# HPCToolkit: Performance Tools for GPU-Accelerated Computing



John Mellor-Crummey Rice University

December 14, 2022











#### HPCToolkit Funding Acknowledgments

- Government
  - Exascale Computing Project 17-SC-20-SC
  - Lawrence Livermore National Laboratory Subcontract B645220
  - Argonne National Laboratory Subcontract 9F-60073
- Corporate
  - Advanced Micro Devices
  - Intel Corporation
  - TotalEnergies EP Research & Technology USA, LLC.



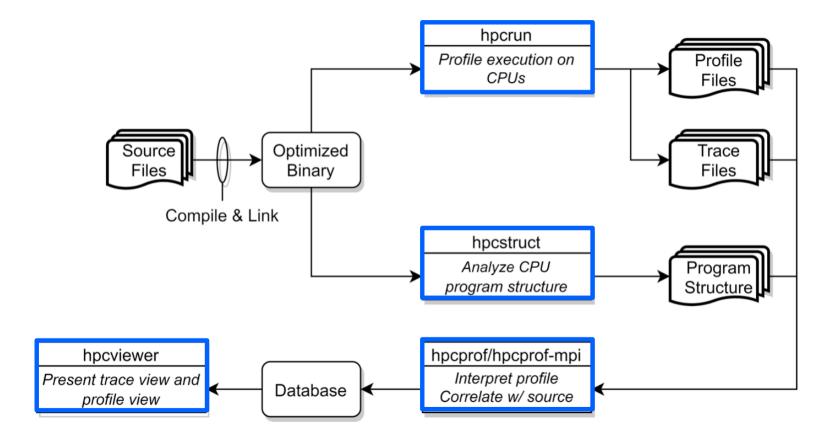
#### Rice University's HPCToolkit Performance Tools

#### Measure and analyze performance of CPU and GPU-accelerated applications

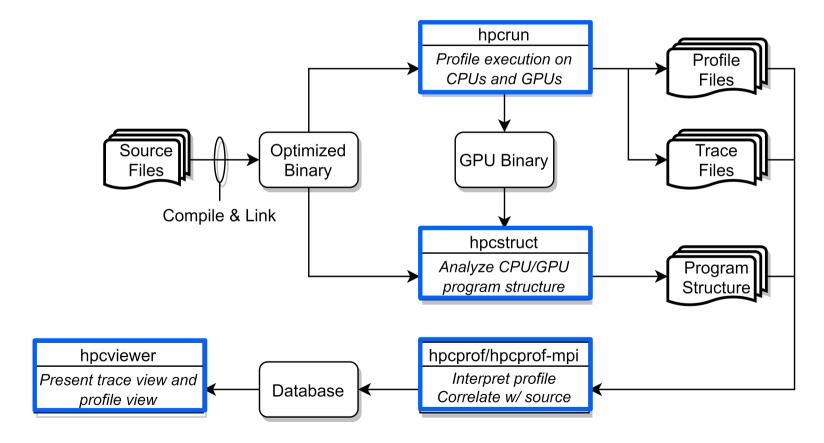
- Easy: profile unmodified application binaries
- Fast: low-overhead measurement
- Informative: understand where an application spends its time and why
  - call path profiles associate metrics with application source code contexts
  - optional hierarchical traces to understand execution dynamics
- Broad audience
  - application developers
  - framework developers
  - runtime and tool developers



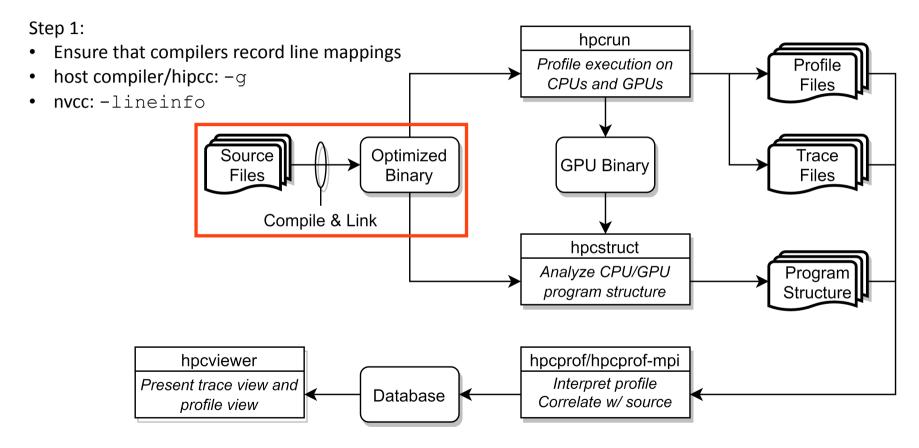
# HPCToolkit's Workflow for CPU Applications













#### Step 2: hpcrun hpcrun collects call path profiles (and Profile execution on **Profile** optionally, traces) of events of interest CPUs and GPUs **Files** Optimized Source Trace **GPU Binary** Files Binary **Files** Compile & Link hpcstruct Analyze CPU/GPU Program program structure Structure hpcprof/hpcprof-mpi hpcviewer Interpret profile Present trace view and Database Correlate w/ source profile view



#### Measurement of CPU and GPU-accelerated Applications

- Sampling using timers and hardware counter overflow on the CPU
- Callbacks when GPU operations are launched and (sometimes) completed
- GPU event stream for GPU operations; PC Samples (NVIDIA)



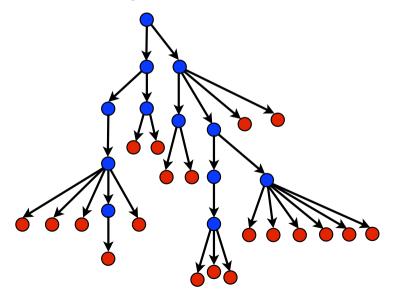
#### Call Stack Unwinding to Attribute Costs in Context

- Unwind when timer or hardware counter overflows
  - measurement overhead proportional to sampling frequency rather than call frequency
- Unwind to capture context for events such as GPU kernel launches

