# Thicket:
# Seeing the performance experiment forest for the individual run trees

Scalable Tools Workshop & HPDC

CASC
Center for Applied
Scientific Computing

19-23 June 2023

**Stephanie Brink**, Michael McKinsey, David Boehme, Connor Scully Allison

Ian Lumsden, Daryl Hawkins, Treece Burgess, Vanessa Lama

Jakob Luettgau, Kate Isaacs, Michela Taufer, Olga Pearce

HPDC
"Open Source
Tools and Data"
Session
Thurs, 6/22,
9am EST

Lawrence Livermore
National Laboratory

# Challenge: Performance analysis in complex HPC ecosystem

- HPC software and hardware are increasingly complex. Need to understand:
  - Strong scaling and weak scaling of applications
  - Impact of application parameters on performance
  - Impact of choice of compilers and optimization levels
  - Performance on different hardware architectures (e.g., CPUs, GPUs)
  - Different tools to measure different aspects of application performance



① Run Code with Measurement Tools

② Call Tree Profiles Produced from Multiple Studies

Data A   Metadata

Data B   Metadata

**Goal: Analyze and visualize performance data from different sources and types**

# Our big picture solution for analyzing and visualizing performance data from different sources and type



① Run Code with Measurement Tools

② Call Tree Profiles Produced from Multiple Studies

Data A
Metadata

Data B
Metadata

③ Load Data Into Thicket Object

Thicket

Performance Data

Metadata

Aggregate Statistics

④ Exploratory Data Analysis (EDA)

Examine

Clustering

Analyze, Model

Manipulate

Documentation: **thicket.readthedocs.io**

https://github.com/LLNL/thicket
https://github.com/LLNL/thicket-tutorial

# What do profiling tools collect per run?         e.g., Caliper

## 1) Call Tree



MAIN

FOO

BAR

BAZ

## 2) Performance data

| Node | Cache Misses |
|------|--------------|
| MAIN | |
| FOO | |
| BAR | |
| BAZ | |

- Time, FLOPS

- Cache misses

- Memory accesses

## 3) Metadata per run

| User | Platform |
|------|----------|
| | |

- Batch submission (user, launch date)

- Hardware info (platform)

- Build info (compiler versions/flags)

- Runtime info (problem parameters, number of MPI ranks used)

# Use Thicket to *compose* performance profiles in Python

**P1**

Metadata

| User | Platform |
|------|----------|
|      |          |

## Performance metrics

| Node | Cache Misses |
|------|--------------|
| MAIN | 24 |
| FOO  |    |
| BAR  |    |
| BAZ  |    |

MAIN
FOO
BAR
BAZ

**P2**

Metadata

| User | Platform |
|------|----------|
|      |          |

## Performance metrics

| Node | Cache Misses |
|------|--------------|
| MAIN | 16 |
| FOO  |    |
| BAR  |    |
| BAZ  |    |

MAIN
FOO
BAR
BAZ

# Use Thicket to *compose* performance profiles in Python

**P1**

Metadata

| User | Platform |
|------|----------|
|      |          |

MAIN
FOO
BAR
BAZ

## Performance metrics

| Node | Cache Misses |
|------|--------------|
| MAIN | 24 |
| FOO |  |
| BAR |  |
| BAZ |  |

**P2**

Metadata

| User | Platform |
|------|----------|
|      |          |

MAIN
FOO
BAR
BAZ

## Performance metrics

| Node | Cache Misses |
|------|--------------|
| MAIN | 16 |
| FOO |  |
| BAR |  |
| BAZ |  |

① Compose functions w/matching call trees

①

## Performance metrics

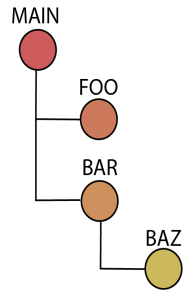| Node | Profile | Cache Misses |
|------|---------|--------------|
| MAIN | P1 | 24 |
| MAIN | P2 | 16 |
| FOO | P1 |  |
| FOO | P2 |  |
| BAR | P1 |  |
| BAR | P2 |  |
| BAZ | P1 |  |
| BAZ | P2 |  |

# Use Thicket to *compose* performance profiles in Python

P1

Metadata

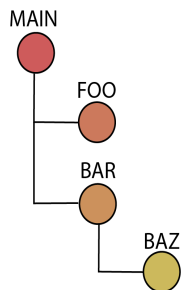| User | Platform |
|------|----------|
|      |          |

