# *Human-Oriented Discussions on Performance Analysis - Outbrief*

## Scalable Tools Workshop '23
## Connor Scully-Allison

**https://bit.ly/stw-humans**

# Participants

- Connor Scully-Allison (The University of Utah)
- Kate Isaacs (The University of Utah)
- Wileam Phan (Rice University)
- Sameer Shende (University of Oregon)
- Jan Laukemann (I'm University of Erlangen-Nuremberg / NHR@FAU)
- Shadmaan Hye (University of Utah)
- David Boehme (LLNL)
- William Jalby (UVSQ/UPSaclay)
- Stephanie Brink (LLNL)
- Kathleen Shoga (LLNL)
- Rahat Zaman (University of Utah)
- David Montoya (Trenza)
- Mahesh Rajan (Trenza)

**https://bit.ly/stw-humans**

# Working Group Goals

- What groups of users exist in the performance tools space?
- Who are they?
  - What do we know about them now?
  - What do we still need to know?
  - How do they differ from us?
- What can we do to understand them better?

**https://bit.ly/stw-humans**

# Exercise: Who is your user?

- User Distinction
  - Code developers are the people we expect to use the tools.
  - Code runners are the people who generate the configurations and performance problems.
    - Domain expert code runners need to relate performance data back to the domain execution
- Code developers are encouraged to add features over performance.
- Code developer teams need a performance champion — or tool developer champions need to come in and set them up and create internal experts.
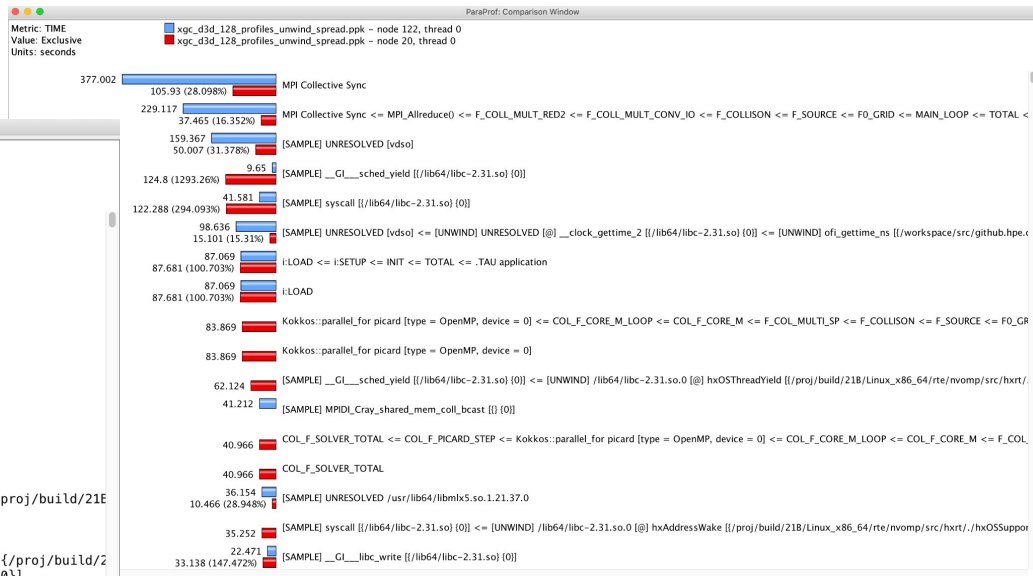
**https://bit.ly/stw-humans**

# Exercise: Who is your user? - Part 2

- People don't want to work with performance tools unless they absolutely need them
    - Then there's a learning curve.
    - Community looking into encouraging them to look at performance continuously.
- There are probably 5-6 kinds of users.
    - Which should we prioritize?
- Users in private sectors can be very secretive about source code access and could benefit from tools that perform analysis without any annotation.

**https://bit.ly/stw-humans**

# Exercise: How do you differ?

- Tool developers have comparatively structured workflows for performance analysis compared to their users
  - With significant foundation understanding of where performance issues commonly lie
- Tool developers are frustrated by the performance and usability of visualization in tools
  - "Zooming on a trace should be like google maps"
- Legacy data formats hold us back from performant GUI development

**https://bit.ly/stw-humans**

# TAU performance data explains why MPI collectives take as long as they do!



Time spent in MPI, MPI_Barrier inside the collective operation (MPI Collective Sync), load imbalance, time spent in Kokkos, in application functions, I/O, sampling, system calls.

