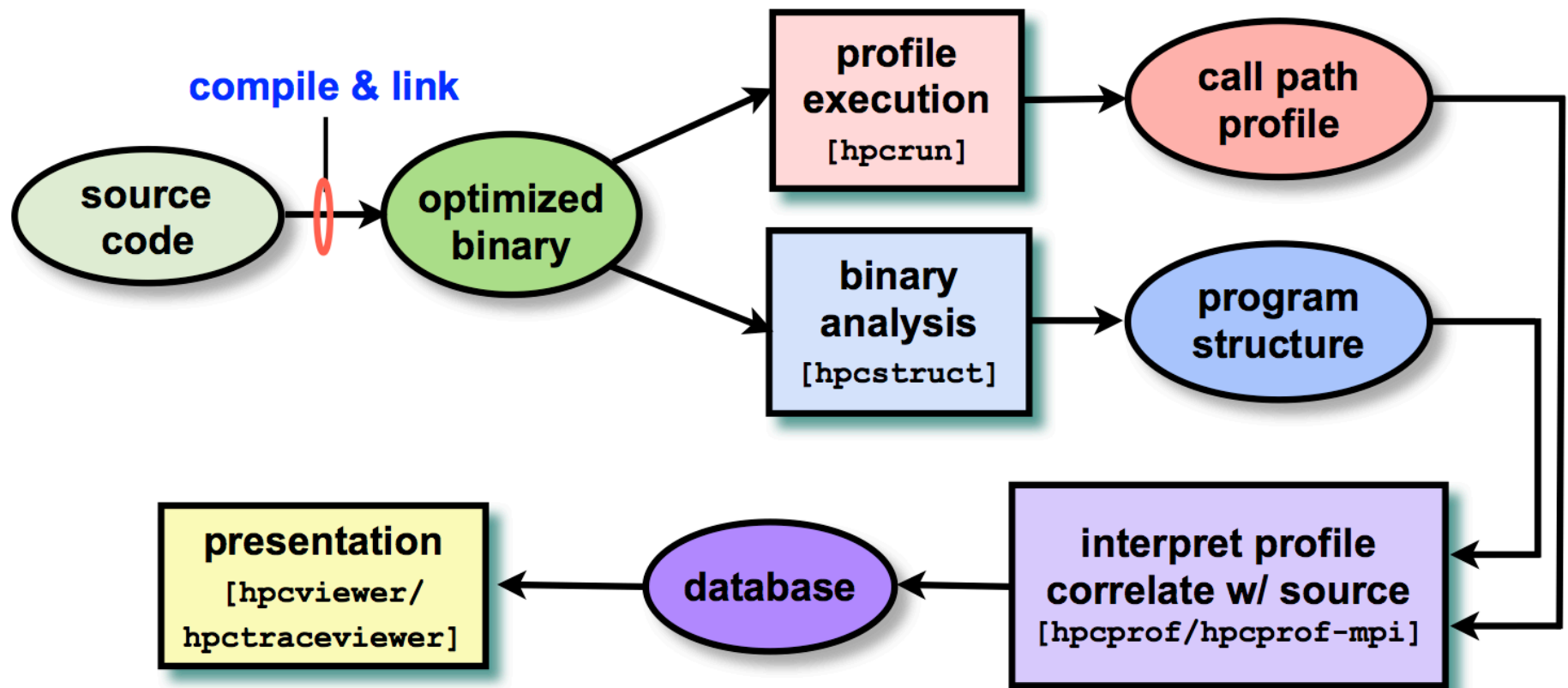
Hardware Counters

- Operating systems support both non-derived and derived PAPI presets: `papi_avail -a` for listing
- A list of some possible hardware counter combinations

For Xeon processors:	
PAPI_FP_INS, PAPI_LD_INS, PAPI_SR_INS	Load store info, memory bandwidth needs
PAPI_L1_DCM, PAPI_L1_TCA	L1 cache hit/miss ratios
PAPI_L2_DCM, PAPI_L2_TCA	L2 cache hit/miss ratios
LAST_LEVEL_CACHE_MISSES, LAST_LEVEL_CACHE_REFERENCES	L3 cache info
MEM_UNCORE_RETIRED:REMOTE_DRAM, MEM_UNCORE_RETIRED:LOCAL_DRAM	Local/nonlocal memory access
For Opteron processors:	
PAPI_FAD_INS, PAPI_FML_INS	Floating point add multiply
PAPI_FDV_INS, PAPI_FSQ_INS	Square root and divisions
PAPI_FP_OPS, PAPI_VEC_INS	Floating point and vector instructions
READ_REQUEST_TO_L3_CACHE:ALL_CORES, L3_CACHE_MISSES:ALL_CORES	L3 cache

Example of Sampling Tools and Workflow: HPCToolkit

- Typically unmodified binary and call stack analysis
- Code centric view, GUI and text-based flat profile



Example of Integrated Tools: TAU (Tuning Analysis Utilities)

- Dynamic, compiler based, source based Instrumentation
- Analysis tools
 - ParaProf
 - PerfExplorer
 - Tracer (Jumpshot, vampir)
- Various built-in graphical tools
- Dis/advantage: compiler/source-based instrumentation

Understanding Sampling Profilers

- Q3: What are the disadvantages of sampling profilers?
 1. Sampling uncertainty (→ statistical analysis with repetition)
 2. Non-Reproducibility (→ use definitive tools)
 3. Interference from other processes (→ standalone nodes)
 4. Jumping between cores (→ setting NUMA affinity, pinning)
 5. All of above
- Strategies
 - Understand your program first (intensity: arithmetic vs. memory)
 - Overview with sampling experiments
 - Focus on critical parts of code: Rule of 80:20
 - Detailed optimization with hardware counter experiments
 - Benchmarking and monitoring of every minor/major update