**Performance metrics**

| Node | Cache Misses |
|------|--------------|
| MAIN | 24           |
| FOO  |              |
| BAR  |              |
| BAZ  |              |

P2

Metadata

| User | Platform |
|------|----------|
|      |          |

**Performance metrics**

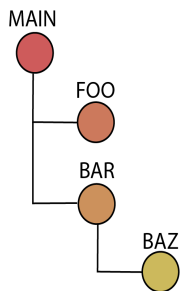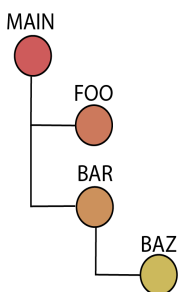| Node | Cache Misses |
|------|--------------|
| MAIN | 16           |
| FOO  |              |
| BAR  |              |
| BAZ  |              |

## Metadata

| Profile | User | Platform |
|---------|------|----------|
| P1      |      |          |
| P2      |      |          |

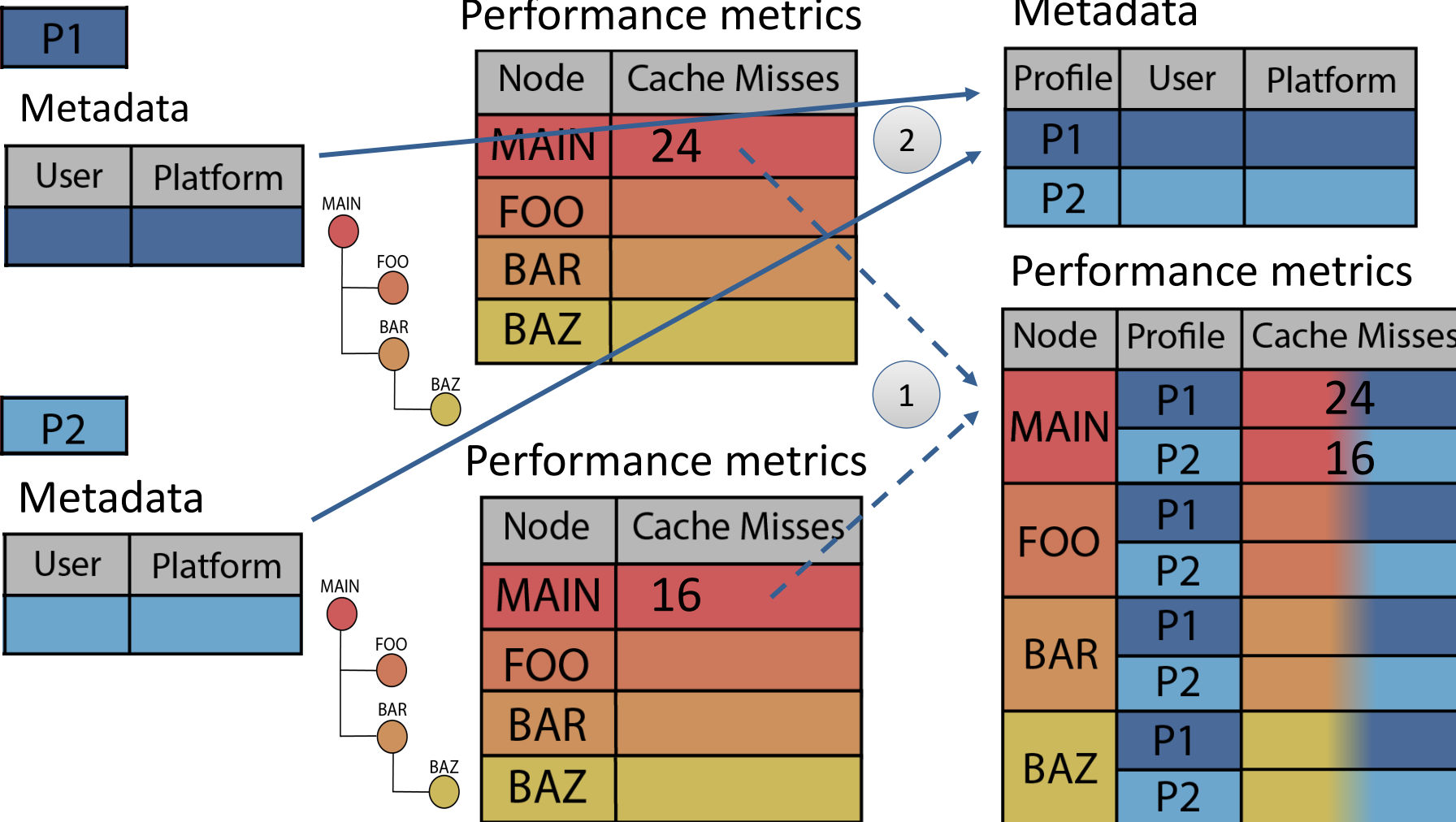## Performance metrics

| Node | Profile | Cache Misses |
|------|---------|--------------|
| MAIN | P1      | 24           |
|      | P2      | 16           |
| FOO  | P1      |              |
|      | P2      |              |