#### Call path sample

return address
return address
return address
instruction pointer

#### Calling context tree





#### 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 qpu=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



#### Step 3: hpcrun hpcstruct recovers program structure Profile execution on **Profile** about lines, loops, and inlined functions CPUs and GPUs **Files** Optimized Trace **GPU Binary** Files Binary Files Compile & Link hpcstruct Analyze CPU/GPU Program program structure Structure hpcprof/hpcprof-mpi hpcviewer Interpret profile Present trace view and Database Correlate w/ source



profile view

## hpcstruct: Analyze CPU and GPU Binaries Using Multiple Threads

Usage

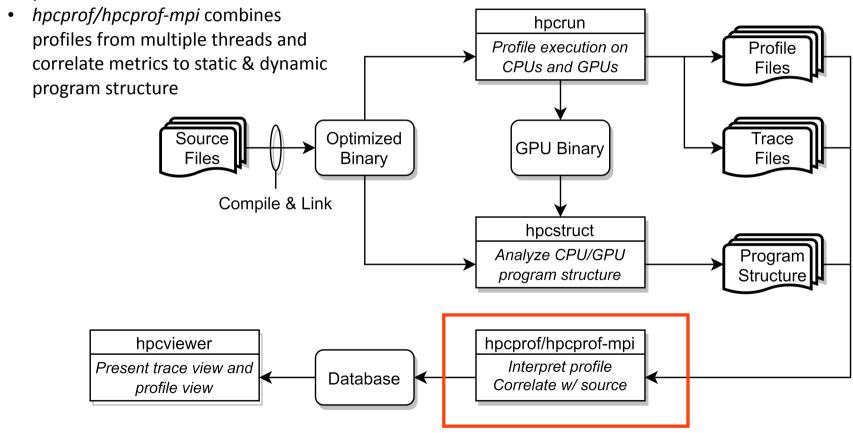
```
hpcstruct [--gpucfg yes] <measurement-directory>
```

- 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 size(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



#### Step 4:





#### hpcprof/hpcprof-mpi: Associate Measurements with Program Structure

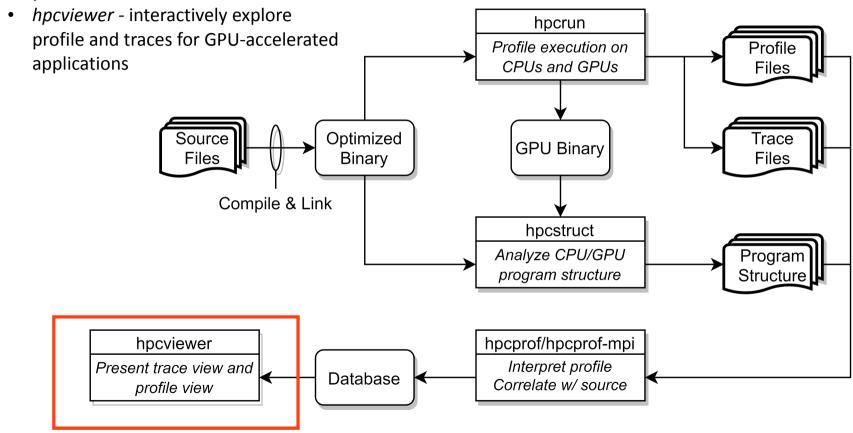
Analyze data from modest executions sequentially
 hpcprof <measurement-directory>

Analyze data from large executions in parallel

```
jsrun -n 2 -a 1 -c 22 -b packed hpcprof-mpi <measurement-directory> srun -N 2 -n 2 -c 126 hpcprof-mpi <measurement-directory>
```

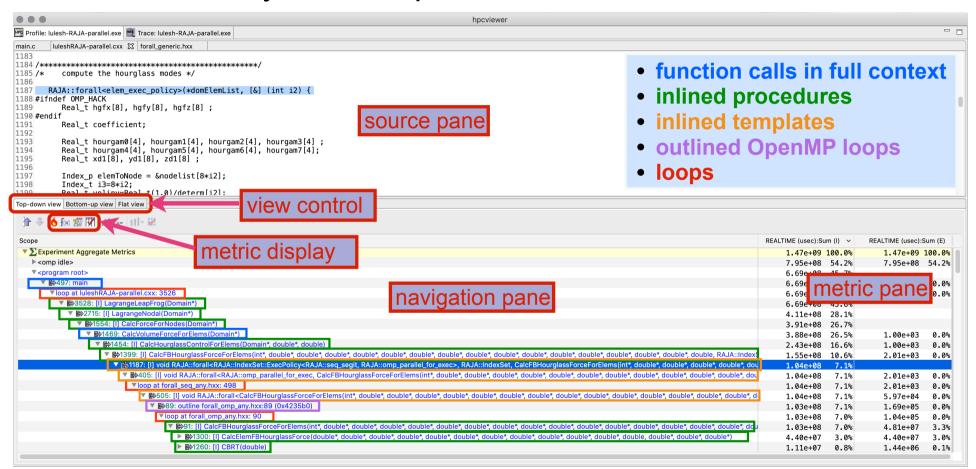


#### Step 4:





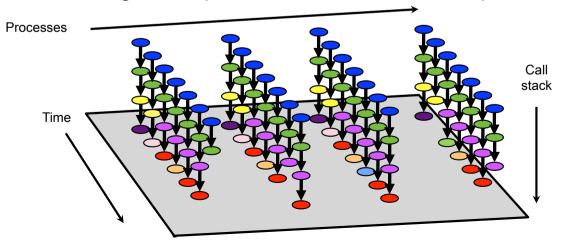
#### Code-centric Analysis with hpcviewer





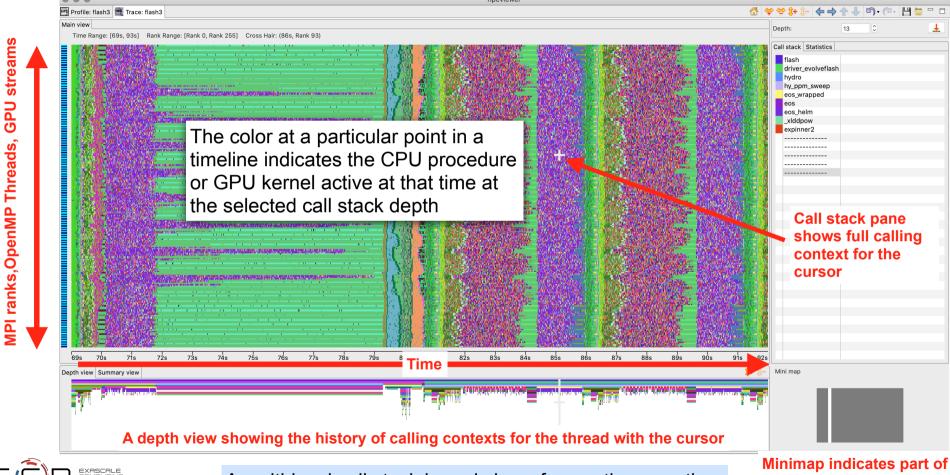
#### **Understanding Temporal Behavior**

