

ARGONNE
ATPESC2023
EXTREME - SCALE COMPUTING

HPCToolkit Performance Tools

Performance analysis of CPU and GPU-accelerated applications at Scale

HPCToolkit at Exascale on Frontier: 8K nodes, 64K MPI ranks + GPU Tiles

John Mellor-Crummey
Professor, Rice University

HPCToolkit Funding Acknowledgments

Government

Exascale Computing Project 17-SC-20-SC

Lawrence Livermore National Laboratory Subcontract B658833

Argonne National Laboratory Subcontract 9F-60073

Corporate

Advanced Micro Devices

TotalEnergies EP Research & Technology USA, LLC.

Rice University's HPCToolkit Performance Tools

Measure and analyze performance of CPU and GPU-accelerated applications

Easy: profile unmodified application binaries

Fast: low-overhead measurement

Informative: understand where an application spends its time and why

- call path profiles associate metrics with application source code contexts

- optional hierarchical traces to understand execution dynamics

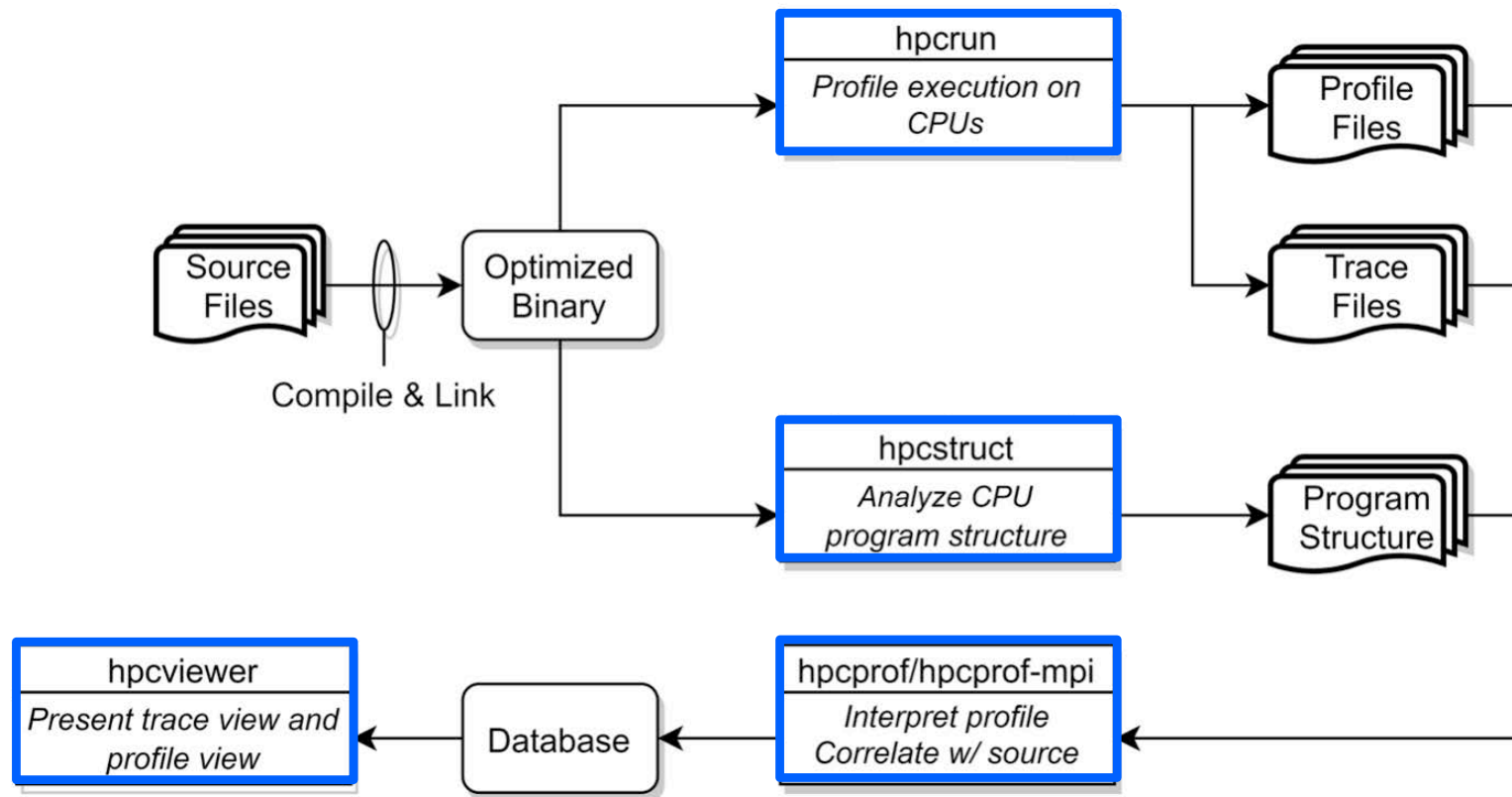
Broad audience

- application developers

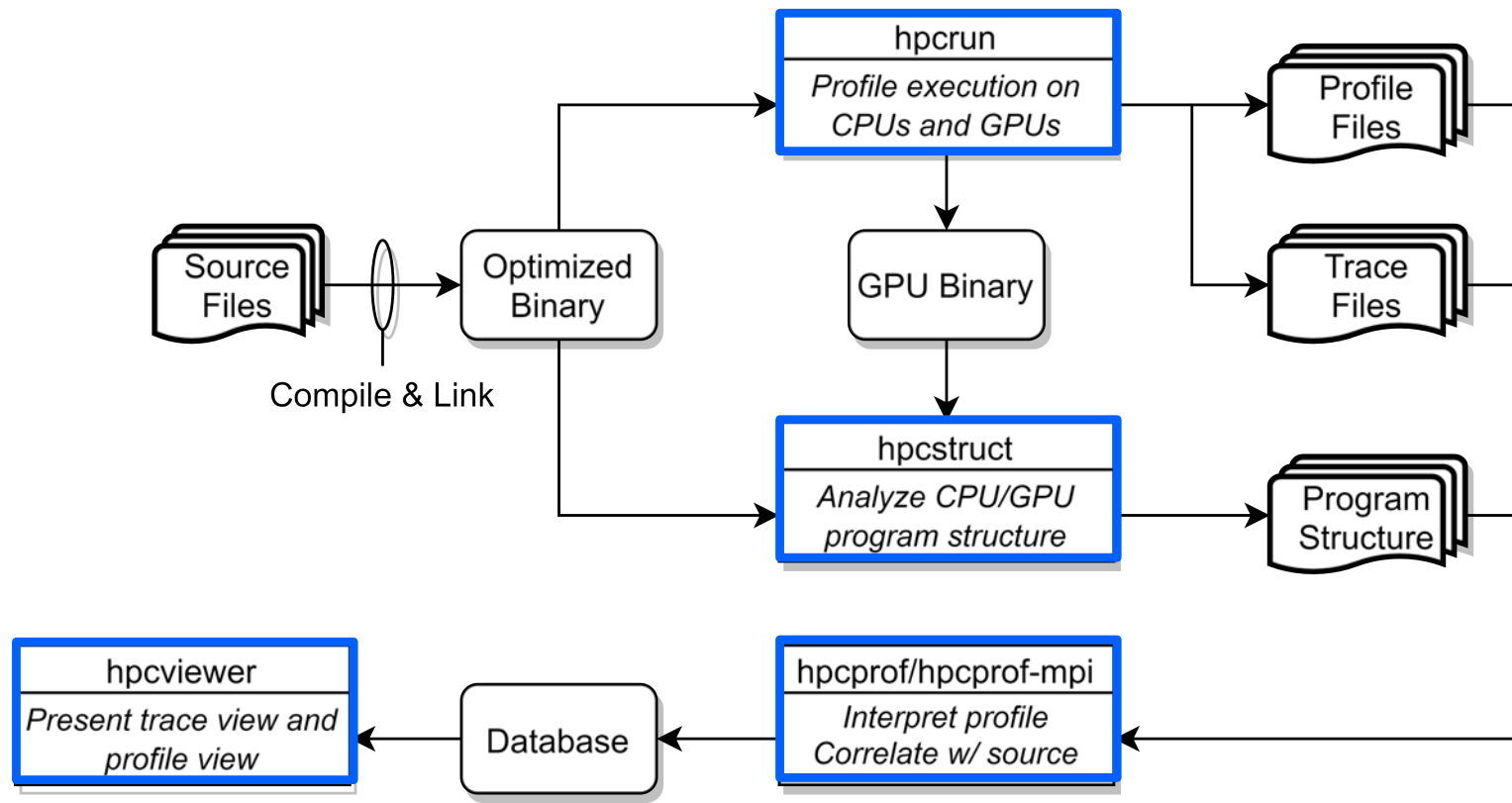
- framework developers

- runtime and tool developers

HPCToolkit's Workflow for CPU Applications



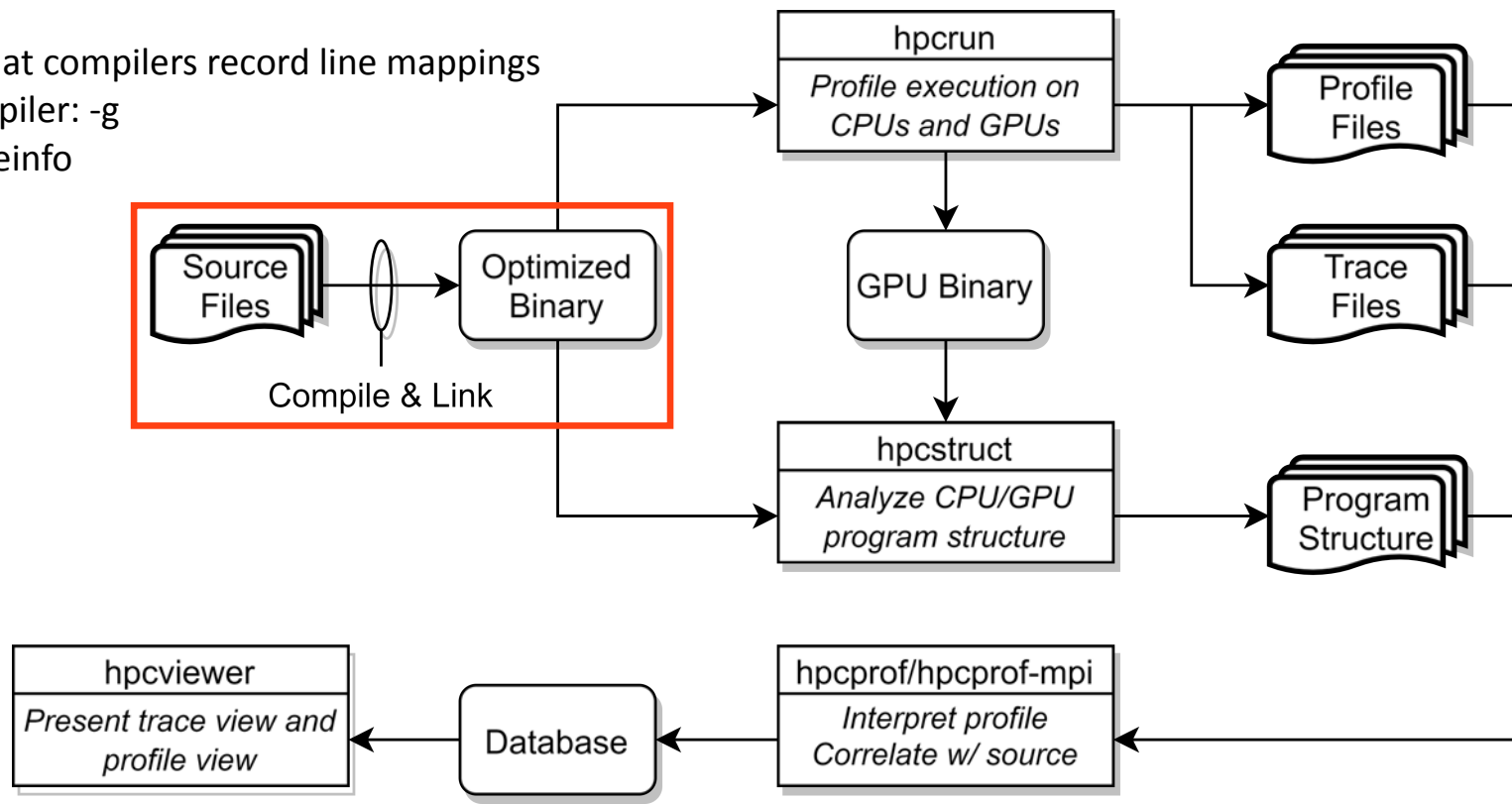
HPCToolkit's Workflow for GPU-accelerated Applications



HPCToolkit's Workflow for GPU-accelerated Applications

Step 1:

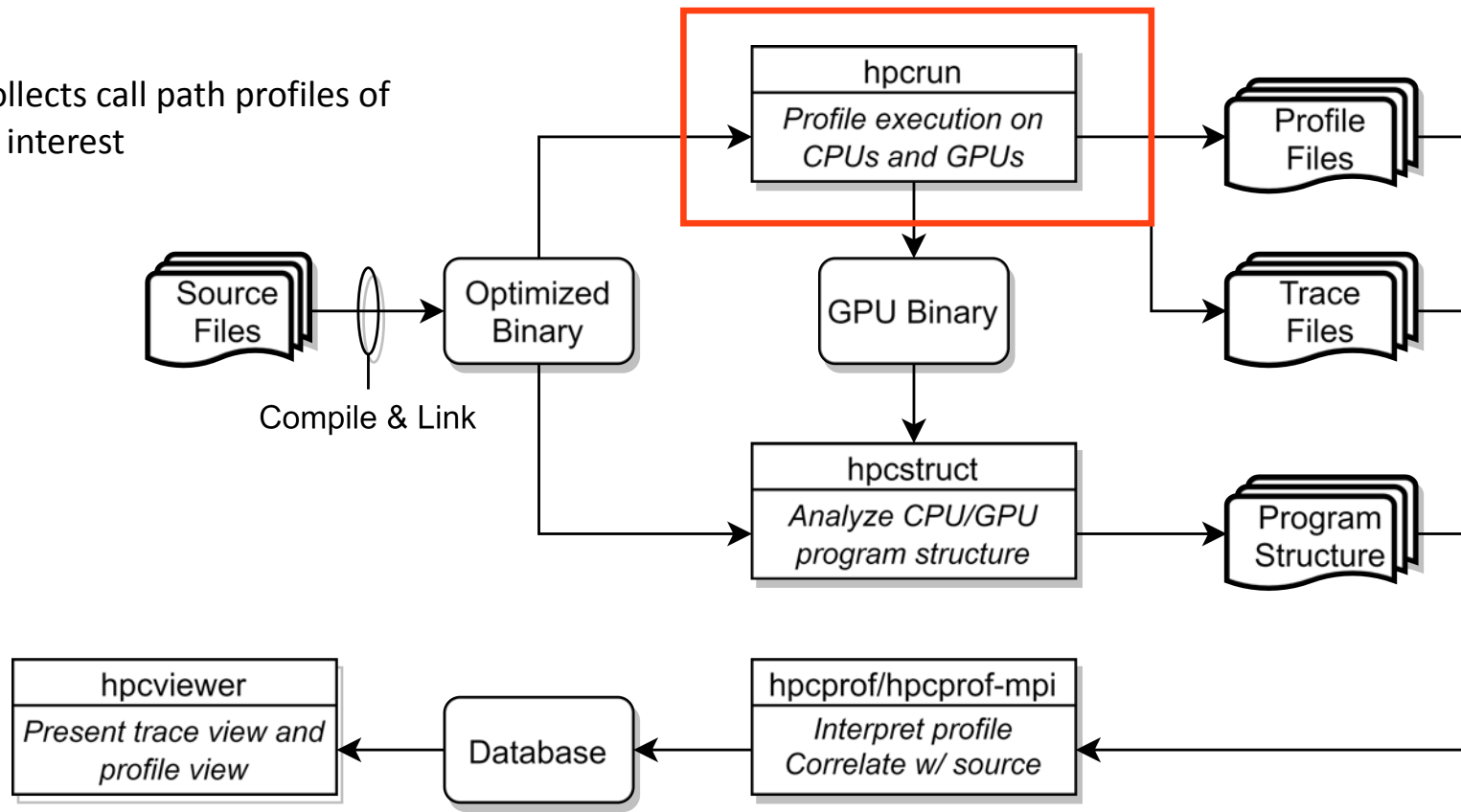
- Ensure that compilers record line mappings
- host compiler: `-g`
- `nvcc`: `-lineinfo`