| BAR  | P1      |              |
|      | P2      |              |
| BAZ  | P1      |              |
|      | P2      |              |

① Compose functions w/matching call trees

② Compose metadata with all fields

# Use Thicket to *compose* performance profiles in Python

# Thicket components are *interconnected*

## Metadata

| Profile | User | Platform |
|---------|------|----------|
| P1 | Jon | lassen |
| P2 | Bob | lassen |

## Filtered Metadata

| Profile | User | Platform |
|---------|------|----------|
| P2 | Bob | lassen |

## Performance metrics

| Node | Profile | Cache Misses |
|------|---------|--------------|
| MAIN | P1 | |
| | P2 | |
| FOO | P1 | |
| | P2 | |
| BAR | P1 | |
| | P2 | |
| BAZ | P1 | |
| | P2 | |

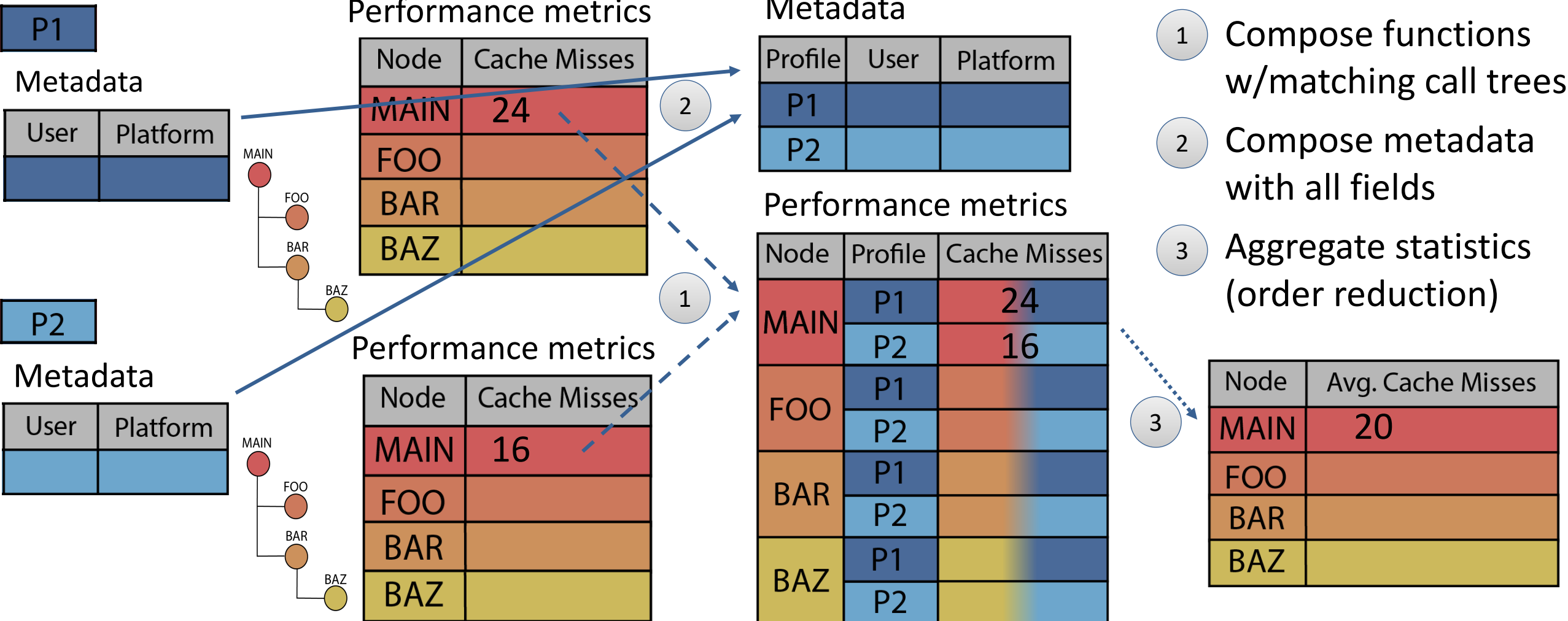Filter on metadata:
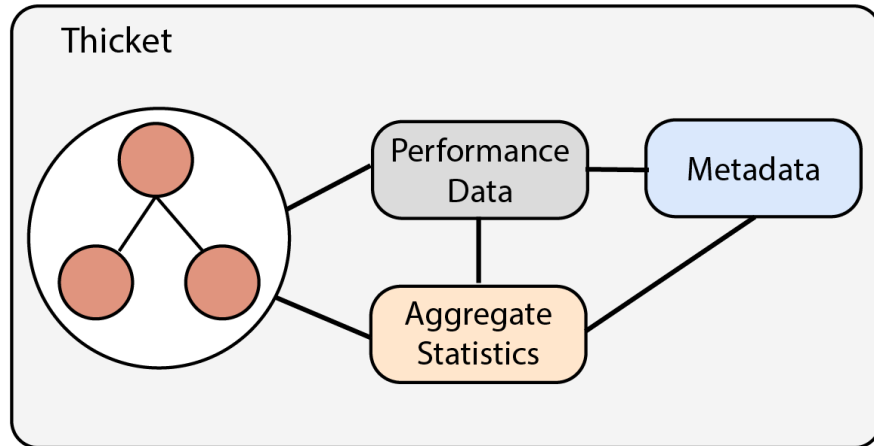platform=="lassen" &&
user=="Bob"

