# Outline

- Introduction to HPCToolkit performance tools
    - Overview of HPCToolkit components and their workflow
    - HPCToolkit's graphical user interfaces
- Analyzing the performance of GPU-accelerated codes with HPCToolkit
    - Slides: Exawind (AMReX)
    - Slides: LAMMPS at Exascale (Kokkos)
    - Demo: GAMESS (OpenMP)
    - Hands-on:
        - Explore available performance databases for a set of applications
        - Collect and explore some measurements for some prepared examples

# Linux Foundation's HPCToolkit Performance Tools

Collect profiles and traces of unmodified parallel CPU and GPU-accelerated applications

Understand where an application spends its time and why

    call path profiles associate metrics with application source code contexts

        analyze instruction-level performance within GPU kernels and attribute it to your source code

    hierarchical traces to understand execution dynamics

Parallel programming models

    across nodes: MPI, SHMEM, UPC++, …

    within nodes: OpenMP, Kokkos, RAJA, HIP, DPC++, Sycl, CUDA, OpenACC, …

Languages

    C, C++, Fortran, Python, …

Hardware

    CPU cores and GPUs within a node

        CPU: x86_64, Power, ARM

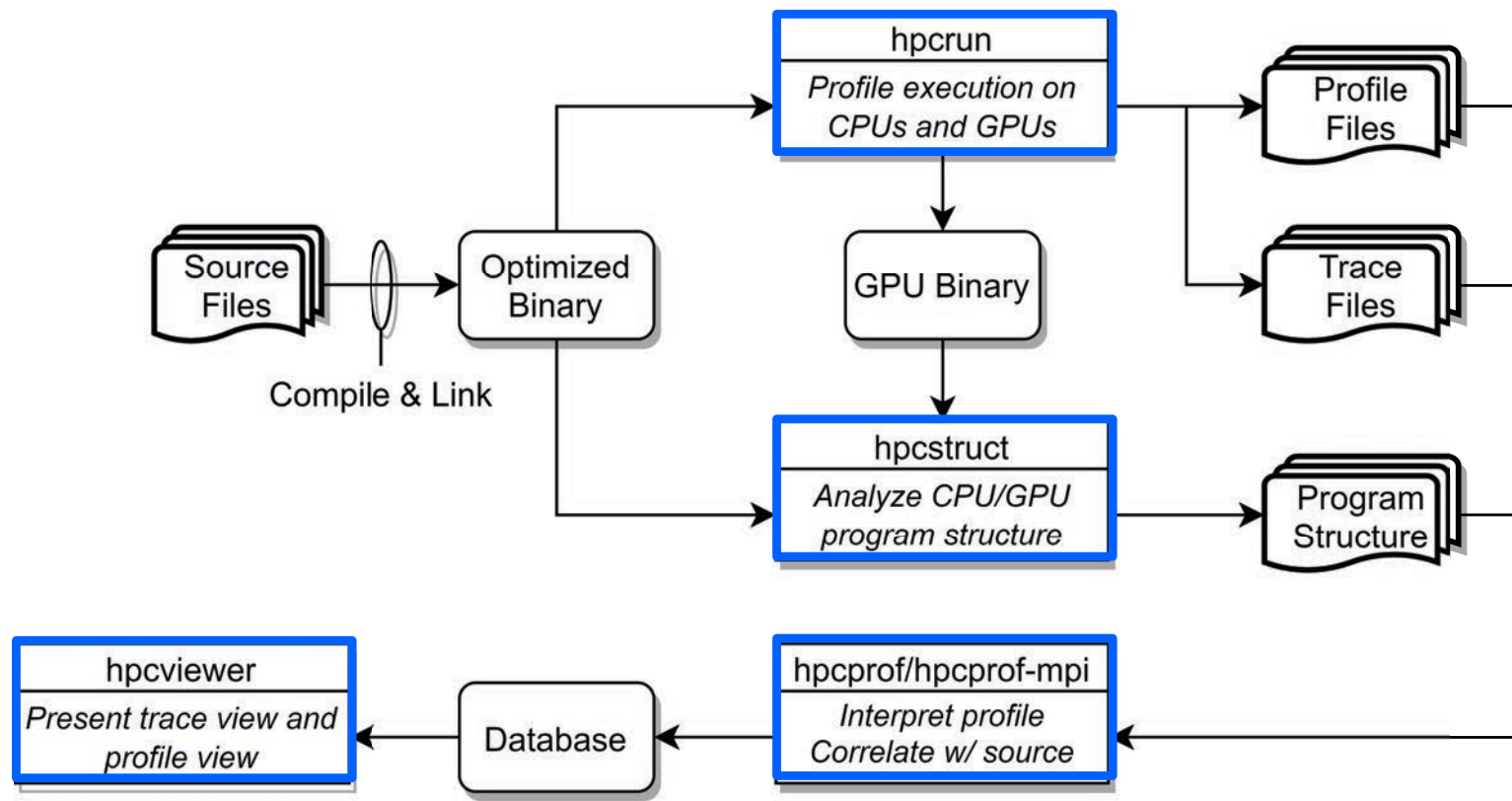        GPU: NVIDIA, AMD, Intel

# Why HPCToolkit?

- 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

- Unlike vendor tools works with a wide range of CPUs and GPUs

# How does HPCToolkit Differ from Vendor Tools?

- NVIDIA NSight Systems
  - —tracing of CPU and GPU streams
  - —analyze traces when you open them with the GUI
    - long running traces are huge and thus extremely slow to analyze, limiting scalability
  - —designed for measurement and analysis within a node
- NVIDIA NSight Compute
  - —detailed measurement of kernels with counters and execution replay
  - —very slow measurement
  - —flat display of measurements within GPU kernels
- Intel VTune: designed for analysis of performance on a single node
- AMD Omtitrace: designed for analysis of performance on a single node
- HPCToolkit
  - —supports more scalable tracing than vendor tools
    - measure exascale executions across many GPUs and nodes
  - —scalable, parallel post-mortem analysis vs. non-scalable in-GUI analysis
  - —detailed reconstruction of estimates for calling context profiles within GPU kernels