HPCToolkit's Workflow for GPU-accelerated Applications

Step 2:

- *hpcrun* collects call path profiles of events of interest



Measurement of CPU and GPU-accelerated Applications

CPU

Sampling on timer interrupts and hardware counter overflows on the CPU

GPU

Callbacks when GPU operations are launched/completed

GPU event stream for GPU operations

PC Samples in GPU kernels (NVIDIA)

Instruction-level instrumentation (Intel)

hpcrun: Measure CPU and/or GPU activity

GPU profiling

```
hpcrun -e gpu=xxx <app> ...
```

xxx ∈ {nvidia,amd,opencl,level0}

GPU instrumentation (Intel GPU only)

```
hpcrun -e gpu=level0,inst=count,latency <app>
```

GPU PC sampling (NVIDIA GPU only)

```
hpcrun -e gpu=nvidia,pc <app>
```

CPU and GPU Tracing (in addition to profiling)

```
hpcrun -e CPUTIME -e gpu=xxx -t <app>
```

Use hpcrun with job launchers

```
jsrun -n 32 -g 1 -a 1 hpcrun -e gpu=xxx <app>
```

```
srun -n 1 -G 1 hpcrun -e gpu=xxx <app>
```

```
aprun -n 16 -N 8 -d 8 hpcrun -e gpu=xxx <app>
```

Profiles: aggregated on the fly

- a calling context tree per thread
- a calling context tree per GPU stream
- instruction level measurements

CPU traces

- trace of call stack samples

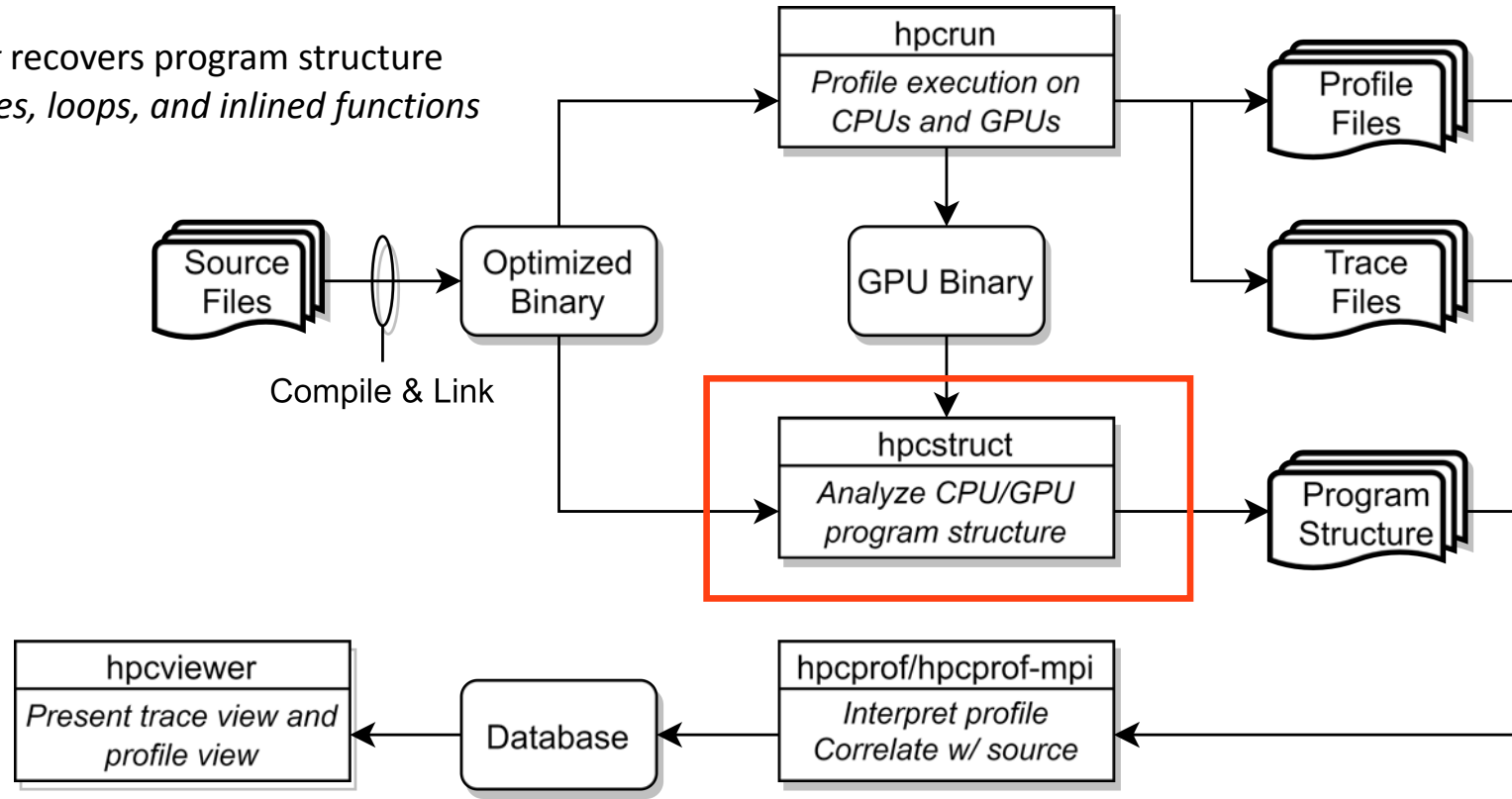
GPU traces

- trace of call stacks that initiate GPU operations

HPCToolkit's Workflow for GPU-accelerated Applications

Step 3:

- *hpcstruct* recovers program structure about lines, loops, and inlined functions



hpcstruct: Analyze CPU and GPU Binaries Using Multiple Threads

Usage

```
hpcstruct [--gpucfg yes] <measurement-directory>
```

What it does

Recover program structure information

Files, functions, inlined templates or functions, loops, source lines

In parallel, analyze all CPU and GPU binaries that were measured by HPCToolkit

default: use $\text{size}(\text{CPU set})/2$ threads

analyze large application binaries with 16 threads

analyze multiple small application binaries concurrently with 2 threads each

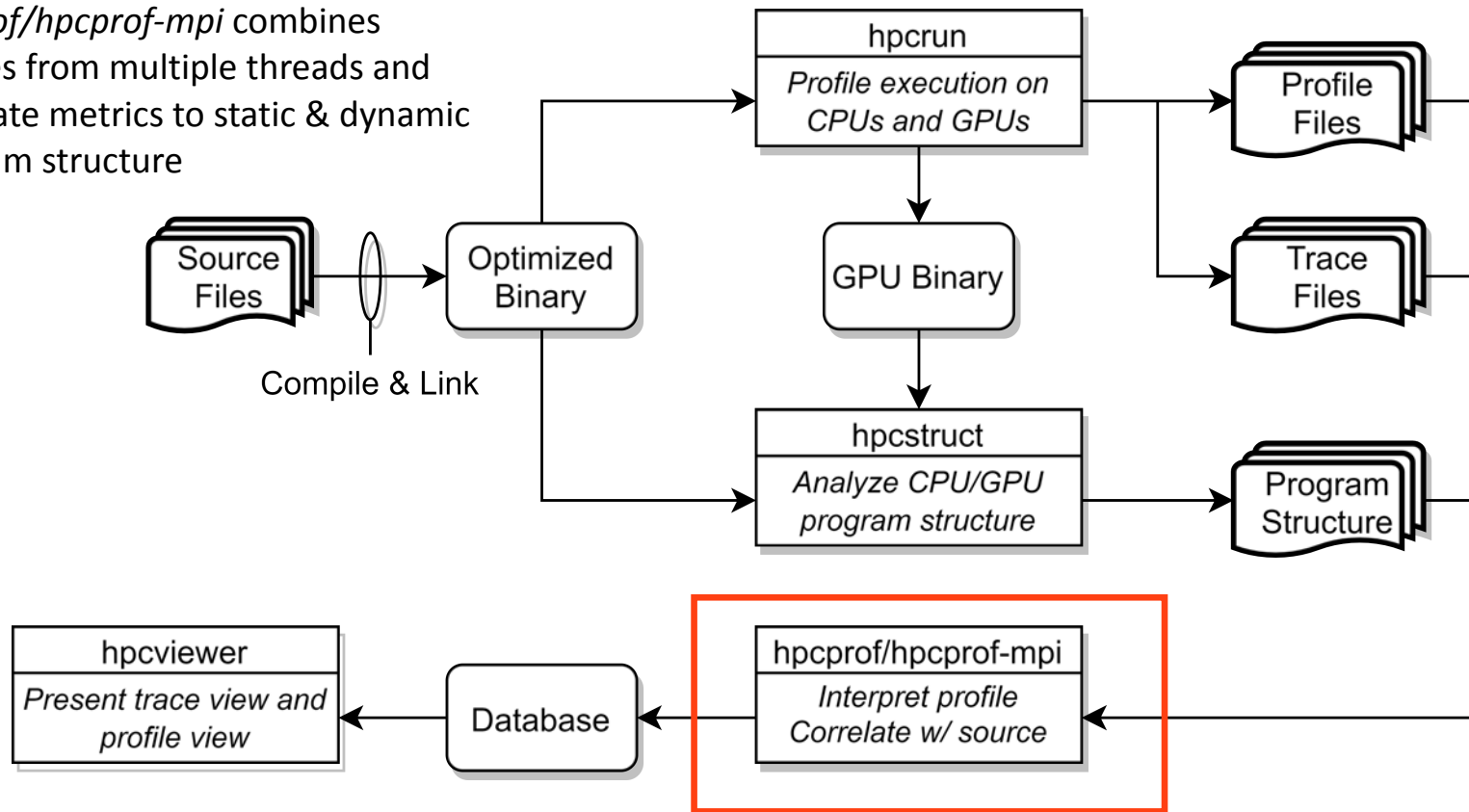
Cache binary analysis results for reuse when analyzing other executions

NOTE: `--gpucfg yes` needed only for analysis of GPU binaries when NVIDIA PC samples were collected

HPCToolkit's Workflow for GPU-accelerated Applications

Step 4:

- *hpcprof/hpcprof-mpi* combines profiles from multiple threads and correlate metrics to static & dynamic program structure



hpcprof/hpcprof-mpi: Associate Measurements with Program Structure

Analyze data from modest executions with threaded parallelism

```
hpcprof <measurement-directory>
```

Analyze data from large executions using both distributed-memory and shared-memory parallelism

```
jsrun -n 32 -a 1 hpcprof-mpi <measurement-directory>
```

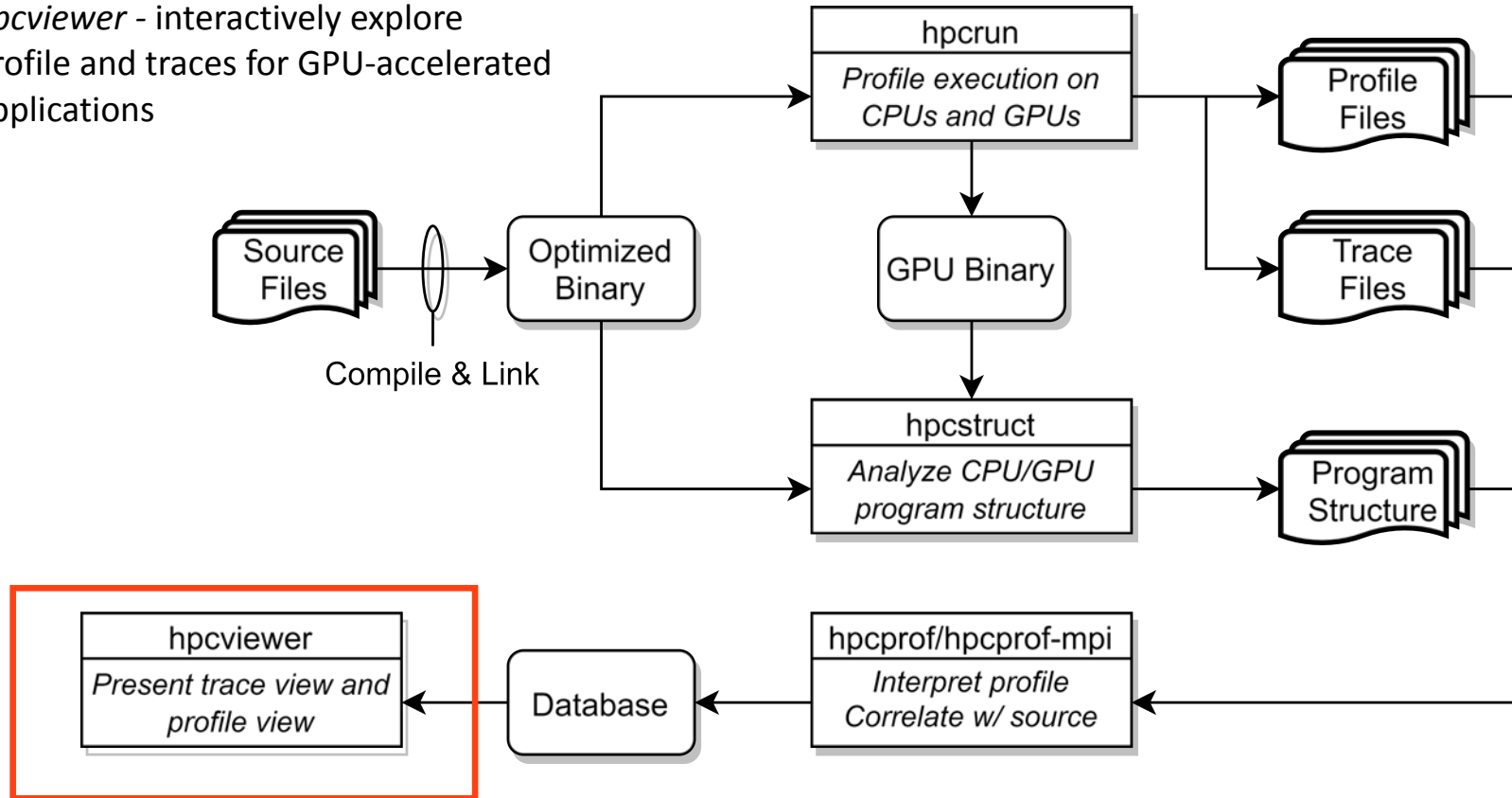
```
srun -n 32 hpcprof-mpi <measurement-directory>
```

```
aprun -n 128 -N 8 hpcprof-mpi <measurement-directory>
```