## Filtered Performance metrics

| Node | Profile | Cache Misses |
|------|---------|--------------|
| MAIN | P2 | |
| FOO | P2 | |
| BAR | P2 | |
| BAZ | P2 | |

Metadata fields useful for understanding and manipulating thicket object!

# Thicket enables exploratory data analysis of multi-run data

③ Load Data Into Thicket Object

Thicket

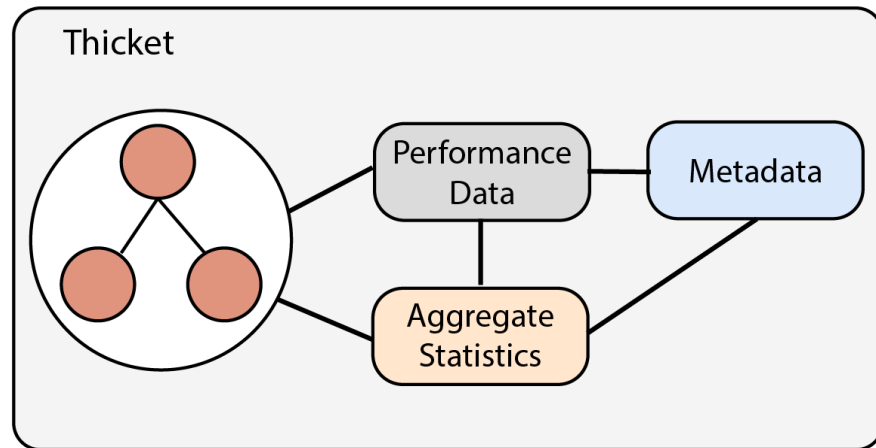Performance Data

Metadata

Aggregate Statistics

- Compose data from diff. sources and types
  - Different scaling (e.g., strong, weak)
  - Different application parameters
  - Different compilers and optimization levels
  - Different hardware types (e.g., CPUs, GPUs)
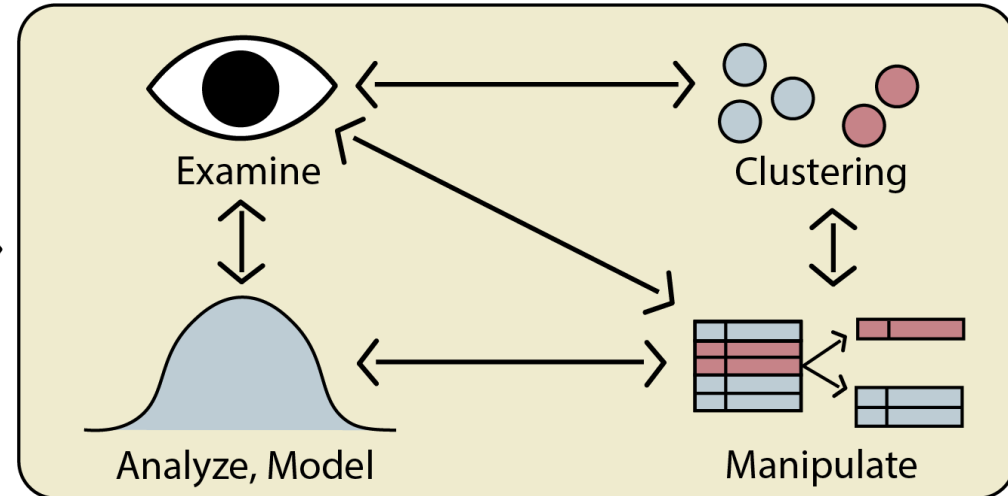  - Different performance tools

# Thicket enables exploratory data analysis of multi-run data



**3** Load Data Into Thicket Object

Thicket

Performance Data — Metadata

Aggregate Statistics

**4** Exploratory Data Analysis (EDA)

Examine

Clustering

Analyze, Model

Manipulate

- Compose data from diff. sources and types
  - Different scaling (e.g., strong, weak)
  - Different application parameters
  - Different compilers and optimization levels
  - Different hardware types (e.g., CPUs, GPUs)
  - Different performance tools

- Perform analysis on the thicket of runs
  - Manipulate the set of data
  - Visualize the dataset
  - Perform analysis on the data
  - Model data
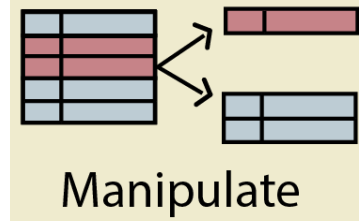  - Leverage third-party tools in the Python ecosystem

# Case Study 1: RAJA Performance Suite

- Open-source suite of loop-based kernels commonly found in HPC applications showcasing performance of different programming models on different hardware