# HPCToolkit's Workflow for GPU-accelerated Applications

# HPCToolkit's Workflow for GPU-accelerated Applications
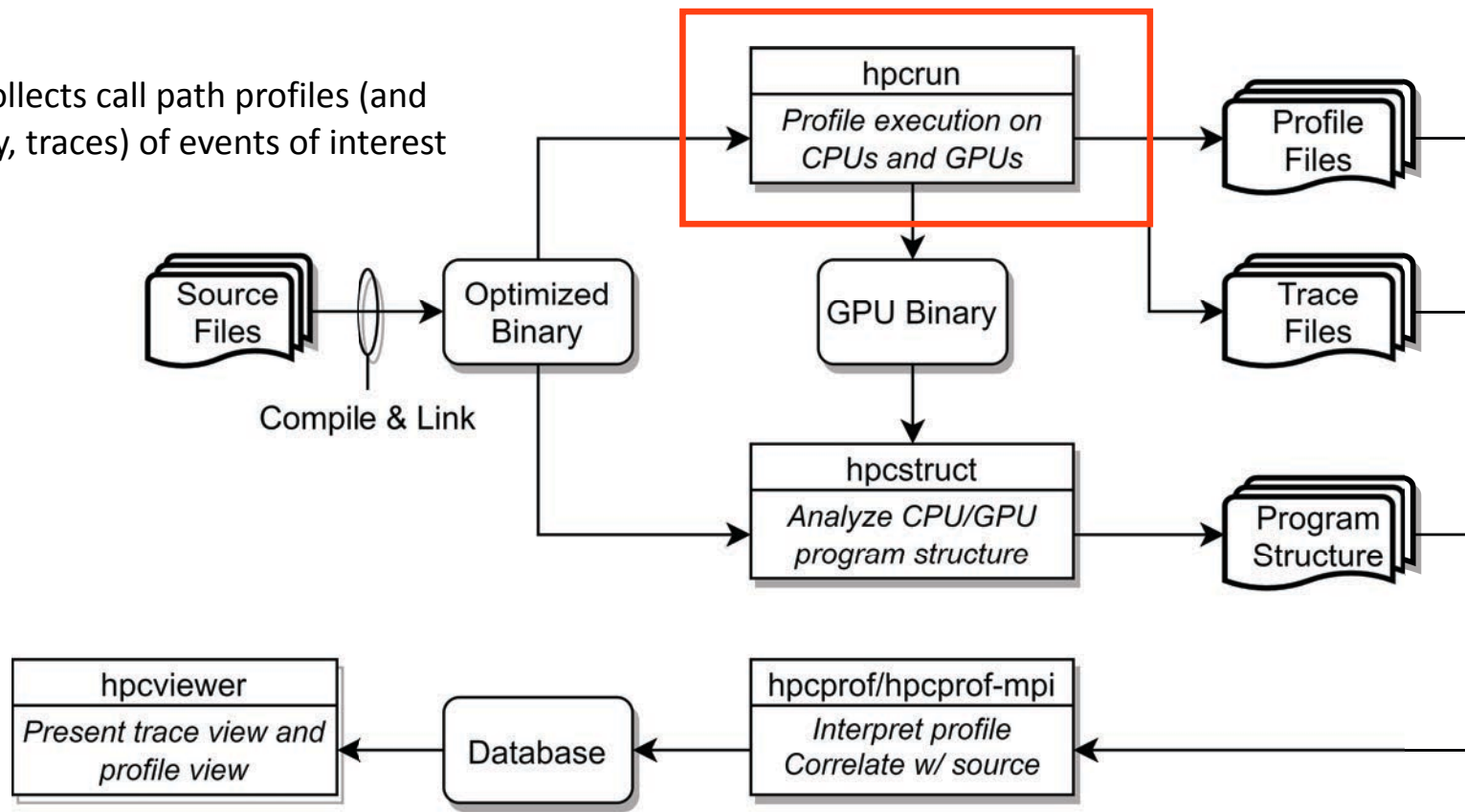
Step 1:
- Ensure that compilers record line mappings
- host compiler: `-g`
- nvcc: `-lineinfo`

# HPCToolkit's Workflow for GPU-accelerated Applications

Step 2:
- *hpcrun* collects call path profiles (and optionally, traces) of events of interest

# Measurement of CPU and GPU-accelerated Applications

- Sampling using Linux timers and hardware counter overflows on the CPU
- Callbacks when GPU operations are launched and (sometimes) completed
- Event stream for GPU operations
- PC Samples: NVIDIA (in progress: AMD, Intel)
- Binary instrumentation of GPU kernels on Intel GPUs for fine-grain measurement

# 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

Calling context tree

Call path sample

● return address

● return address

● return address

● instruction pointer

# hpcrun: Measure CPU and/or GPU activity

- GPU profiling
  - —hpcrun -e gpu=xxx <app> ….                                    xxx ∈ *{nvidia,amd,opencl,level0}*

- 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 -tt <app>

- Use hpcrun with MPI on Polaris or Aurora
  - —mpiexec -n <ranks> … hpcrun -e gpu=xxx <app>

RICE

# 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
    - —typically analyze large application binaries with 16 threads
    - —typically 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 for interpreting PC samples on NVIDIA GPUs

# 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 multithreading  (moderate scale)

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

- Analyze data from large executions with distributed-memory parallelism + multithreading (large scale)

```
mpiexec –n ${NODES} --ppn 1 —depth=128 \
    hpcprof–mpi <measurement-directory>
```

# HPCToolkit's Workflow for GPU-accelerated Applications

Step 4:
* *hpcviewer* - interactively explore profile and traces for GPU-accelerated applications

# 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



A multi-level call stack based view of execution over time

# Case Studies

- ExaWind
- GAMESS  (OpenMP)
- Quicksilver (CUDA)
- LAMMPS (Kokkos) at exascale

# ExaWind: Wakes from Three Turbines over Time



Figure credit: Jon Rood, NREL

# ExaWind: Visualization of a Wind Farm Simulation



Figure credit: Jon Rood, NREL

# ExaWind: Execution Traces on Frontier Collected with HPCToolkit

Traces on roughly ~70K MPI ranks for ~17minutes

Before: MPI waiting (bad), shown in red                    After: MPI overhead negligible*



Figure credits: Jon Rood, NREL                    *replaced non-blocking send/recv with ialltoallv

# ExaWind Testimonials for HPCToolkit

*I just wanted to mention we've been using HPCToolkit a lot for our ExaWind application on Frontier, which is a hugely complicated code, and your profiler is one of the only ones we've found that really lets us easily instrument and then browse what our application is doing at runtime including GPUs. As an example, during a recent hackathon we had, we improved our large scale performance by 24x by understanding our code better with HPCToolkit and running it on 1000s of nodes while profiling. We also recently improved upon this by 10% for our total runtime.*

*- Jon Rood NREL (5/31/2024)*

*One big thing for us is that we can't overstate how complicated ExaWind is in general, and how complicated it is to build, so finding out that HPCToolkit could easily profile our entire application without a ton of instrumentation during the build process, and be able to profile it on a huge amount of Frontier with line numbers and visualizing the trace was really amazing to us.*

*- Jon Rood NREL (6/3/2024)*

# LAMMPS on Frontier: Executions with Kernel Duration of Milliseconds

# LAMMPS on Frontier: Executions with Kernel Duration of Milliseconds

# LAMMPS on Frontier: Executions with Kernel Duration of Milliseconds

# LAMMPS on Frontier: Executions with Kernel Duration of Milliseconds

# LAMMPS on Frontier: Executions with Kernel Duration of Milliseconds

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

## Kernel duration of microseconds

# 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 prior Perlmutter hackathon
    - Single node with 4 GPUs
    - Five nodes with 20 GPUs