PART-II

Profiling Results of LArTest with IgProf and OpenSpeedshop

Application for this Tutorial: LArTest

- A standalone Geant4 application (developed by H. Wenzel)
 - Cubic (5mx5mx5m) LAr fiducial volume
 - GDML to assign step limits and sensitive detector to volumes
 - Optical (scintillation) photons produced in sensitive detector
- Computing performance monitoring features
 - Event time
 - Memory (IgProf, statm)
 - Statistics of the number for tracks/steps per particle type
- Profiling examples with IgProf and OpenSpeedshop

Installing and Running IgProf

- Installation: <http://igprof.org/install.html>
- Implementation

```
if (void *sym = dlsym(0, "igprof_dump_now")) {  
    dump_ = __extension__ (void (*)(const char *)) sym;  
} else { /* message */ ; }
```

- Running igprof on your application (-mp: memory profiling)
<http://igprof.org/running.html>

```
igprof -d -mp -z -o ${IG_OUT} $exe {args...}
```

- Analysis (web-navigable version of the report, -r for \$mode)

```
cmd="igprof-analyse --sqlite -d -v -g -r"  
$cmd ${mode} ${IG_OUT} | sqlite3 out.sql3
```

– \${mode} = MEM_LIVE, MEM_MAX, MEM_TOTAL

IgProf

- Snapshot live memory on the heap (for every N-events)

```
if ( dump_ && evt->GetEventID() % 25 == 0) {  
    sprintf(outfile,"| gzip -9c > IgProf.%d.gz",evt->GetEventID()+1);  
    dump_(outfile);  
}
```

```
cmd="igprof-analyse --sqlite -d -v -g -r"  
$cmd ${mode} -b out1.gz --diff-mode out2.gz| sqlite3 diff.sql3
```

- Performance report formats
 - ascii text (flat file)
 - sqlite database files
- Demo for the web-navigable report

https://g4cpt.fnal.gov/g4p/oss_10.3.r04_lArTest_01/index_igprof.html

Hint from IgProf

- Q4: How to improve performance of this function? Avoid ...

```
MaterialPropertyVector*
MaterialPropertiesTable::GetProperty(const char *key)
{
    // Returns a Material Property Vector corresponding to a key
    if (string(key) == "GROUPVEL") return SetGROUPVEL();

    MPTiterator i;
    i = MPT.find(string(key));
    if ( i != MPT.end() ) return i->second;
    return nullptr;
}
```

1. String comparison
2. String conversion
3. String search (find)
4. Race condition (map)

Introduction to OpenSpeedshop (OSS)

- Comprehensive performance analysis of sequential, multithreaded, and MPI applications
- Open source (the Krell institute, <https://openspeedshop.org>) and one of ASCR profiling tools
- The base functionality includes
 - Sampling experiment (light-weighted)
 - Support call stack analysis
 - Hardware performance (PAPI) counters
 - Multi-threaded, MPI profiling and tracing
 - Memory function tracing, I/O profiling and tracing, etc...
- Tested on a variety of Linux clusters and supports parallel hardware architectures (Intel MIC, NVIDIA CUDA) as well as HPC systems (Cray, Blue Gene)

OSS: Installation and Performance Measurement

- Installation: a typical build (with the version 2.2)
 - `./install-tool --build-krell-root`
 - `--krell-root-prefix ${install_dir}/krellroot_v2.2`
 - `--with-openmpi /usr/local/openmpi-1.8.1`
 - `./install-tool --build-offline`
 - `--openss-prefix ${install_dir}/openspeedshop2.2`
 - `--krell-root-prefix ${install_dir}/krellroot_v2.2`
 - `--with-openmpi /usr/local/openmpi-1.8.1`
- Running an experiment: unmodified binary instrumentation
 - `osspcsamp "lArTest lArBox.gdml profile.pi-5GeV" [frequency]`
- Performance analysis: command-line (-cli) or GUI (-f)
 - `openss -cli lArTest-pcsamp.openss`

Demo: OSS Command-line Analysis (-cli)

> openss -cli	Open the CLI.
openss>> expcreate -f "mutatee 2000" pcsamp	Create an experiment using pcsamp with this application.
openss>> expgo	Run the experiment and create the database
openss>> expview	Display the default view of the performance data.

help or help commands	list -v obj
expview	list -v ranks
expview -v statements	list -v hosts
expview -v loops	expview -m <metric>
expview -v linkedobjects	expview -v calltrees,fullstack
expview -v calltrees,fullstack	<experiment type> <number>
expview -m loadbalance	expview -v calltrees,fullstack usertime2
expview -r <rank_num>	expview <experiment-name><number>
expcompare -r 1 -r 2 -m time	expview pcsamp2
list -v metrics	expview -v statements
list -v src	<experiment-name><number>

OSS (GUI): Default View and Statistical Panel

Toolbars

Top Functions

Showing Functions Report:

% of Total Exclusive CPU Time	Exclusive CPU time	Inclusive CPU time	% of Total Exclusive CPU Time	Function (defining location)
11.216171	105.742855	105.742855	11.216171	__tls_get_addr (/lib64/ld-2.12.so)
10.128193	95.485712	95.485712	10.128193	__ieee754_log (/lib64/libm-2.12.so)
6.012668	56.685713	56.685713	6.012668	__ieee754_exp (/lib64/libm-2.12.so)
4.748917	44.771428	62.885713	4.748917	G4HadronCrossSections::CalcScatteringCrossSections
3.279086	30.914285	156.999997	3.279086	G4CrossSectionDataStore::GetCrossSection (/home/s
	20.400000	64.571427	2.163833	G4ElasticHadrNucleusHE::HadrNucDifferCrSec (/hor
	16.028571	47.514285	1.700155	G4Navigator::LocateGlobalPointAndSetup (/home/sy
	15.400000	15.400000	1.633482	cmsExpMagneticField::GetFieldValue (/home/syjun
	12.571428	62.085713	1.333455	G4hPairProductionModel::ComputeDMicroscopicCro
	11.885714	14.371428	1.260721	G4ProductionCutsTable::ScanAndSetCouple (/home/
other				