- 560 runs/profiles:
  - 2 clusters (CPU, CPU+GPU)
  - 4 problem sizes
  - 3 compilers, 4 optimizations
  - 3 programming models (sequential, OpenMP, CUDA)
  - 3 performance tools (Caliper, PAPI, Nsight Compute)

| | cluster | systype build | problem size | compiler | compiler optimizations | omp num threads | cuda compiler | block sizes | RAJA variant | #profiles |
|---|---------|---------------|--------------|----------|------------------------|-----------------|---------------|-------------|--------------|-----------|
| 0 | quartz | toss_3_x86_64_ib | [1M, 2M, 4M, 8M] | clang++-9.0.0 | [-O0, -O1, -O2, -O3] | 1 | N/A | N/A | Sequential | 160 |
| 1 | quartz | toss_3_x86_64_ib | [1M, 2M, 4M, 8M] | g++-8.3.1 | [-O0, -O1, -O2, -O3] | 1 | N/A | N/A | Sequential | 160 |
| 2 | quartz | toss_3_x86_64_ib | [1M, 2M, 4M, 8M] | clang++-9.0.0 | -O0 | 72 | N/A | N/A | OpenMP | 40 |
| 3 | quartz | toss_3_x86_64_ib | [1M, 2M, 4M, 8M] | g++-8.3.1 | -O0 | 72 | N/A | N/A | OpenMP | 40 |
| 4 | lassen | blueos_3_ppc64le_ib_p9 | [1M, 2M, 4M, 8M] | xlc++_r-16.1.1.12 | -O0 | 1 | nvcc-11.2.152 | [128, 256, 512, 1024] | CUDA | 160 |

# Use Thicket to *compose* multi-platform, multi-tool data


Manipulate

Thicket object composed of 2 profiles run on CPU

|  |  | time (exc) | Reps | Retiring | Backend bound |
|---|---|---|---|---|---|
| **node** | **problem_size** |  |  |  |  |
| **Apps_NODAL_ACCUMULATION_3D** | 1M | 0.204583 | 100 | 0.144928 | 0.783786 |
|  | 4M | 0.795511 | 100 | 0.139002 | 0.788017 |
| **Apps_VOL3D** | 1M | 0.067061 | 100 | 0.402238 | 0.510525 |
|  | 4M | 0.241508 | 100 | 0.400775 | 0.515976 |

