



TOOLS FOR GPU COMPUTING

With focus on NVIDIA GPUs

03.12.2019 | MICHAEL KNOBLOCH

MOTIVATION

Make it work,
make it right,
make it fast.

Kent Beck

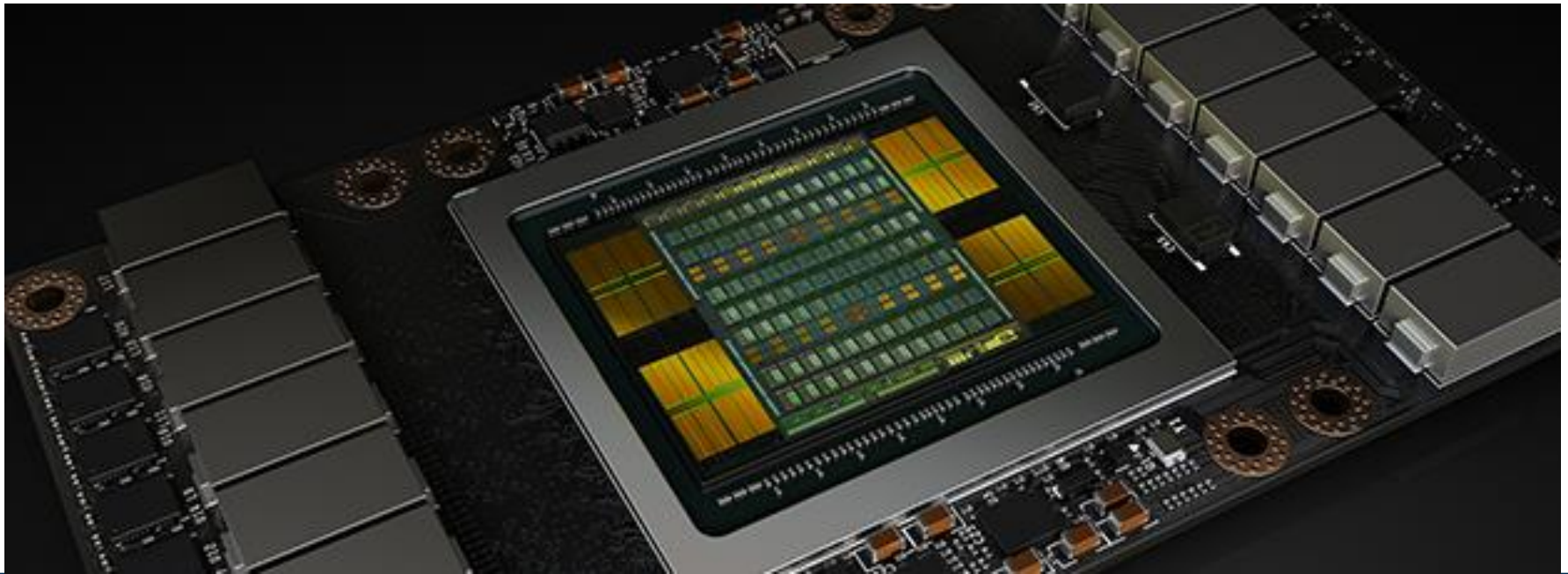
- IDEs
- Porting
- Libraries

Debugger:

- CUDA-MEMCHECK
- CUDA-GDB
- TotalView
- DDT

Performance Tools:

- NVIDIA Visual Profiler
- NVIDIA Nsight System
- NVIDIA Nsight Compute
- Score-P
- Vampir
- Performance Reports
- TAU
- HPCToolkit



GPU PROGRAMMING MODELS

CURRENT STATE OF THE MESS

TRADITIONAL HPC

- **Inter-node:**

- MPI

- **Intra-node:**

- OpenMP
- Pthreads

- C/C++ and Fortran

MODERN HPC

- **Inter-node:**

- MPI
- PGAS (SHMEM, GASPI, ...)

- **Intra-node:**

- OpenMP
 - Pthreads
 - Tasking, C++11 threads, TBB, ...
- C/C++, Fortran and Python

GPU PROGRAMMING

- **Low-level:**

- CUDA (NVIDIA), ROCm (AMD)
- OpenCL

- **Pragma-based:**

- OpenACC
- OpenMP target

- On top: SYCL, oneAPI, HIP, KOKKOS, ...

```
uint32_t maxIndex = firstElementIndex;
uint32_t nextElement;
uint32_t i = firstElementIndex + threadsCount;

for (; i < ARRAY_SIZE; i += threadsCount) {
    nextElement = array[i];
    if (nextElement > max) {
        max = nextElement;
        maxIndex = i;
    }
}
threadMax[threadIdx.x] = max;
threadMaxIdx[threadIdx.x] = maxIndex;

reduce(threadMax, threadMaxIdx);

if (!threadIdx.x) { // After reduce max will be in thread 0
    array[blockIdx.x] = threadMax[0];
    array[blockIdx.x + BLOCKS] = threadMaxIdx[0];
}
}

uint32_t hostFindMax(const uint32_t array[], uint32_t *index, const uint32_t arrayLength) {
    uint32_t i, max = 0;
    for (i = 0; i < arrayLength; i++) {
        if (array[i] > max) {
            *index = i;
            max = array[i];
        }
    }
}
```

inttypes.h
stdio.h
stdlib.h
string.h
ARRAY_SIZE
BLOCKS
THREADS
MEMORY_BANKS
CUDA_CHECK_RETURN()
• findMaxSingleThread(uint32_t*) : void
• reduce(uint32_t*, uint32_t*) : void
• cudaFindGlobalMax(uint32_t*) : void
• cudaFindMax(uint32_t*) : void
• hostFindMax(const uint32_t*, uint32_t*) : void
• deviceFindMax(const uint32_t*, uint32_t*) : void
• initArray(uint32_t* const) : void
• verifyResult(uint32_t, uint32_t) : void
• main(int, char**) : int

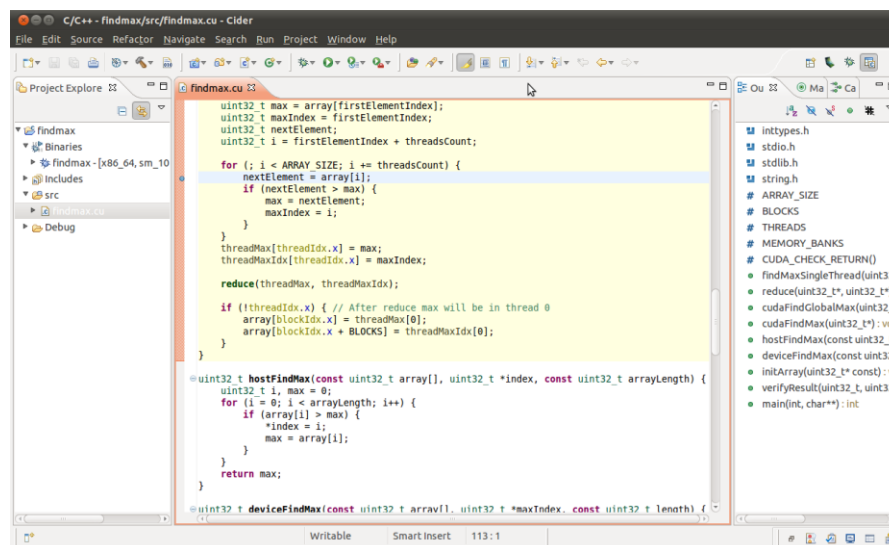
MAKE IT WORK

DEVELOPMENT OF GPU APPS

IDE

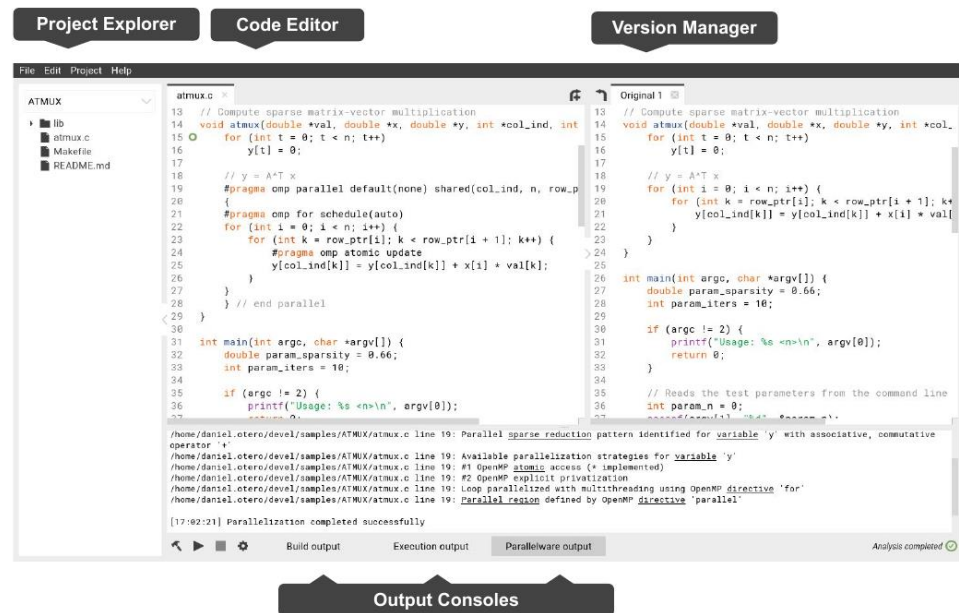
Integrated Development Environment

- Integrates Editor, Build system, Debugger, and Profiler
- NVIDIA Nsight (Linux: Eclipse, Windows: Visual Studio)
- Nsight Code Editor
 - CUDA aware code completion and inline help
 - CUDA code highlighting
 - CUDA aware refactoring