- 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

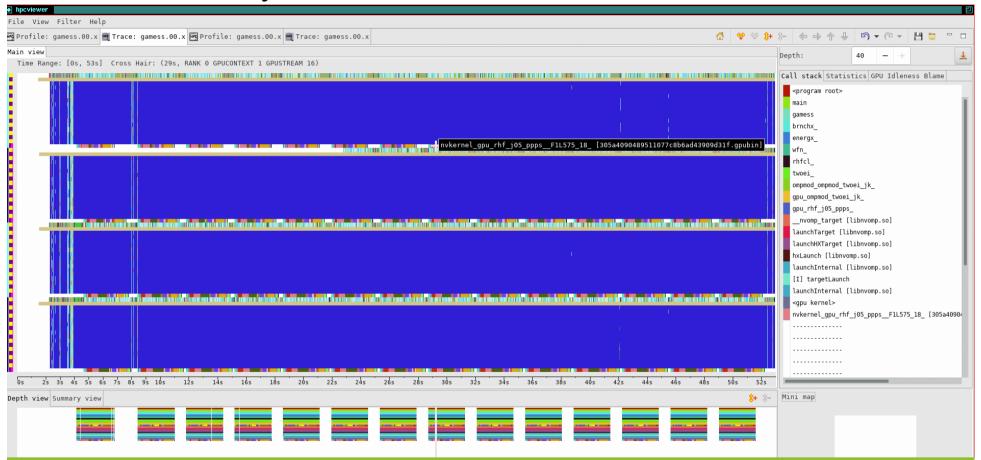




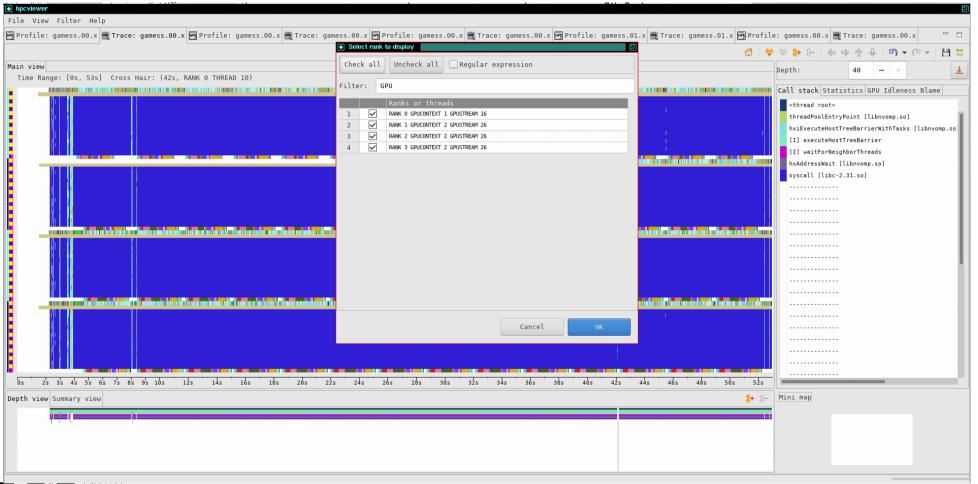
## Case Study: GAMESS

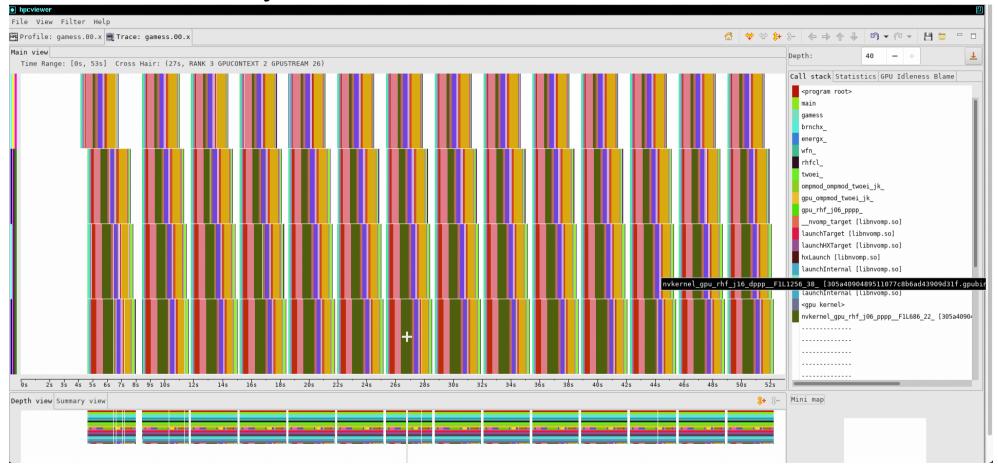
- General Atomic and Molecular Electronic Structure System (GAMESS)
  - general ab initio quantum chemistry package
- Calculates the energies, structures, and properties of a wide range of chemical systems
- Experiments
  - GPU-accelerated nodes at a Perlmutter hackathon
    - Single node with 4 GPUs
    - Five nodes with 20 GPUs
  - GPU-accelerated node on Crusher