**Perlmutter node at a glance**

AMD Milan CPU
4 NVIDIA A100 GPUs
256 GB memory

GAMESS original    All CPU threads and GPU streams

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



GAMESS original · All CPU threads and GPU streams

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



GAMESS original    All GPU streams, whole execution

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



GAMESS original

GPU streams: 1 iteration

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



GAMESS original

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



GAMESS improved          CPU Threads and GPU Streams

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



GAMESS improved

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



GAMESS improved with better manual distribution of work in input

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



GAMESS improved adding Rank 0 Thread 0 to GPU streams

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



1 CPU Stream, 2 GPU Streams: 6 Iterations

# 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

# Case Study: Quicksilver

- Proxy application that represents some elements of LLNL's Mercury workload
- Solves a simplified dynamic Monte Carlo particle transport problem
  - Attempts to replicate memory access patterns, communication patterns, and branching or divergence of Mercury for problems using multigroup cross sections
- Parallelization: MPI, OpenMP, and CUDA
- Performance Issues
  - load imbalance (for canned example)
  - latency bound table look-ups
  - a highly branchy/divergent code path
  - poor vectorization potential

# Quicksilver: Detailed analysis within a Kernel using PC Sampling

# Quicksilver: Detailed analysis within a Kernel using PC Sampling

# HPCToolkit Resources

- Documentation
  - User manual for HPCToolkit: http://hpctoolkit.org/manual/HPCToolkit-users-manual.pdf
  - Cheat sheet: https://gitlab.com/hpctoolkit/hpctoolkit/-/wikis/HPCToolkit-cheat-sheet
  - User manual for hpcviewer: https://hpctoolkit.gitlab.io/hpcviewer
  - Tutorial videos
    - http://hpctoolkit.org/training.html
    - recorded demo of GPU analysis of Quicksilver: https://youtu.be/vixa3hGDuGg
    - recorded tutorial presentation including demo with GPU analysis of GAMESS: https://vimeo.com/781264043
- 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

# Some Hpcviewer Tips

# Information for Using Hpcviewer

- Filtering GPU traces
  - Can use the filter menu to select what execution traces you want to see
    - cpu only, gpu, a mix
    - type a string or a regular expression in the chooser select or unselect the new set
    - only traces that exceed a minimum number of samples
- Filtering GPU calling context tree nodes to hide clutter
  - hide individual CCT nodes: e.g. lines that have no source code mapping library@0x0f450
  - hide subtrees: MPI implementation, implementation of CUDA primitives
- When inspecting GPU activity, be aware that hpcviewer has two modes
  - expose GPU traces or not
    - means: when displaying GPU trace lines, don't just show GPU activity if the time in the middle of a pixel is in a GPU operation. instead, show the first (if any) GPU operation between the time in the middle of the pixel and the middle of the next pixel
    - why? GPU activity is so short, it may be hard to find if we don't "expose" where it is
    - downside: makes the GPU appear more active than it is
      - can correct the statistics by turning the mode off
  - mode can be selected from <File>:<Preferences>:<Traces>

# Hands-on Examples

# Two Kinds of Hands-on Examples

- **Pre-collected databases to explore**
  - **gain experience using hpctoolkit's hpcviewer graphical user interface to analyze performance data**
- **Hands-on examples**
  - **build, run, and view several codes to get the full experience**
    - **hpcrun: measure an application as it executes**
    - **hpcstruct: recover program structure information for mapping measurements to source code**
    - **hpcprof: combine measurements with program structure information**
    - **hpcviewer: explore profiles and traces**

# Performance Databases to Explore

- **On an aurora login node**

```
% qsub -I -l select=2,walltime=1:00:00,place=scatter -l
   filesystems=flare -A gpu_hack -q gpu_hack_prio -X
```

- **On an aurora compute node**

```
% module use /soft/perftools/hpctoolkit/modulefiles
% module load hpctoolkit
% cd /flare/gpu_hack/hpctoolkit/data
% ls
       arborx   gamess   minitest   pelelmex   quicksilver
```

NOTE: all of the databases in these directories end with the suffix ".d". For the gamess examples, the hpctoolkit databases are one directory deeper, i.e. in subdirectories that begin with a number.

# More about the Available Performance Databases

See /flare/gpu_hack/hpctoolkit/data for each of the following

- quicksilver: Monte Carlo particle transport proxy application **(C++ + CUDA)**
  - hpctoolkit-qs-gpu-cuda.d - profile and trace on 4 CPUs + 4 GPUs
  - hpctoolkit-qs-gpu-cuda-pc.d - instruction-level measurements within kernels using PC sampling
  - EXERCISES
- pelelmex: Adaptive mesh hydrodynamics simulation code for low Mach number reacting flows **(C++ + AMReX)**
  - pelelmex.db -a large  trace with load imbalance from 2025 NERSC hackathon run on 16 CPUs + 16 GPU
- gamess: General Atomic and Molecular Electronic Structure System **(Fortran + OpenMP)**
  - 1.singlegroup-unbalanced/hpctoolkit-gamess-1n-chol-noDS.d
  - 2.singlegroup-balanced/hpctoolkit-gamess-1n-chol-fix_load_balance_noDS.d/
  - 3.multigroup-unbalanced-mtarbr/hpctoolkit-gamess-5n.d/
  - 4.multigroup-balanced/hpctoolkit-gamess-5n-manualbalance.d/
  - 5.multigroup-unbalanced-pc/hpctoolkit-gamess-5n-pc.d/
  - 6.scale/hpctoolkit-gamess-22n-test.d/
- arborx **(C++ + Kokkos)**
- minitest **(SYCL and OpenMP TARGET)**