PORTING

- Several tools exist helping expose parallelism
- Example: Appentra Parallelware Trainer
 - Identifies parallelizable sections in sequential applications
 - Supports OpenMP and OpenACC
 - Supports versioning of changes
 - Start program directly from GUI



The screenshot displays a multi-pane debugger interface. The left pane shows a list of threads with columns for Thread State, TID, and Name. The middle pane shows a table of variables with columns for Name, Type, and Value. The right pane shows C++ source code with line numbers. The bottom pane shows a call stack with columns for ID, Type, Stop, Location, and Line.

Thread State	TID	Name
Breakpoint	1.1	tensorflow::SoftmaxXentWithLoss
Stopped	1.2	pthread_cond_wait
Stopped	1.3	pthread_cond_wait
Stopped	1.4	pthread_cond_wait

Name	Type	Value
_	int	0x0000000000000000...
nstar	int	0x000000006 (6)
grap...	int	0x000000015 (21)

```

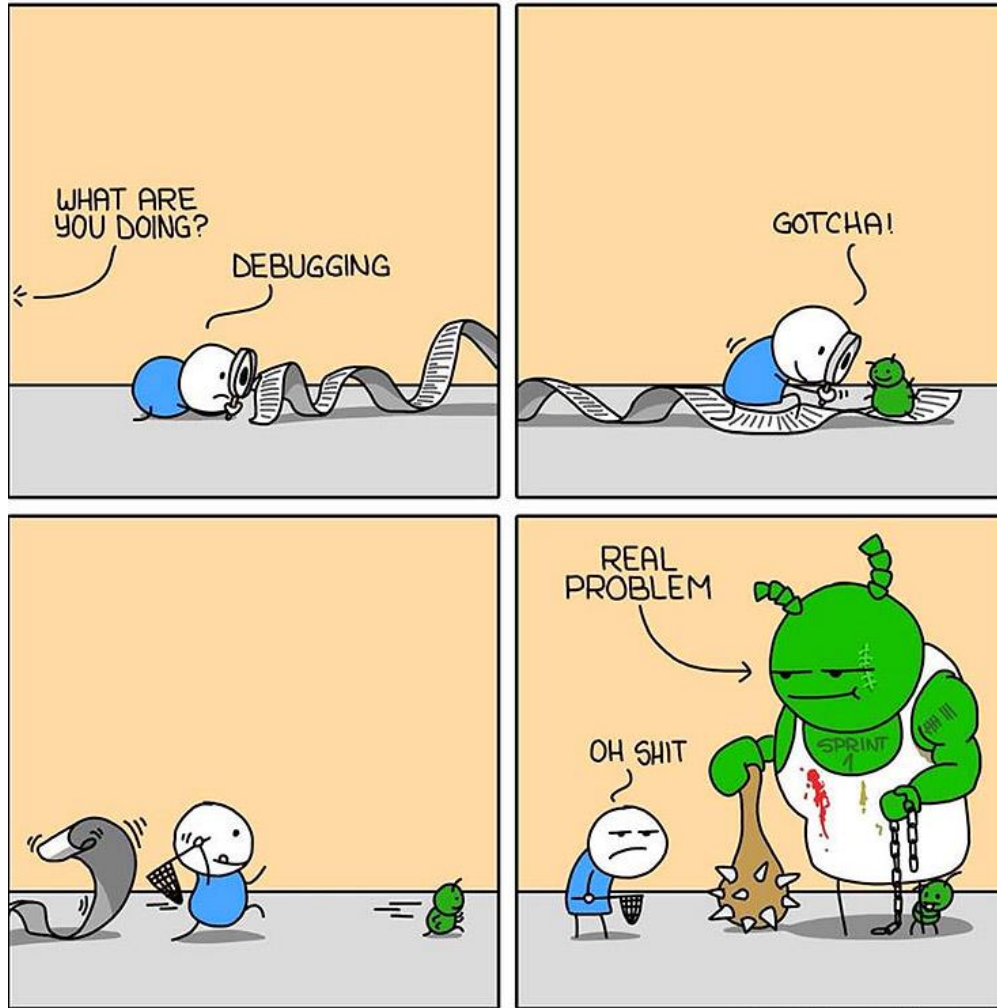
601 // Target nodes
602 const char** c_target_oper_names, int ntargets,
603 TF_Buffer* run_metadata, TF_Status* status) {
604 TF_Run_Setup(noutputs, c_outputs, status);
605 std::vector<std::pair<tensorflow::string, Tensor>> input_pairs(ntargets);
606 if (!TF_Run_Inputs(c_inputs, &input_pairs, status)) return;
607 for (int i = 0; i < ninputs; ++i) {
608   input_pairs[i].first = c_input_names[i];
609 }
610 std::vector<tensorflow::string> output_names(noutputs);
611 for (int i = 0; i < noutputs; ++i) {
612   output_names[i] = c_output_names[i];
613 }
614 std::vector<tensorflow::string> target_oper_names(ntargets);
615 for (int i = 0; i < ntargets; ++i) {
616   target_oper_names[i] = c_target_oper_names[i];
617 }
618 TF_Run_Helper(s->session, nullptr, run_options, input_pairs, output_names,
619               c_outputs, target_oper_names, run_metadata, status);
620 }
621
622 void TF_PRUNSetup(TF_DeprecatedSession* s,
623                  // Input names
624                  const char** c_input_names, int ninputs,
625                  // Output names
626                  const char** c_output_names, int noutputs,
627                  // Target nodes
628                  const char** c_target_oper_names, int ntargets,
629                  const char** handle, TF_Status* status) {
630   status->status = Status::OK();
631 }
632 std::vector<tensorflow::string> input_names(ninputs);
633 std::vector<tensorflow::string> output_names(noutputs);
634 std::vector<tensorflow::string> target_oper_names(ntargets);
  
```

ID	Type	Stop	Location	Line
			tensorflow::FunctionLibraryRuntime::GetOrLoadFunctionLibrary	
			tensorflow::DirectSession::GetOrLoadFunctionLibrary	
			std::FunctionHandler<tensorflow::Status>::GetOrLoadFunctionLibrary	
			std::function<tensorflow::Status>::GetOrLoadFunctionLibrary	
			tensorflow::Sunnamed_namespaces::GetOrLoadFunctionLibrary	
			tensorflow::NewLocalExecutor	
			tensorflow::DirectSession::GetOrLoadFunctionLibrary	
			tensorflow::DirectSession::Run	
			TF_Run_Helper	
			TF_Run	
			tensorflow::TF_Run_wrapper_helper	
			tensorflow::TF_Run_wrapper	
			_run_fn	
			ext_do_call	
			_do_call	
			_do_run	

MAKE IT RIGHT DEBUGGER AND MEMORY ANALYZER

DEBUGGING – A PAINFUL PROCESS

ROOT CAUSE



STARECAT.COM

MONKEYUSER.COM

TOOL COMPATIBILITY MATRIX

Tool	CUDA	OpenACC	OMPD	OpenCL
CUDA-MEMCHECK	✓	(x)	x	x
CUDA-GDB	✓	(x)	x	x
TotalView	✓	✓	(x)	x
DDT	✓	✓	x	x