Thicket object composed of 2 profiles run on GPU

|  |  | time (gpu) | gpu__compute_memory_throughput | gpu__dram_throughput | sm__throughput |
|---|---|---|---|---|---|
| **node** | **problem_size** |  |  |  |  |
| **Apps_NODAL_ACCUMULATION_3D** | 1M | 0.007478 | 70.689752 | 46.724767 | 7.330745 |
|  | 4M | 0.026951 | 74.275834 | 51.257993 | 7.688628 |
| **Apps_VOL3D** | 1M | 0.006028 | 81.012826 | 67.751194 | 35.676942 |
|  | 4M | 0.021422 | 91.929933 | 70.122011 | 35.386470 |

CPU                                           GPU

|  |  | time (exc) | Reps | Retiring | Backend bound | time (gpu) | gpu__compute_memory_throughput | gpu__dram_throughput | sm__throughput |
|---|---|---|---|---|---|---|---|---|---|
| **node** | **problem_size** |  |  |  |  |  |  |  |  |
| **Apps_NODAL_ACCUMULATION_3D** | 1M | 0.204583 | 100 | 0.144928 | 0.783786 | 0.007478 | 70.689752 | 46.724767 | 7.330745 |
|  | 4M | 0.795511 | 100 | 0.139002 | 0.788017 | 0.026951 | 74.275834 | 51.257993 | 7.688628 |
| **Apps_VOL3D** | 1M | 0.067061 | 100 | 0.402238 | 0.510525 | 0.006028 | 81.012826 | 67.751194 | 35.676942 |
|  | 4M | 0.241508 | 100 | 0.400775 | 0.515976 | 0.021422 | 91.929933 | 70.122011 | 35.386470 |

# Analyze multi-architecture/multi-tool data

- Dataset: 4 types of profiles side-by-side to compare CPU to GPU performance
  - (1) Basic CPU metrics from Caliper
  - (2) Top-down metrics from Caliper/PAPI
  - (3) GPU runtime from Caliper
  - (4) GPU metrics from Nsight Compute

- Examples of analysis:
  - Compute CPU/GPU speedup
  - Correlate memory and compute usage on the CPU vs. GPU
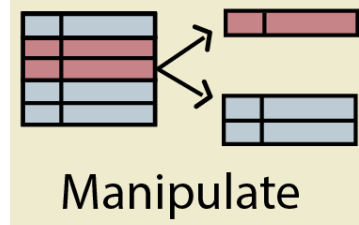
| | | (1) CPU | | | (2) CPU top-down | | (3) GPU | (4) GPU Nsight Compute | | | | Derived |
| | | | | | | | | | | | | |
| Node | Problem size | time (exc) | Bytes/Rep | Flops/Rep | Retiring | Backend bound | time (gpu) | gpu__compute_memory_throughput | gpu__dram_throughput | sm__throughput | sm__warps_active | speedup |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Apps_VOL3D | 8M | 0.498815 | 282109496 | 632421288 | 0.377843 | 0.540604 | 0.040761 | 93.742058 | 72.140428 | 36.206767 | 54.459589 | 12.237556 |
| Lcals_HYDRO_1D | 8M | 2.077556 | 201326600 | 41943040 | 0.032965 | 0.909545 | 0.242928 | 92.944968 | 92.944968 | 6.595714 | 95.266148 | 8.552147 |