# Hands-on Tutorial Examples on Aurora

```
% git clone https://gitlab.com/hpctoolkit/tutorial-examples
% cd hpctoolkit-tutorial-examples/gpu/intel
% ls
    arborx.kokkos.sycl    minitest.sycl    minitest.omp


  For your chosen example
  1. cd to the example directory
  2. source setup-env/aurora.sh  # custom for each example
  3. make build # build the code
  4. make run # submit to the batch queue
      wait until the hpctoolkit database ending in .d appears and you see the file
      log.run.done appear
  5. make view # launch hpcviewer to explore the performance data collected
```
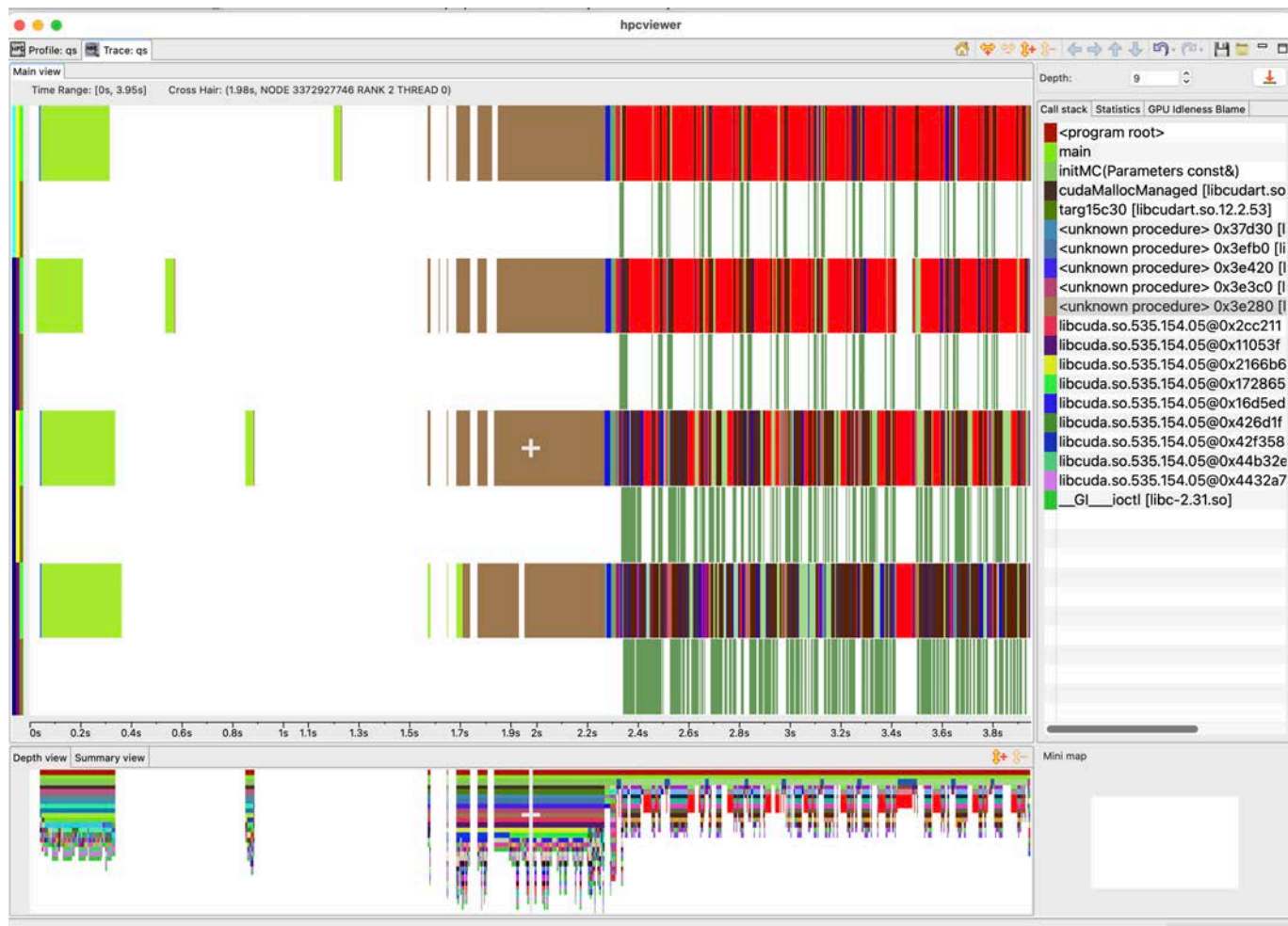
# Hands-on Tutorial Example on Polaris: Quicksilver

% git clone https://gitlab.com/hpctoolkit/tutorial-examples

% cd hpctoolkit-tutorial-examples/gpu/nvidia/quicksilver.cuda

1. source setup-env/polaris.sh  # custom for each example

2. make build # build the code

3. make run # submit to the batch queue

   wait until the hpctoolkit database hpctoolkit-qs.d appears and you see
   the file log.run.done appear

4. make view # launch hpcviewer to explore the profiles and traces

5. make run-pc # collect instruction-level kernel measurements

   wait until the hpctoolkit database hpctoolkit-qs-pc.d appears and you
   see the file log.run-pc.done appear

6. make view-pc # launch hpcviewer to explore the pc sampling data

# Inspecting the Pre-collected Quicksilver Data

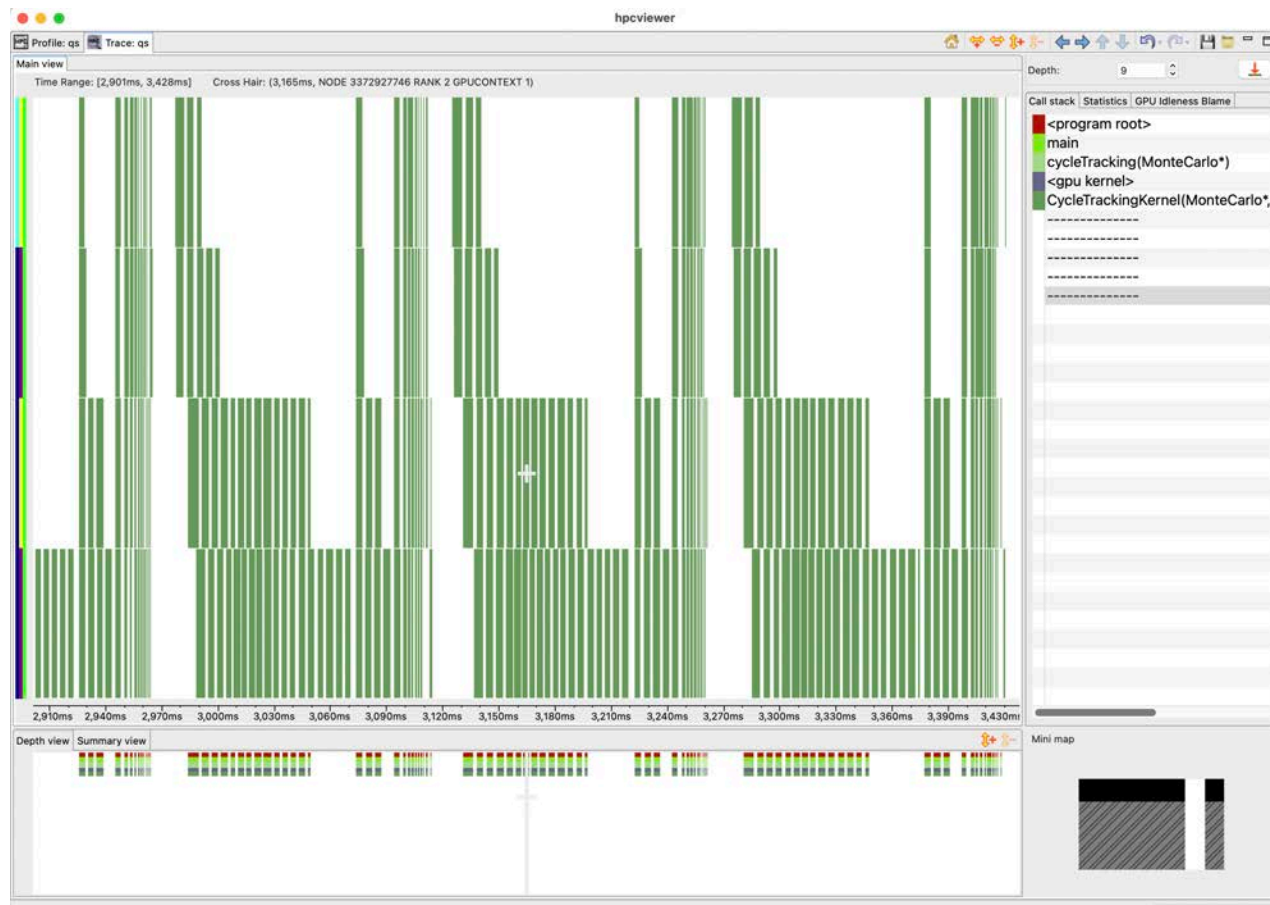# Select the Tab "Trace: qs"