CUDA-MEMCHECK



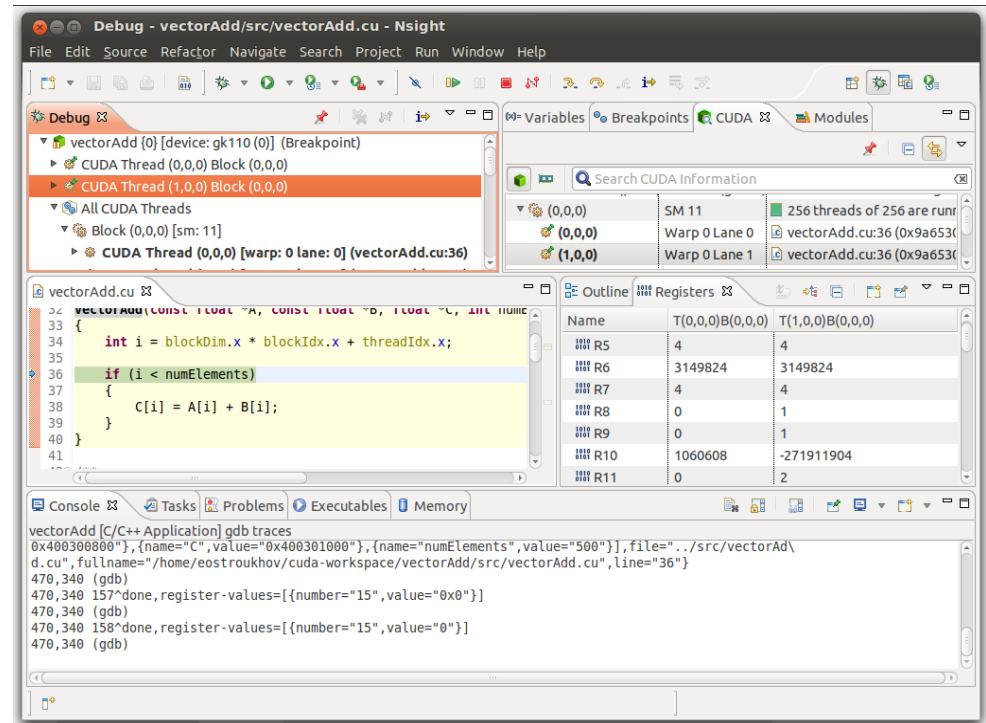
- Valgrind for GPUs
- Monitors hundreds of thousands of threads running concurrently on each GPU
- Reports detailed information about global, local, and shared memory access errors (e.g. out-of-bounds, misaligned memory accesses)
- Reports runtime executions errors (e.g. stack overflows, illegal instructions)
- Reports detailed information about potential race conditions
- Displays stack back-traces on host and device for errors
- And much more
- Included in the CUDA Toolkit

```
Applications Places System
File Edit View Terminal Help
linux64:~/demo2010$ ./ptrchecktest
unspecified launch failure : 79
linux64:~/demo2010$ cuda-memcheck ./ptrchecktest
===== CUDA-MEMCHECK
unspecified launch failure : 79
===== Invalid __global__ read of size 4
===== at 0x00000158 in ptrchecktest.cu:27:kernel2
===== by thread (0,0,0) in block (0,0)
===== Address 0xfd0000001 is misaligned
=====
===== ERROR SUMMARY: 1 error
linux64:~/demo2010$ cuda-memcheck --continue ./ptrchecktest
===== CUDA-MEMCHECK
Checking...
Done
Checking...
Error: 3 (0)
Done
Checking...
Error: 1 (0)
Error: 3 (0)
Error: 5 (0)
Error: 7 (0)
Done
===== Invalid __global__ read of size 4
===== at 0x00000158 in ptrchecktest.cu:27:kernel2
===== by thread (0,0,0) in block (0,0)
===== Address 0xfd0000001 is misaligned
=====
===== Invalid __global__ read of size 4
===== at 0x00000198 in ptrchecktest.cu:18:kernel1
===== by thread (3,0,0) in block (5,0)
===== Address 0xfd00000028 is out of bounds
=====
===== Invalid __global__ write of size 8
===== at 0x000001d0 in ptrchecktest.cu:38:kernel3
===== by thread (1,0,0) in block (8,0)
===== Address 0xfd00000204 is misaligned
=====
===== Invalid __global__ write of size 4
===== at 0x000000f0 in ptrchecktest.cu:44:kernel4
===== by thread (63,0,0) in block (22,0)
===== Address 0x00000000 is out of bounds
=====
===== ERROR SUMMARY: 4 errors
```

CUDA-GDB



- Extension to gdb
- CLI and GUI (Nsight)
- Simultaneously debug on the CPU and multiple GPUs
- Use conditional breakpoints or break automatically on every kernel launch
- Can examine variables, read/write memory and registers and inspect the GPU state when the application is suspended
- Identify memory access violations
 - Run CUDA-MEMCHECK in integrated mode to detect precise exceptions.



- UNIX Symbolic Debugger
for C/C++, Fortran, Python, PGI HPF, assembler programs
- JSC's "standard" debugger
- Special, non-traditional features
 - Multi-process and multi-threaded
 - Multi-dimensional array data visualization
 - Support for **parallel debugging** (MPI: automatic attach, message queues, OpenMP, Pthreads)
 - Scripting and **batch debugging**
 - Advanced memory debugging
 - **CUDA** and **OpenACC** support
- <http://www.roguewave.com>
- **NOTE:** JSC license limited to 2048 processes (shared between all users)

TOTALVIEW: MAIN WINDOW

The screenshot displays the TOTALVIEW IDE interface with several callouts highlighting key features:

- Toolbar for common options:** Located at the top of the window, containing icons for file operations, execution, and debugging.
- Thread control:** A panel on the left side of the window, used to manage and control the execution of threads.
- Break points:** A panel at the bottom left, used to set and manage breakpoints in the code.
- Stack trace:** A panel on the right side, showing the current call stack and the sequence of function calls.
- Source code window:** The central area displaying the source code of the program being debugged.
- Local variables for selected stack frame:** A panel on the right side, showing the local variables and their values for the currently selected stack frame.

The main window shows the following components:

- Processes & Threads:** A table listing the processes and threads running in the application.
- Source code window:** Displays the source code of the program, with the current line of execution highlighted.
- Call Stack:** A panel showing the sequence of function calls.
- Action Points:** A panel showing the action points (breakpoints) set in the code.
- Data View:** A panel showing the data view of the program, including variables and their values.

- UNIX Graphical Debugger for C/C++, Fortran, and Python programs
- Modern, easy-to-use debugger
- Special, non-traditional features
 - Multi-process and multi-threaded
 - Multi-dimensional array data visualization
 - Support for **MPI parallel debugging** (automatic attach, message queues)
 - Support for **OpenMP** (Version 2.x and later)
 - Support for **CUDA** and **OpenACC**
 - Job submission from within debugger
- <https://developer.arm.com>
- **NOTE:** JSC license limited to 64 processes (shared between all users)