# Manipulate: Filter using call path query

```
0.001 Base_CUDA
└─ 0.000 Algorithm
   ├─ 0.000 Algorithm_MEMCPY
   │  ├─ 0.002 Algorithm_MEMCPY.block_128
   │  ├─ 0.009 Algorithm_MEMCPY.block_256
   │  └─ 0.006 Algorithm_MEMCPY.library
   ├─ 0.000 Algorithm_MEMSET
   │  ├─ 0.001 Algorithm_MEMSET.block_128
   │  ├─ 0.004 Algorithm_MEMSET.block_256
   │  └─ 0.003 Algorithm_MEMSET.library
   ├─ 0.000 Algorithm_REDUCE_SUM
   │  ├─ 0.003 Algorithm_REDUCE_SUM.block_128
   │  ├─ 0.004 Algorithm_REDUCE_SUM.block_256
   │  └─ 0.002 Algorithm_REDUCE_SUM.cub
   └─ 0.000 Algorithm_SCAN
      └─ 0.006 Algorithm_SCAN.default
```
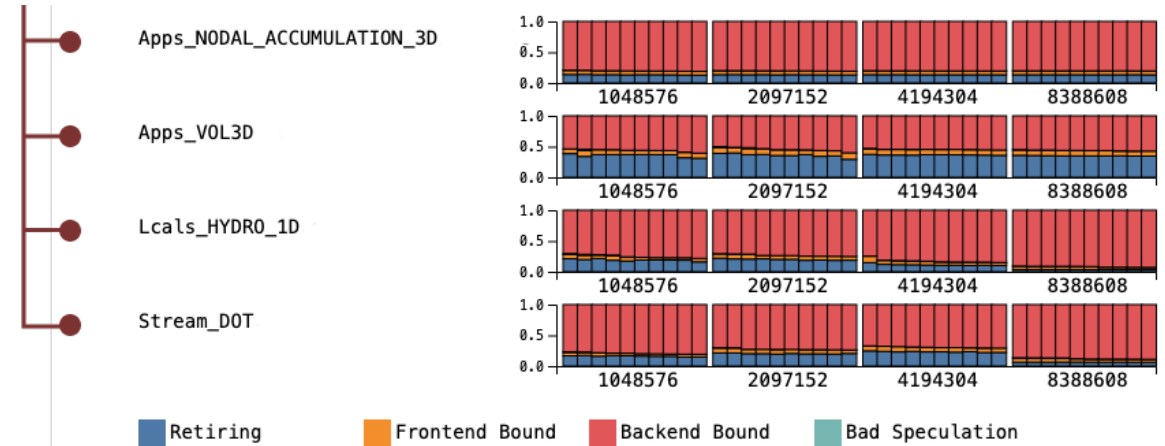
Filter on call path:
(1) Node named
    "Base_CUDA"

```
0.001 Base_CUDA
```

Input call tree

Output call tree

I Lumsden et al. "Enabling Call Path Querying in Hatchet to Identify Performance Bottlenecks in Scientific Applications", e-Science 2022

# Manipulate: Filter using call path query

```
0.001 Base_CUDA
└─ 0.000 Algorithm
   ├─ 0.000 Algorithm_MEMCPY
   │  ├─ 0.002 Algorithm_MEMCPY.block_128
   │  ├─ 0.009 Algorithm_MEMCPY.block_256
   │  └─ 0.006 Algorithm_MEMCPY.library
   ├─ 0.000 Algorithm_MEMSET
   │  ├─ 0.001 Algorithm_MEMSET.block_128
   │  ├─ 0.004 Algorithm_MEMSET.block_256
   │  └─ 0.003 Algorithm_MEMSET.library
   ├─ 0.000 Algorithm_REDUCE_SUM
   │  ├─ 0.003 Algorithm_REDUCE_SUM.block_128
   │  ├─ 0.004 Algorithm_REDUCE_SUM.block_256
   │  └─ 0.002 Algorithm_REDUCE_SUM.cub
   └─ 0.000 Algorithm_SCAN
      └─ 0.006 Algorithm_SCAN.default
```

Input call tree

**Filter on call path:**
(1) Node named "Base_CUDA"
(2) Node with "block_128" in name (and any nodes in between)

```
0.001 Base_CUDA
└─ 0.000 Algorithm
   ├─ 0.000 Algorithm_MEMCPY
   │  └─ 0.002 Algorithm_MEMCPY.block_128
   ├─ 0.000 Algorithm_MEMSET
   │  └─ 0.001 Algorithm_MEMSET.block_128
   └─ 0.000 Algorithm_REDUCE_SUM
      └─ 0.003 Algorithm_REDUCE_SUM.block_128
```

Output call tree

I Lumsden et al. "Enabling Call Path Querying in Hatchet to Identify Performance Bottlenecks in Scientific Applications", e-Science 2022

# Visualize: Intel CPU top-down analysis

- *Top-down analysis* uses HW counters in a hierarchy to identify bottlenecks*

- Use Caliper's top-down module to derive top-down metrics for call-tree regions

- Thicket's *tree+table* visualization shows top-down metrics as stacked bar charts, each bar is a profile
  - Apps_VOL3D has the highest retiring rates
  - Lcals_HYDRO and Stream_DOT become more backend bound as problem size grows

* Yasin, A.: A Top-Down Method for Performance Analysis and Counters Architecture. In: 2014 IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS). pp. 35–44. IEEE, CA, USA (Mar 2014).