Sampling Experiments in OSS

- pcsamp (periodic sampling of program counters)
 - low overhead overview of time distribution
- usertime (call path profiling)
 - inclusive and exclusive timing data
 - call paths, caller and callee relationships
- hwcsamp (periodic sampling hardware counters)
 - profiling of hardware counter events (PAPI events)
- pthreads (POSIX thread tracing)
- mem (memory tracing)
 - call paths of memory related function call events
 - aggregate and individual rank, thread, or processing times
- io (I/O tracing)
- Many other useful experiments

OSS: Measurement Overheads and Output Size

- pcsamp: exclusive time - insensitive to sampling frequency (default 100Hz)

Frequency	Time(sec)	OverHead(%)	DB size(MB)
base :	52.20	-	
50 Hz:	52.27	0.13	0.376832
100 Hz:	52.62	0.80	0.486400
200 Hz	52.36	0.31	0.607232
500 Hz	52.98	1.49	0.811008
1000 Hz:	52.65	0.86	0.971776
10000 Hz:	52.76	1.07	1.012736

- usertime: inclusive time and call paths – large overhead (default 35): similar overhead for hwcsamp

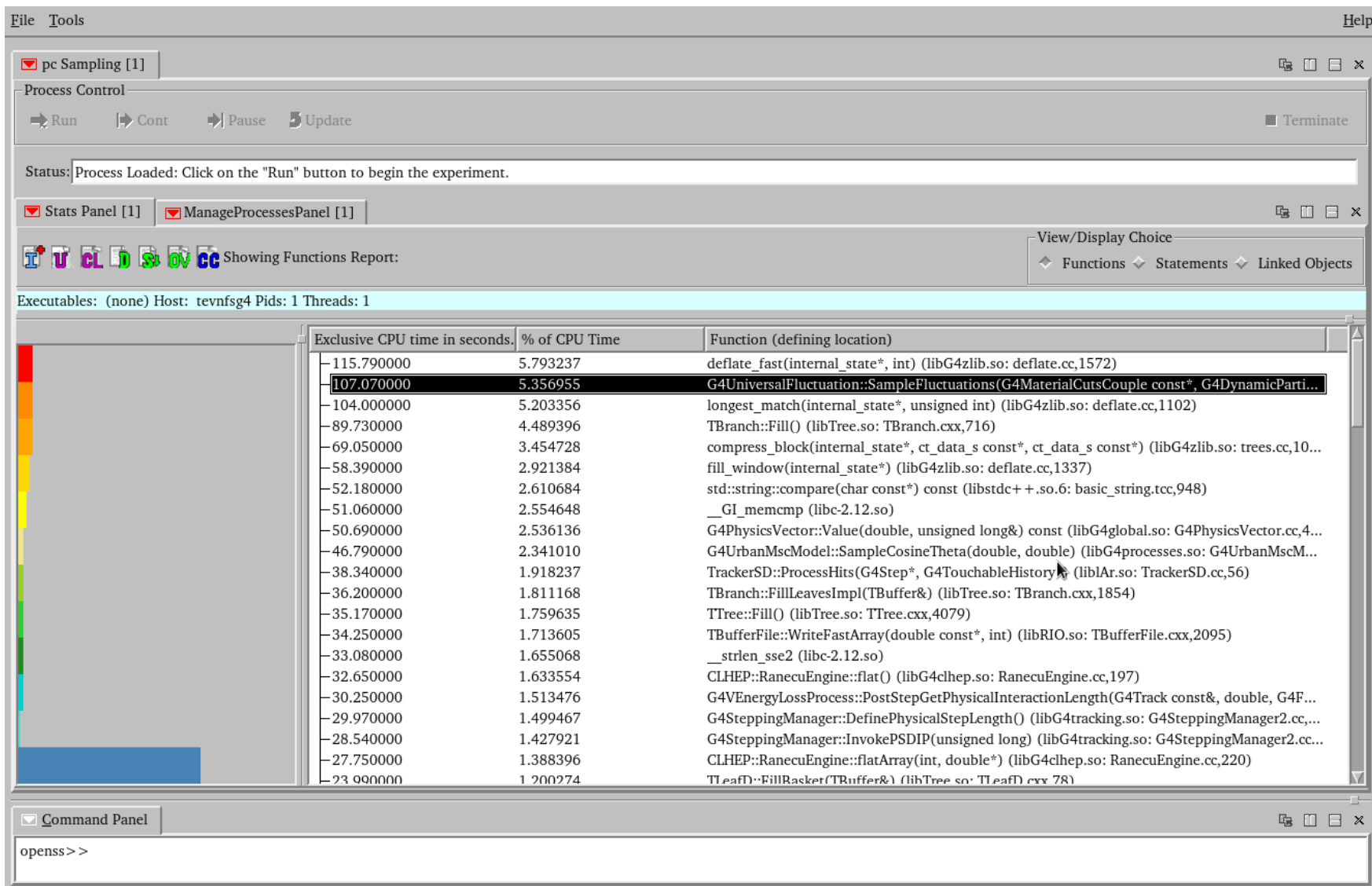
Frequency	Time(sec)	OverHead(%)	DB size(MB)
base :	52.80	-	
35 Hz:	53.89	2.06	1.087488
50 Hz:	54.33	2.90	1.430528
100 Hz:	56.21	6.46	2.355200
200 Hz	60.25	14.11	4.208640
1000 Hz:	92.84	75.83	18.725888