DDT: MAIN WINDOW

Process controls

CUDA Thread stepping

Variables

CUDA Thread control

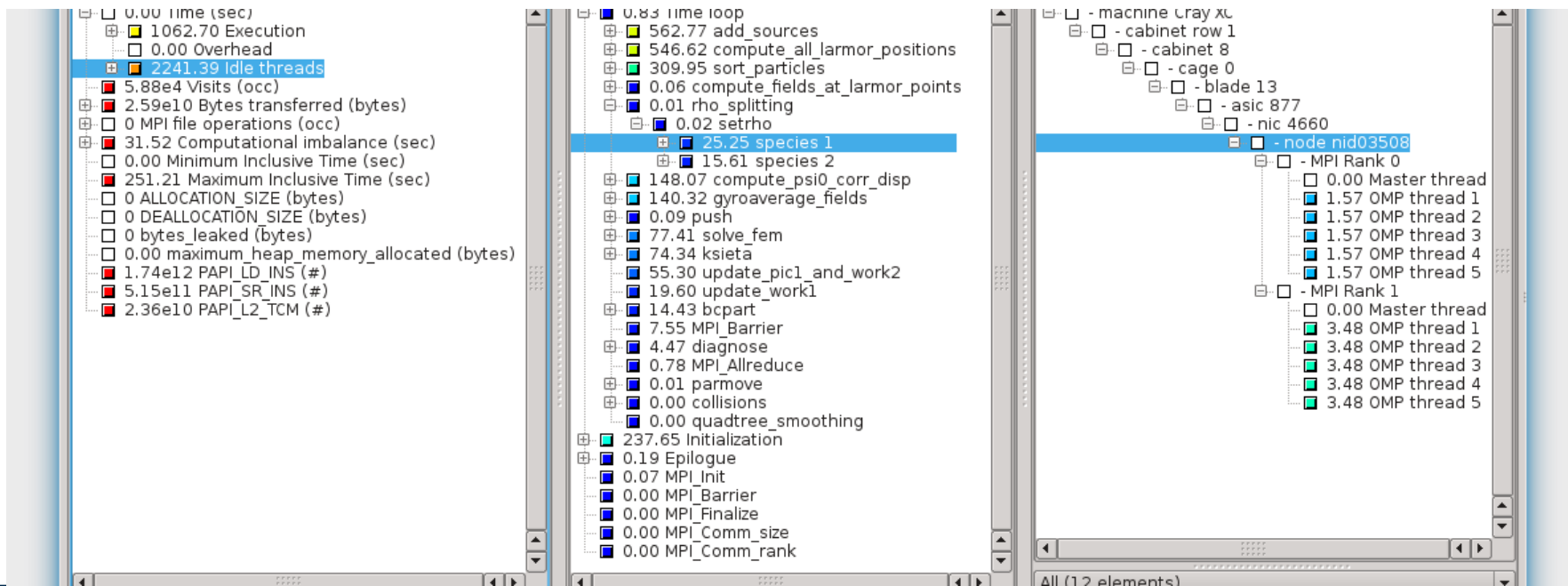
Source code

GPU Device information

Expression evaluator

Stack trace

The screenshot displays the Arm DDT - Arm Forge 19.1.1 interface. At the top, there are process and thread control buttons. Below this is a 'Threads' panel showing a list of threads, with thread 992 selected. The main area shows the source code for 'edge.cu', which contains a 2D convolution filter implementation. To the right of the source code is a 'Locals' panel showing variables like 'cx', 'cy', 'output', 'x', 'y', and 'index'. Below the source code is a 'Stacks' panel showing a stack trace with frames for 'conv2d_global' and 'main'. At the bottom right is an 'Evaluate' panel for the expression evaluator. The bottom status bar indicates 'Ready Connected to: root@jet...'



MAKE IT FAST

PERFORMANCE ANALYSIS TOOLS

TOOL COMPATIBILITY MATRIX

Tool	CUDA	OpenACC	OMPT	OpenCL
Score-P	✓	✓	(✗)	✓
NVIDIA Tools	✓	✓	✗	✗
Perf. Reports	✓	(✗)	✗	✗
TAU	✓	✓	(✗)	✓
HPCToolkit	✓	✗	(✗)	✗

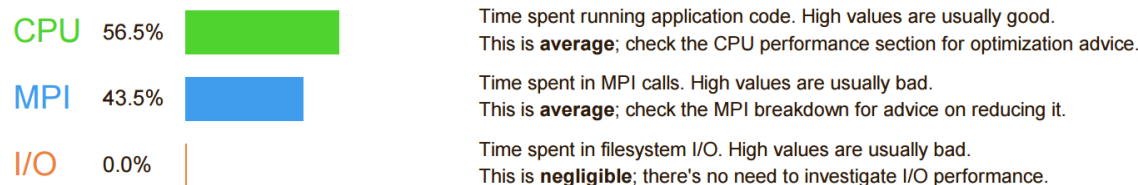
ARM PERFORMANCE REPORTS

- **Single page** report provides quick overview of performance issues
- Works on unmodified, optimized executables
- Shows CPU, GPU, memory, network and I/O utilization
- Supports MPI, multi-threading and accelerators
- Saves data in HTML, CVS or text form
- <https://www.arm.com/products/development-tools/server-and-hpc/performance-reports>
- **Note:** JSC license limited to 512 processes (with unlimited number of threads)

EXAMPLE PERFORMANCE REPORTS

Summary: cp2k.popt is **CPU-bound** in this configuration

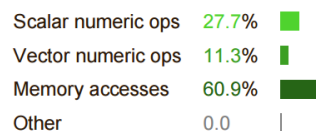
The total wallclock time was spent as follows:



This application run was **CPU-bound**. A breakdown of this time and advice for investigating further is in the **CPU** section below.

CPU

A breakdown of how the **56.5%** total CPU time was spent:

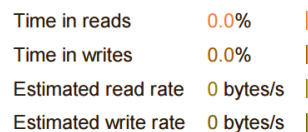


The per-core performance is **memory-bound**. Use a profiler to identify time-consuming loops and check their cache performance.

Little time is spent in **vectorized instructions**. Check the compiler's vectorization advice to see why key loops could not be vectorized.

I/O

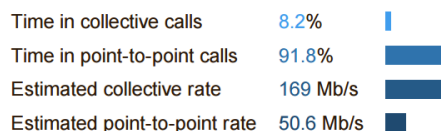
A breakdown of how the **0.0%** total I/O time was spent:



No time is spent in **I/O operations**. There's nothing to optimize here!

MPI

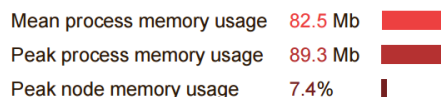
Of the **43.5%** total time spent in MPI calls:



The **point-to-point** transfer rate is low. This can be caused by inefficient message sizes, such as many small messages, or by imbalanced workloads causing processes to wait. Use an MPI profiler to identify the problematic calls and ranks.

Memory

Per-process memory usage may also affect scaling:




The **peak node memory usage** is low. You may be able to reduce the total number of CPU hours used by running with fewer MPI processes and more data on each process.

PERFORMANCE REPORTS

ACCERLERATOR

Accelerators

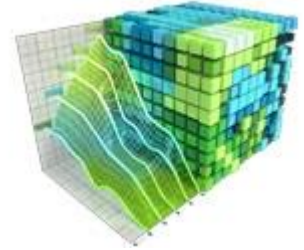
A breakdown of how accelerators were used:

GPU utilization	47.8%	
Global memory accesses	1.6%	
Mean GPU memory usage	0.8%	
Peak GPU memory usage	0.8%	

GPU utilization is low; identify CPU bottlenecks with a profiler and offload them to the accelerator.

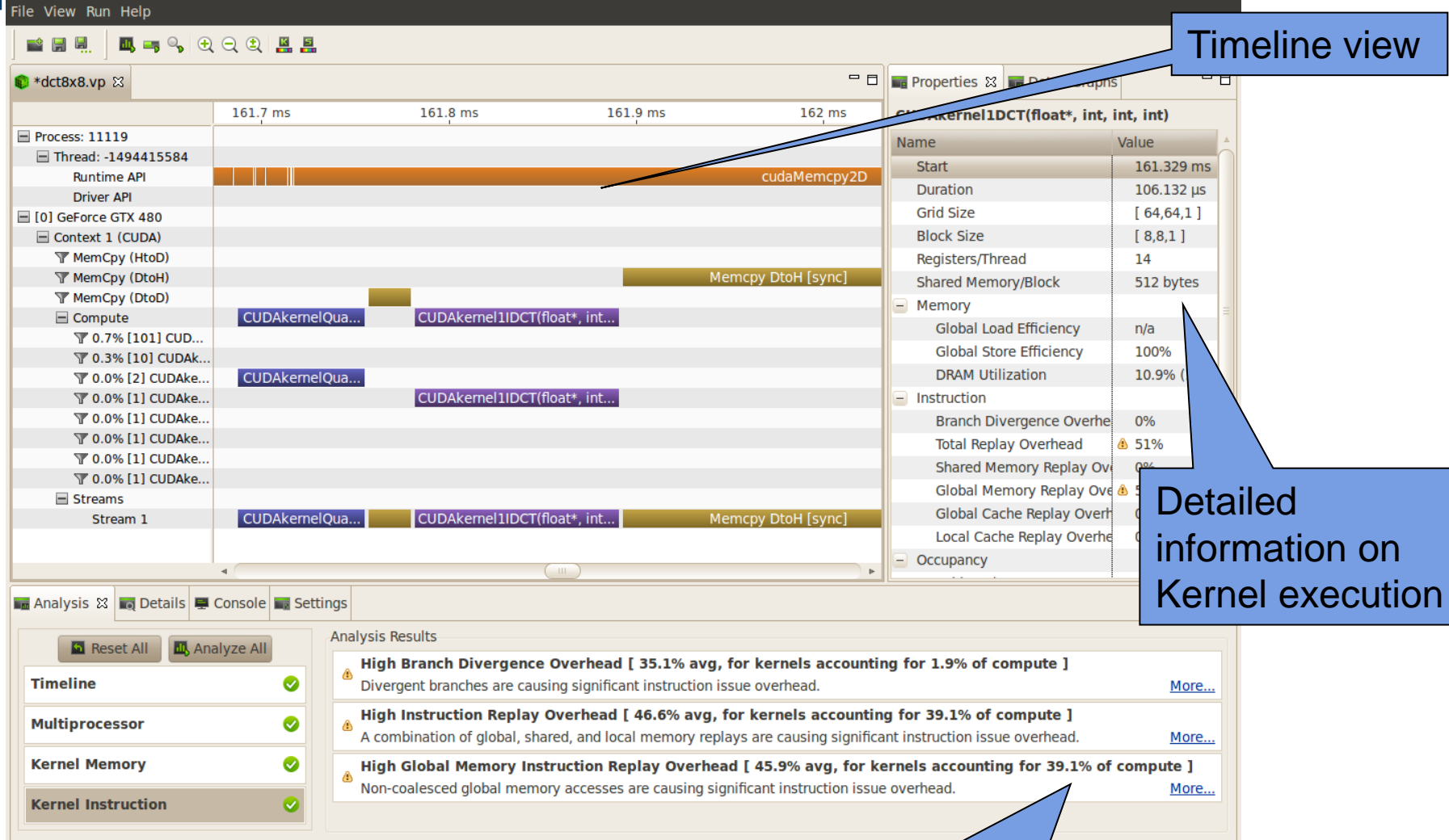
The **peak GPU memory usage** is low. It may be more efficient to offload a larger portion of the dataset to each device.