# Use the Filter to "Uncheck all" and Check "GPU" streams

# See Load Imbalance Across the Four GPUs

# Analyzing Quicksilver Traces

**Using a measurement database with profiles and traces**
- Select the Trace tab "Trace: qs"
- Identifying the traces
    - Select a pixel on a trace line
    - Look at legend on the top of the display, which reports the location of the "cross hair"
    - Is this a CPU or GPU trace line?
    - Repeat this a few times to identify what each of the trace lines represents
- Notice that each time you select a colored pixel on a trace line, you will be shown the function call stack in the rightmost pane
- At the top of the pane is a "depth" indicator, that indicates what level in the call stack you are viewing. The selected level will also be highlighted
- You can change the depth of your view by using the depth up/down, typing a depth, or simply selecting a frame in the call stack at the desired depth
- You can select ↧ above the call stack frame to show the call stacks at the deepest depth
    - If a sample doesn't have an entry at the selected depth, its deepest frame will be shown

# Analyzing Quicksilver Traces

**Using a measurement database with profiles and traces**

- Zoom in on a region in a trace by selecting it in the trace display

- Use the back button ↶ to undo a zoom

- Use the control buttons ⇔ ⇔ ⬍+ ⬍− ← → ↑ ↓ at the top of the trace pane to
    - expand or contract the pane
    - move left, right, up, or down
- Keep an eye on the minimap in the lower right corner of the display to know what part of the trace you are viewing

- Use the home button 🏠 to reset the trace view to show the whole trace

# Analyzing Quicksilver Traces

**Using a measurement database with profiles and traces**
- Select the Trace tab "Trace: qs"
- Configure filtering
  - Use the Filter menu to select Filter Execution Contexts
  - In the filtering menu, select "Uncheck all"
  - Now, in the empty box preceded by "Filter:", type "GPU" and then click "Check all"
  - Select "OK".
  - Now, the Trace View will show only trace lines for the GPUs.

- Inspect the trace data
  - Is the work load balanced across the GPUs? How can you tell?
  - Bring up the filter menu again. Select "Uncheck all". Type in "RANK 3" in the Filter box. Select thread 0 and the GPU context. Select "OK".
  - Move the call stack to depth 2
    - What CPU function is Rank 3 thread 0 executing when the GPU is idle?
    - Does this suggest any optimization opportunities?

# Analyzing the Quicksilver Summary Profile

**Using a measurement database with profiles and traces**

- Select the Profile Tab "Profile: qs"
- Use the column selector 📊 to deselect and hide the two REALTIME columns
- Select the GPU OPS column, which represents time spent in all GPU operations
- Select the 🔥 button to show the "hot path" according to the selected column
  - the hot path of parent will continue into a child as long as the child accounts for 50% or more of the parent's cost
- The hot path will select "CycleTrackingKernel" — a GPU kernel that consumes 100% of the GPU cost in this profile
- Use the 📊 button to graph "GPU OPS (I)" — inclusive GPU operations across the profiles
  - Are the GPU operations balanced or not across the execution contexts (ranks)?

# Analyzing the Quicksilver Summary Profile

- You will notice that for quicksilver, HPCToolkit doesn't report any data copies between the host and device

  - The quicksilver code uses "unified memory" so that all of the data movement occurs between CPU and GPU using page faults rather than explicit copies

  - Today's GPU hardware doesn't support attribution of page faults to individual instructions

    - We could profile them, but not attribute them to code

# The Profile View in the other "PC Sampling" Database

# Analyzing Quicksilver PC Samples

**Using a measurement database with traces that was collected \*with\* PC sampling enabled**

Using the default top-down view of the profile

- Select the column "GINS (I)" to focus on the measurement of inclusive GPU Instructions
- Select use the flame button to look at where the instructions are executed
- In the call stack revealed, you will <gpu kernel> placeholder that separates CPU activity (above) from GPU kernel activity (below)
- Below the <gpu kernel> placeholder you will see the function calls, inlined functions, loops and statements in HPCToolkit's reconstruction of calling contexts within the CycleTrackingKernel

- Using the bottom-up view of the profile
    - Select the bottom-up tab of above the control pane
    - Select the GINS STL_ANY (E) column, which will sort the functions by the exclusive GPU instruction stalls within that function
    - Scroll right to see which of the types of contributing types of stalls accounts for most of the STL_ANY amount
    - Select the function that has the most exclusive stalls
    - Select the the hot path to see where this function is called from.
        - Where do the calls to the costly function come from?
        - Does there appear to be an opportunity to reduce the number of calls to this function?

# Filtering Tips to Hide Unwanted Implementation Details

- Filter "descendants-only" of CCT nodes with names *MPI* to hide the details of MPI implementation in profiles and traces

- Filter internal details of RAJA and SYCL templates to suppress unwanted detail using a "self-only" filter

# Downloading, Installing, and Using Hpcviewer Graphical User Interface on Your Laptop

# Hpcviewer Graphical User Interface on Your Laptop

<span style="color:red">Prepare to explore performance data on your laptop</span>

- Download and install hpcviewer

  - See https://hpctoolkit.org/download.html

    Select the right one for your laptop: MacOS (Apple Silicon, Intel), Windows, Linux

    Install a recent Java (17 or 21) if you don't have one, using the directions on the page

- User manual for hpcviewer: https://hpctoolkit.gitlab.io/hpcviewer

# Viewing Performance Data

- Copy a performance database directory to your laptop and open it locally

- Open a performance database on a remote system

Note: using a HPCViewer with a remote system presumes that hpcserver has already been installed on the remote system

— hpcserver has been installed on Aurora

— you can download and install hpcserver on your local cluster as well (ask in Slack for directions)

# Configuring Hpcviewer Remote Access

Run hpcviewer

From the file menu, select "Open remote database"

