# Validating the Performance of GPU-Offloading with Differential Performance Models

Department of Computer Science | Laboratory for Parallel Programming | Alexander Geiß

# Motivation

- When offloading computation to the GPU we expect

    - Things to get faster

    - More energy efficient


- Are these expectations for a specific code fulfilled?

- What about individual application parts?

- What happens when scaling up?

# Workflow

**1** Obtain hardware capabilities

**2** Collect baseline for the CPU-only version
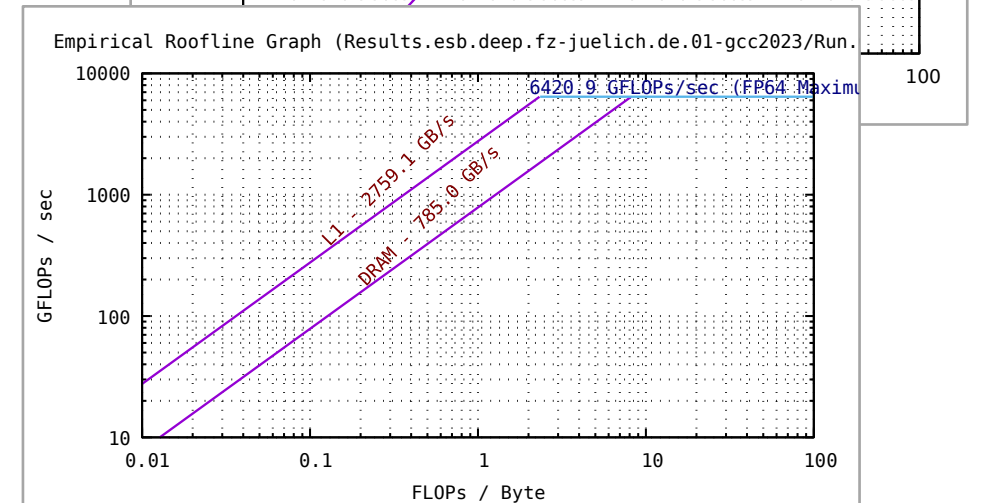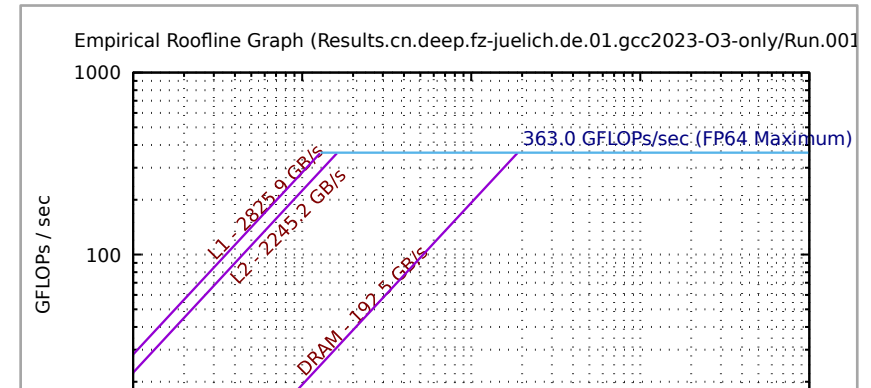
**3** Benchmark the GPU port