NVIDIA VISUAL PROFILER



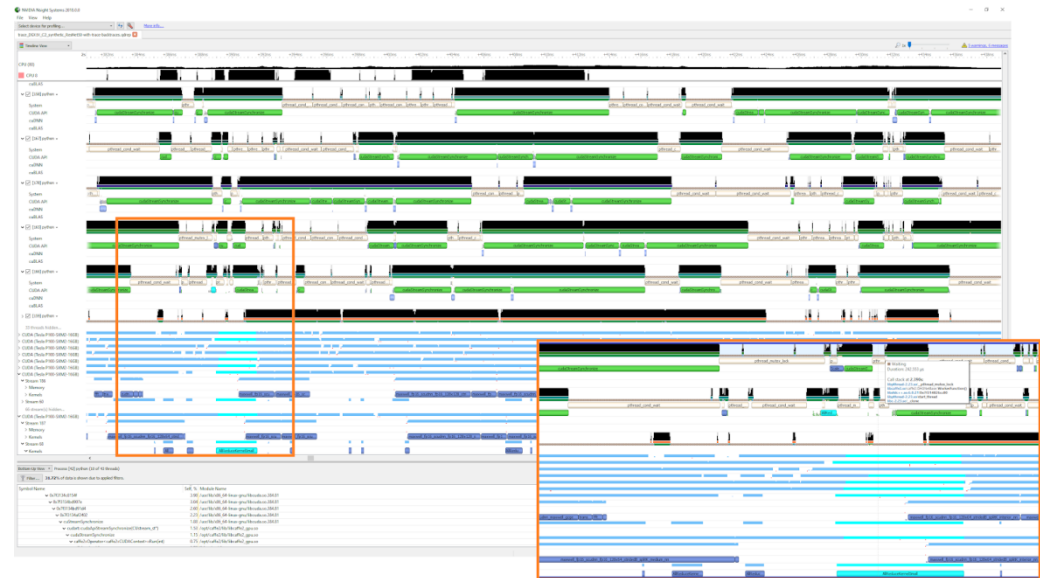
- Part of the CUDA Toolkit
- Supports all CUDA enabled GPUs
- Supports CUDA and OpenACC on Windows, OS X and Linux
- Unified CPU and GPU Timeline
- CUDA API trace
 - Memory transfers, kernel launches, and other API functions
- Automated performance analysis
 - Identify performance bottlenecks and get optimization suggestions
- Guided Application Analysis
- Power, thermal, and clock profiling

NVIDIA VISUAL PROFILER - EXAMPLE



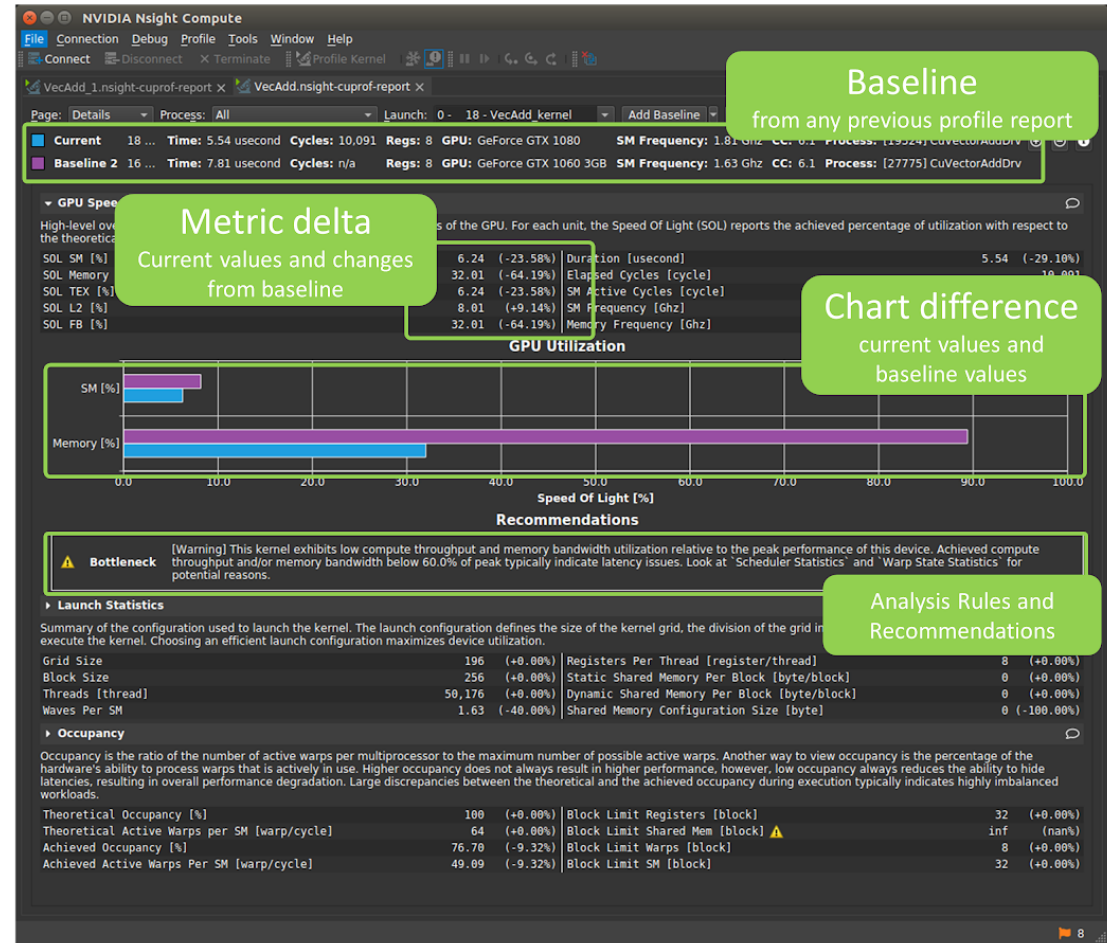
NVIDIA NSIGHT SYSTEMS

- System wide performance analysis tool
- High-level, low overhead
- Similar functionality as NVVP
 - No automated/guided analysis
 - Can launch Nsight Compute for in-depth kernel analysis
- CLI and GUI



NVIDIA NSIGHT COMPUTE

- Interactive kernel profiler
- Detailed performance metrics
- Guided analysis
- Baseline feature to compare versions
- Customizable and data-driven UI
- Supports analysis scripts for post-processing results
- CLI and GUI



SCORE-P

- Community instrumentation and measurement infrastructure
- Developed by a consortium of performance tool groups