Preliminary Performance Experiments with LArTest

- LArTest configuration
 - Beam: 5 GeV pi-
 - Step limit: 0.01 cm
 - Physics list: FTFP_BERT (uses standard EM)
 - 1000 events
- osspcsamp (100 Hz)
 - I/O (digitization) ON
 - Analysis ON
- ossusertime and osshwcsamp (35 Hz)
 - I/O (digitization) OFF
 - Analysis OFF

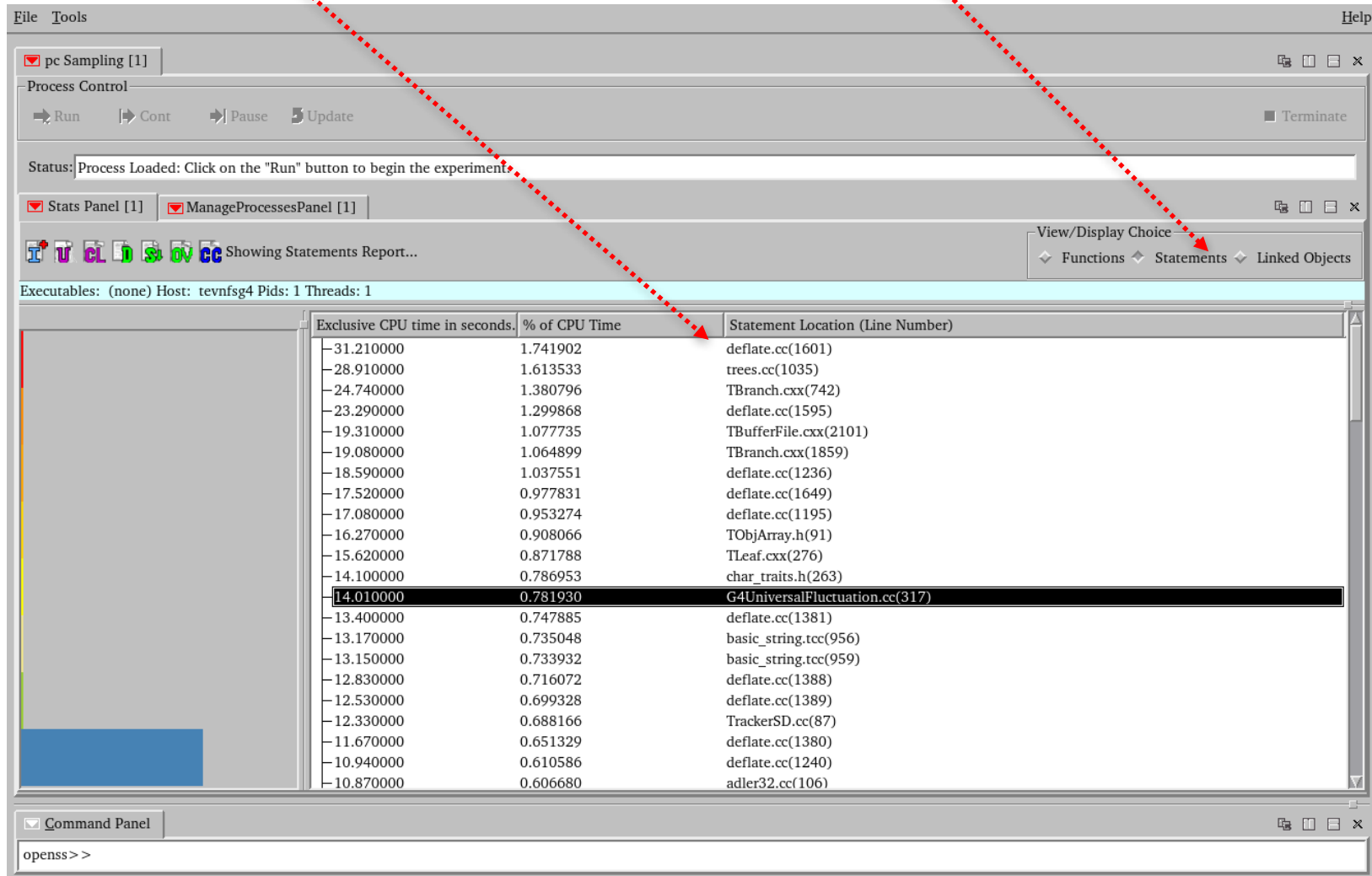
osspcsamp : Functions