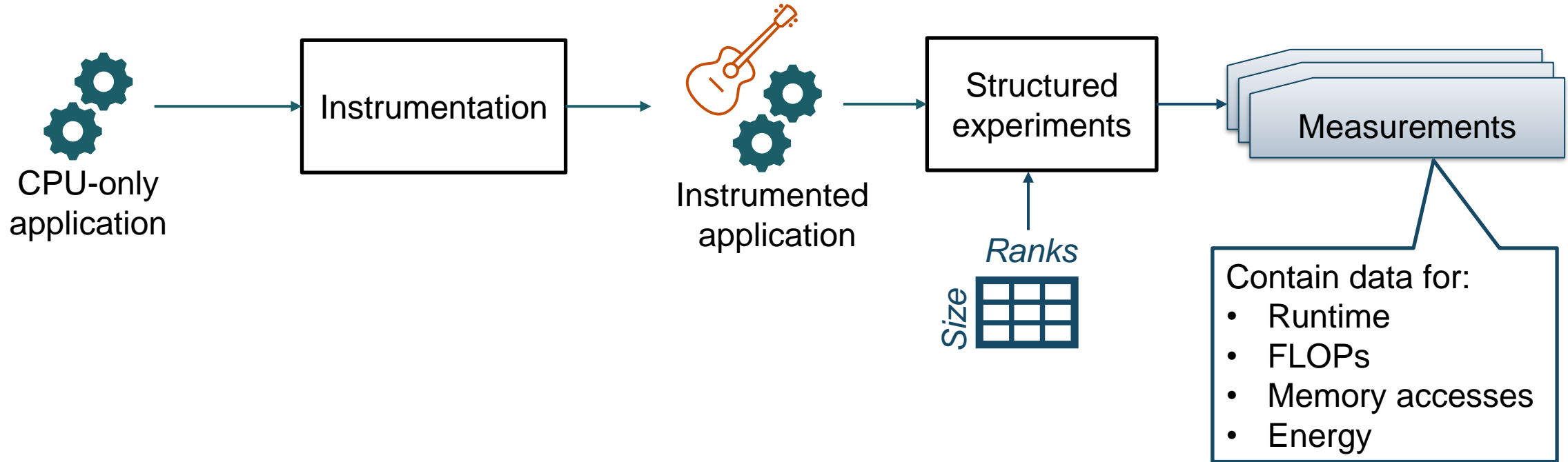**4** Generate performance models and explore them
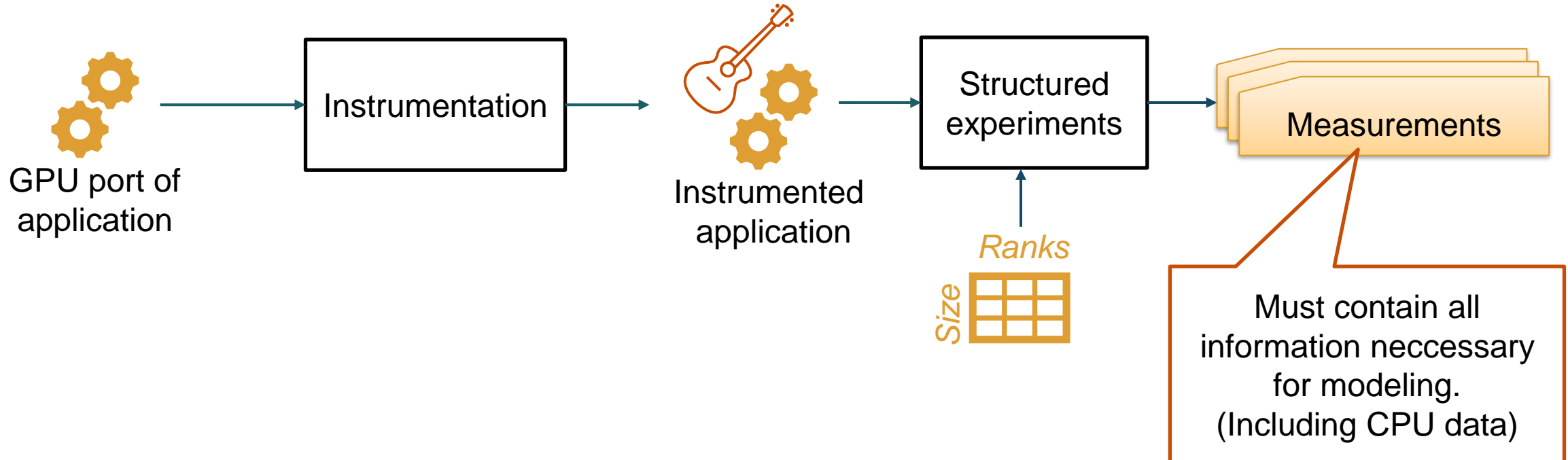
# 1. Obtaining Hardware Capabilities