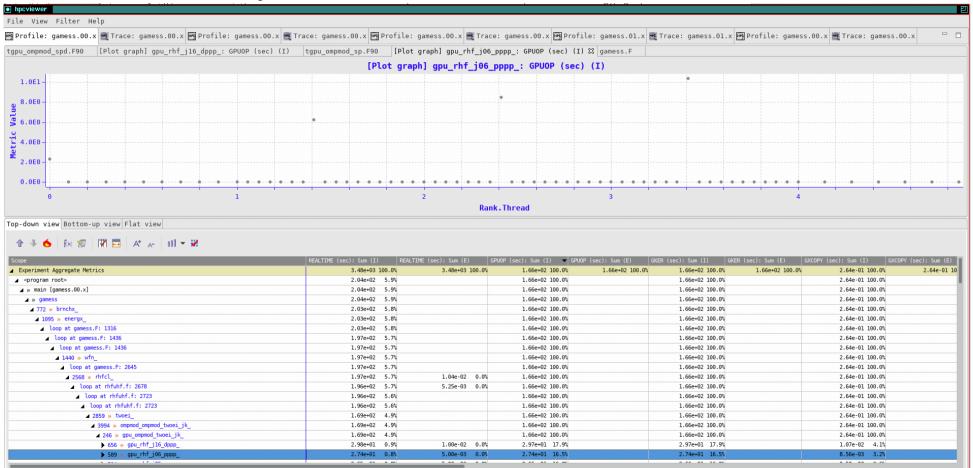




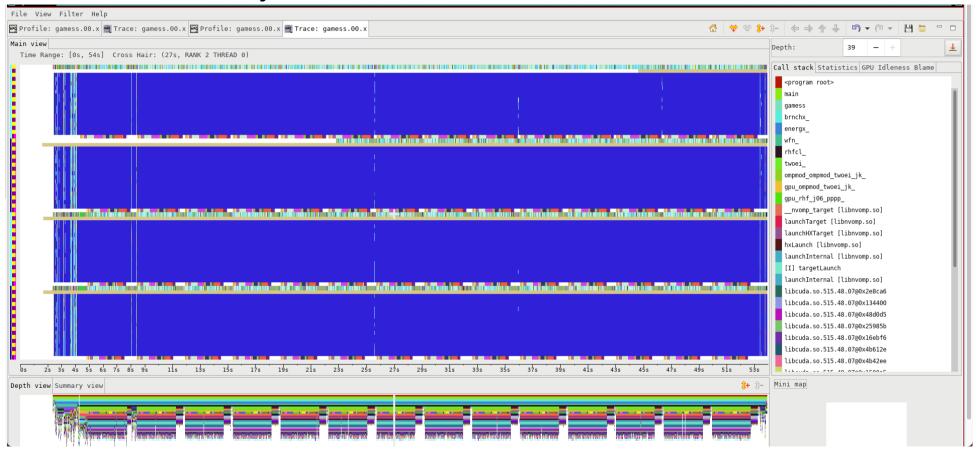




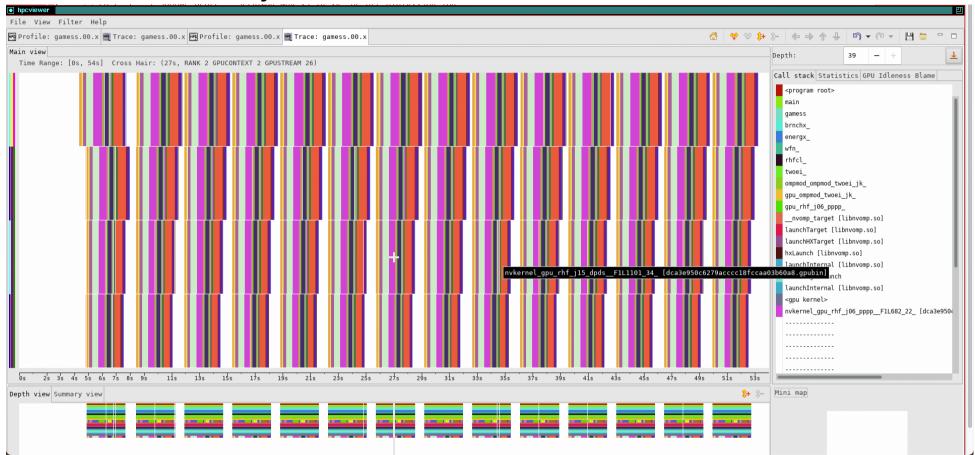




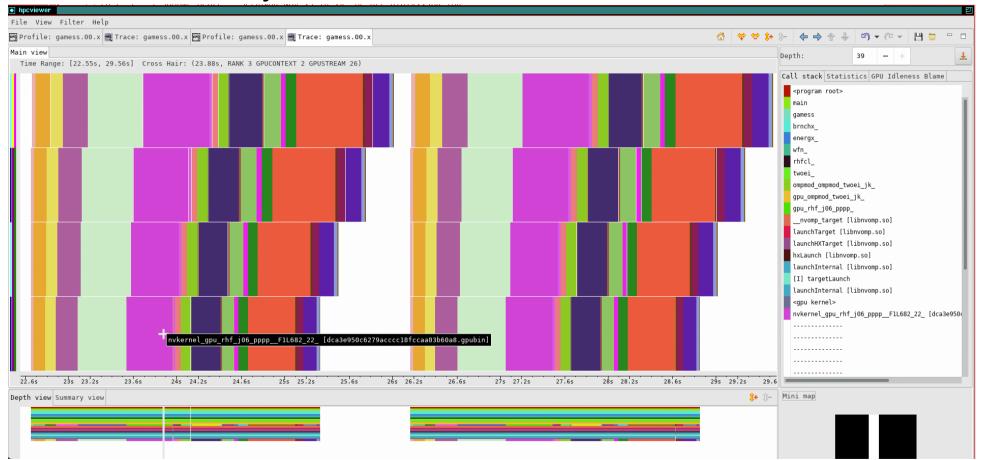




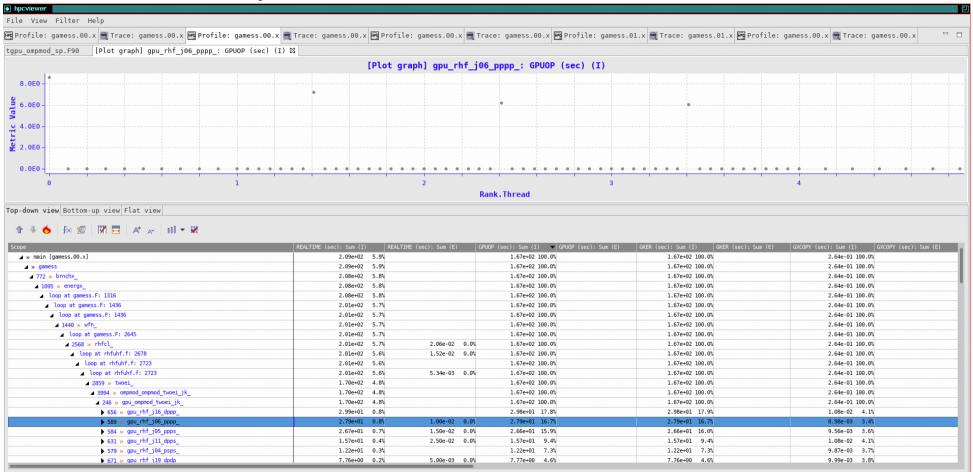




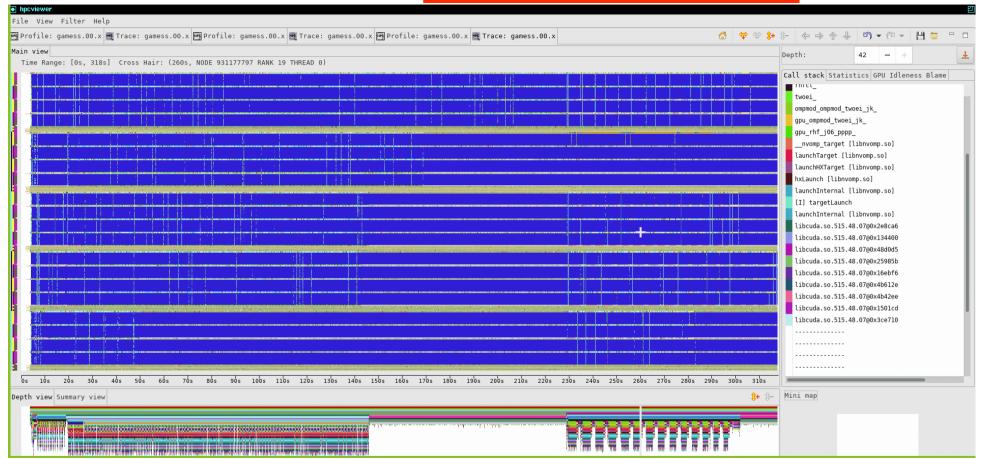




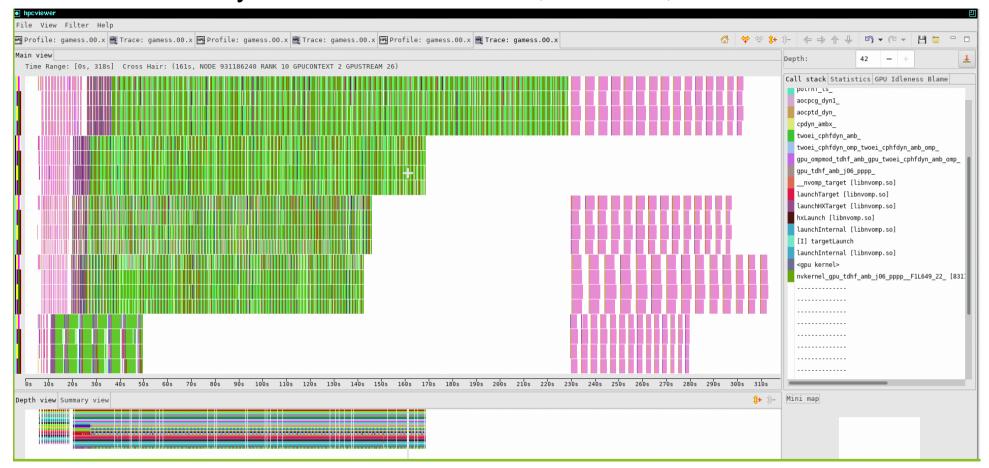




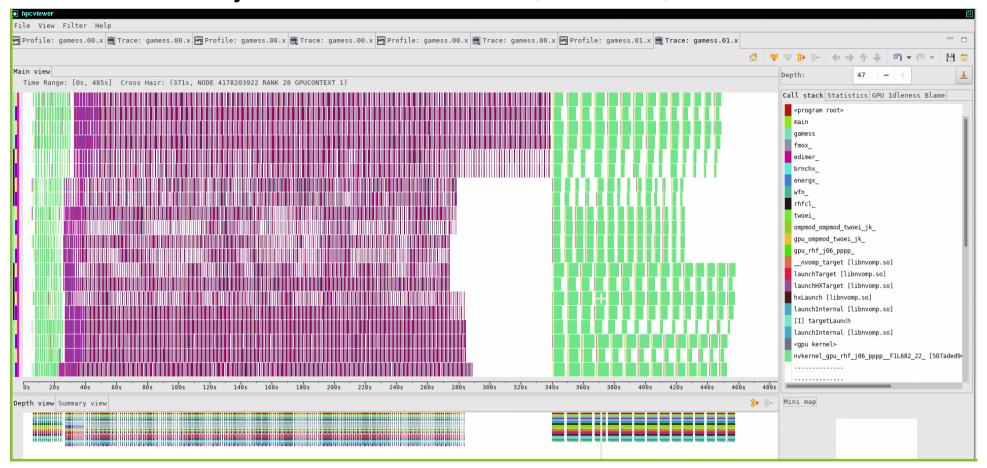




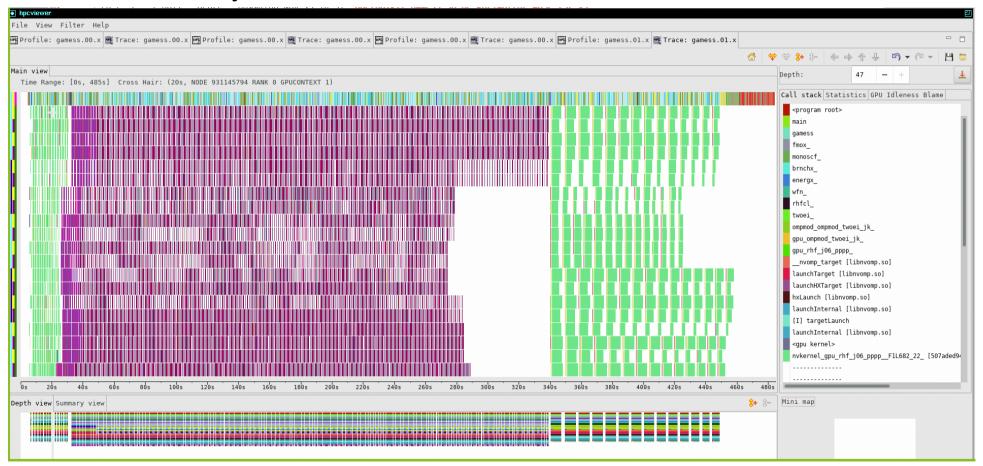




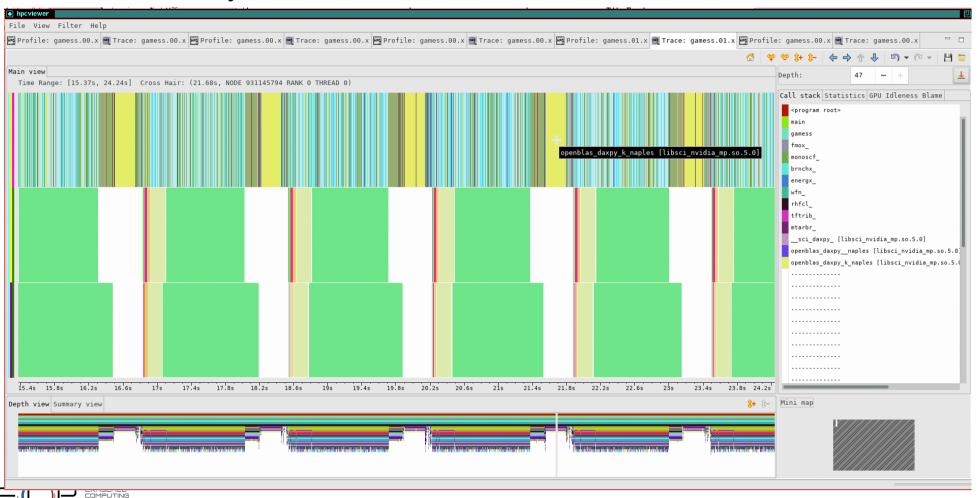










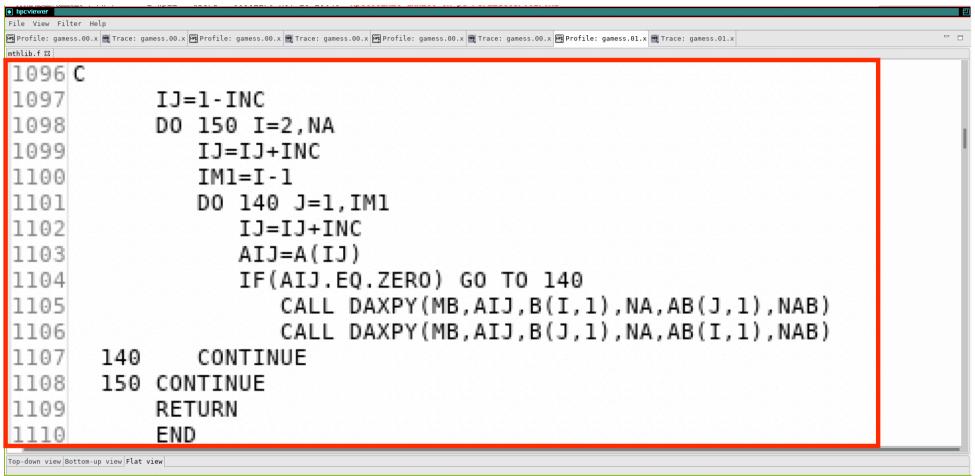


```
File View Filter Help
🕝 Profile: gamess.00.x 🖻 Trace: gamess.00.x 🖻 Profile: gamess.00.x 🖻 Profile: gamess.00.x 🖻 Trace: gamess.00.x 🖻 Trace: gamess.00.x
mthlib.f ₩
 1055 C*MODULE MTHLIB *DECK MTARBR
 1060
          IMPLICIT DOUBLE PRECISION(A-H, 0-Z)