**https://bit.ly/stw-humans**

# HPCToolkit analysis of GAMESS 5-node run on Perlmutter



1. Look at the profile, specifically the total time spent on GPU (GPUOP), then compare time spent in kernels (GKER) vs time spent in data copies (GXCOPY).
This one spends most of the time in kernels (**good!**)

**https://bit.ly/stw-humans**

# HPCToolkit analysis of GAMESS 5-node run on Perlmutter



2. Sort by GKER and click the "flame" button to find the most expensive GPU kernel
(This looks the same as their 1-node runs)

**https://bit.ly/stw-humans**

# HPCToolkit analysis of GAMESS 5-node run on Perlmutter



3. Look at the trace, filter to show only thread 0 on CPU and all the GPU activity. Here we see the problem is **load imbalance** across different compute nodes

**https://bit.ly/stw-humans**

# What else should we be asking? What do we want to know about them?

- How often is analysis being done?
- Granularity of what they are looking for/ their interest.
  - "See everything at absolutely no overhead"
- Where does a user lie on the interactive non interactive analysis spectrum
- Users are interested in seeing data related to abstractions that they are familiar with
  - Annotations that they have inserted into the code
  - "Kokkos parallel for picard"
  - They want to relate the performance to the domain application function

# What actionable steps can we take to understand our users better?

- Provide clear use cases/samples alongside API Documentation
  - **Quickstart for <class of user>**
    - This makes it transferable between different users on the same team
    - Provide a template github repo alongside documentation
    - Keep tutorials short
- Acquaint new users with the diversity of commonly used tools
  - There is few opportunities to try new tools and they are often driven by problems
  - **A reference document of some variety that guides tool developers**
    - **Provide pros and cons of the tools**
- We need to measure user workflows and integrate them better.
  - **HCI/HCC experts can help here**

**https://bit.ly/stw-humans**

# Raw Notes Slides

# Exercise: Who is **your** "user"? - Stage Right Group (High Level)

- Code developers are the people we expect to use the tools. Code runners are the people who generate the configurations and performance problems. There isn't always a good channel for the configurations/data to pass back to the developers.
- Code developers are encouraged to add features over performance.
- Code developer teams need a performance champion — or tool developer champions need to come in and set them up and create internal experts.
- People don't want to work with performance tools unless they absolutely need them but then there's a learning curve. Currently community looking into encouraging them to look at performance continuously.
- There are probably 5-6 kinds of users. Which should we prioritize.

**https://bit.ly/stw-humans**

# What we hope to answer:

- What groups of users exist in the performance tools space?
- Who are they?
  - What do we know about them now?
  - What do we still need to know?
  - How do they differ from us?
- What can we do to understand them better?

**https://bit.ly/stw-humans**

# Exercise: Who is **your** "user"? - Stage Left

What is their expected performance analysis workflow (with or without your tool)?
- Scientific simulations application builder
- Two types (domain scientists and coders)
- PThread based performance analysis
  - Async io programming
  - With python
  - Financial Users are very closed off
- Straightforward/traditional performance analysts
- Slurm Based-Job Launchers/MPI
- TAU traces visualized in Vampir/Perfetto
- Nvidia Nsight -> only for NVIDIA GPUs
- Proprietary software does not scale on a desktop

What is their biggest frustration?
- Scalability of scaling tools
- AMD is in early stages
- They cannot even try the big runs
- HPCToolkit -> Learning curve is difficult
  - Needs to be handheld with how to analyze the data
  - It's not evident how the data relates to runtime performance issues
  - They do not have the time or resources to learn it in depth
- They don't want to change the binary
- Users don't want to be overwhelmed with unnecessary information
- Not enough visibility into the main problems
- They need to understand why their collective operations are taking so long
- Performance data needs to be interpretable in an abstractions they understand

How much time do you think they spend analyzing performance in their job? How long per session?

- Only when needed

How often do they try new tools?

- Not often at all
- There is a focus on desktop-class systems with a GPU

**https://bit.ly/stw-humans**

# Exercise: Who is **your** "user"? - Stage Right Group (Raw)

What is their expected performance analysis workflow (with or without your tool)?
- They are scientists, less experienced in CS, physicists, chemists, etc. Others have those more used to coding.
- Many students use tools for class in CS and computational engineering.
- They would like to use tools but they don't have time to use tools. They know it's important but as long as things are working they don't do things. (Guilty?) ← Might be a bit of a disconnect
- Teams have their own cultures and it's hard to get into it unless someone from the inside wants to do something (and adds money). Some will create their own and some will go find tool developers. So there's some Not-Invented-Here. You need humans in the teams to help adoption.
- We tell them run a profile, find hotspot, focus on them, go a level deeper with markers and annotations to actually see what's going on. Use LIKWID or PAPI to see memory bound or compute bound... OR compare to the expected, is there space for optimization if you're at 90% of peak... but in practice people don't want to spend that much time. They want to run a command and get an answer.
- This leads to teaching them to do automated performance regression on the daily with Caliper or other tools. Some teams are using it.
- Tools like TotalView are for when you absolutely need it. Performance is the same thing but then there's a learning curve. You don't know how to use it and you don't understand the data. So how do we build performance as a continuous process. If they see something every day that alerts them there's a level of performance in different categories, then they can see the expected. Teams are in different shape with preparedness.
- User group is so vast – the application USERS and the application DEVELOPERS. The developers don't necessarily run at scale but the users/scientists do. Developers are doing features. At Sandia there used to be a separate performance team.This is sometimes an 'application readiness' team.
  - There is some moving tools from Developers to Code Runners problem. Developers in theory can read a profile and optimize their code.The Runner is the one who runs into the performance bottlenecks.
  - It's hard to run multiple performance runs because of the bank limitations
- We should still develop for the code developers and not as much the code users… the users could generate data but it's hard to get them to use it. It would be nice to have their input decks for things like application parameters.