- Once per system/hardware

- Determining hardware capabilities

  - We use roofline models

    - Easy to create & understand

  - Generated with empirical roofline toolkit

# 2. Baseline Measurements



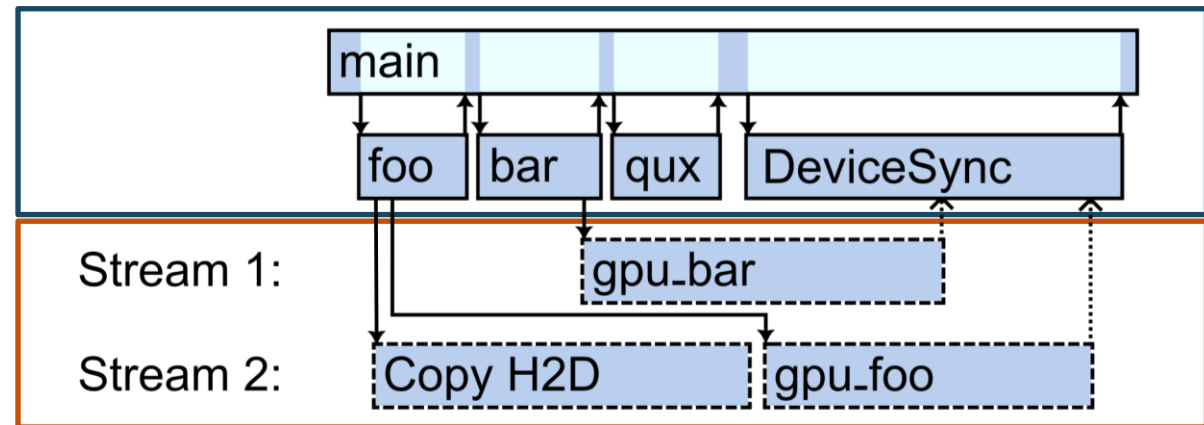CPU-only application → Instrumentation → Instrumented application → Structured experiments

*Ranks*

*Size*

Measurements

Contain data for:
- Runtime
- FLOPs
- Memory accesses
- Energy

# 3. Benchmark the GPU port



GPU port of application → Instrumentation → Instrumented application → Structured experiments → Measurements

*Size* *Ranks*

Must contain all information neccessary for modeling.
(Including CPU data)

# Joint CPU–GPU-Profiling

```
cudaStream_t stream1, stream2;
void foo(double* a_gpu, ...) { ...
  cudaMemcpyAsync(a_gpu, ..., H2D, stream2);
  gpu_foo<<<..., stream2>>>(a_gpu, ...); ...
}


void bar(...) { ...
  gpu_bar<<<..., stream1>>>(...); ...
}


int main() { ...
  foo(a_gpu, ...); ...
  bar(...); ...
  qux(...); ...
  cudaDeviceSynchronize(); ...
}
```

# Joint CPU–GPU-Profiling II

- We collect timestamps on start and end of each function and kernel

  - Using compiler instrumentation on CPU

  - CUPTI Callback API for GPU

- We profile the CPU and trace the GPU events

  - We store only limited trace data

  - Calling context is tracked

  - During post-processing we convert to a unified profile

# Joint CPU–GPU-Profiling III

- This allows recording of

  - Synchronization of GPU actions with CPU functions

  - Overlap of concurrent GPU activities

- Small result files with data necessary for modeling

- Full call path up to kernel

# 4. Differential Performance Modeling

# 4. Differential performance modeling

Call-tree mapping

Defining expectations

Calculating differential models

Checking expectations

# Call-Tree Mapping

- Prefix matching with aggregation

Remember we are using exclusive values

# Defining Expectations

**Faster than CPU-only implementation**

- Runtime of GPU port is lower than runtime of CPU-only version

**Uses less energy**

- Energy usage of GPU port is lower than runtime of CPU-only version

**Uses the hardware well**

- Achieves same or higher hardware efficiency

# Performance Modeling with Extra-P

## Performance model

Performance model:
$f(p)$ = runtime
$f(p) = c \times \sqrt[3]{p}$
where p is the number of processes