 1061 C
          DIMENSION A(*),B(NA,MB),AB(NAB,MB)
 1063 C
 1064
          PARAMETER (ZERO=0.0D+00)
 1066 C* 31 OCT 1979
 1067 C*
 L068 C*FUNCTION
                - TO MULTIPLY SYMMETRIC MATRIX A
                  TIMES RECTANGULAR MATRIX B AND GET RECTANGULAR MATRIX AB
 1069 C*
 1070 C*
 1071 C*PARAMETERS
 1072 C* A - THE INPUT REAL SYMMETRIC MATRIX OF ORDER NA
                 STORED IN SYMMETRIC STOAGE MODE.
               - THE INPUT REAL NA BY MB RECTANGULAR MATRIX
              - THE ORDER OF MATRIX A
              - THE COLUMN DIMENSION OF MATRICES B AND AB
              - THE OUTPUT PRODUCT NA BY MB MATRIX
         NAB - THE INPUT ROW DIMENSION OF MATRIX AB
         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
          DO 120 I=1.NA
            IJ=IJ+I*INC
 1088
             AIJ=A(IJ)
            DO 110 K=1.MB
               AB(I,K)=AIJ*B(I,K)
 1091 110 CONTINUE
 1092 120 CONTINUE
         IF(NA.EQ.1) RETURN
 1094 C
            PROCESS OFF-DIAGONAL ELEMENTS OF INPUT MATRIX A
          IJ=1-INC
          DO 150 I=2,NA
            IJ=IJ+INC
             IM1=I-1
            DO 140 J=1.IM1
               IJ=IJ+INC
               IF(AIJ.EO.ZERO) GO TO 140
                  CALL DAXPY(MB, AIJ, B(I, 1), NA, AB(J, 1), NAB)
                  CALL DAXPY(MB, AIJ, B(J, 1), NA, AB(I, 1), NAB)
      140
          CONTINUE
 1108 150 CONTINUE
          RETURN
          END
Top-down view Bottom-up view Flat view
```