Fill in the hostname/IP address: aurora.alcf.anl.gov

Fill in your username on Aurora

Fill in the remote installation directory for hpcviewer's server: `/soft/perftools/hpctoolkit/hpcserver`
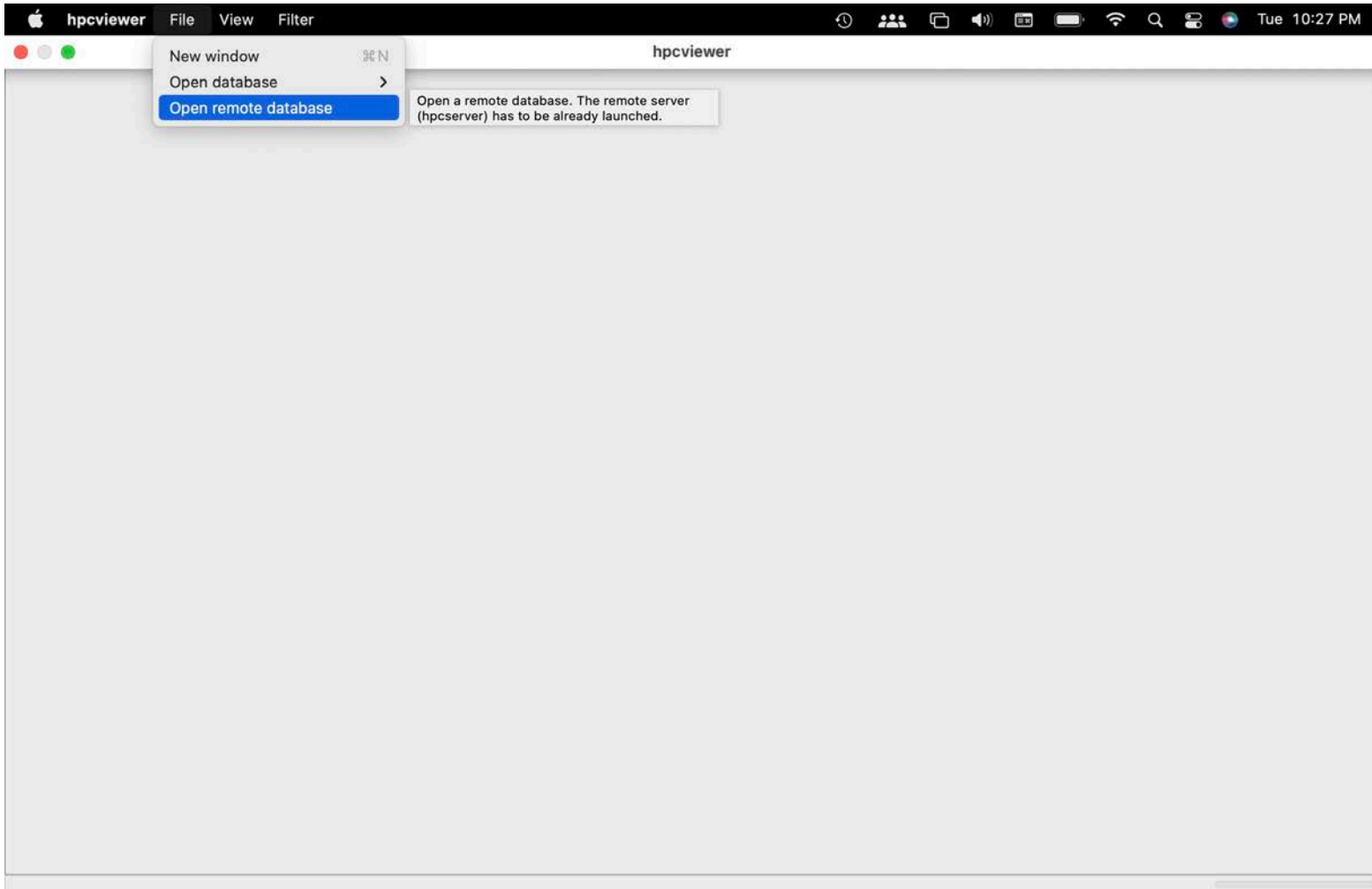
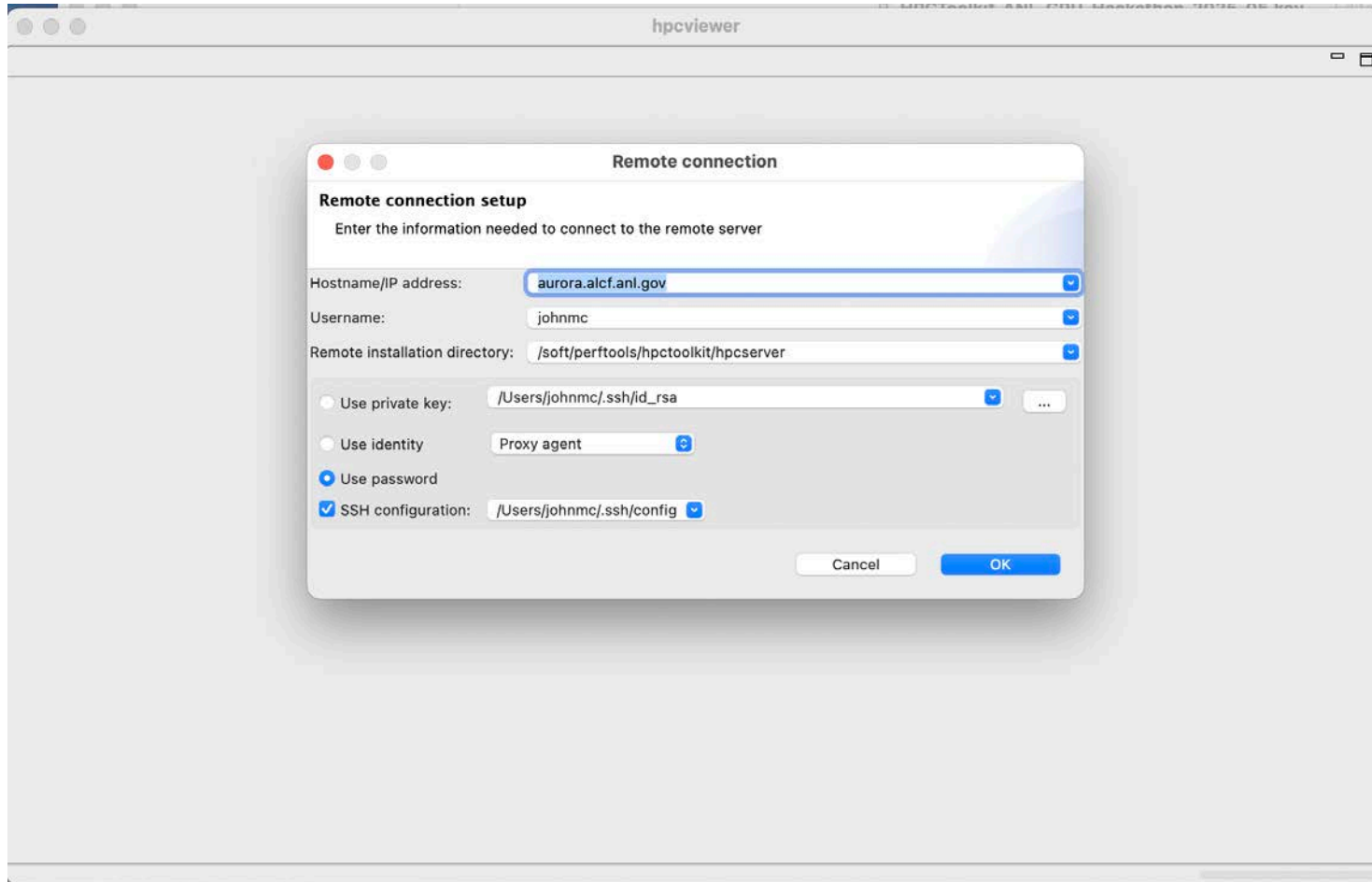Select the authentication method: "Use password"