- Exclusive CPU time - an overall performance view



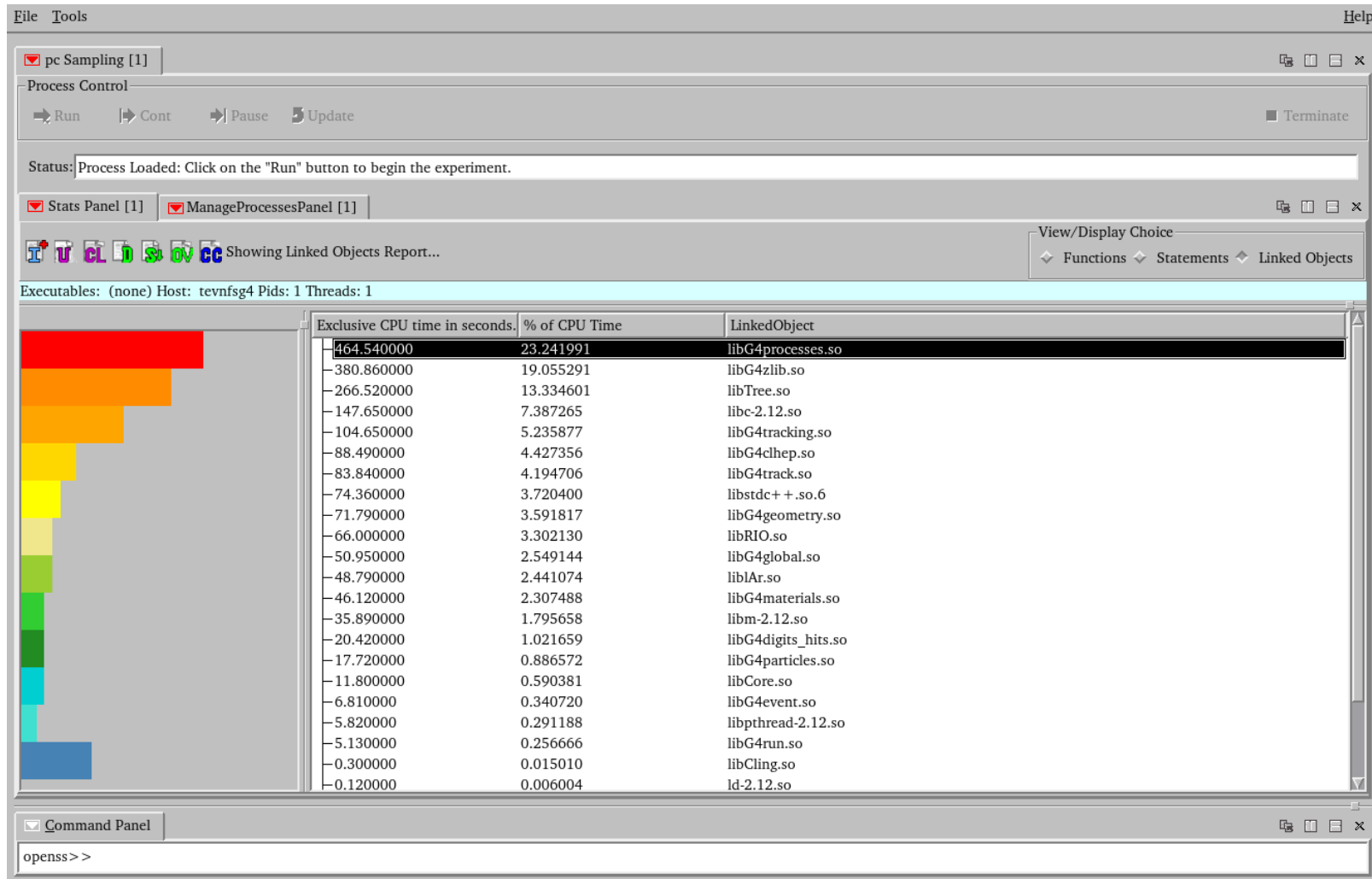
osspsamp: Statements (Line Numbers)

- Select statement level granularity
- List line numbers in program that took most of time



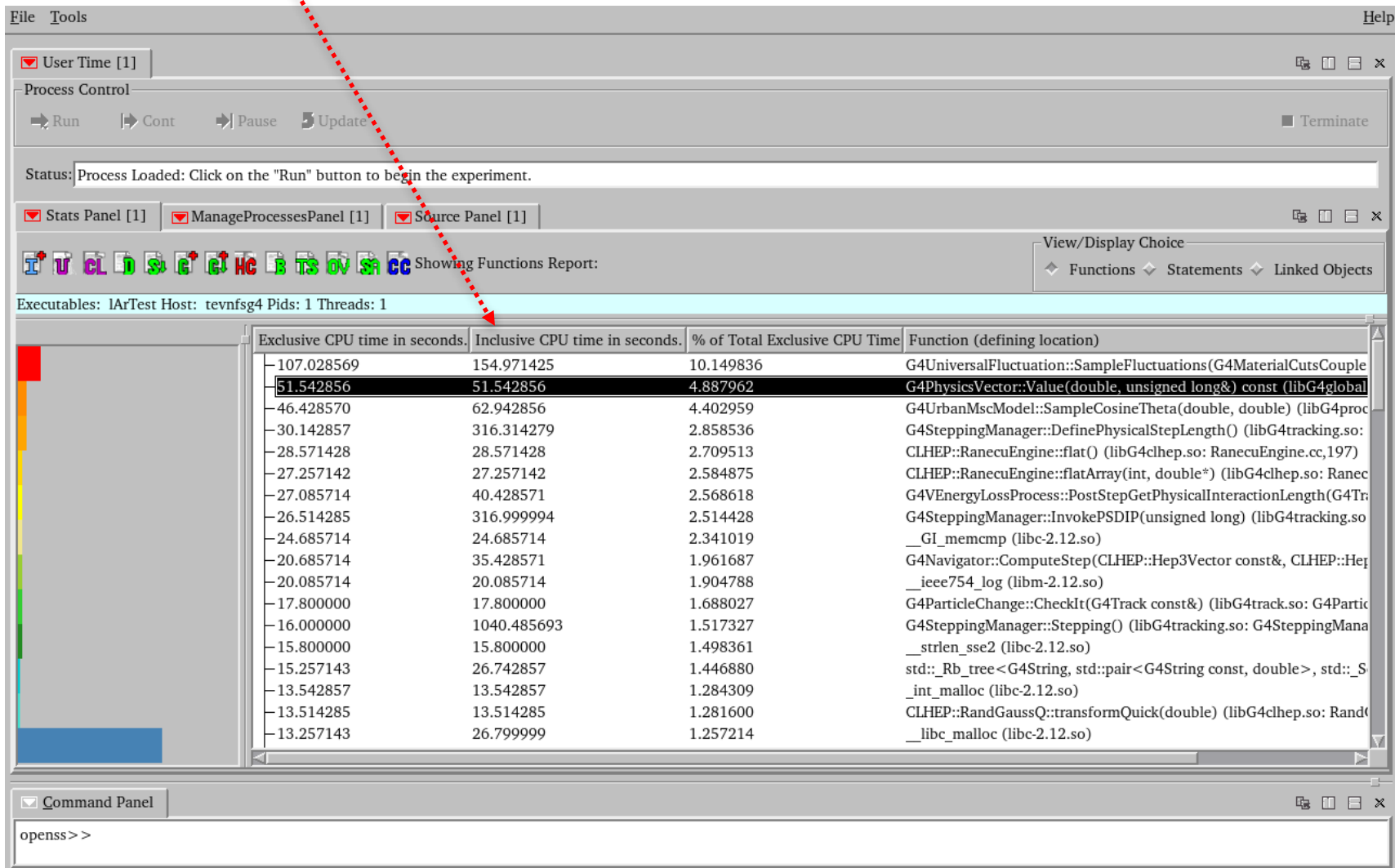
osspsamp: Linked Objects

- The library in which the associated function is located (aggregated by shared objects)