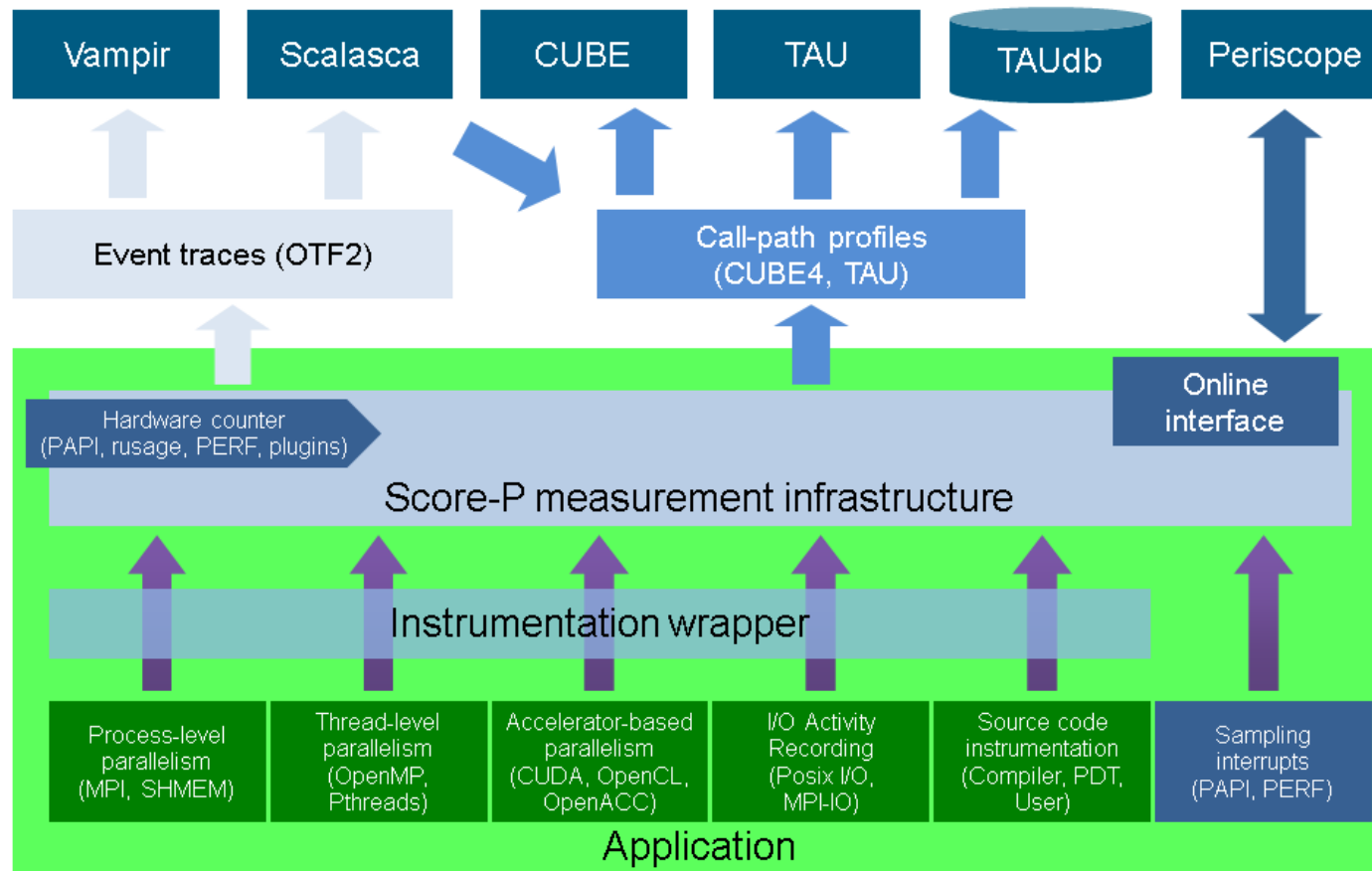


UNIVERSITY OF OREGON



- Next generation measurement system of
 - Scalasca 2.x
 - Vampir
 - TAU
 - Periscope
- Common data formats improve tool interoperability
- <http://www.score-p.org>

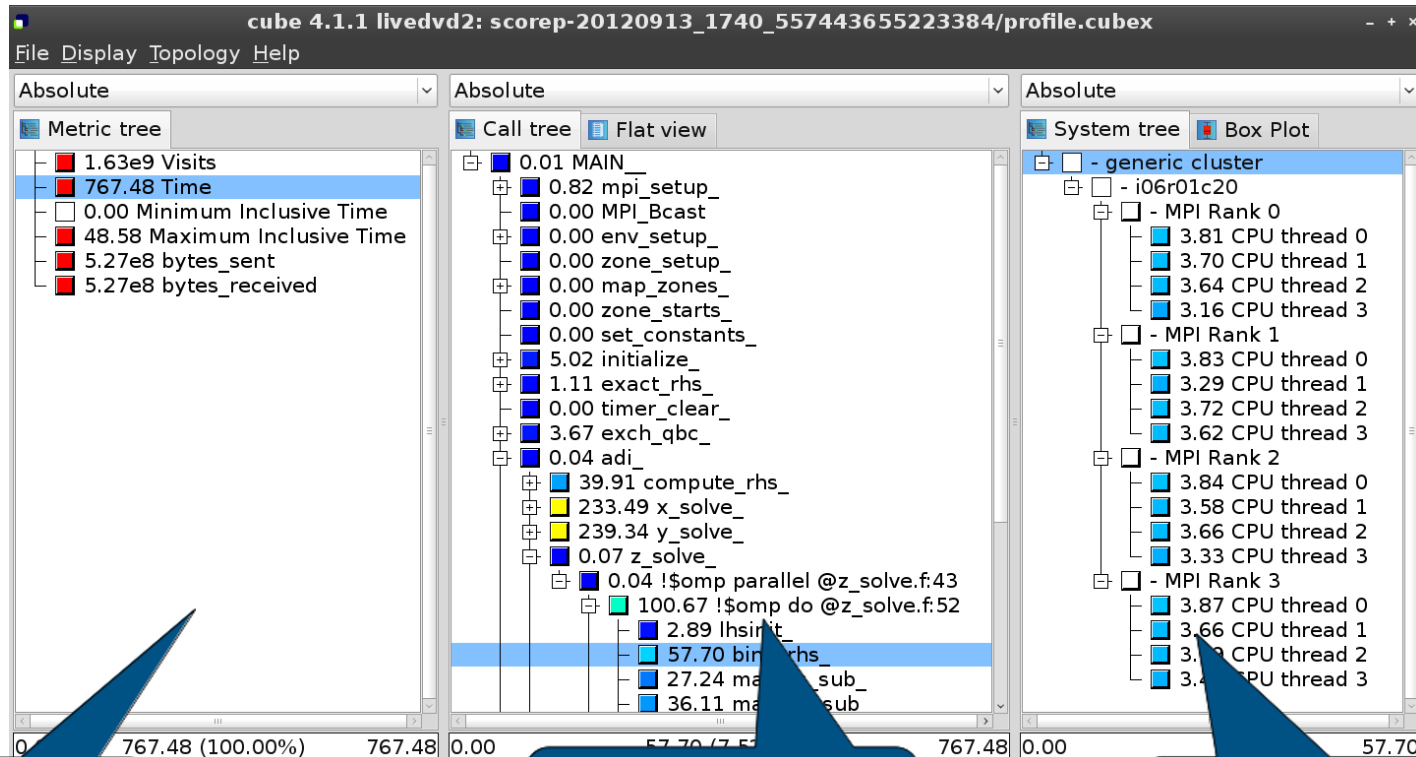
SCORE-P OVERVIEW



SCORE-P GPU MEASUREMENTS

- OpenACC
 - Prefix compiler and linker command with `scorep --openacc`
 - `export ACC_PROFLIB=$SCOREP_ROOT/lib/libscorep_adapter_openacc_event.so`
 - `export SCOREP_OPENACC_ENABLE=yes`
 - yes refers to: regions, wait, enqueue
 - Full list of options in User Guide
- CUDA
 - Prefix compiler and linker command with `scorep --cuda`
 - `export SCOREP_CUDA_ENABLE=yes`
 - yes refers to: runtime, kernel, memcpy
 - Full list of options in User Guide
- OpenCL similar (use `SCOREP_OPENCL_ENABLE=yes`)

CUBE OVERVIEW



What kind of performance metric?

Where is it in the source code?
In what context?

How is it distributed across the processes/threads?