Click "OK"

Authenticate using your token as you normally do

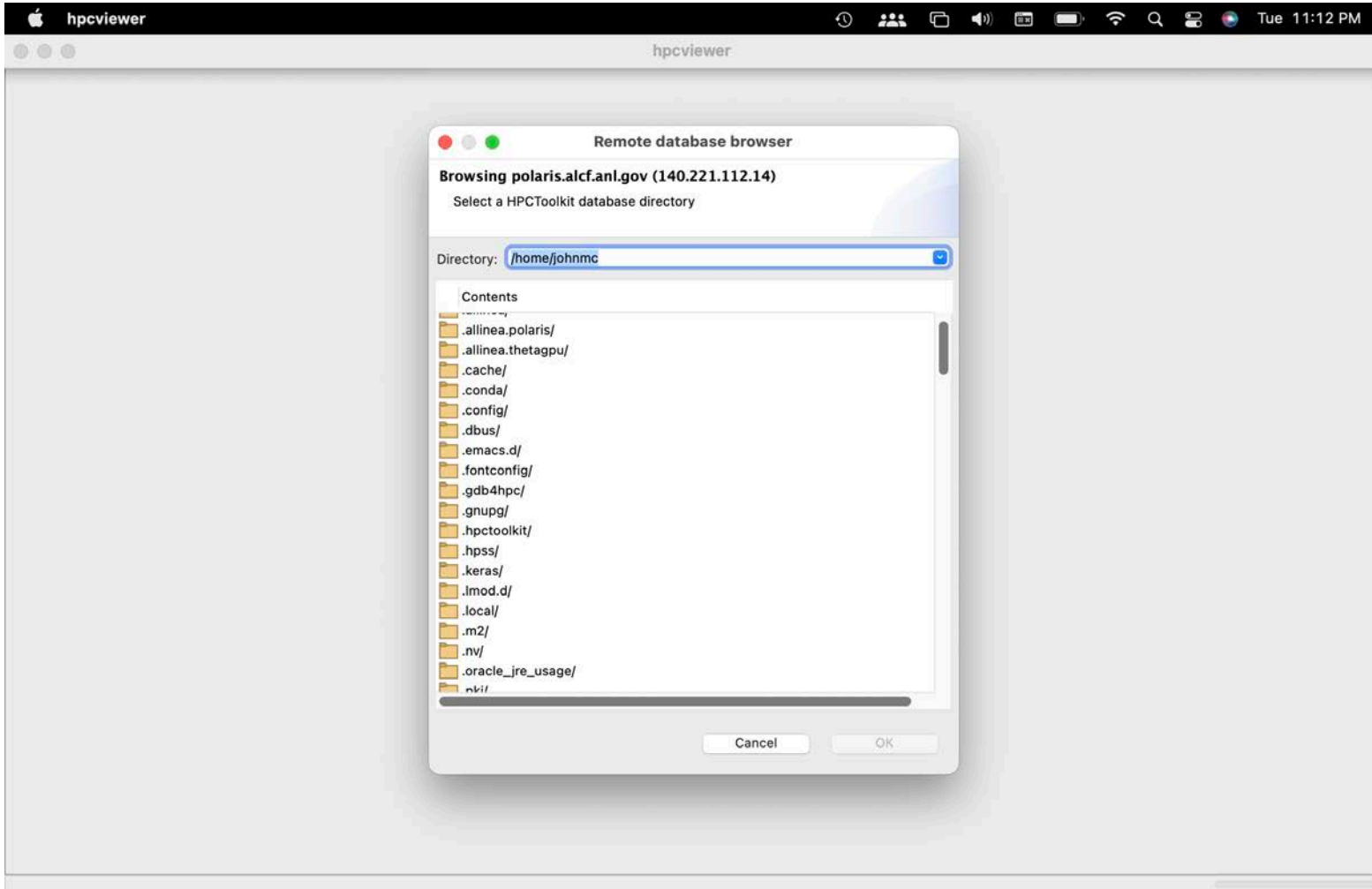Navigate to a database with the file chooser in /flare/gpu_hack/hpctoolkit/data