HPCToolkit's Workflow for GPU-accelerated Applications

Step 4:

- *hpcviewer* - interactively explore profile and traces for GPU-accelerated applications



Understanding Temporal Behavior

Profiling compresses out the temporal dimension

Temporal patterns, e.g. serial sections and dynamic load imbalance are invisible in profiles

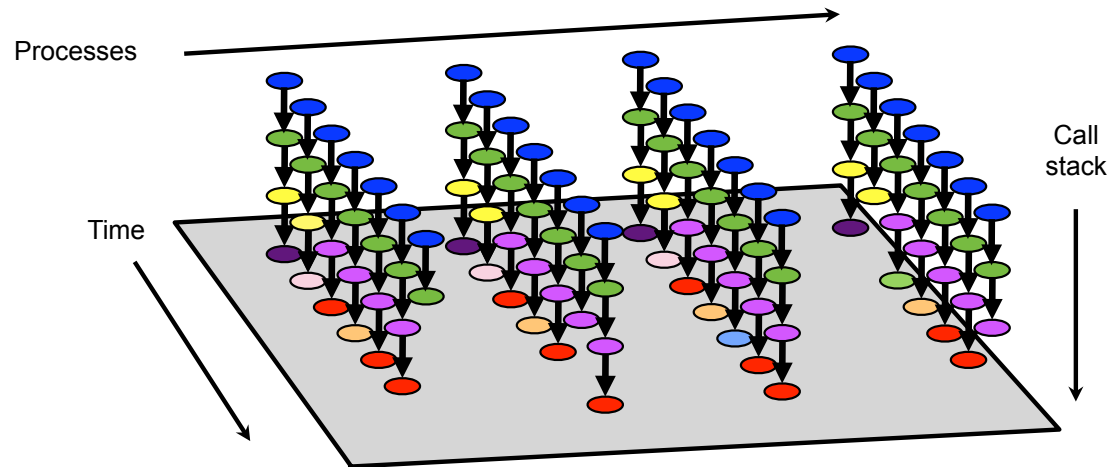
What can we do? Trace call path samples

N times per second, take a call path sample of each thread

Organize the samples for each thread along a time line

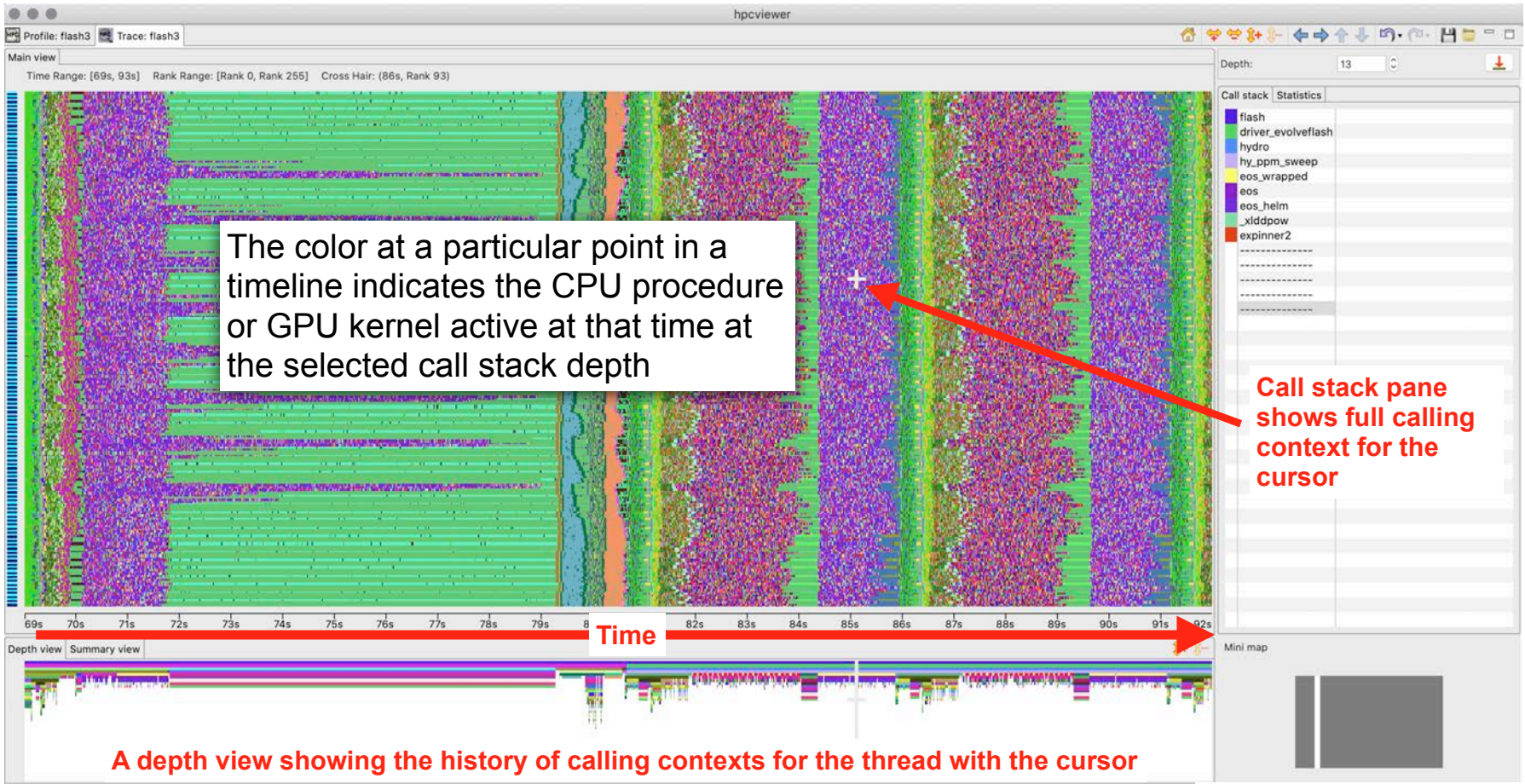
View how the execution evolves left to right

What do we view? assign each procedure a color; view a depth slice of an execution