Represents performance metric (e.g., execution time or energy consumption) as a function of one or more execution parameters (e.g., the number of processes or problem size)

Watch Extra-P overview video

# Automatic Performance Modeling

```
main() {
    foo()
    bar()
    compute()

}
```

**Instrumentation**

- All functions

Performance measurements

$M_i$     $M_j$



**Extra-P**

**Input**

**Output**

Human-readable
performance models
of all functions
(e.g., $t(p) = c_1 \cdot \log(p) + c_2$)

# Performance Model Normal Form

$$f(x) = \sum_{k=1}^{n} c_k \cdot x^{i_k} \cdot log_2^{j_k}(x)$$

$$\begin{aligned} n &\in \mathbb{N} \\ i_k &\in I \\ j_k &\in J \\ I, J &\subset \mathbb{Q} \end{aligned}$$

$$n = 1$$
$$I = \{0, 1, 2\}$$
$$J = \{0, 1\}$$

$c_1$   $c_1 \cdot \log x$

$c_1 \cdot x$   $c_1 \cdot x \cdot \log x$

$c_1 \cdot x^2$   $c_1 \cdot x^2 \cdot \log x$

# Calculating Differential Models

- All models $f_c^{\mathrm{GPU}}$ and $f_c^{\mathrm{CPU}}$ are mathematical expressions

  - We can calculate with them

**Differential models express the difference**

- $\Delta_c(p_1, \dots) = f_c^{\mathrm{GPU}}(p_1, \dots) - f_c^{\mathrm{CPU}}(p_1, \dots)$

  - For a specific call tree entry $c$

# Hardware Efficiency

$$\text{Hardware efficiency} = \frac{\text{achieved FLOPs per second}}{\text{achievable FLOPs per second}}$$