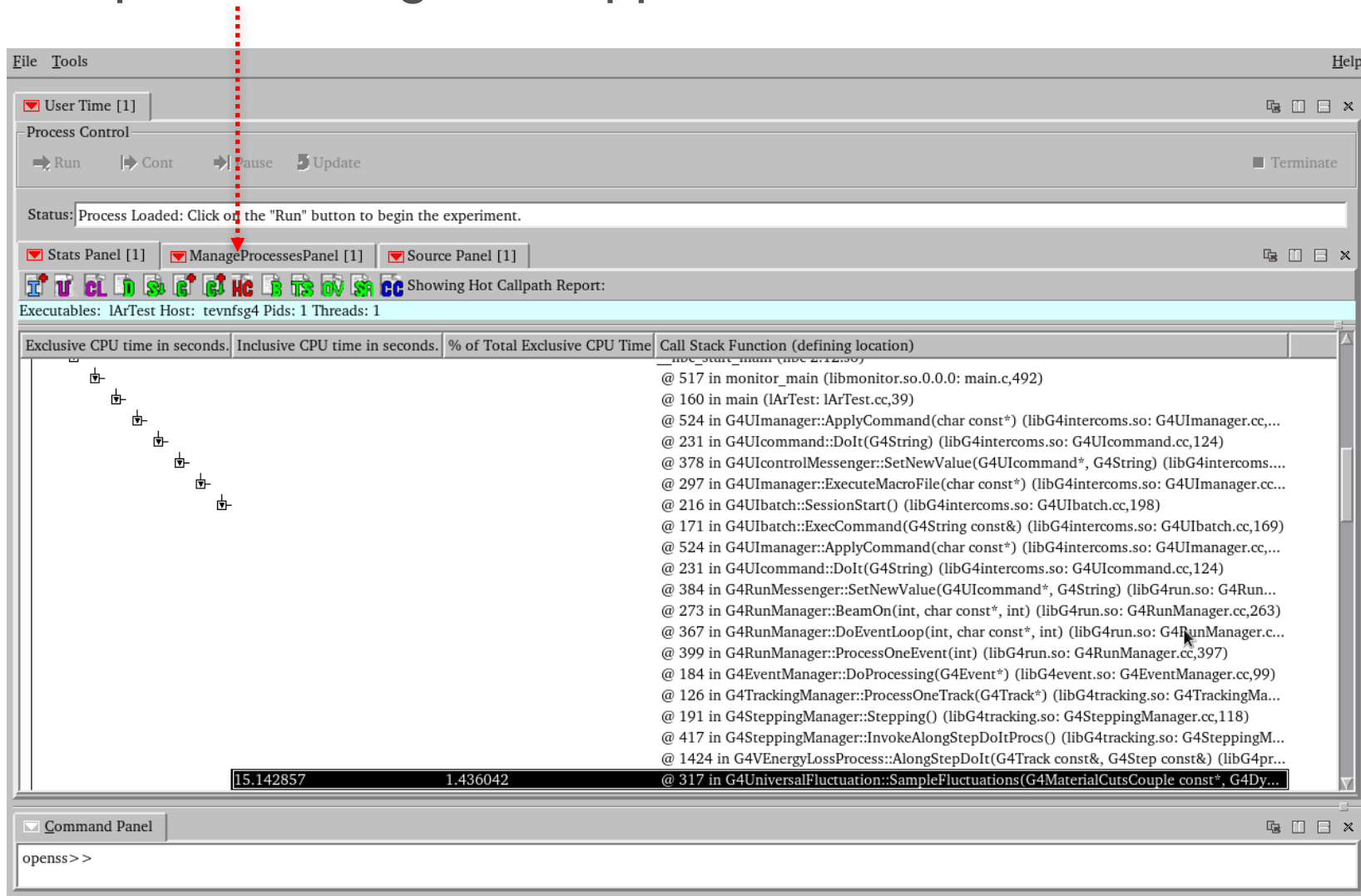
ossusertime: Call Path (Functions)

- Function calls observed anywhere in the stack
- The inclusive time taken by the function and all its callees