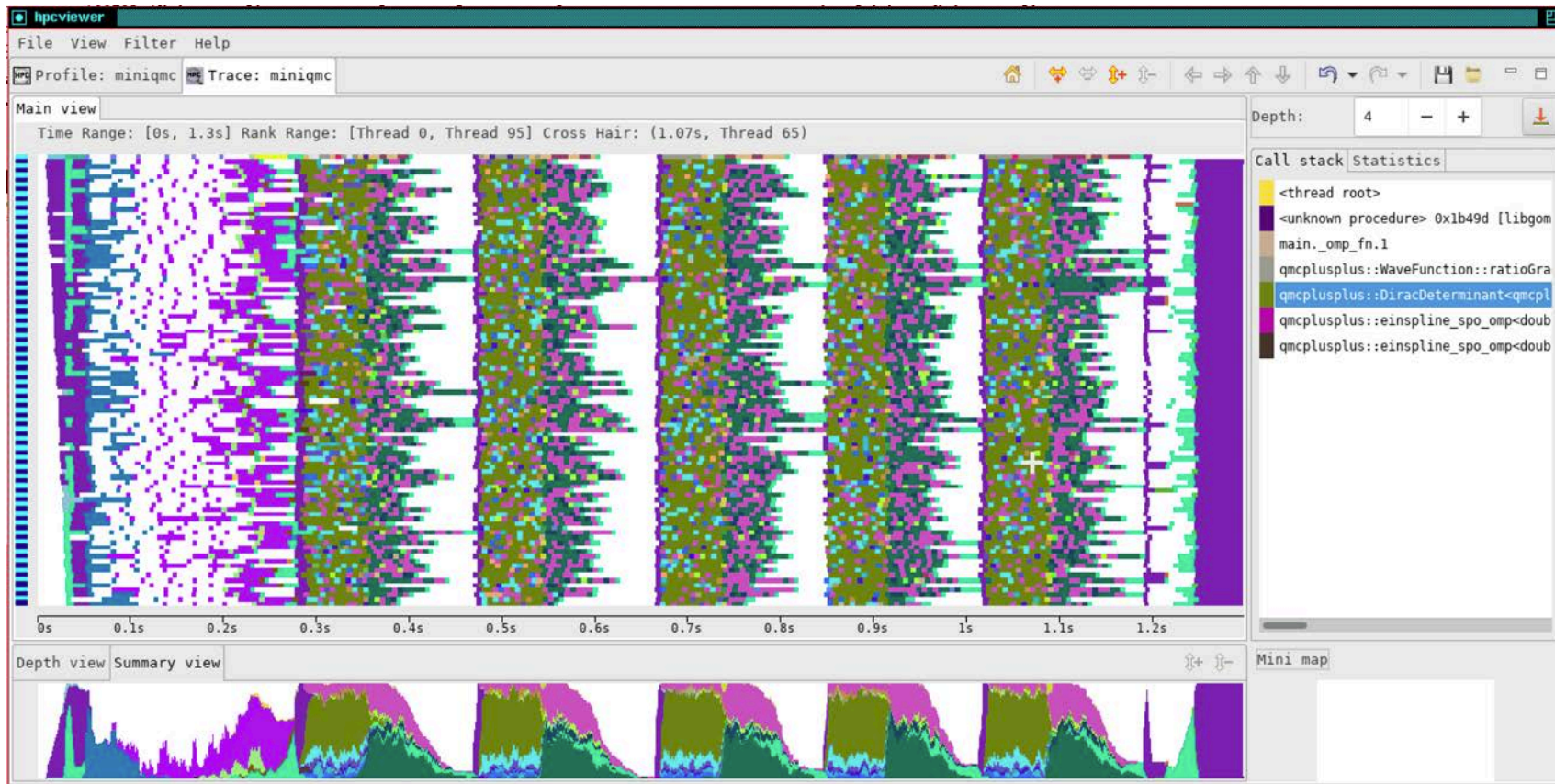
Time-centric Analysis with hpcviewer

MPI ranks, OpenMP Threads, GPU streams

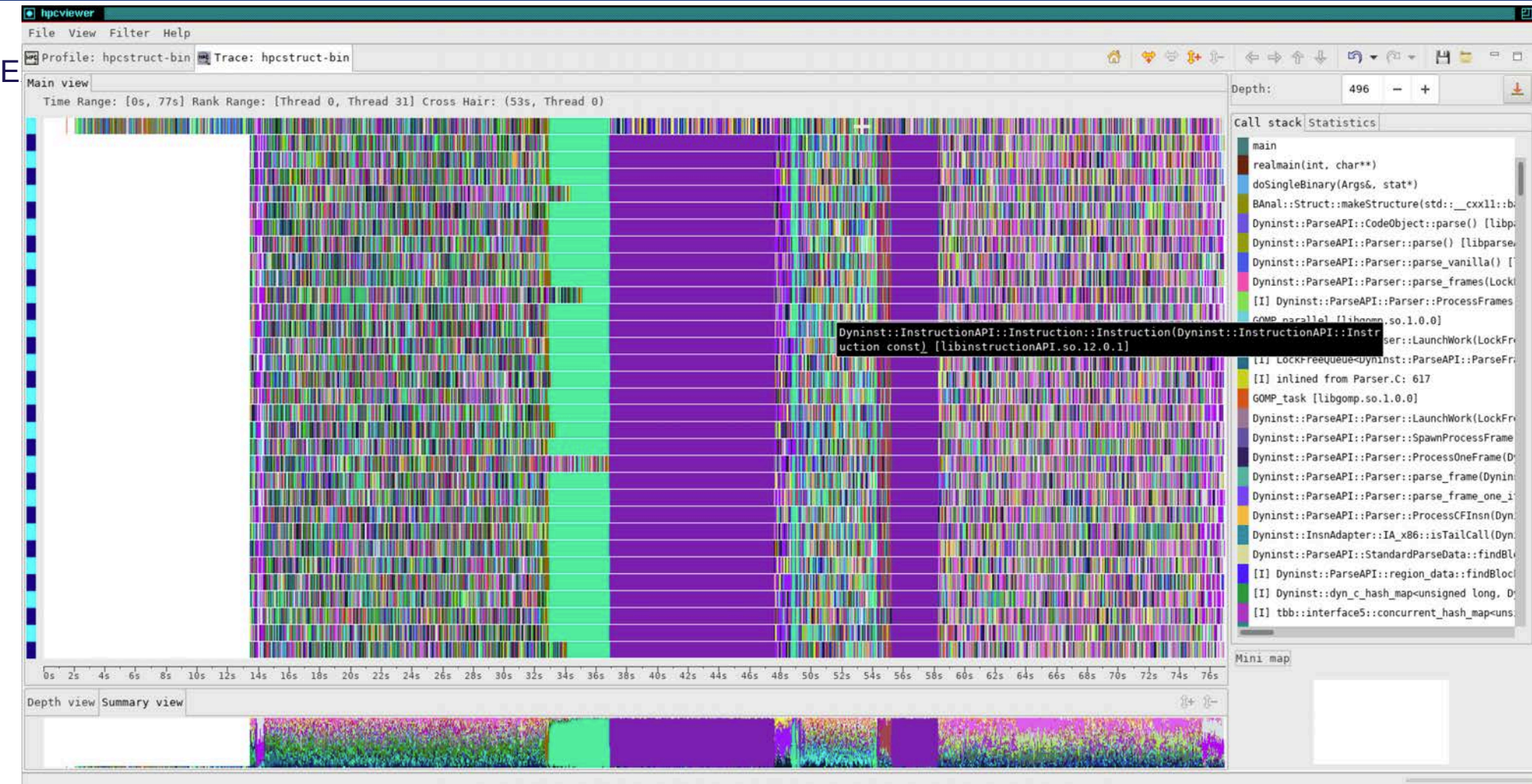


White Intervals in Traces Indicate Blocking of CPU Threads or GPU Streams

Miniqmc: OpenMP on 32 CPU threads

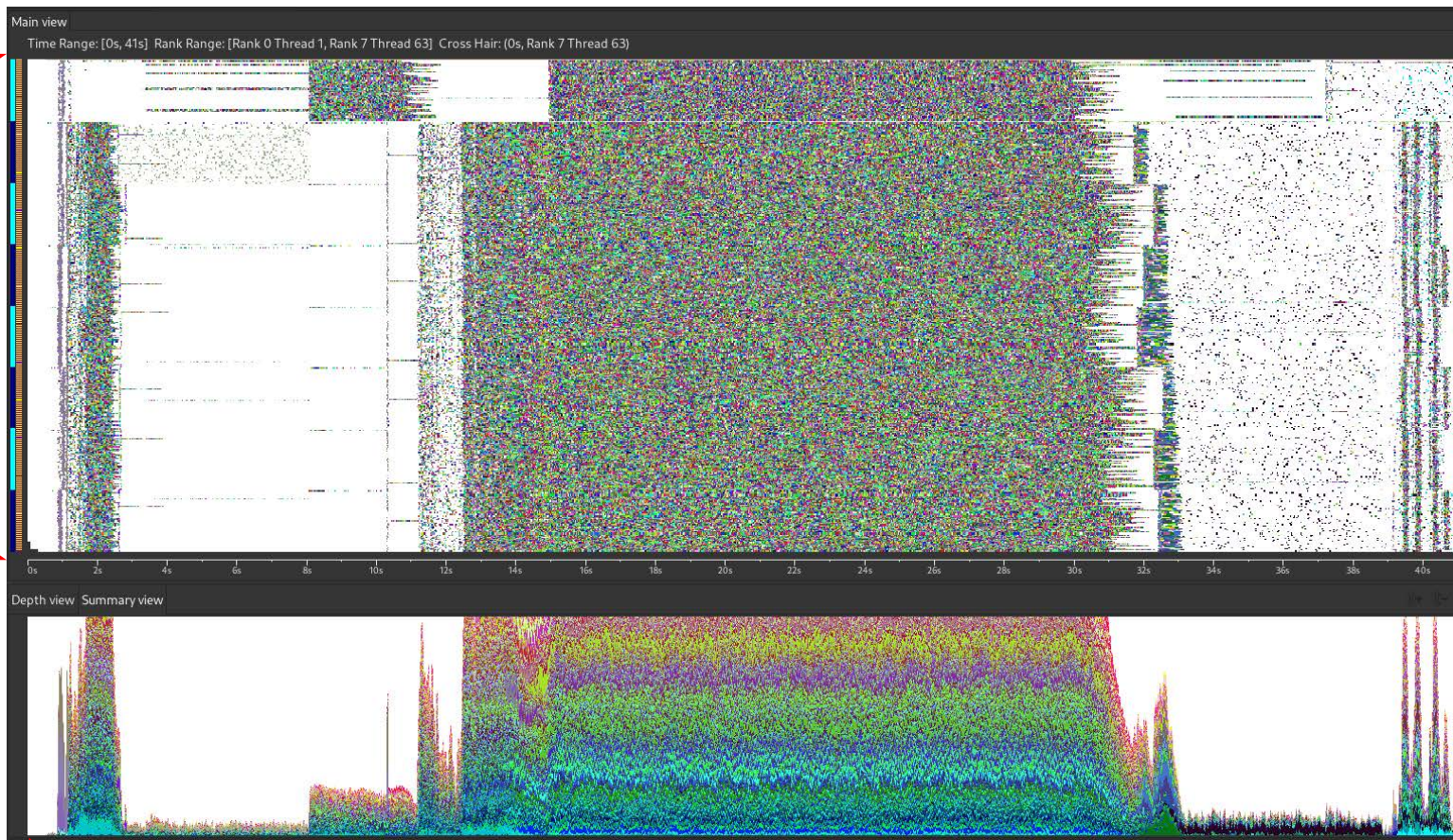


hpcstruct Example: Analyze 7.7GB TensorFlow library (170MB text) in 77s



hpcprof-mpi: Analyze Measurements of LAMMPS @ 2K threads + 2K GPUs

Analysis on 8 nodes
using 504 threads!



Completes in 41s!

Coarse- and Fine-grain Measurement on NVIDIA GPUs: LLNL's Quicksilver

Compute Node

-2xPower9 + 4xNVIDIA GPUs

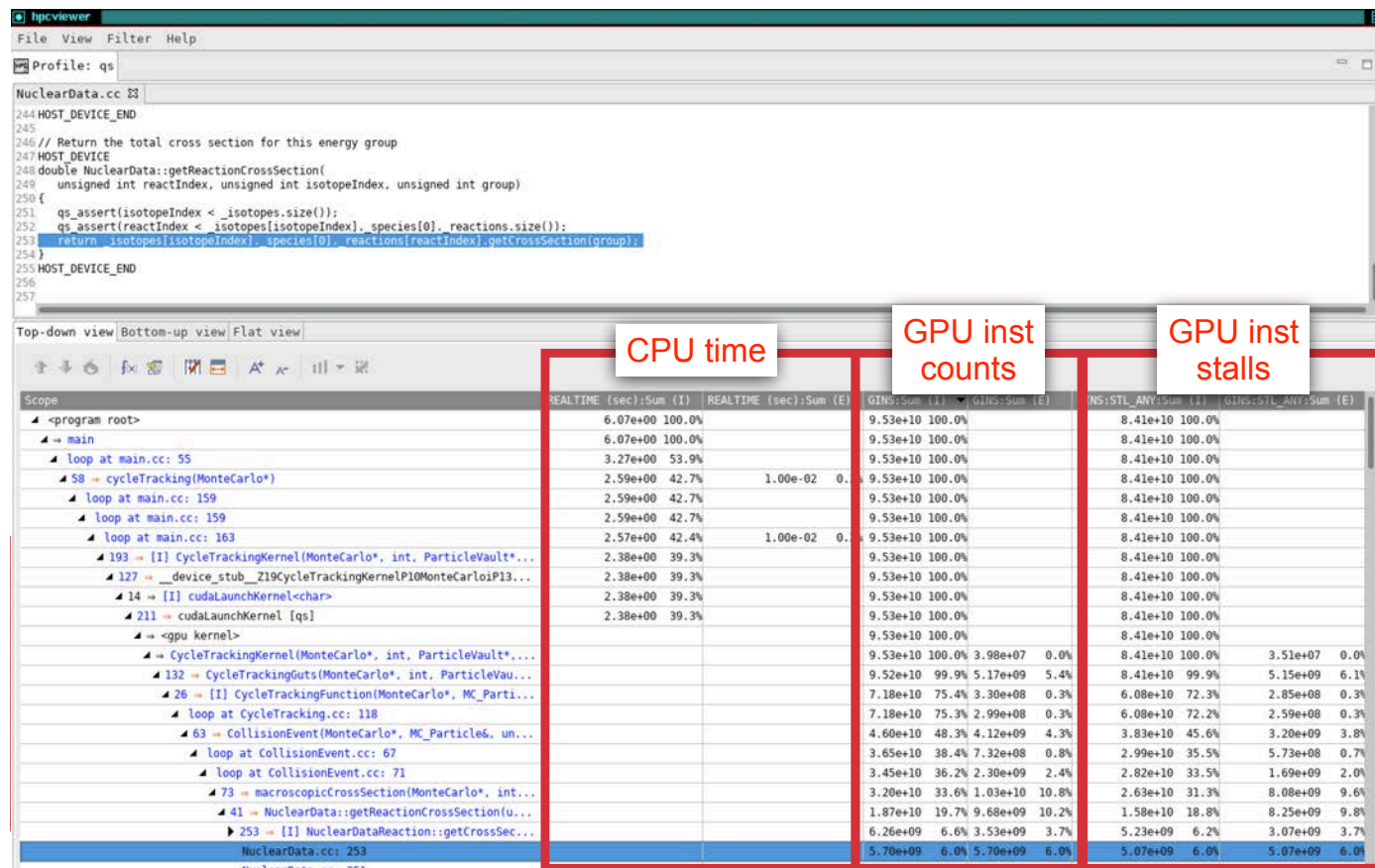
Optimized (-O2) compilation with nvcc

Detailed measurement and attribution using PC sampling

Attribute information to heterogeneous calling context

Key Metrics

- instructions executed
- instruction stalls and reasons
- GPU utilization



K. Zhou, M. W. Krentel, and J. Mellor-Crummey. Tools for top-down performance analysis of GPU-accelerated applications. International Conference on Supercomputing. ACM, NY, NY, USA, June, 2020.

Analysis of PeleC using PC Sampling on an NVIDIA GPU

```

438 UserData udata = static_cast<ARKODEUserData*>(user_data);
439 udata->dt_save = t;
440
441 #ifdef AMREX_USE_GPU
442 const auto ec = amrex::Gpu::ExecutionConfig(udata->ncells_d);
443 amrex::launch_global<<<
444 udata->nbBlocks, udata->nbThreads, ec.sharedMem, udata->stream>>>{
445 [=] AMREX_GPU_DEVICE() noexcept {
446     for (int icell = blockDim.x * blockDim.x + threadIdx.x,
447         stride = blockDim.x * gridDim.x;
448         icell < udata->ncells_d; icell += stride) {
449         fKernelSpec(
450             icell, udata->dt_save, udata->ireactor_type, yvec_d, ydot_d,
451             udata->rho_init_d, udata->rhoescr_ext_d, udata->rYsrc_d);
452     }
453 };
454 #else
455 for (int icell = 0; icell < udata->ncells_d; icell++) {
456     fKernelSpec(

```

Cause:
passed udata structure pointer to lambda capture

Improvement:
pass udata components as scalars
<https://github.com/AMReX-Combustion/PelePhysics/pull/192>
4% speedup on PeleC PMF drm19 test case

Scope	GINS:[0,0] (I)	GINS:[0,0] (E)	GINS:STL_ANY:[0,0] (I)	GINS:STL_ANY:[0,0] (E)	GINS:STL_GMEM:[0,0] (I)	GINS:STL_GMEM:[0,0] (E)
loop at AMReX_Amr.cpp: 2061	1.24e+13	88.6%	1.05e+13	88.7%	5.58e+12	89.3%
loop at Advance.cpp: 302	1.24e+13	88.5%	1.05e+13	88.6%	5.57e+12	89.2%
loop at React.cpp: 109	9.43e+12	67.2%	8.14e+12	68.7%	4.06e+12	65.0%
loop at React.cpp: 446	5.14e+12	36.6%	5.35e+10	0.4%	4.62e+10	0.4%
AMReX_Gpu.launchGlobal.H: 12	1.75e+10	0.1%	1.75e+10	0.1%	1.70e+10	0.1%