How much time do you think they spend analyzing performance in their job? How long per session?

How often do they try new tools?

- Tools are not super-intuitive. Need learning curve so it scares away
- Users may not have time for that
- Tools that are too simple and easy to use but that takes away some of their potential power. So it's a trade off of simple results versus more power.
- Each tool has a different interfaces, and they remember having struggles and then there's a wall to get into things. The benefit has to be really obvious.
- Having human champions to do it for them has been great. So people from the Caliper team come in and put things in.
- Social dynamics may get in the way of valuing performance expertise versus science expertise. The champions teach the scientists until they have expertise.
- Tools have help visualize in an ensemble are beautiful and so nice

What are barriers/frustrations?

- Documentation limitations – why is it hard to find info on Caliper? Is there no quickstart guide? Again a learning curve.
- What's the Apple design philosophy? How can you make this so easy? First compParmentalize, what are our 5-6 categories of users.
  - Funding/focus tends to be on a few major applications and there's some rickling down over that.
  - We are identifying the groups, figuring out what data they need, figuring out how to get the data, and then finding ways to build up from there. However, we have tech issues at making sue we can help people do things at scale
- We tend to focus on individual parts of the system (compilers, I/O) but there's a need to look at the whole workflow. This can be ensembles or it could be heterogeneous. Workflow managers are starting to collect data. Monitoring is involved. Application-specific not enough.
- We have difficulty getting enough traction to have someone look at it.
- Users want to be able to turn on and off instrumentation. Once the instrumentation is there,then we can show them how to use the data.Meet with them occasionally has been somewhat successful.
- Silo'd institutions make it difficult for application users to talk to the developers. Institutional problem between application users not having the data going directly back to the developers

**https://bit.ly/stw-humans**

# Exercise 2: How do **you** differ? - Stage Left

What is your perf analysis workflow?
- Sameer -> Run application with `tau_exec` -> look at hardware performance counters -> look at traces through vampir or other
  - This workflow is more structured than a users workflow, maybe
- Wil -> Collect CPU+GPU profile+trace w/ HPCToolkit, open database in HPCViewer
  - HPCViewer profile view is a tree+table paradigm
  - Can sort by inclusive vs exclusive metrics
  - Looking at GPU kernels before anything else
  - Looking at GPU kernels vs. GPU copies
    - To understand execution vs. transportation overhead
  - A case with 30% time on the GPU is a very very good case
    - Some apps can be as low as 1%
  - Select a metric and use the flame button to highlight a hotpath immediately
  - Look at the trace view
  - Load imbalance is a commonly identified issue

How much time do you spend analyzing performance in your job? How long per session?

What is your biggest frustration?
- Would like to have visualization that scales with levels of the runtime.
  - Google maps as a metaphor
  - No overview to zoom
  - Formats the support hierarchical views
  - There is no performance visualization with different scales
- Our data formats are legacy and do not adhere to hierarchical views "like a microscope"

How often do you try new tools?

**https://bit.ly/stw-humans**

18

# What actionable steps can we take to understand our users better?

- Is there merit in socializing users to multiple tools?
  - There is few opportunities to try new tools
  - **A reference document of some variety that guides tool developers**
    - **Provide pros and cons of the tools**
- Did we all learn performance and debugging under duress of an issue?
  - time() as the "ultimate performance tool"
    - Why is this the case: is it a usability issue, discoverability of tools issue, learability issue?
- We need to measure the workflow and integrate it better.
  - We start from assumptions that the space of performance problems are pretty limited
- Complexity of tool installation is very difficult
  - Containers, Different Configurations
  - A user space package manager system
    - Spack fills this role to a certain extent

# What actionable steps can we take to understand our users better?

- Provide clear use cases/samples alongside API Documentation
  - Sometimes these tools are hard to move out of their silos
  - There are often very clear and simple "5 steps" we do when we use a tool
  - Should documentation be more task based first?
    - Is that an easier step or too constraining
    - **Quickstart for <class of user>**
      - This makes it transferable between different users on the same team
      - Will users get too impatient for tutorials
    - Provide a template github repo alongside documentation
    - Are we different from our users in this learning workflow?
  - Are there academic users? Are they any different? Should we do anything for them?
    - Do we need to focus down our documentation to avoid struggling with meeting too many diverse needs?

- 

**https://bit.ly/stw-humans**