ATPESC2024
EXTREME-SCALE COMPUTING

Argonne
NATIONAL LABORATORY

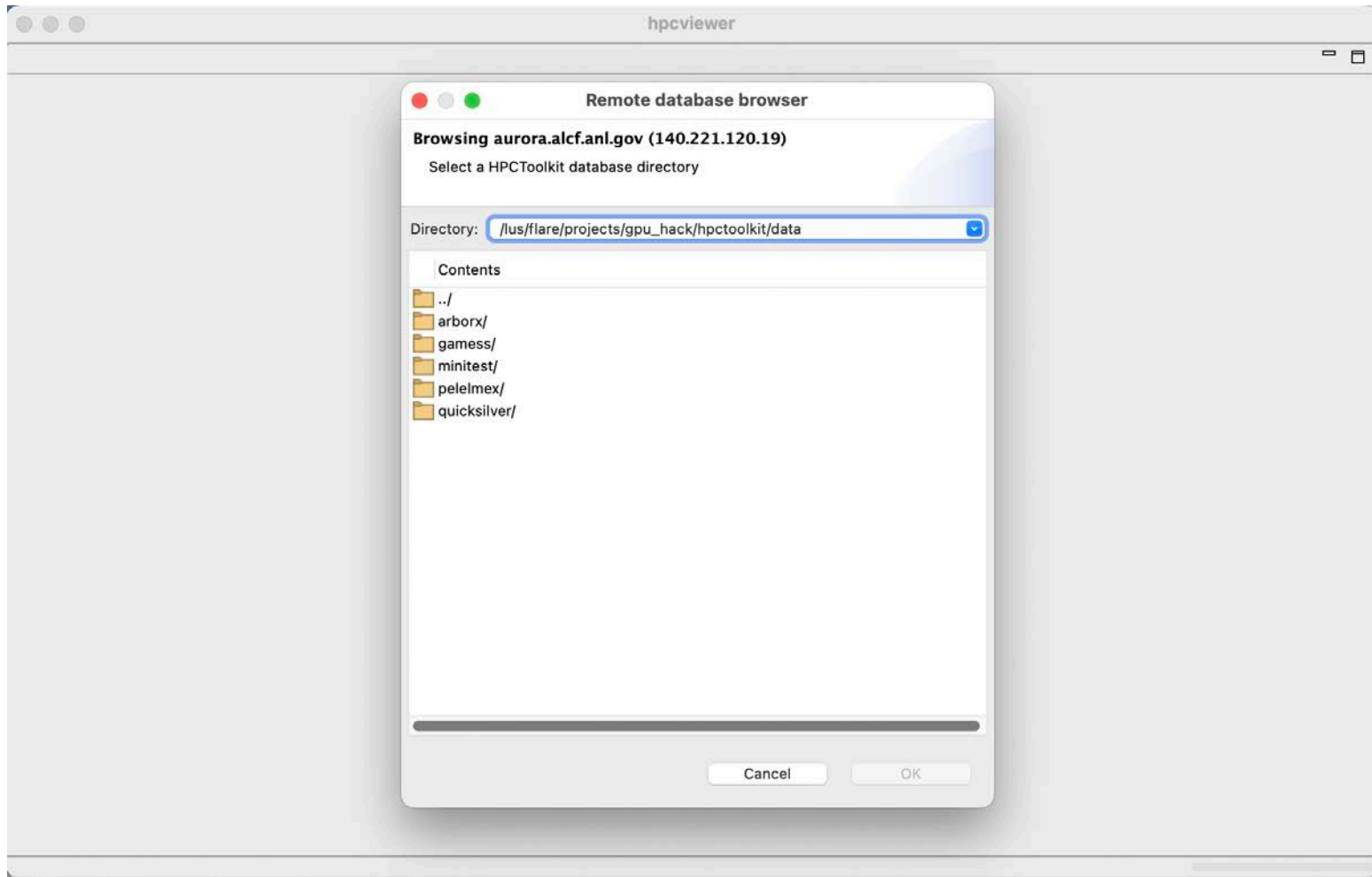# Opening a Remote Database

# Configuring for remote access to Aurora using hpcserver

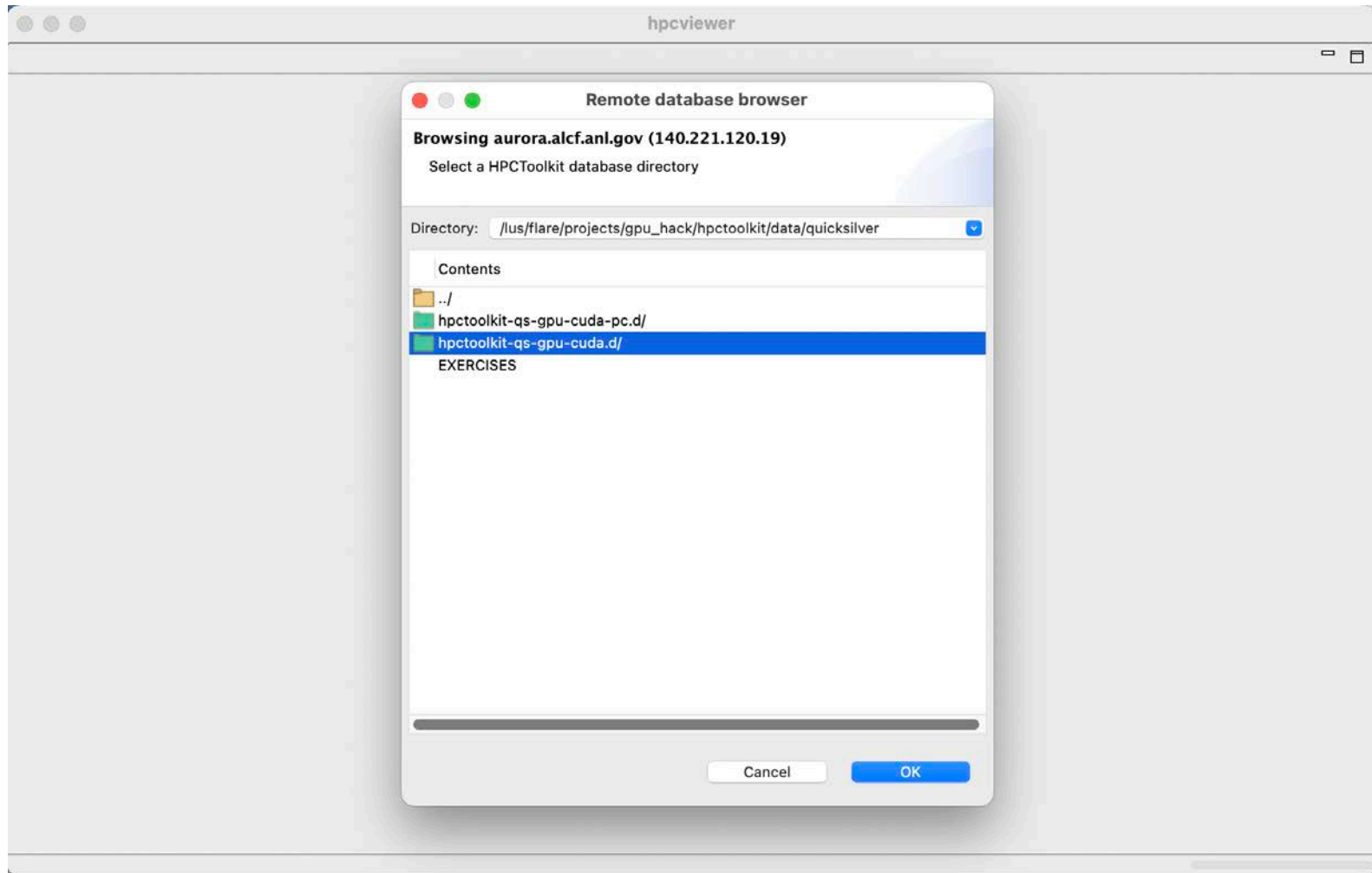# First View of Aurora: Your Home Directory

# Navigate to Example Databases

# Select a Quicksilver Database with Traces

# After Selecting hpctoolkit-qs-gpu-cuda.d

# Troubleshooting Tips

# Why can't I see my Source Code?

- To relate performance measurements to your application source code, the code must be compiled with a "-g" option in addition to your optimization flags. Otherwise, there is no information for any tool to map performance to anything finer grain than at the procedure level
  - For instance, if you are building with make, you will want to build **RelWithDebInfo** rather than **Release** for the best experience