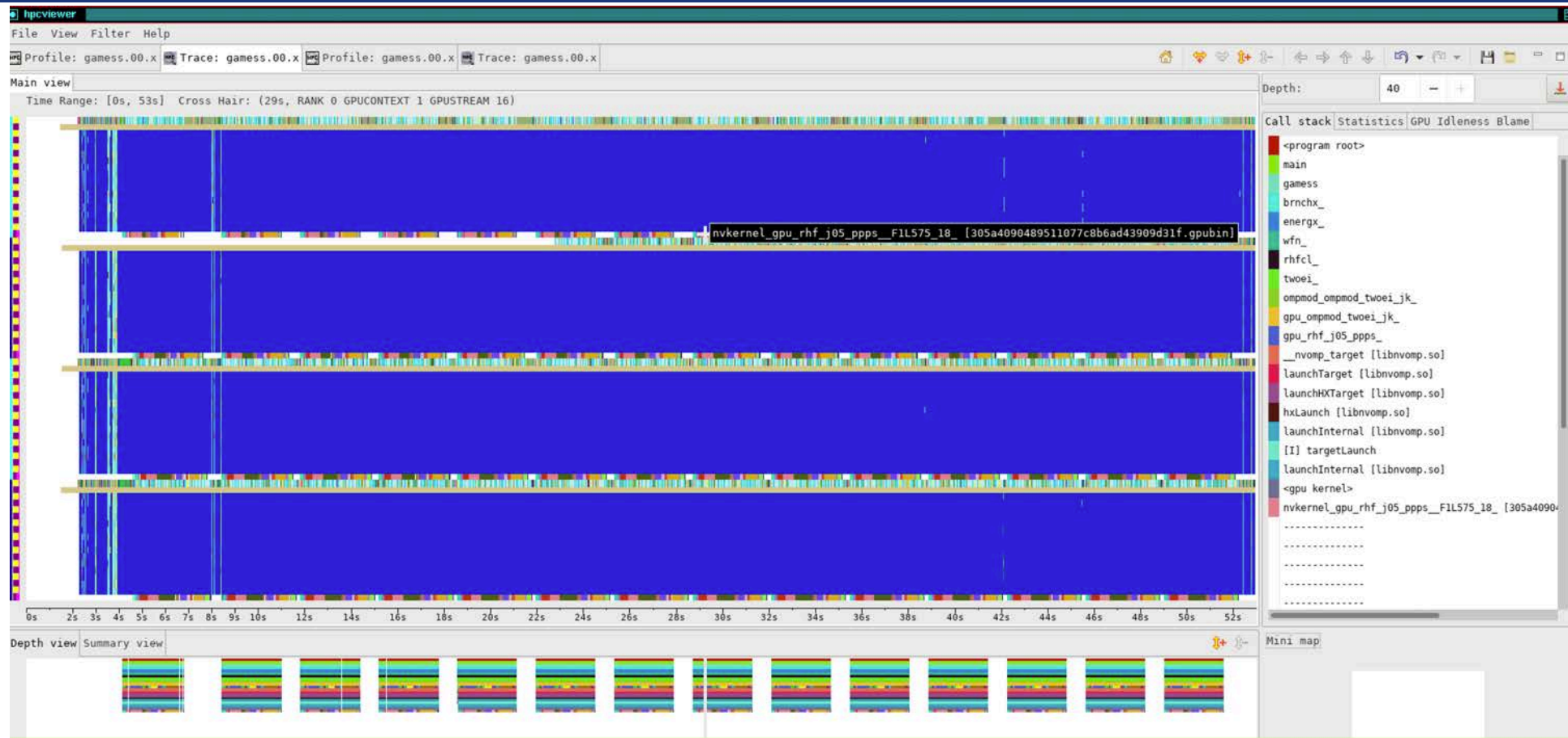
CPU context

GPU context

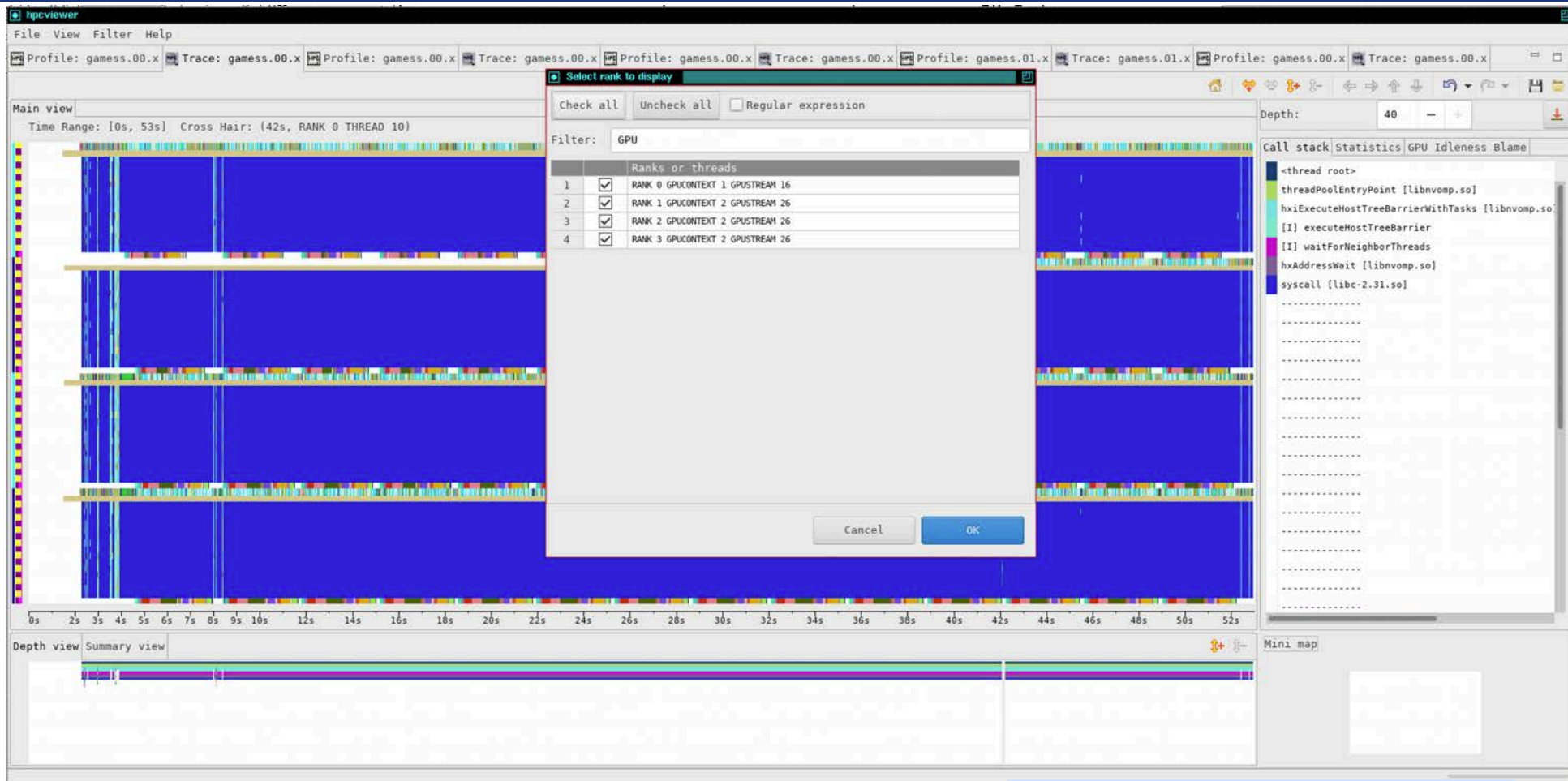
9.4% GPU stalls outside the loop

mostly memory stalls

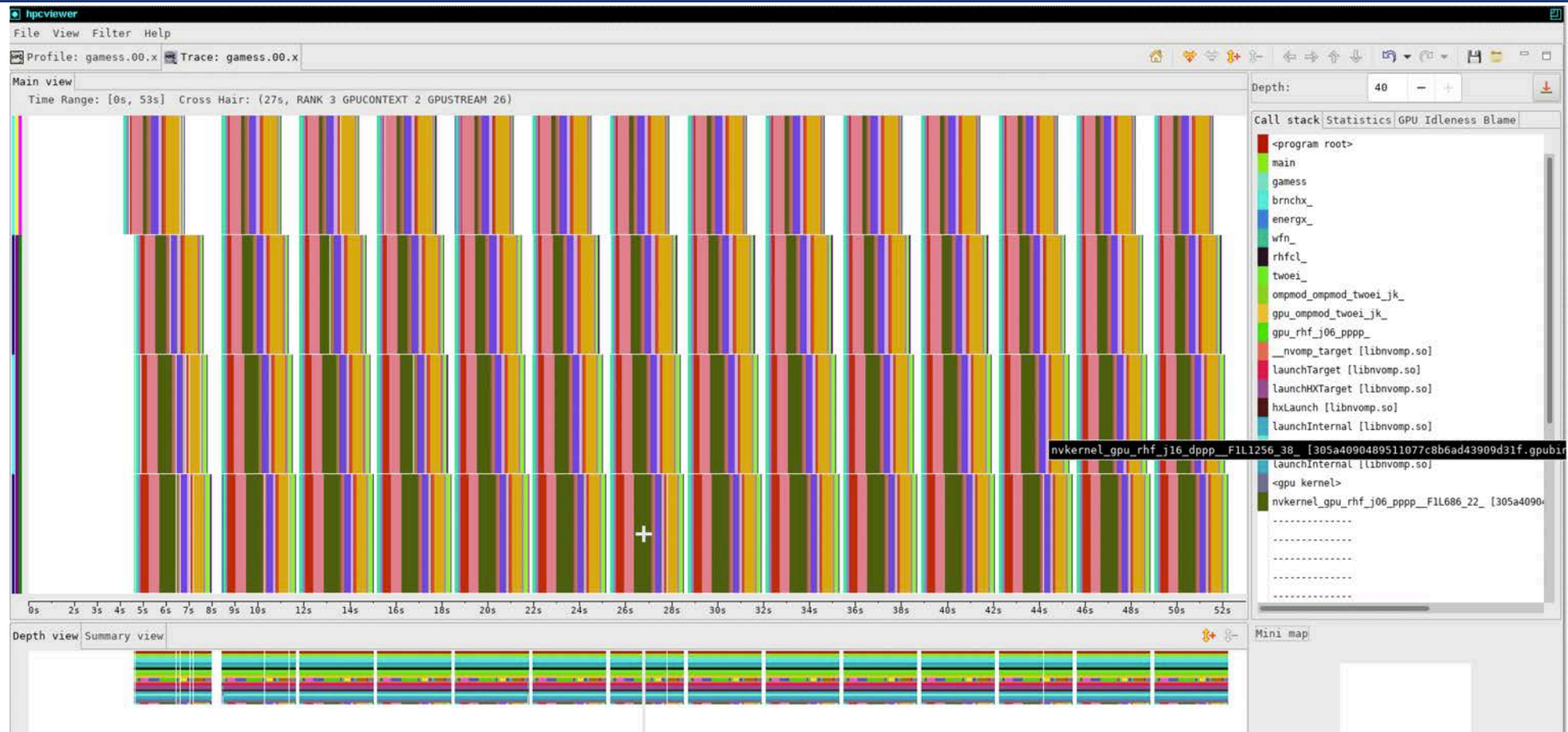
Time-centric Analysis: GAMESS 4 ranks, 4 GPUs on Perlmutter



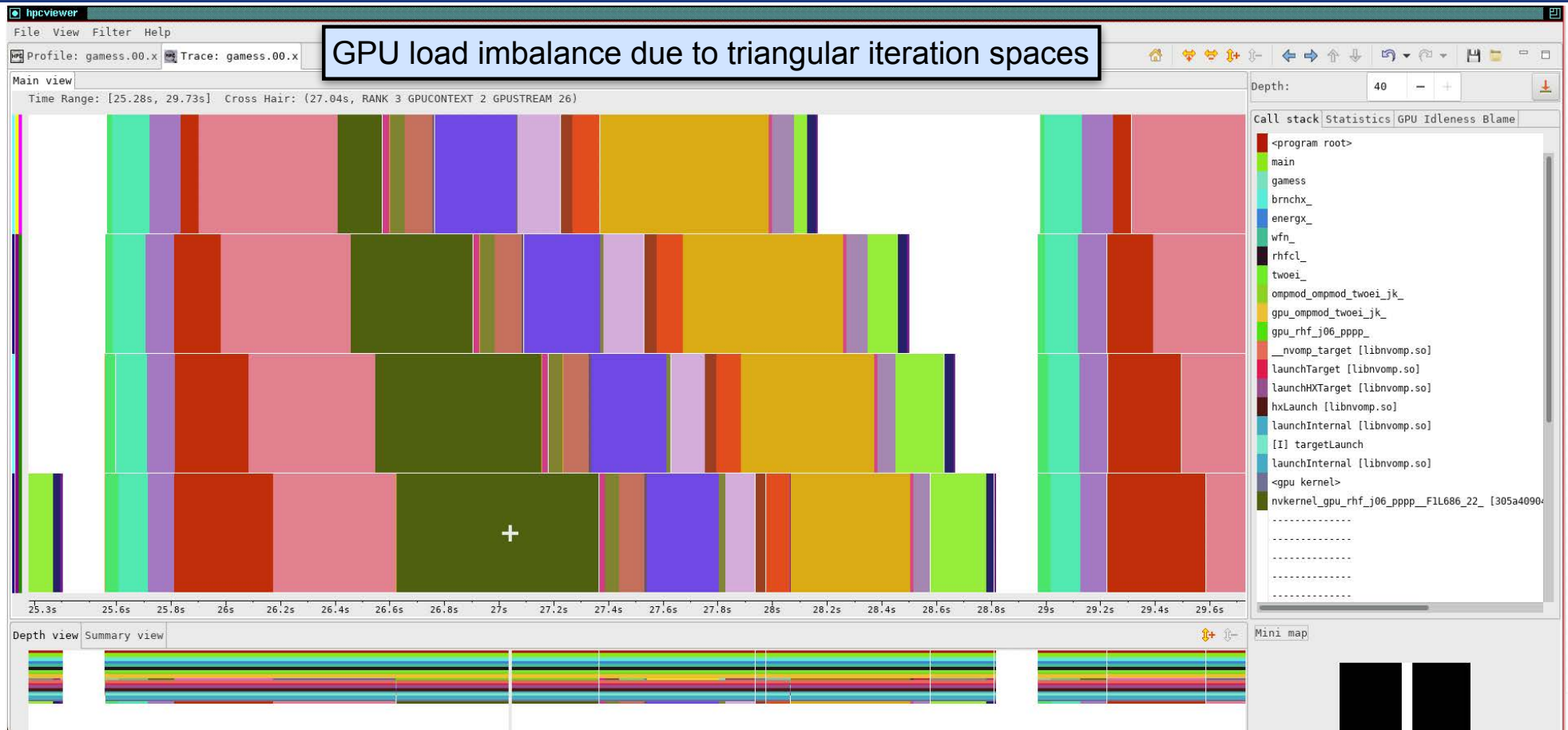
Time-centric Analysis: GAMESS 4 ranks, 4 GPUs on Perlmutter



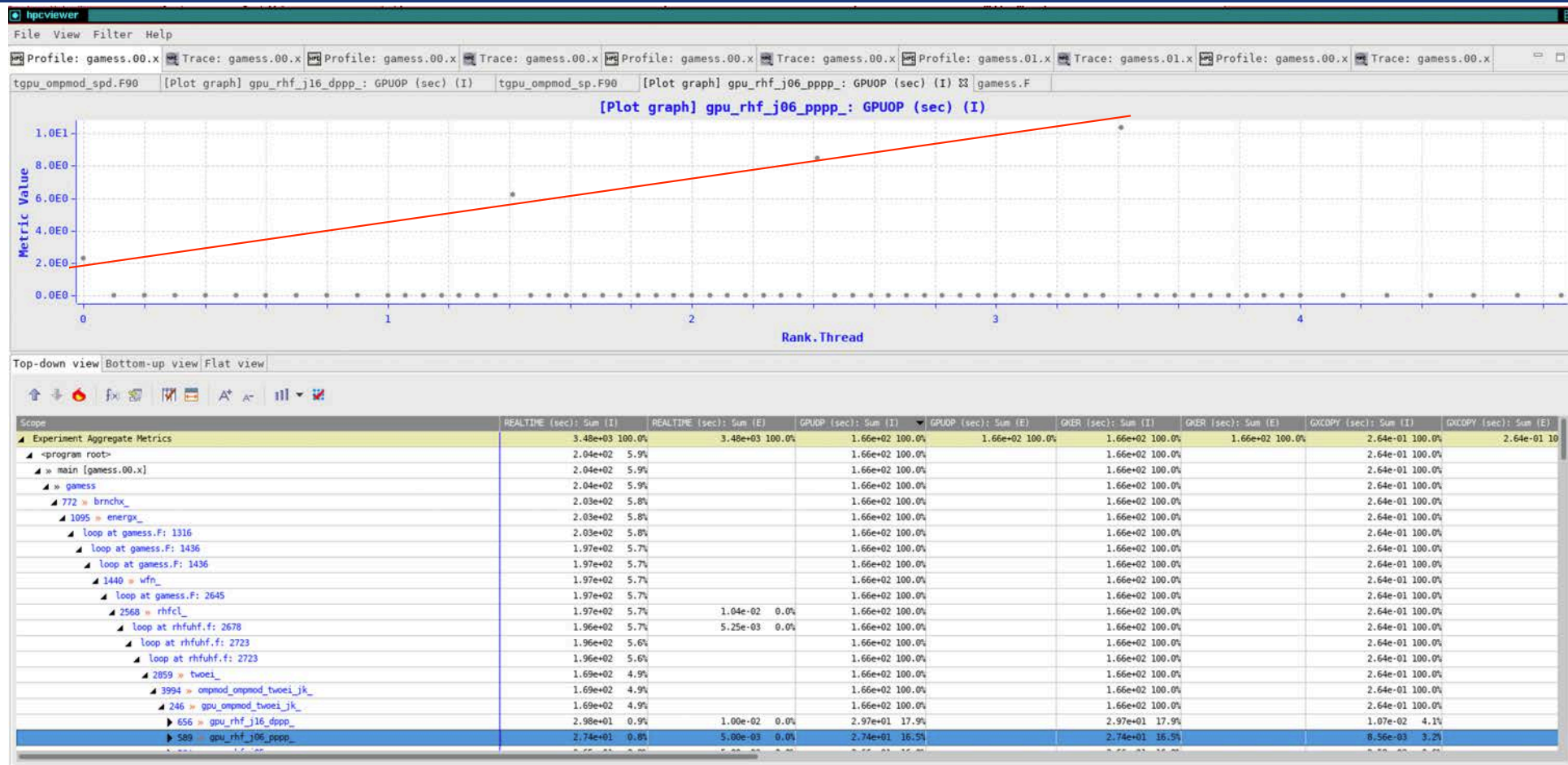
Time-centric Analysis: GAMESS 4 ranks, 4 GPUs on Perlmutter



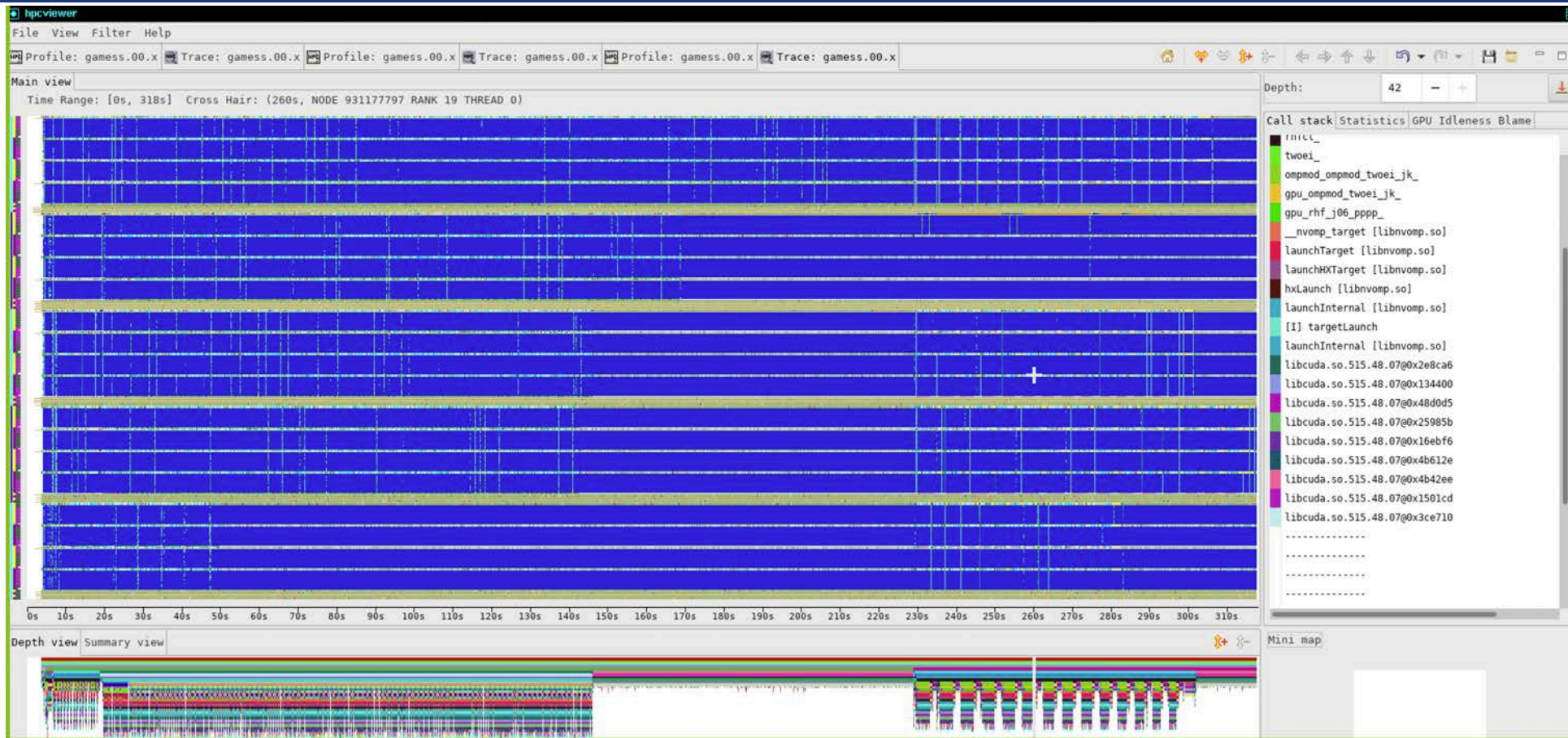
Time-centric Analysis: GAMESS 4 ranks, 4 GPUs on Perlmutter