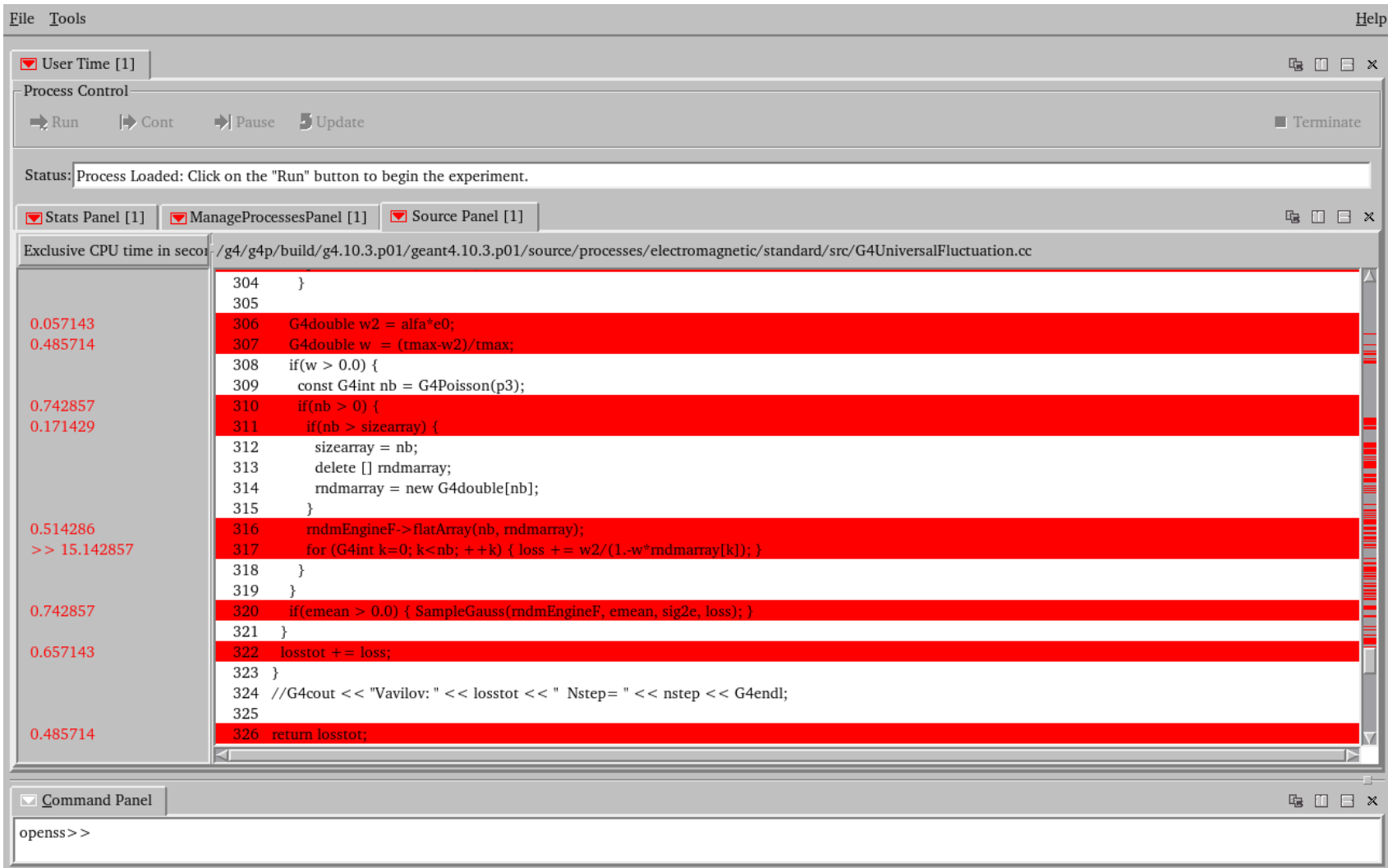
ossusertime: Hot Call Path

- Relationship between caller and callee
- The paths through the application that take the most time



ossusertime: Hot Call (Source)

- Exclusive time on highlighted lines that indicate relatively high CPU times



Experiments with Hardware Counters

- Periodic sampling of hardware counters (hwcsamp)
- Supports both derived and non-derived PAPI presets
- Metrics for instructions, FLOPS, memory and resource patterns

The screenshot shows the HWCSamp Panel [1] interface. At the top, there's a 'Process Control' section with buttons for Run, Cont, Pause, Update, and Terminate. Below this is a 'Status' section with a message: 'Process Loaded: Click on the "Run" button to begin the experiment.' The main section is titled 'Stats Panel [1]' and 'ManageProcessesPanel [1]'. It features a 'View/Display Choice' dropdown menu with options: Functions, Statements, and Linked Objects. Below this is a table of hardware counter metrics. The table has columns: Exclusive CPU time in seconds, % of CPU Time, papi_tot_cyc, papi_tot_ins, papi_fp_ops, and Function (defining location). The table lists various functions and their corresponding hardware counter values. A red box labeled 'PAPI hwc' points to the table. A red circle highlights the 'View/Display Choice' dropdown menu. A red circle also highlights the first row of the table.

Exclusive CPU time in seconds	% of CPU Time	papi_tot_cyc	papi_tot_ins	papi_fp_ops	Function (defining location)
103.844443	9.880537	238072205847	103028210581	16883847741	G4UniversalFluctuation::SampleFluctuations(G4MaterialCutsCou
55.111111	5.243683	126318620429	102749596516	8946035679	G4PhysicsVector::Value(double, unsigned long&) const (libG4glol
46.555555	4.429644	106775342756	86543030788	7919089722	G4UrbanMscModel::SampleCosineTheta(double, double) (libG4p
31.444444	2.991860	72048169097	58654838134	5042279306	G4SteppingManager::DefinePhysicalStepLength() (libG4tracking.s
28.955555	2.755048	66430726872	54158585775	4639678425	G4VEnergyLossProcess::PostStepGetPhysicalInteractionLength(G4
28.488889	2.710646	65348816099	53202548312	4609499319	CLHEP::RanecuEngine::flat() (libG4clhep.so: RanecuEngine.cc,197
27.955555	2.659901	64028514071	52208146862	4480292193	CLHEP::RanecuEngine::flatArray(int, double*) (libG4clhep.so: Ran
26.622222	2.533037	61110837179	49772120291	4327842442	G4SteppingManager::InvokePSDIP(unsigned long) (libG4tracking.s
24.444444	2.325827	56020565509	45584419721	3972397216	__GI_memcmp (libc-2.12.so)
19.733333	1.877577	45257897510	36798636053	3148575513	__ieee754_log (libm-2.12.so)
19.355555	1.841632	44392061033	36175379636	3100280785	G4Navigator::ComputeStep(CLHEP::Hep3Vector const&, CLHEP::I
16.400000	1.560419	37599864927	30652840846	2624429457	G4ParticleChange::CheckIt(G4Track const&) (libG4track.so: G4Pa
15.644444	1.488529	35836916031	29148246541	2526679670	__strlen_sse2 (libc-2.12.so)
15.466667	1.471614	35422081311	28829396851	2495986647	G4SteppingManager::Stepping() (libG4tracking.so: G4SteppingMa