# Use third-party Python libraries, e.g., Scikit-learn clustering

Clustering

1. Select data of interest

— Filter 8M problem size

— Use query language to extract all implementations of the Stream kernel

2. (optional) Normalize data

3. Apply scikit-learn clustering to top-down analysis metrics of runs with different compiler optimization levels



**Optimization Level**
- 🔵 -O0
- 🟢 -O1
- 🟡 -O2
- 🔴 -O3

**K-Means Clusters**
- ● 0
- ✖ 1
- ■ 2

**Kernels**
- —— Stream_ADD
- – – – Stream_COPY
- ········ Stream_DOT
- – · – Stream_MUL
- —·· Stream_TRIAD

# Case Study 2: MARBL multi-physics code

- MARBL is a next-generation multi-physics code developed at LLNL

- 60 runs/profiles:
  - 2 clusters (rztopaz, AWS ParallelCluster)
  - 2 MPI libraries (impi, openmpi)
  - 6 node/rank counts
  - 5 repeat runs per config

| | cluster | ccompiler | mpi | version | numhosts | mpi.world.size | #profiles |
|---|---------|-----------|-----|---------|----------|----------------|-----------|
| 0 | ip---- | /usr/tce/packages/clang/clang-9.0.0 | impi | v1.1.0-203-gcb0efb3 | [1, 2, 4, 8, 16, 32] | [36, 72, 144, 288, 576, 1152] | 30 |
| 1 | rztopaz | /usr/tce/packages/clang/clang-9.0.0 | openmpi | v1.1.0-201-g891eaf1 | [1, 2, 4, 8, 16, 32] | [36, 72, 144, 288, 576, 1152] | 30 |

# Manipulate: Compute noise and scaling



1. Use groupby(mpi.world.size) to generate unique subsets of data which are repeated runs; compute noise

2. Compose runs on different platforms and at different scales

3. Generate strong scaling plot with matplotlib
   — Deviation shown in shaded region, dots are average of 5 runs

impi and OpenMPI scale well up to 16 nodes

# Model: Use third-party Python library, Extra-P

Extra-P derives an analytical performance model from an ensemble of profiles covering one or more modeling parameters                 http://github.com/extra-p/extrap

- Select functions of interest
- Call Extra-P to model scaling on different hardware types



**AWS Extra-P Model**

— 154.8848323145599 + -14.012557071778664 * p^(1/3)
● M_solver->Mult

**CTS Extra-P Model**

— 200.23124269331294 + -18.278533682209932 * p^(1/3)
● M_solver->Mult

# Visualize metadata with parallel coordinates plot

- Thicket's interactive parallel coordinates plot shows relationships between metadata variables, and between metadata and performance data
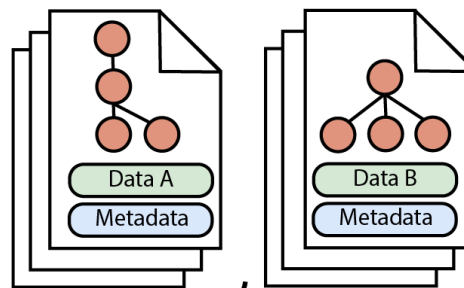
# Thicket is a toolkit for exploratory data analysis of multi-run data



① Run Code with Measurement Tools

② Call Tree Profiles Produced from Multiple Studies

Data A
Metadata
Data B
Metadata

③ Load Data Into Thicket Object

Thicket

Performance Data — Metadata — Aggregate Statistics

④ Exploratory Data Analysis (EDA)

Examine — Clustering
Analyze, Model — Manipulate

Lawrence Livermore National Laboratory
LLNL-PRES-850268

Documentation: **thicket.readthedocs.io**

https://github.com/LLNL/thicket
https://github.com/LLNL/thicket-tutorial

# CASC

Center for Applied
Scientific Computing

# Lawrence Livermore
National Laboratory