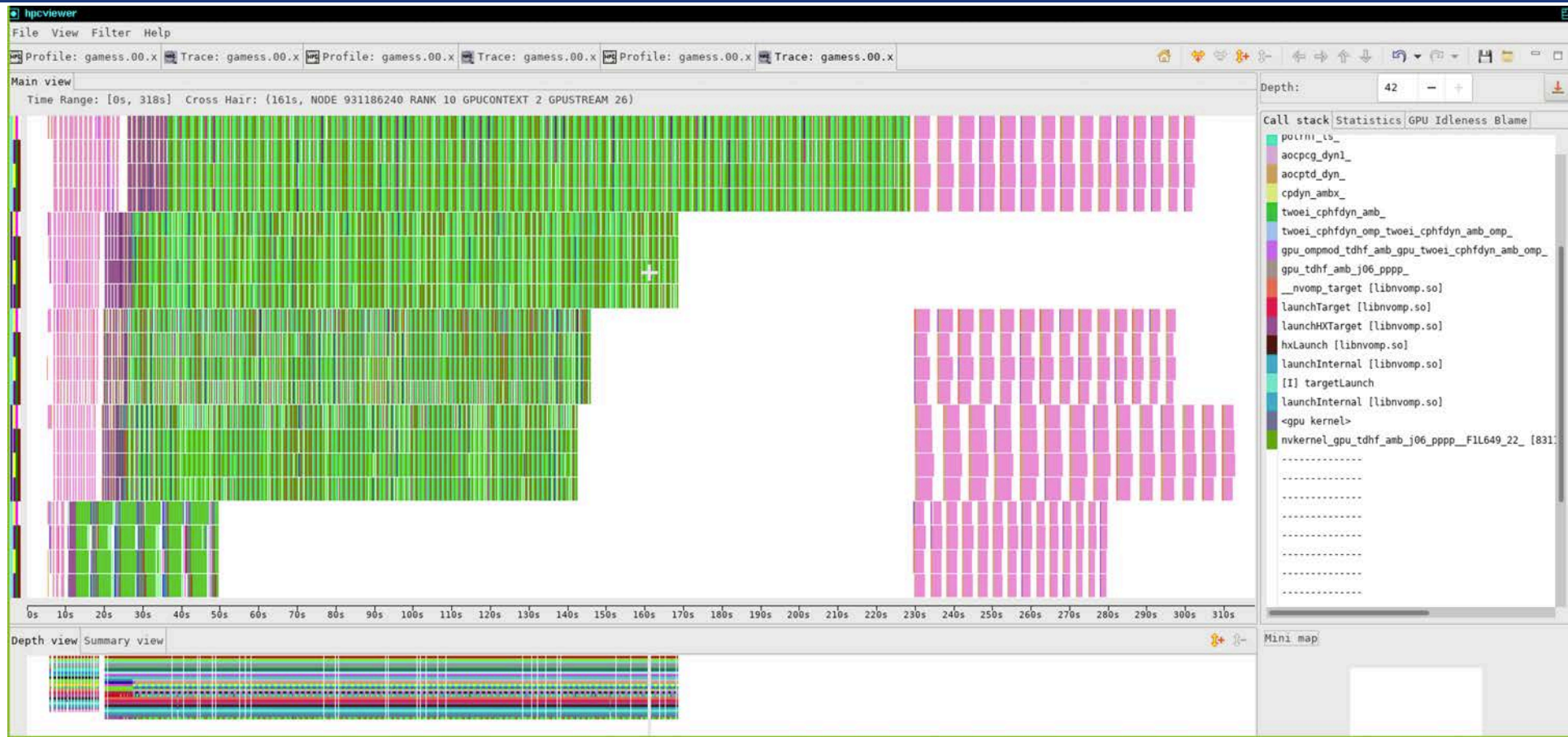
Time-centric Analysis: GAMESS 4 ranks, 4 GPUs on Perlmutter



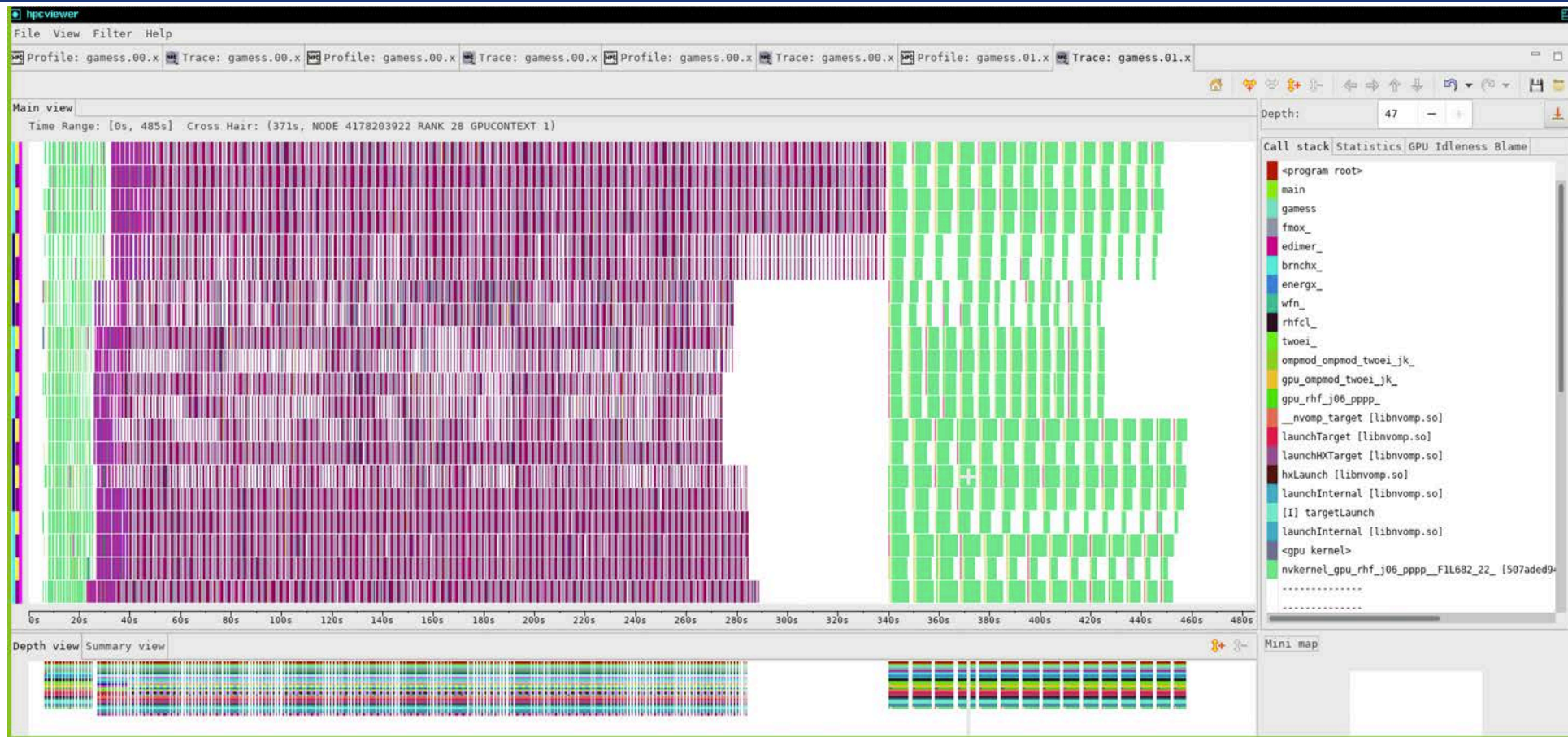
Time-centric Analysis: GAMESS 5 nodes, 40 ranks, 20 GPUs on Perlmutter



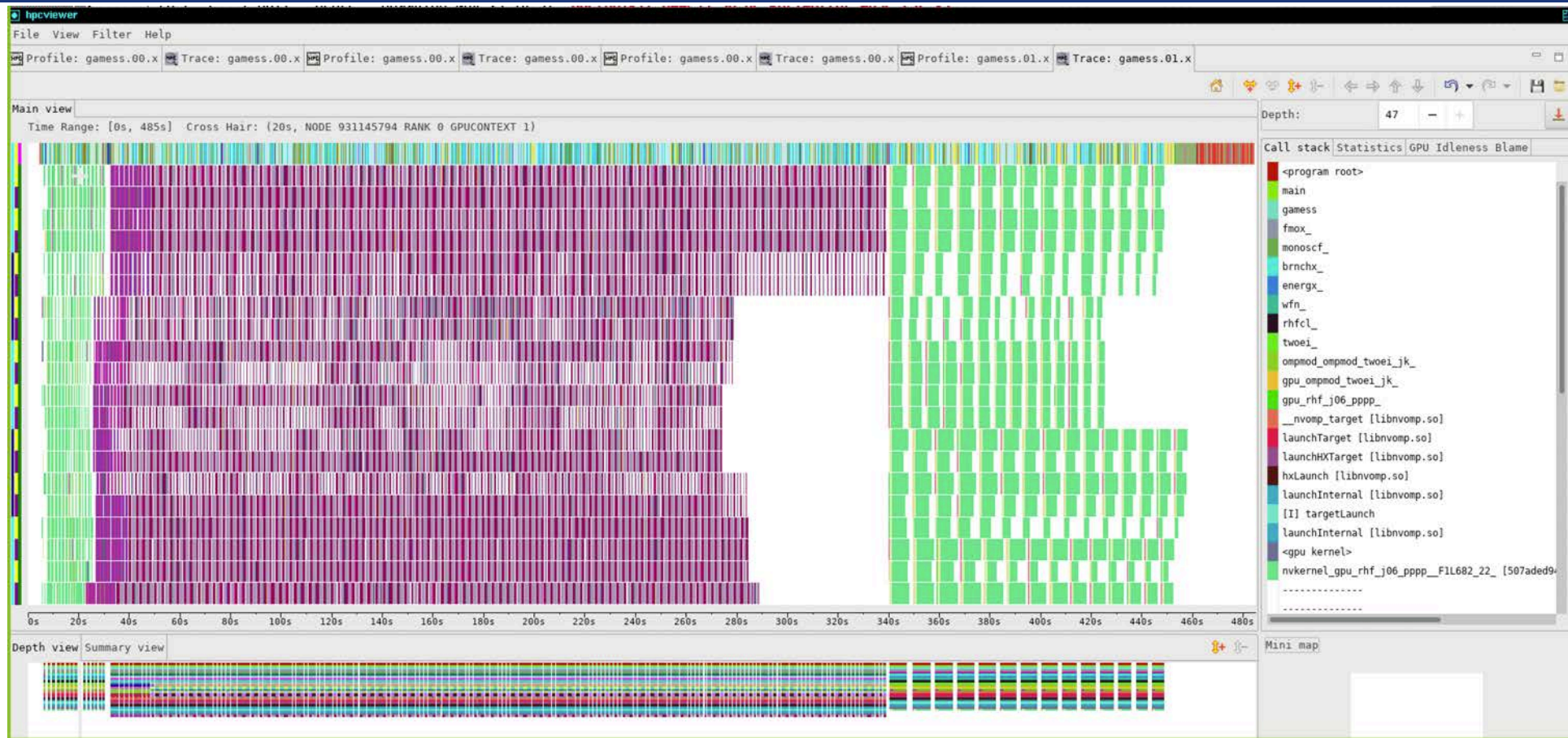
Time-centric Analysis: GAMESS 5 nodes, 40 ranks, 20 GPUs on Perlmutter



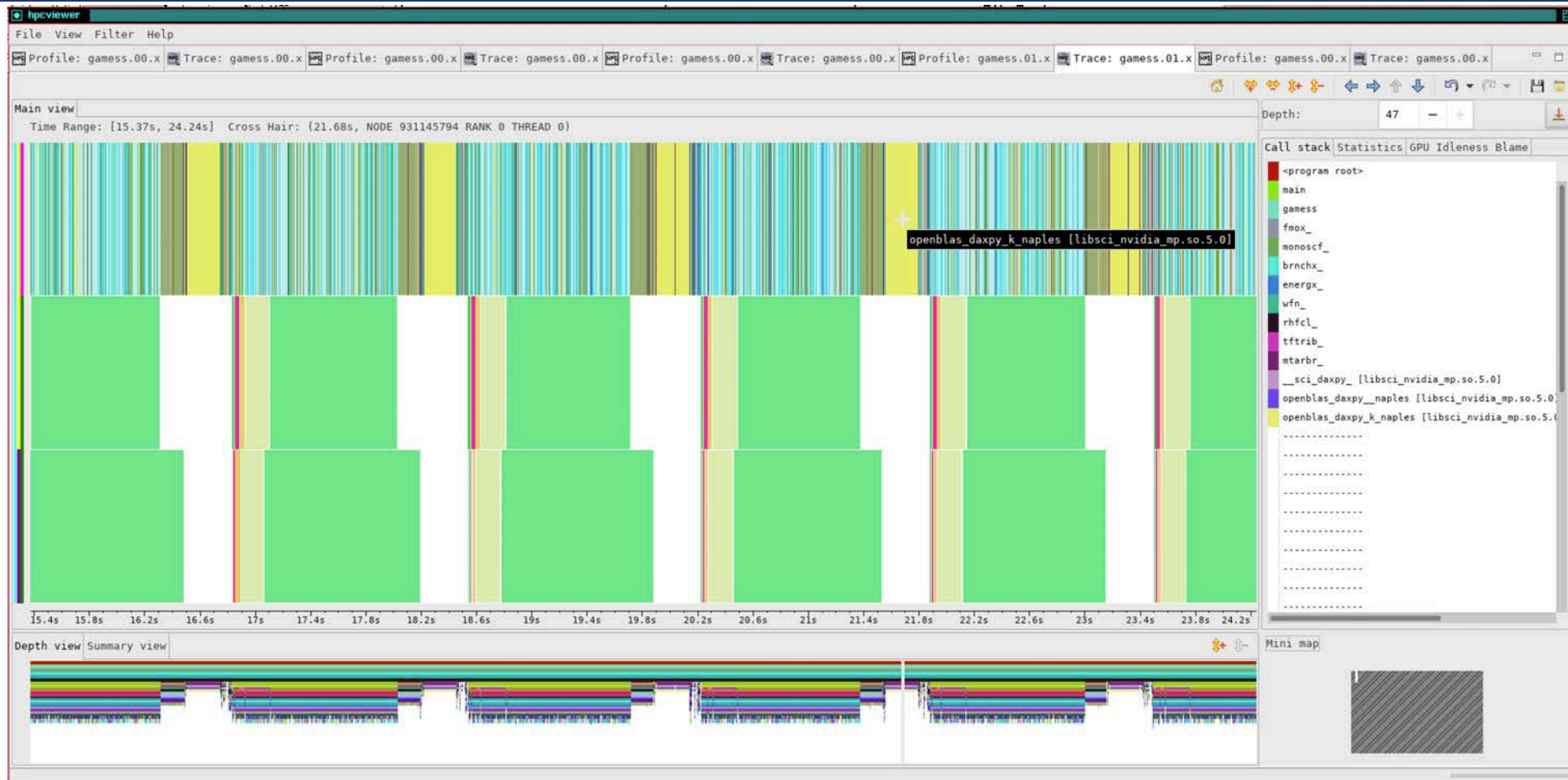
Time-centric Analysis: GAMESS 5 nodes, 40 ranks, 20 GPUs on Perlmutter