EXAMPLE: OPENACC

File Display Plugins Help

Synchronize state of ... Restore Setting Save Settings Delete Settings

Absolute

Metric tree

- 1.00e4 Visits (occ)
- 6.61 Time (sec)
- 0.00 Minimum Inclusive Time (sec)
- 6.61 Maximum Inclusive Time (sec)
- 0 bytes_put (bytes)
- 0 bytes_get (bytes)
- 0 ALLOCATION_SIZE (bytes)
- 0 DEALLOCATION_SIZE (bytes)
- 0 bytes_leaked (bytes)
- 0.00 maximum_heap_memory_allocate

Absolute

Call tree Flat view

- 0.04 laplace2d_acc
 - 0.00 ..acc_cuda_funcreg_constructor_1
 - 0.05 main
 - 0.05 acc_init
 - 0.00 StartTimer
 - 0.05 acc_data_enter@laplace2d.c:85
 - 4.16 acc_compute@laplace2d.c:91
 - 2.21 acc_compute@laplace2d.c:103
 - 0.05 acc_data_exit@laplace2d.c:85
 - 0.00 GetTimer

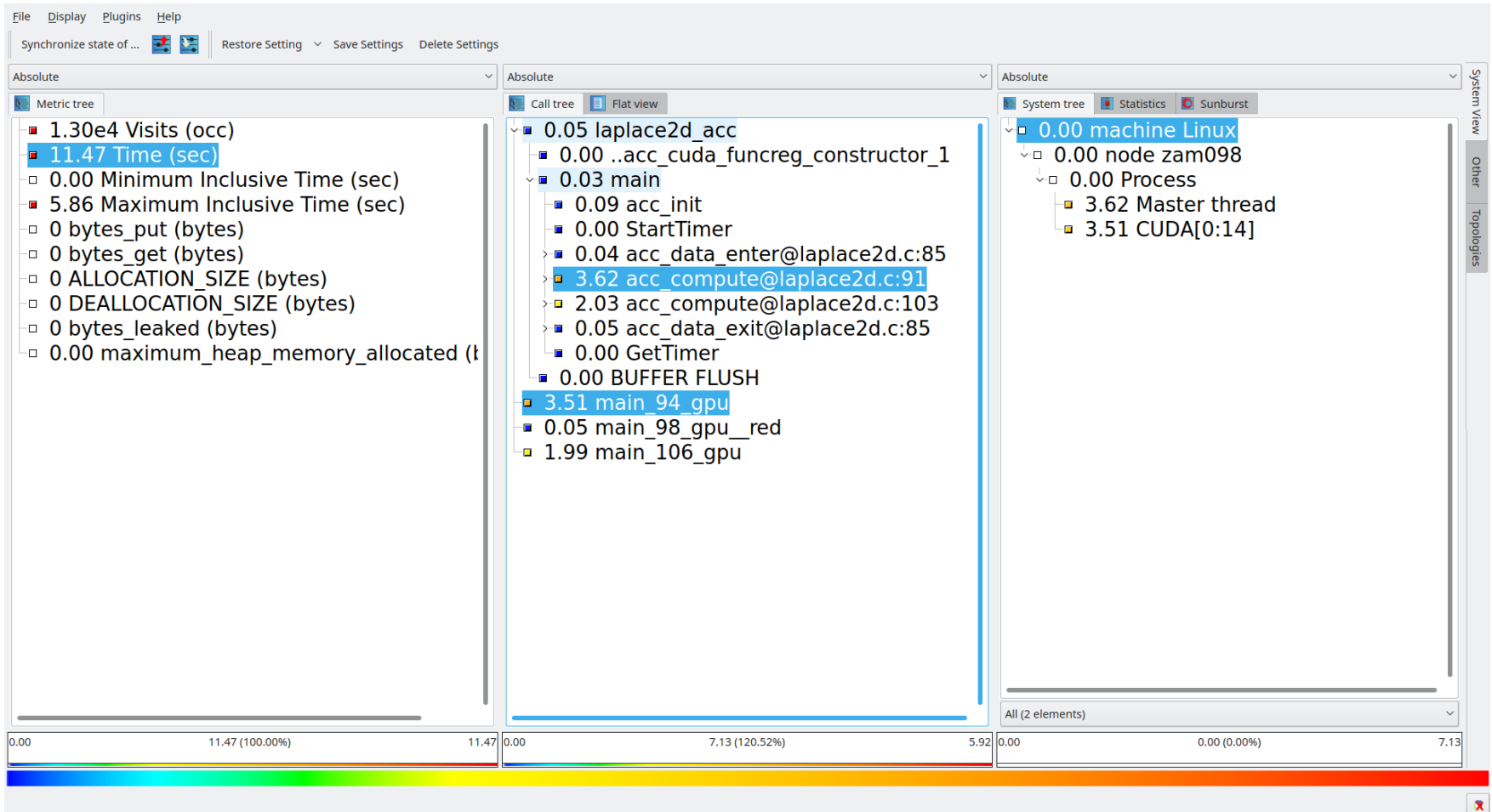
Where are my kernels?

Advisor Score-P Configuration ScorePion Source

```
79 for (int j = 1; j < n; j++)
80 {
81     Anew[j][0] = y0[j];
82     Anew[j][m-1] = y0[j]*exp(-pi);
83 }
84
85 #pragma acc data copy(A, Anew)
86 while ( error > tol && iter < iter_max )
87 {
88     error = 0.f;
89
90 #pragma omp parallel for shared(m, n, Anew, A)
91 #pragma acc kernels
92 for( int j = 1; j < n-1; j++)
93 {
94     for( int i = 1; i < m-1; i++)
95     {
96         Anew[j][i] = 0.25f * ( A[j][i+1] + A[j][i-1]
97                               + A[j-1][i] + A[j+1][i]);
98         error = fmaxf( error, fabsf(Anew[j][i]-A[j][i]));
99     }
100 }
101
102 #pragma omp parallel for shared(m, n, Anew, A)
103 #pragma acc kernels
104 for( int j = 1; j < n-1; j++)
105 {
106     for( int i = 1; i < m-1; i++)
107     {
108         A[j][i] = Anew[j][i];
109     }
110 }
```

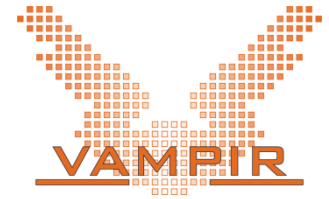
Pure OpenACC measurements give host-side events only

EXAMPLE: OPENACC + CUDA



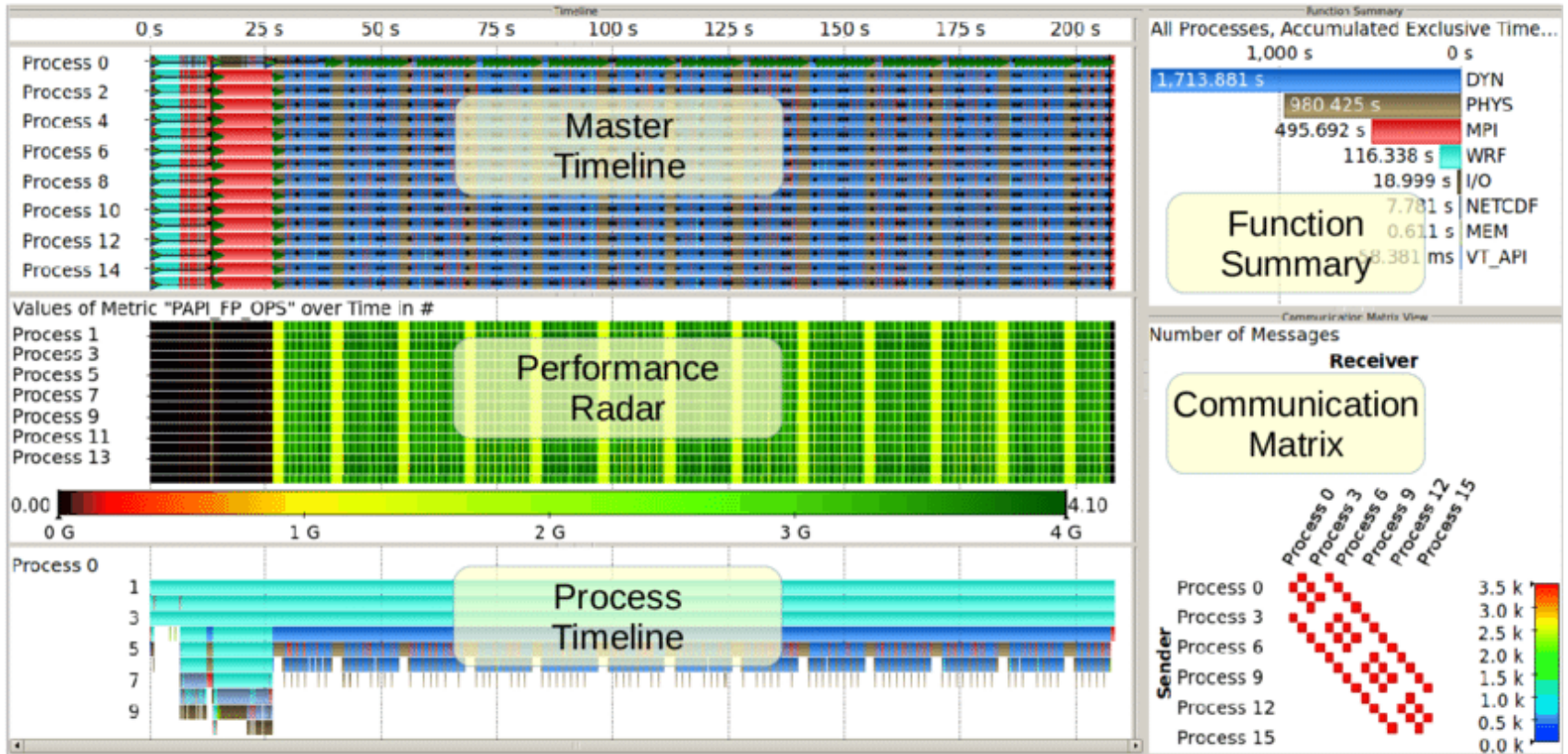
Enabling CUDA also shows kernels on the GPU