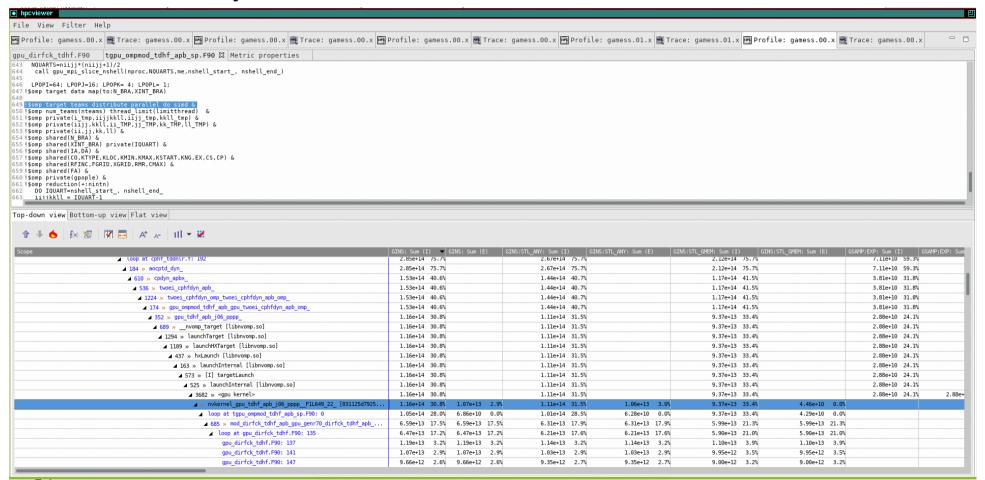


```
File View Filter Help
🕝 Profile: gamess.00.x 🖻 Trace: gamess.00.x 🖻 Profile: gamess.00.x 🖻 Profile: gamess.00.x 🖻 Trace: gamess.00.x 🖻 Trace: gamess.00.x
mthlib.f ₩
 1055 C*MODULE MTHLIB *DECK MTARBR
 1060
          IMPLICIT DOUBLE PRECISION(A-H, 0-Z)
 1061 C
          DIMENSION A(*),B(NA,MB),AB(NAB,MB)
 1063 C
 1064
          PARAMETER (ZERO=0.0D+00)
 1066 C* 31 OCT 1979
 1067 C*
 L068 C*FUNCTION
                - TO MULTIPLY SYMMETRIC MATRIX A
                  TIMES RECTANGULAR MATRIX B AND GET RECTANGULAR MATRIX AB
 1069 C*
 1070 C*
 1071 C*PARAMETERS
 1072 C* A - THE INPUT REAL SYMMETRIC MATRIX OF ORDER NA
 1073 C*
                 STORED IN SYMMETRIC STOAGE MODE.
               - THE INPUT REAL NA BY MB RECTANGULAR MATRIX
              - THE ORDER OF MATRIX A
              - THE COLUMN DIMENSION OF MATRICES B AND AB
              - THE OUTPUT PRODUCT NA BY MB MATRIX
         NAB - THE INPUT ROW DIMENSION OF MATRIX AB
         INCA - ADDRESS DIFFERENCE OF TWO ADJACENT ELEMENTS OF A
 1080 C*
          INC=INCA
 1082 C
 1083 C
            PROCESS DIAGONAL ELEMENTS OF INPUT MATRIX A
 1084 C
          IJ=1-INC
          DO 120 I=1.NA
            IJ=IJ+I*INC
             AIJ=A(IJ)
            DO 110 K=1.MB
               AB(I,K)=AIJ*B(I,K)
 1091 110 CONTINUE
 1092 120 CONTINUE
         IF(NA.EQ.1) RETURN
             DDOCESS DEE-DIAGONAL ELEMENTS DE INDUIT MATDIY /
          IJ=1-INC
          DO 150 I=2,NA
            IJ=IJ+INC
            IM1=I-1
            DO 140 J=1.IM1
               IJ=IJ+INC
               AIJ=A(IJ)
               IF(AIJ.EO.ZERO) GO TO 140
                  CALL DAXPY(MB, AIJ, B(I, 1), NA, AB(J, 1), NAB)
                  CALL DAXPY(MB, AIJ, B(J, 1), NA, AB(I, 1), NAB)
          CONTINUE
      150 CONTINUE
          RETURN
 Top-down view Bottom-up view Flat view
```