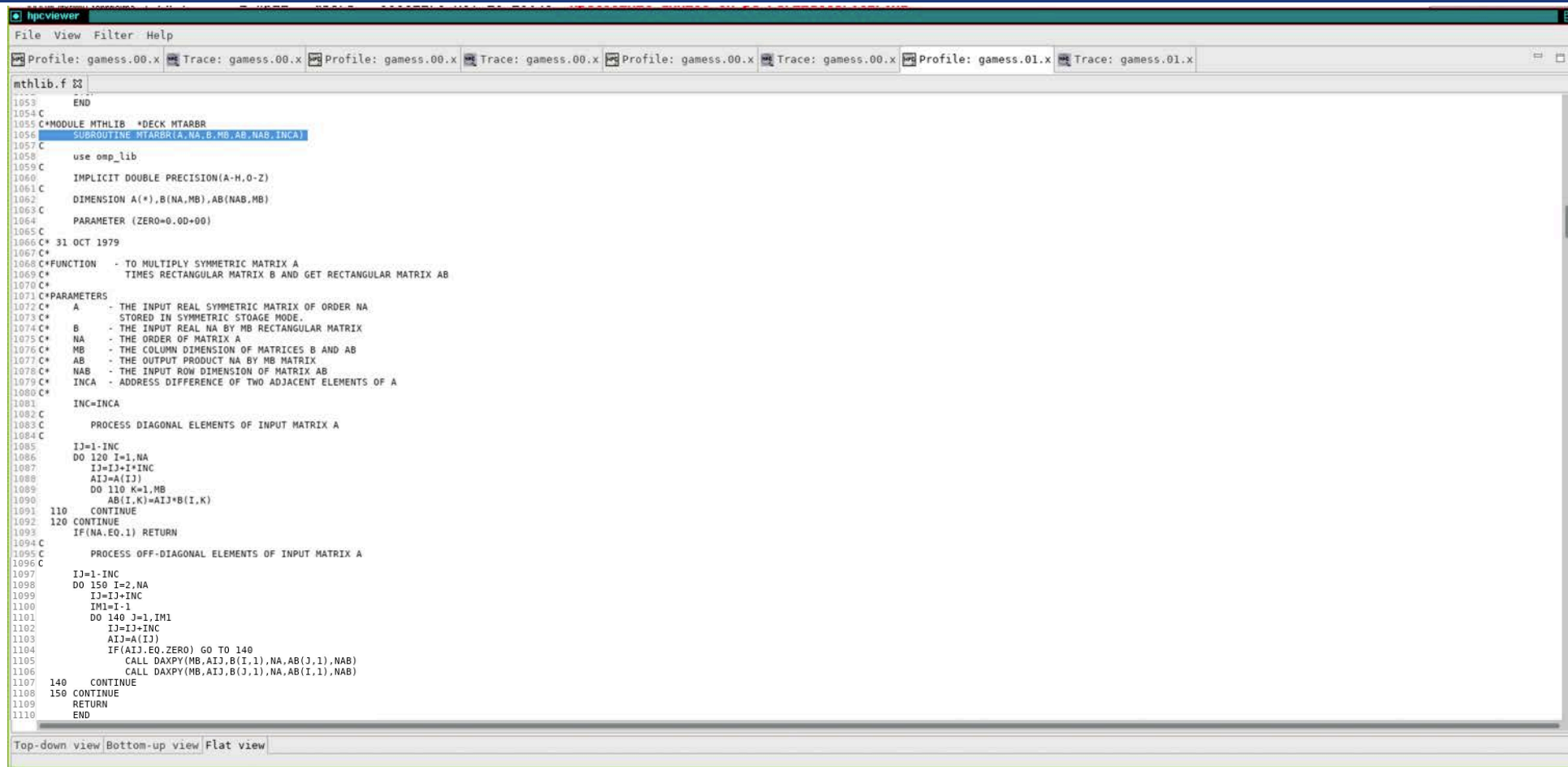
Time-centric Analysis: GAMESS 5 nodes, 40 ranks, 20 GPUs on Perlmutter



Time-centric Analysis: GAMESS 5 nodes, 40 ranks, 20 GPUs on Perlmutter

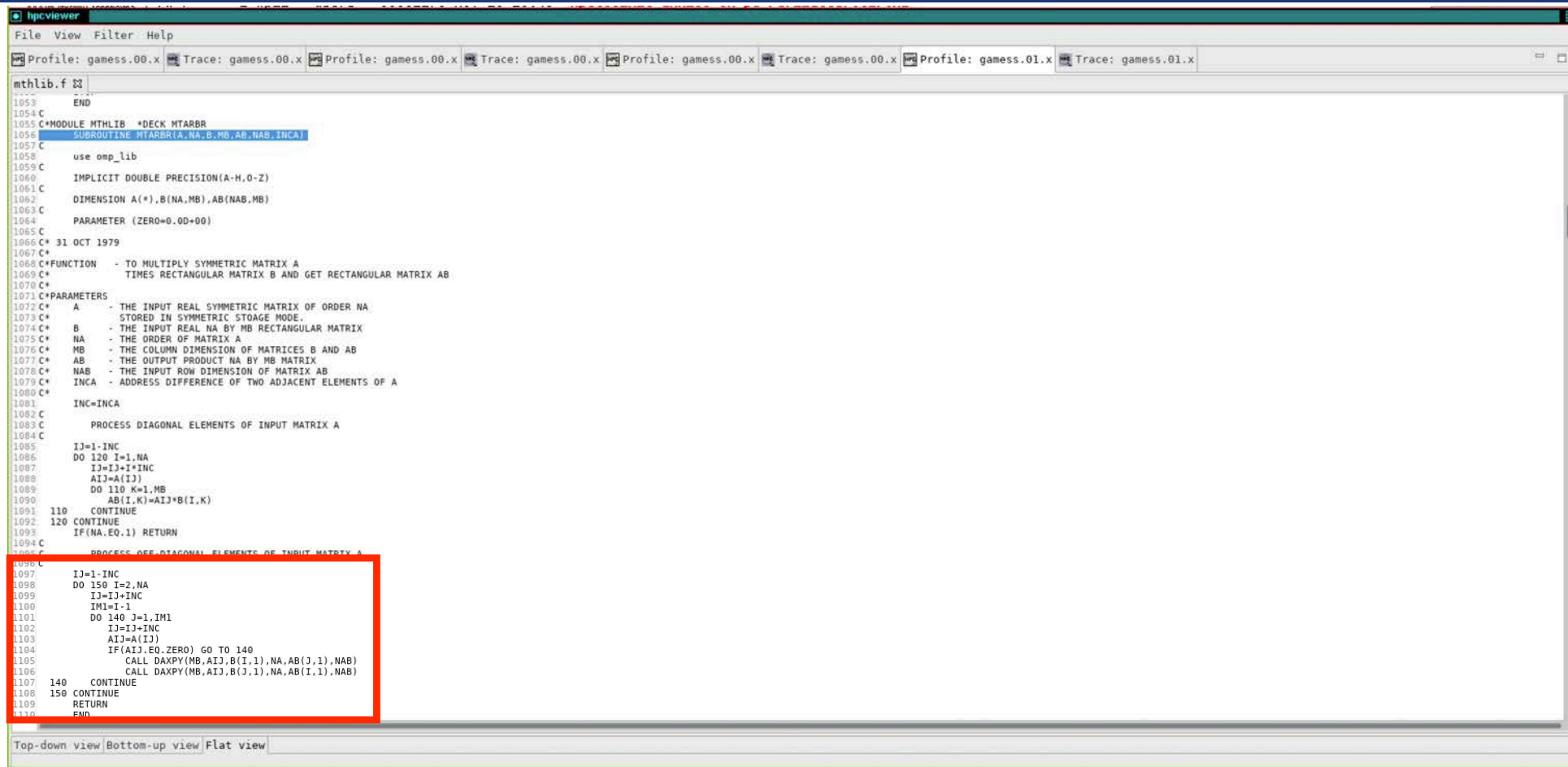


Time-centric Analysis: GAMESS 5 nodes, 40 ranks, 20 GPUs on Perlmutter



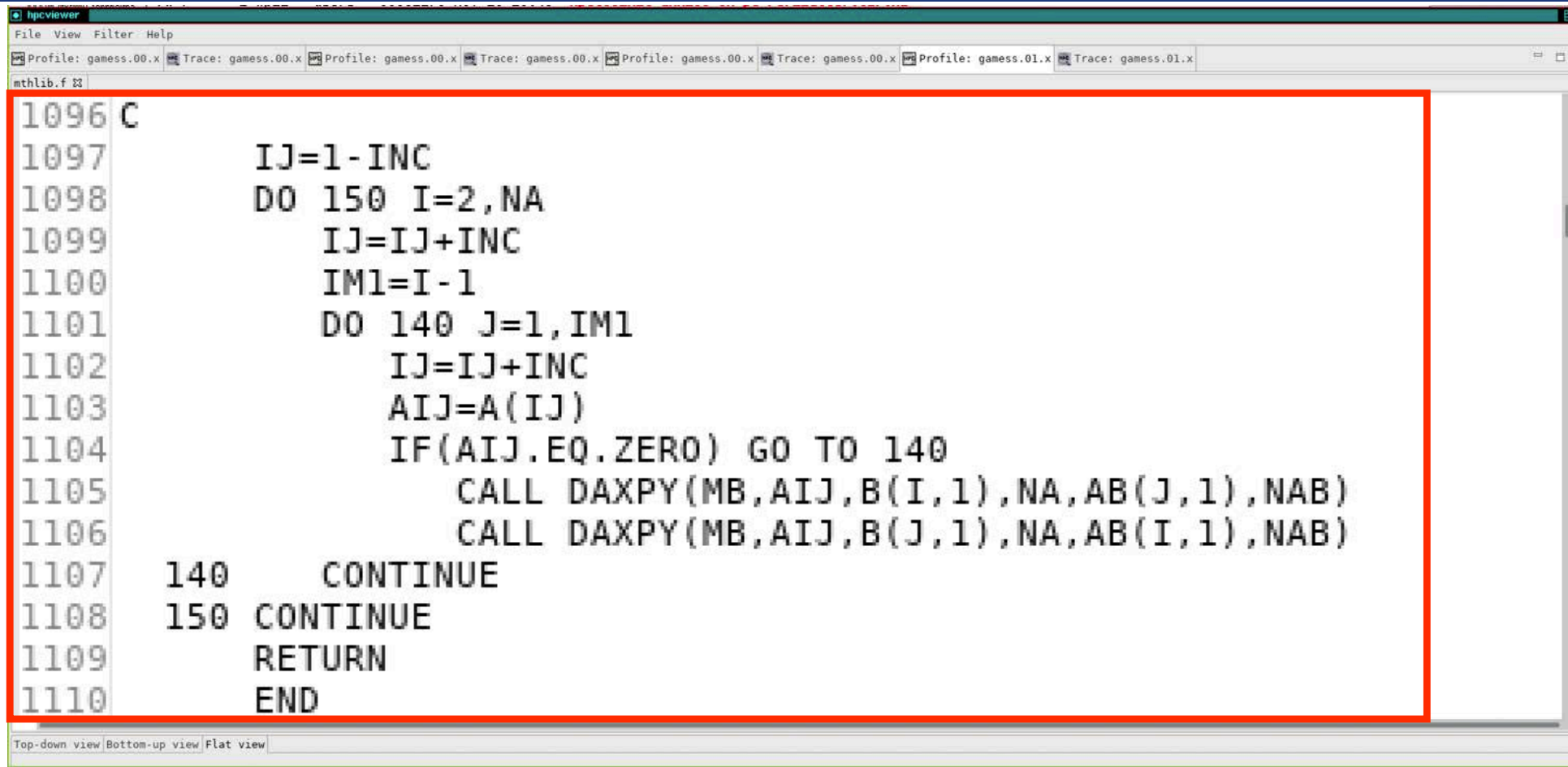
```
hpcviewer
File View Filter Help
Profile: gamess.00.x Trace: gamess.00.x Profile: gamess.00.x Trace: gamess.00.x Profile: gamess.00.x Trace: gamess.00.x Profile: gamess.01.x Trace: gamess.01.x
mthlib.f
1053 END
1054 C
1055 C*MODULE MTHLIB *DECK MTARBR
1056 SUBROUTINE MTARBR(A,NA,B,MB,AB,NAB,INCA)
1057 C
1058 use omp_lib
1059 C
1060 IMPLICIT DOUBLE PRECISION(A-H,O-Z)
1061 C
1062 DIMENSION A(*),B(NA,MB),AB(NAB,MB)
1063 C
1064 PARAMETER (ZERO=0.0D+00)
1065 C
1066 C* 31 OCT 1979
1067 C*
1068 C*FUNCTION - TO MULTIPLY SYMMETRIC MATRIX A
1069 C* TIMES RECTANGULAR MATRIX B AND GET RECTANGULAR MATRIX AB
1070 C*
1071 C*PARAMETERS
1072 C* A - THE INPUT REAL SYMMETRIC MATRIX OF ORDER NA
1073 C* STORED IN SYMMETRIC STORAGE MODE.
1074 C* B - THE INPUT REAL NA BY MB RECTANGULAR MATRIX
1075 C* NA - THE ORDER OF MATRIX A
1076 C* MB - THE COLUMN DIMENSION OF MATRICES B AND AB
1077 C* AB - THE OUTPUT PRODUCT NA BY MB MATRIX
1078 C* NAB - THE INPUT ROW DIMENSION OF MATRIX AB
1079 C* INCA - ADDRESS DIFFERENCE OF TWO ADJACENT ELEMENTS OF A
1080 C*
1081 INC=INCA
1082 C
1083 C PROCESS DIAGONAL ELEMENTS OF INPUT MATRIX A
1084 C
1085 IJ=1-INC
1086 DO 120 I=1,NA
1087 IJ=IJ+INC
1088 AIJ=A(IJ)
1089 DO 110 K=1,MB
1090 AB(I,K)=AIJ*B(I,K)
1091 110 CONTINUE
1092 120 CONTINUE
1093 IF(NA.EQ.1) RETURN
1094 C
1095 C PROCESS OFF-DIAGONAL ELEMENTS OF INPUT MATRIX A
1096 C
1097 IJ=1-INC
1098 DO 150 I=2,NA
1099 IJ=IJ+INC
1100 IM1=I-1
1101 DO 140 J=1,IM1
1102 IJ=IJ-INC
1103 AIJ=A(IJ)
1104 IF(AIJ.EQ.ZERO) GO TO 140
1105 CALL DAXPY(MB,AIJ,B(I,1),NA,AB(J,1),NAB)
1106 CALL DAXPY(MB,AIJ,B(J,1),NA,AB(I,1),NAB)
1107 140 CONTINUE
1108 150 CONTINUE
1109 RETURN
1110 END
Top-down view Bottom-up view Flat view
```