- Modeled with Extra-P

  - We build models for FLOPs, time, and memory accesses from measurements

    - All models are mathematical expressions

- We present this as a hardware adjusted runtime model to the user

# Interactive Exploration

# Interactive Exploration

# Interactive Exploration



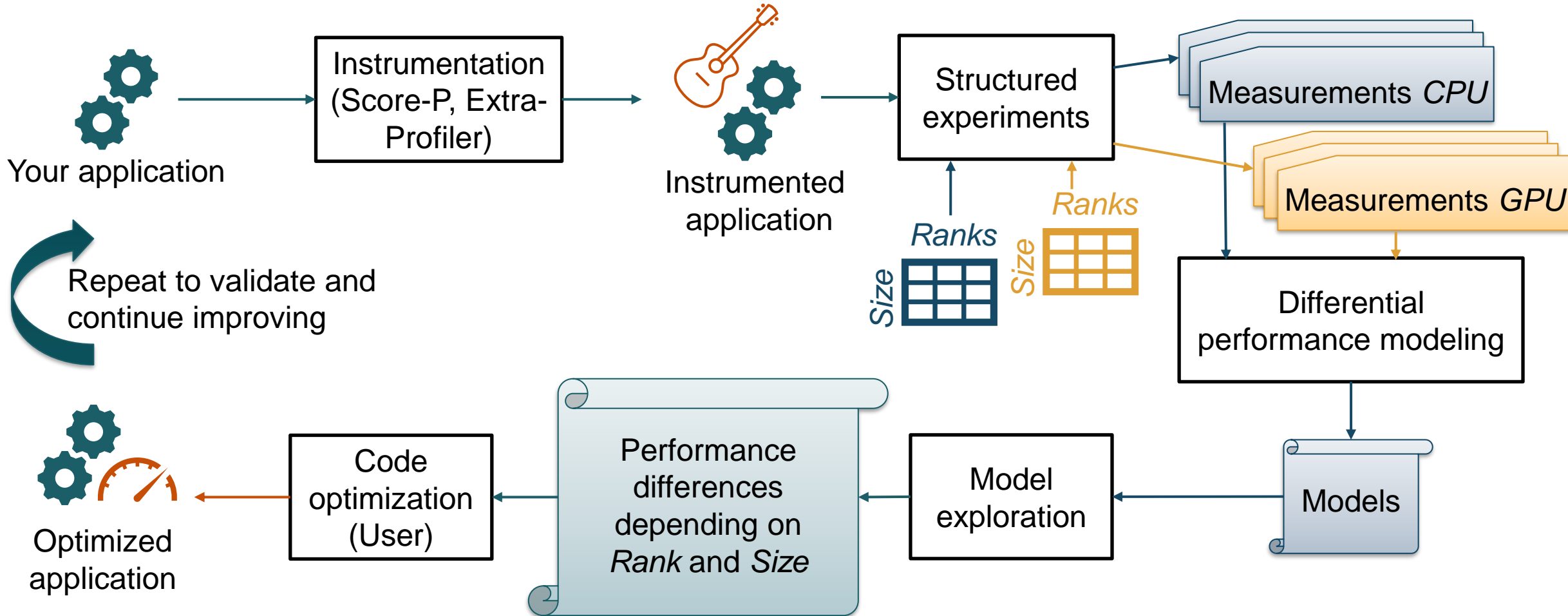| Callpath | Annota | Value |
|---|---|---|
| ▼ main | | $\Delta = 1.079 \times 10^{-20} * n^{7/4} * p^{9/4}$ ... |
| ▶ miniFE::get_parameters | | $\Delta = -1.274 \times 10^{-20} * n^{7/4} - 2.7$... |
| miniFE::initialize_mpi | | $\Delta = 7.598 \times 10^{-4} * p^{1/4} * \log_2$... |
| miniFE::mytimer | | $\Delta = -1.715 \times 10^{-20} * n^{7/4} + 4.$... |
| miniFE::broadcast_para... | | $\Delta = 1.147 \times 10^{-7} * \log_2(n) * \mathsf{l}$... |
| ▶ box_partition | | $\Delta = 3.07 \times 10^{-8} * p * \log_2(p)$ ... |
| YAML_Doc::YAML_Doc | | $\Delta = -4.023 \times 10^{-13} * p^3 * \log_2$... |
| ▶ add_params_to_yaml | | $\Delta = -1.573 \times 10^{-9} * p^{3/4} * \log$... |
| ▶ add_configuration_to_ya... | | $\Delta = 4.832 \times 10^{-8} * p^{4/5} + 6.51$... |
| ▶ add_timestring_to_yaml | | $\Delta = -3.314 \times 10^{-12} * p^3 * \log_2$... |
| ▼ miniFE::driver | | $\Delta = 1.079 \times 10^{-20} * n^{7/4} * p^{9/4}$... |
| ▶ miniFE::compute_im... | | $\Delta = 3.375 \times 10^{-15} * n^{4/3} * \log$... |
| miniFE::mytimer | | $\Delta = 1.485 \times 10^{-12} * n^{1/3} * \log$... |
| ▶ miniFE::create_map_i... | | $\Delta = 1.707 \times 10^{-15} * n^{3/4} * p^{1/4}$... |
| miniFE::find_row_for... | | $\Delta = 1.196 \times 10^{-8} * n^{2/3} * \log_2$... |
| ▶ miniFE::generate_ma... | | $\Delta = -3.91 \times 10^{-19} * n^{7/4} + 1.4$... |
| ▶ miniFE::assemble_FE... | | $\Delta = 4.535 \times 10^{-20} * n^{7/4} + 1.3$... |

Result: GPU port...
- ... is slower than expected
- ... is faster than expected
- ... meets expectation

Result: GPU port is asymptotically...
- ... slower than expected
- ... faster than expected
- ... faster and slower than expected depending on the parameter

Differential performance models

# Summary

# Acknowledgement