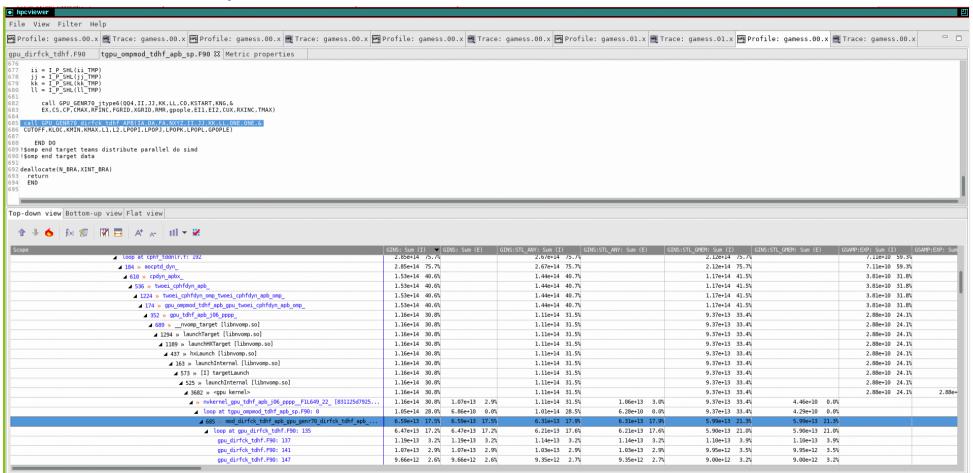




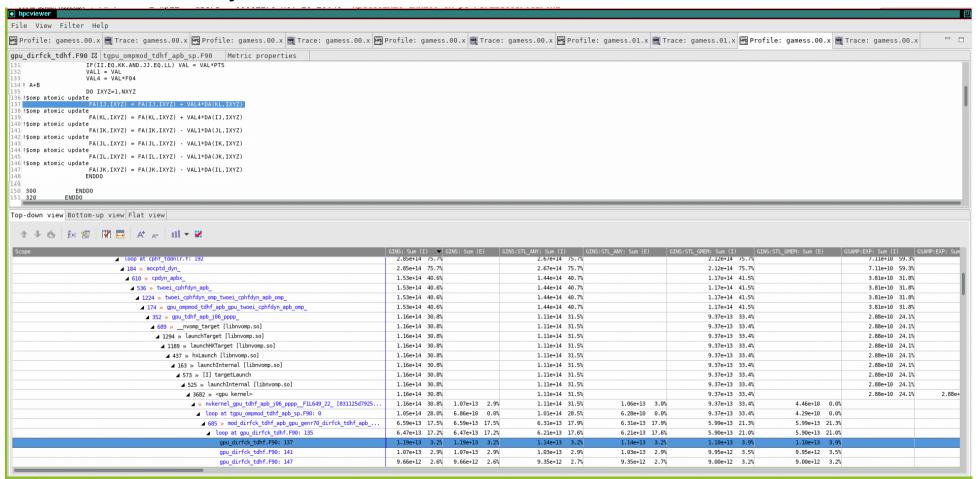




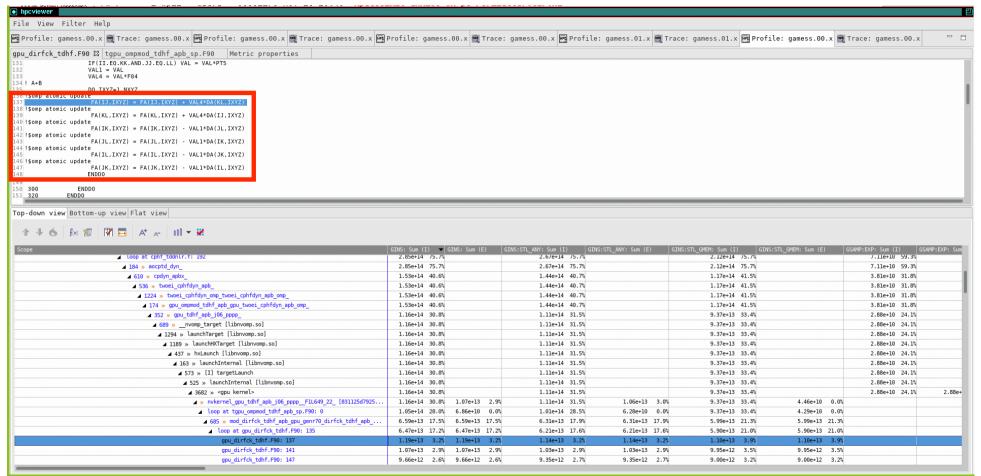




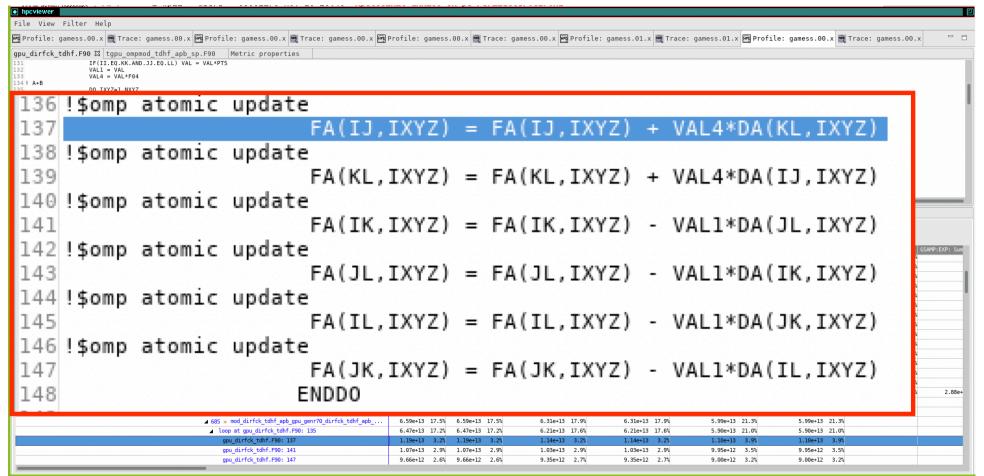














# Time-centric Analysis: GAMESS 16 ranks, 8 GPUs on Crusher





# Time-centric Analysis: GAMESS 16 ranks, 8 GPUs on Crusher



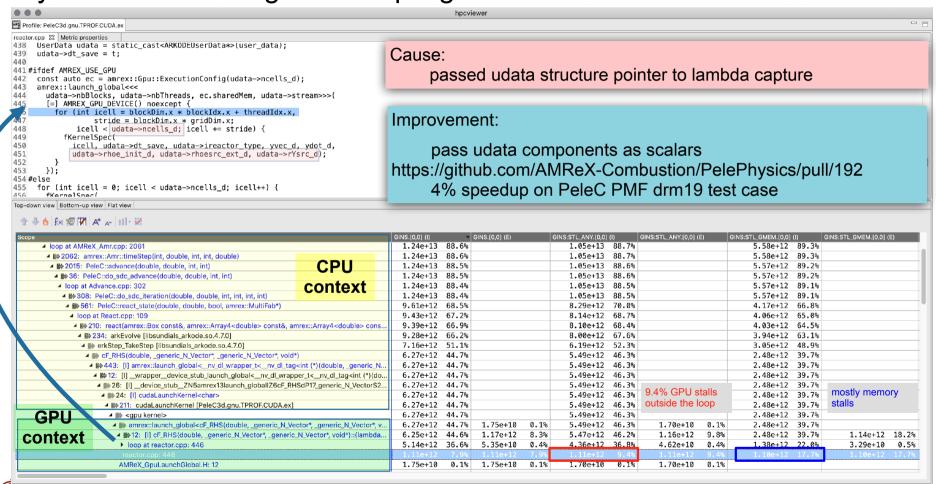


# Case Study: PeleC

- An adaptive mesh refinement solver for compressible reacting flows
- Experiment
  - PC Sampling on Summit node with NVIDIA GPU



## Analysis of PeleC using PC Sampling on an NVIDIA GPU



#### **HPCToolkit Status on GPUs**

- NVIDIA
  - heterogeneous profiles, including 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, including GPU instruction-level measurements with kernel instrumentation and heuristic latency attribution to instructions
  - traces



## Using HPCToolkit at OLCF

- Summit (NVIDIA GPUs)
  - module use /gpfs/alpine/csc322/world-shared/modulefiles/ppc64le
  - module load hpctoolkit
- Crusher (AMD GPUs)
  - module use /gpfs/alpine/csc322/world-shared/modulefiles/x86 64
  - module load hpctoolkit
- Join our ECP Engagement Slack and ask questions
  - https://join.slack.com/t/hpctoolkit-ecp/shared\_invite/zt-1lgrzjt93aEB43o3SVgKFy1yFCjwlyA



#### **HPCToolkit Resources**

- Documentation
  - User manual
    - http://hpctoolkit.org/manual/HPCToolkit-users-manual.pdf
  - Tutorial videos
    - http://hpctoolkit.org/training.html
    - recorded demo of GPU analysis: <a href="https://youtu.be/vixa3hGDuGg">https://youtu.be/vixa3hGDuGg</a>
  - Cheat sheet
    - <a href="https://gitlab.com/hpctoolkit/hpctoolkit/-/wikis/home">https://gitlab.com/hpctoolkit/hpctoolkit/-/wikis/home</a>
- 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