Time-centric Analysis: GAMESS 5 nodes, 40 ranks, 20 GPUs on Perlmutter



```
hpcviewer
File View Filter Help
Profile: gamess.00.x Trace: gamess.00.x Profile: gamess.00.x Trace: gamess.00.x Profile: gamess.00.x Trace: gamess.00.x Profile: gamess.01.x Trace: gamess.01.x
mthlib.f
1053 END
1054 C
1055 C*MODULE MTHLIB *DECK MTARBR
1056 SUBROUTINE MTARBR(A,NA,B,MB,AB,NAB,INCA)
1057 C
1058 use omp_lib
1059 C
1060 IMPLICIT DOUBLE PRECISION(A-H,O-Z)
1061 C
1062 DIMENSION A(*),B(NA,MB),AB(NAB,MB)
1063 C
1064 PARAMETER (ZERO=0.0D+00)
1065 C
1066 C* 31 OCT 1979
1067 C*
1068 C*FUNCTION - TO MULTIPLY SYMMETRIC MATRIX A
1069 C* TIMES RECTANGULAR MATRIX B AND GET RECTANGULAR MATRIX AB
1070 C*
1071 C*PARAMETERS
1072 C* A - THE INPUT REAL SYMMETRIC MATRIX OF ORDER NA
1073 C* STORED IN SYMMETRIC STORAGE MODE.
1074 C* B - THE INPUT REAL NA BY MB RECTANGULAR MATRIX
1075 C* NA - THE ORDER OF MATRIX A
1076 C* MB - THE COLUMN DIMENSION OF MATRICES B AND AB
1077 C* AB - THE OUTPUT PRODUCT NA BY MB MATRIX
1078 C* NAB - THE INPUT ROW DIMENSION OF MATRIX AB
1079 C* INCA - ADDRESS DIFFERENCE OF TWO ADJACENT ELEMENTS OF A
1080 C*
1081 INC=INCA
1082 C
1083 C PROCESS DIAGONAL ELEMENTS OF INPUT MATRIX A
1084 C
1085 IJ=1-INC
1086 DO 120 I=1,NA
1087 IJ=IJ+INC
1088 AIJ=A(IJ)
1089 DO 110 K=1,MB
1090 AB(I,K)=AIJ*B(I,K)
1091 110 CONTINUE
1092 120 CONTINUE
1093 IF(NA.EQ.1) RETURN
1094 C
1095 C PROCESS OFF-DIAGONAL ELEMENTS OF INPUT MATRIX A
1096 C
1097 IJ=1-INC
1098 DO 150 I=2,NA
1099 IJ=IJ+INC
1100 IM1=I-1
1101 DO 140 J=1,IM1
1102 IJ=IJ-INC
1103 AIJ=A(IJ)
1104 IF(AIJ.EQ.ZERO) GO TO 140
1105 CALL DAXPY(MB,AIJ,B(I,1),NA,AB(J,1),NAB)
1106 CALL DAXPY(MB,AIJ,B(J,1),NA,AB(I,1),NAB)
1107 140 CONTINUE
1108 150 CONTINUE
1109 RETURN
1110 ENDO
Top-down view Bottom-up view Flat view
```

Time-centric Analysis: GAMESS 5 nodes, 40 ranks, 20 GPUs on Perlmutter



The image shows a screenshot of the hpcviewer application. The main window displays the source code for a Fortran file named 'mthlib.f'. The code is enclosed in a red rectangular border. The code consists of several lines of Fortran, including a loop structure with nested DO loops and a RETURN statement. The code is as follows:

```
1096 C
1097     IJ=1-INC
1098     DO 150 I=2,NA
1099         IJ=IJ+INC
1100         IM1=I-1
1101         DO 140 J=1,IM1
1102             IJ=IJ+INC
1103             AIJ=A(IJ)
1104             IF(AIJ.EQ.ZERO) GO TO 140
1105                 CALL DAXPY(MB,AIJ,B(I,1),NA,AB(J,1),NAB)
1106                 CALL DAXPY(MB,AIJ,B(J,1),NA,AB(I,1),NAB)
1107     140     CONTINUE
1108     150 CONTINUE
1109         RETURN
1110     END
```

At the bottom of the hpcviewer window, there are view options: 'Top-down view', 'Bottom-up view', and 'Flat view'.

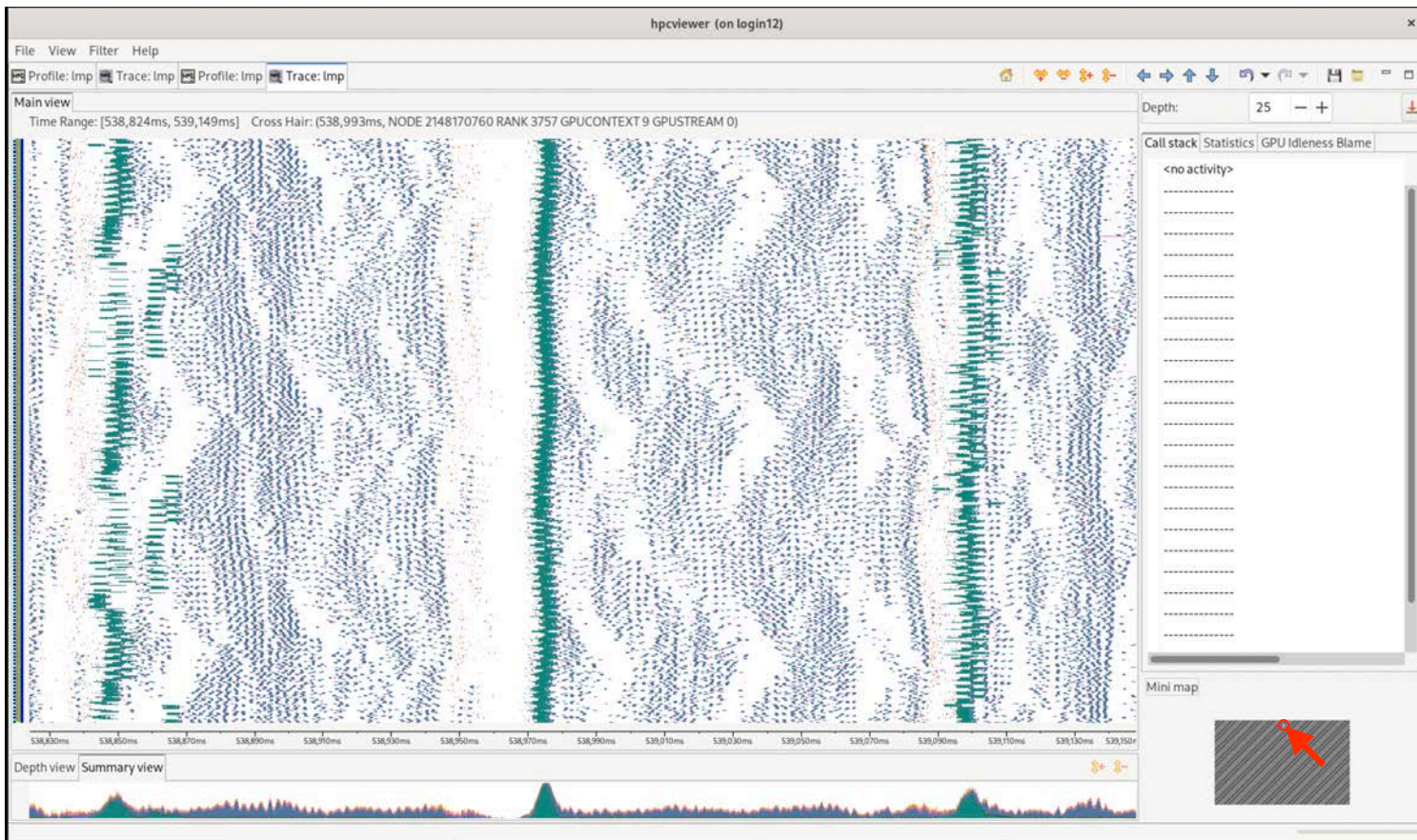
Measure and Attribute OpenMP Offloading

The screenshot shows the hpcviewer interface. The top pane displays C++ code from `einspline_spo_omp.cpp`. Line 162 contains a `PRAGMA OFFLOAD` directive: `PRAGMA OFFLOAD("omp target teams distribute num_teams(NumTeams) thread_limit(ChunkSizePerTeam)")`. The code below it shows a loop over `team_id` with nested calculations for `first` and `last` indices, and a conditional expression involving `nSplinesPerBlock_local`.

The bottom pane shows a performance metrics table with columns for Scope, GPUOP (sec):Sum (I), GPUOP (sec):Sum (E), GKER (sec):Sum (I), GKER (sec):Sum (E), and GMEM (sec):Sum (I). The table is expanded to show nested scopes, including `<program root>`, `main`, and various loops and function calls. The `<omp tgt kernel>` scope is highlighted in blue.

Scope	GPUOP (sec):Sum (I)	GPUOP (sec):Sum (E)	GKER (sec):Sum (I)	GKER (sec):Sum (E)	GMEM (sec):Sum (I)
Experiment Aggregate Metrics	6.23e+00 100.0%	6.23e+00 100.0%	5.50e+00 100.0%	5.50e+00 100.0%	7.10e-03 100.0%
<program root>	6.23e+00 100.0%		5.50e+00 100.0%		7.10e-03 100.0%
main	6.23e+00 100.0%		5.50e+00 100.0%		7.10e-03 100.0%
loop at miniqmc.cpp: 402	6.02e+00 96.6%		5.33e+00 96.9%		
404 => .omp_outlined..62	6.02e+00 96.6%		5.33e+00 96.9%		
404 => [I] .omp_outlined_debug_.61	6.02e+00 96.6%		5.33e+00 96.9%		
loop at miniqmc.cpp: 405	6.02e+00 96.6%		5.33e+00 96.9%		
loop at miniqmc.cpp: 472	5.01e+00 80.5%		4.49e+00 81.6%		
loop at miniqmc.cpp: 476	5.01e+00 80.5%		4.49e+00 81.6%		
loop at miniqmc.cpp: 478	5.01e+00 80.5%		4.49e+00 81.6%		
485 => qmcplusplus::WaveFunction::ratio(q...	5.01e+00 80.5%		4.49e+00 81.6%		
qmcplusplus::DiracDeterminant<qmcplus...	5.01e+00 80.5%		4.49e+00 81.6%		
198 => qmcplusplus::einspline_spo_omp<...	5.01e+00 80.5%		4.49e+00 81.6%		
187 => qmcplusplus::einspline_spo_omp...	5.01e+00 80.5%		4.49e+00 81.6%		
loop at einspline_spo_omp.cpp: 157	5.01e+00 80.5%		4.49e+00 81.6%		
162 => <omp tgt kernel>	4.49e+00 72.0%		4.49e+00 81.6%		
<gpu kernel>	4.49e+00 72.0%	4.49e+00 72.0%	4.49e+00 81.6%	4.49e+00 81.6%	
unknown files: libhpcrun_solv_0	4.49e+00 72.0%	4.49e+00 72.0%	4.49e+00 81.6%	4.49e+00 81.6%	

LAMMPS on Frontier: 8K nodes, 64K MPI ranks + GPU times



HPCToolkit Status on GPUs

NVIDIA

- heterogeneous profiles

- GPU instruction-level execution and stalls using PC sampling

- traces

AMD

- heterogeneous profiles

- no GPU instruction-level measurements within kernels

- measure OpenMP offloading using OMPT interface

- traces

Intel

- heterogeneous profiles

- GPU instruction-level measurements with instrumentation; heuristic latency attribution to instructions

- measure OpenMP offloading using OMPT interface

- traces

Ongoing Work

Enhancing measurement to identify root causes of scalability losses

- Better measurement of delays caused by GPU and communication

Improving the scalability of hpcprof-mpi

- Avoid unnecessary serialization of I/O

Adding a Python-based interface for analysis of performance results

- Python API supports arbitrary queries and analysis of profiles and traces

- Automatic analysis to identify notable features in executions

 - e.g. load imbalance, trace line equivalence classes

HPCToolkit Resources

Documentation

User manual

<http://hpctoolkit.org/manual/HPCToolkit-users-manual.pdf>

Tutorial videos

<http://hpctoolkit.org/training.html>

Software

Download hpcviewer GUI binaries for your laptop, desktop, cluster, or supercomputer

OS: Linux, Windows, MacOS

Processors: x86_64, aarch64, ppc64le

<http://hpctoolkit.org/download.html>

Install HPCToolkit on your Linux desktop, cluster, or supercomputer using Spack

<http://hpctoolkit.org/software-instructions.html>

ARGONNE
ATPESC2023
EXTREME - SCALE COMPUTING

HPCToolkit Hands-On Directions

Performance analysis of CPU and GPU-accelerated applications at Scale

John Mellor-Crummey
Professor, Rice University

extremecomputingtraining.anl.gov



EXTREME
COMPUTING
PROJECT



Sample Performance Databases for You to Analyze on Polaris

- Setup on polaris
 - module use /soft/perftools/hpctoolkit/polaris/modulefiles
 - module load hpctoolkit/default
- Data on theta and polaris: /grand/ATPESC2023/track-6-tools-hpctoolkit/data
 - CPU
 - QMCPACK - quantum Monte Carlo electronic structure calculations (early experiment)
 - GPU-accelerated
 - GAMESS - ab initio quantum chemistry
 - 1.singlegroup-unbalanced
 - 2.singlegroup-balanced
 - 3.multigroup-unbalanced-mtarbr
 - 4.multigroup-balanced
 - 5.multigroup-unbalanced-pc
 - 6.scale
 - PeleC - AMR Solver (AMREx) for compressible reacting flows
 - Pytorch-deepwave - GPU-accelerated reverse-time migration using Pytorch
 - Quicksilver - proxy application for dynamic Monte Carlo transport

GPU: Profiling Quicksilver with HPCToolkit on Polaris or Perlmutter

- git clone <https://github.com/hpctoolkit/hpctoolkit-tutorial-examples>
- cd hpctoolkit-tutorial-examples/examples/gpu/quicksilver
- polaris:
 - export HPCTOOLKIT_TUTORIAL_PROJECTID=ATPESC2023
 - export HPCTOOLKIT_TUTORIAL_RESERVATION=default
 - source setup-env/polaris.sh
- perlmutter:
 - export HPCTOOLKIT_TUTORIAL_PROJECTID=ntrain5_g
 - export HPCTOOLKIT_TUTORIAL_RESERVATION=default
 - source setup-env/perlmutter.sh
- make build
- make run
- make run-pc
- make view
- make view-pc

CPU: Profiling AMG2013 with HPCToolkit on Theta

- `git clone https://github.com/hpctoolkit/hpctoolkit-tutorial-examples`
- `cd hpctoolkit-tutorial-examples/examples/cpu/mpi+openmp/amg2013`
- `export HPCTOOLKIT_TUTORIAL_PROJECTID=ATPESC2023`
- `export HPCTOOLKIT_TUTORIAL_RESERVATION=debug-cache-quad`
- `source setup-env/theta.sh`
- `make build`
- `make run`
 - `# wait for $COBALT_JOBID.done to appear in your directory`
- `make analyze`
- Alternatives
 - `make view`
 - `hpcviewer hpctoolkit-amg2013.d`