VAMPIR EVENT TRACE VISUALIZER

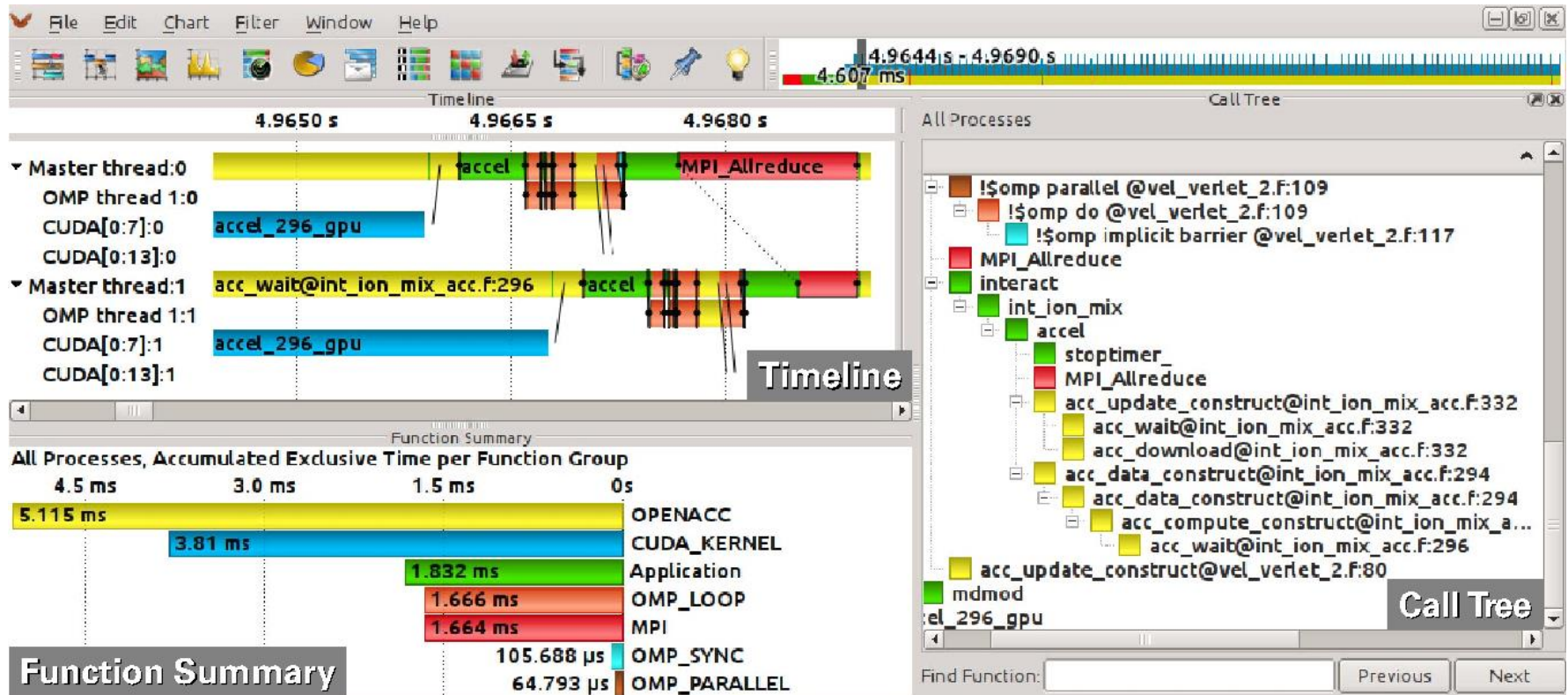


- Offline trace visualization for Score-P's OTF2 trace files
- Visualization of MPI, OpenMP, GPU and application events:
 - All diagrams highly customizable (through context menus)
 - Large variety of displays for ANY part of the trace
- <http://www.vampir.eu>
- Advantage:
 - Detailed view of dynamic application behavior
- Disadvantage:
 - Requires event traces (huge amount of data)
 - Completely manual analysis

VAMPIR DISPLAYS



VAMPIR COMPLEX APPLICATION



REMARK: NO SINGLE SOLUTION IS SUFFICIENT!



☞ *A combination of different methods, tools and techniques is typically needed!*

- Analysis
 - Statistics, visualization, automatic analysis, data mining, ...
- Measurement
 - Sampling / instrumentation, profiling / tracing, ...
- Instrumentation
 - Source code / binary, manual / automatic, ...

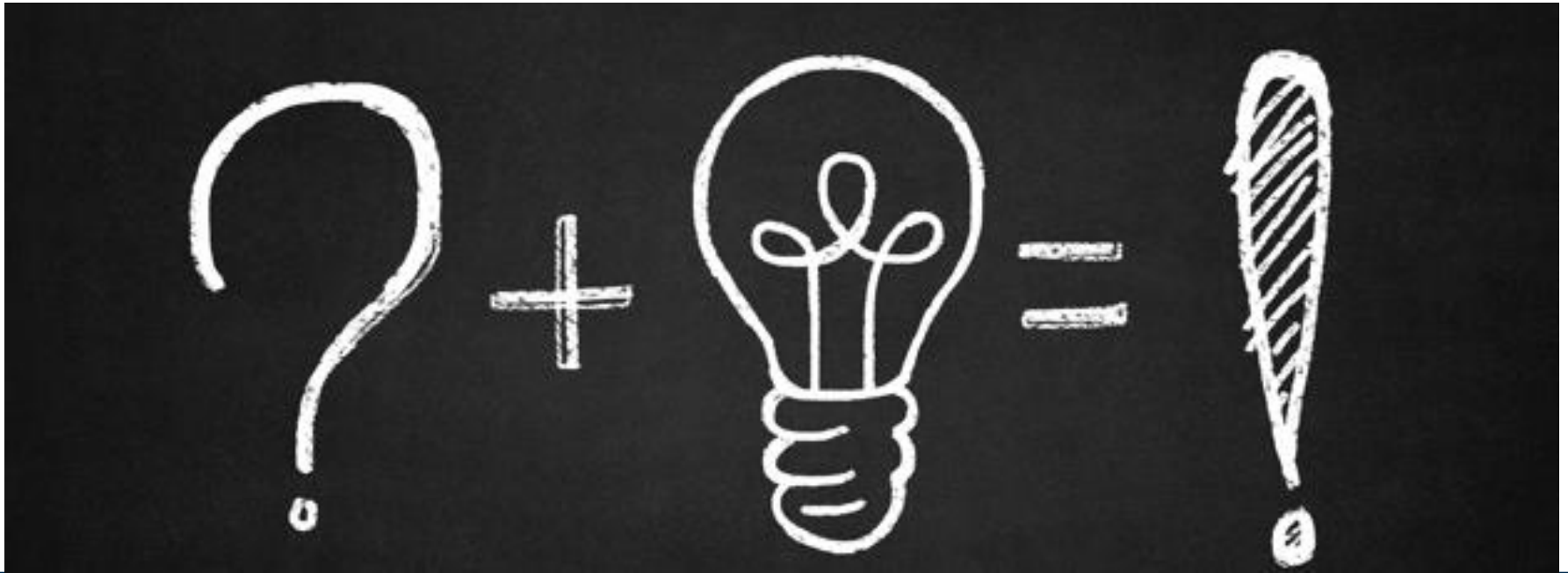
WHAT NOW?

- The tools are there – what now?
- Development phase:
 - Use NVIDIA tools
 - Debug: CUDA-MEMCHECK/CUDA-GDB
 - Performance: Nsight Systems and Compute
- Scaling up:
 - Use 3rd-party tools
 - Debug: TotalView/DDT
 - Performance: Score-P, Vampir

NEED HELP?

- Talk to the experts
 - NVIDIA Application Lab
 - CST Performance Analysis
 - CST Application Optimization
 - Apply for a POP audit

☞ Successful performance engineering often is a collaborative effort



QUESTIONS