Code Performance by Hardware Counter Metrics

- Derivatives: examples

Hardware Counter Metrics Derivatives	Performance
IPC (Instruction/Cycle)	Large values suggest good balance with minimal stalls.
FPC (FLOPS/Cycle)	Large values for floating point intensive code suggests efficient CPU utilization
FMO (FLOPS/Memory Ops)	Good data locality, Computational Intensity
LPC (Loads/Cycle)	Useful for calculating FMO, may indicate good stride through arrays.
SPC (Stores/Cycle)	Useful for calculating FMO, may indicate good stride through arrays.

- LArTest (Overall): 5 GeV pi- (Intel Xeon X5650@2.67GHz)
 - IPC = 0.79 (relatively small)
 - FMO = 0.32

Other useful OSS Features

- Flexible analysis options ([GUI](#), [command line](#), [online](#))
- Export report data in different formats (text, cvs, chart)
- Multi-threading capability
- Compare two experiments (osscompare): examples
 - two releases
 - two experiments with the different numbers of threads
- Call path analysis based on DB
- Experiments for parallel code (MPI tracing)

Demo: Example of Performance Profiling Report

- Monitor Geant4 part of performance changes for LAr-based detectors by
 - Beam energy/Particle type/Physics list
 - Geant4 (reference) release

https://g4cpt.fnal.gov/g4p/oss_10.3.r04_IArTest_01/index_sprof.html

https://g4cpt.fnal.gov/g4p/oss_10.3.r04_IArTest_01/index_igprof.html

OpenSpeedshop

Geant4.10.3.r04 IArTest

Sample	Physics List	B-Field	Energy
e-	FTFP_BERT	OFF (0 T)	1 GeV 5 GeV
e+	FTFP_BERT	OFF (0 T)	1 GeV 5 GeV
mu-	FTFP_BERT	OFF (0 T)	1 GeV 5 GeV
mu+	FTFP_BERT	OFF (0 T)	1 GeV 5 GeV
pi-	FTFP_BERT	OFF (0 T)	1 GeV 5 GeV
pi+	FTFP_BERT	OFF (0 T)	1 GeV 5 GeV
K-	FTFP_BERT	OFF (0 T)	1 GeV 5 GeV
K+	FTFP_BERT	OFF (0 T)	1 GeV 5 GeV
p	FTFP_BERT	OFF (0 T)	1 GeV 5 GeV

Memory Profiler/IgProf

Memory profiling reports

- MEM_LIVE: memory that has not been freed - snapshot of the heap, i.e. a heap profile.
- MEM_MAX: the largest single allocation by any function
- MEM_TOTAL: the total amount of memory allocated by any function - a snapshot of poor memory locality
- N: memory snapshot at the end of N-th event
- Diff(N-M): memory difference between N-th and M-th event - direct memory leakage
- End of Run: memory snapshot at the End of Run

Geant4.10.3.r04 IArTest B=4.0T

Sample	Physics List	Energy	MEM_LIVE	MEM_MAX	MEM_TOTAL
e-	FTFP_BERT	1 GeV	1 Diff(1001-1) 1001 End of Run	1 Diff(1001-1) 1001 End of Run	1 Diff(1001-1) 1001 End of Run
		5 GeV	1 Diff(1001-1) 1001 End of Run	1 Diff(1001-1) 1001 End of Run	1 Diff(1001-1) 1001 End of Run
e+	FTFP_BERT	1 GeV	1 Diff(1001-1) 1001 End of Run	1 Diff(1001-1) 1001 End of Run	1 Diff(1001-1) 1001 End of Run
		5 GeV	1 Diff(1001-1) 1001 End of Run	1 Diff(1001-1) 1001 End of Run	1 Diff(1001-1) 1001 End of Run
mu-	FTFP_BERT	1 GeV	1 Diff(1001-1) 1001 End of Run	1 Diff(1001-1) 1001 End of Run	1 Diff(1001-1) 1001 End of Run
		5 GeV	1 Diff(1001-1) 1001 End of Run	1 Diff(1001-1) 1001 End of Run	1 Diff(1001-1) 1001 End of Run

Summary

- Performance profiling and analysis is an essential part of the software development cycle
 - Modern hardware architectures are demanding (parallelism)
 - HEP applications are big and complex
 - Profilers will help to identify critical parts of code, monitor changes of performance and provide opportunities of optimization
- Where you can start:
 - Try profiling your programs with basic tools
 - IgProf: <http://igprof.org/index.html>
 - OpenSpeedshop: <https://openspeedshop.org/>
 - HPCToolkits: <http://hpctoolkit.org/index.html>
 - TAU: <http://www.cs.uoregon.edu/research/tau/home.php>
- Above tools are quite suitable for continuous integration tests

Acronym

- CISC: Complicated instruction set computer)
- RISC: Reduced instruction set computer)
- ARM: Advanced RISC Machines
- CPI: Cycles per instruction
- IPC: Instructions per cycle
- MIPS: Million instructions per second
- FMO: Floating point operations per memory operation
- DRAM: Dynamic random-access memory
- ASCR: Advanced Scientific Computing Research
- MuMMI: Multiple Metrics Modeling Infrastructure

Answer sheet

- Q1:
- Q2:
- Q3:
- Q4: