

Performance Analysis and Automatic Tuning of Hash Aggregation on GPUs

Viktor Rosenfeld
viktor.rosenfeld@dfki.de

Sebastian Breß
sebastian.bress@tu-berlin.de

Steffen Zeuch
steffen.zeuch@dfki.de

Tilmann Rabl
tilmann.rabl@hpi.de

Volker Markl
volker.markl@tu-berlin.de

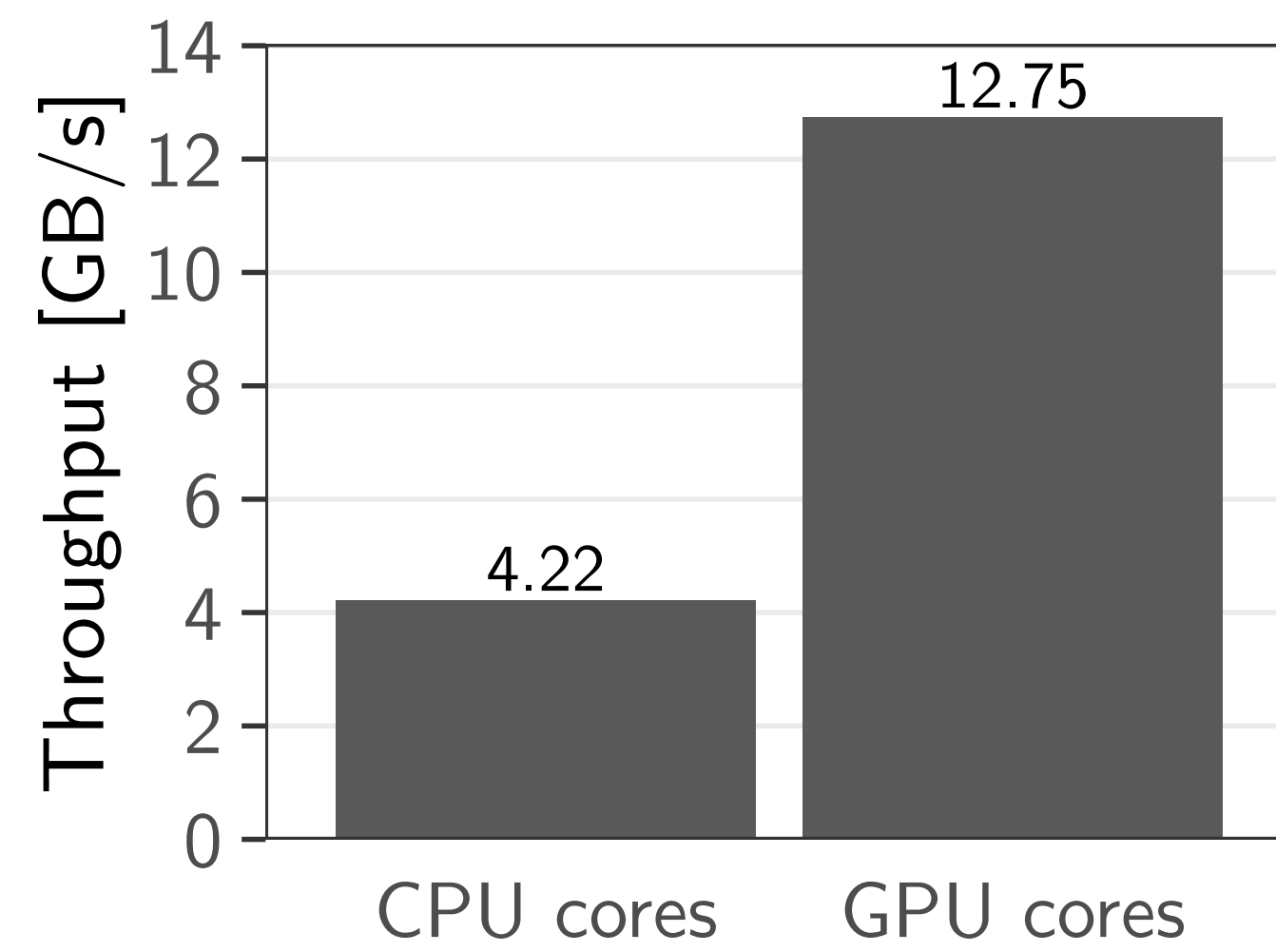
Introduction

Hash aggregation:

- Used to implement GROUP BY and DISTINCT.
- Can be significantly accelerated on GPUs.

Example:

- Query: `SELECT G, sum(V) FROM R GROUP BY G`
- Processor: AMD A10-7850K APU (CPU and GPU cores on same die).
- GPU cores 1.6 – 4.8 times faster.



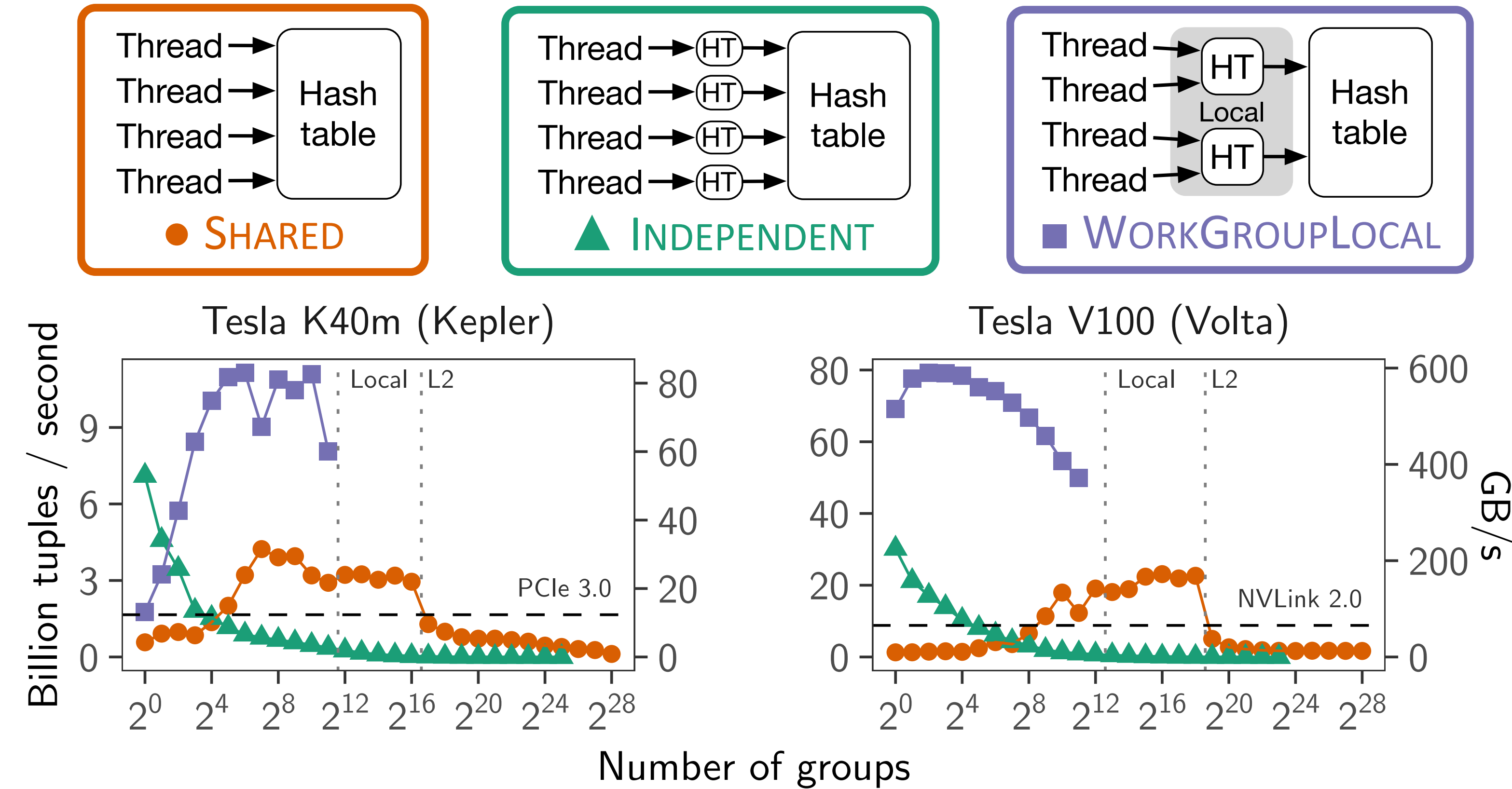
Previous work:

- Analyzed influence of parallelization strategy and thread configuration on aggregation performance.
- Formulated heuristics based on analysis of single NVIDIA Kepler GPU.

Our contributions:

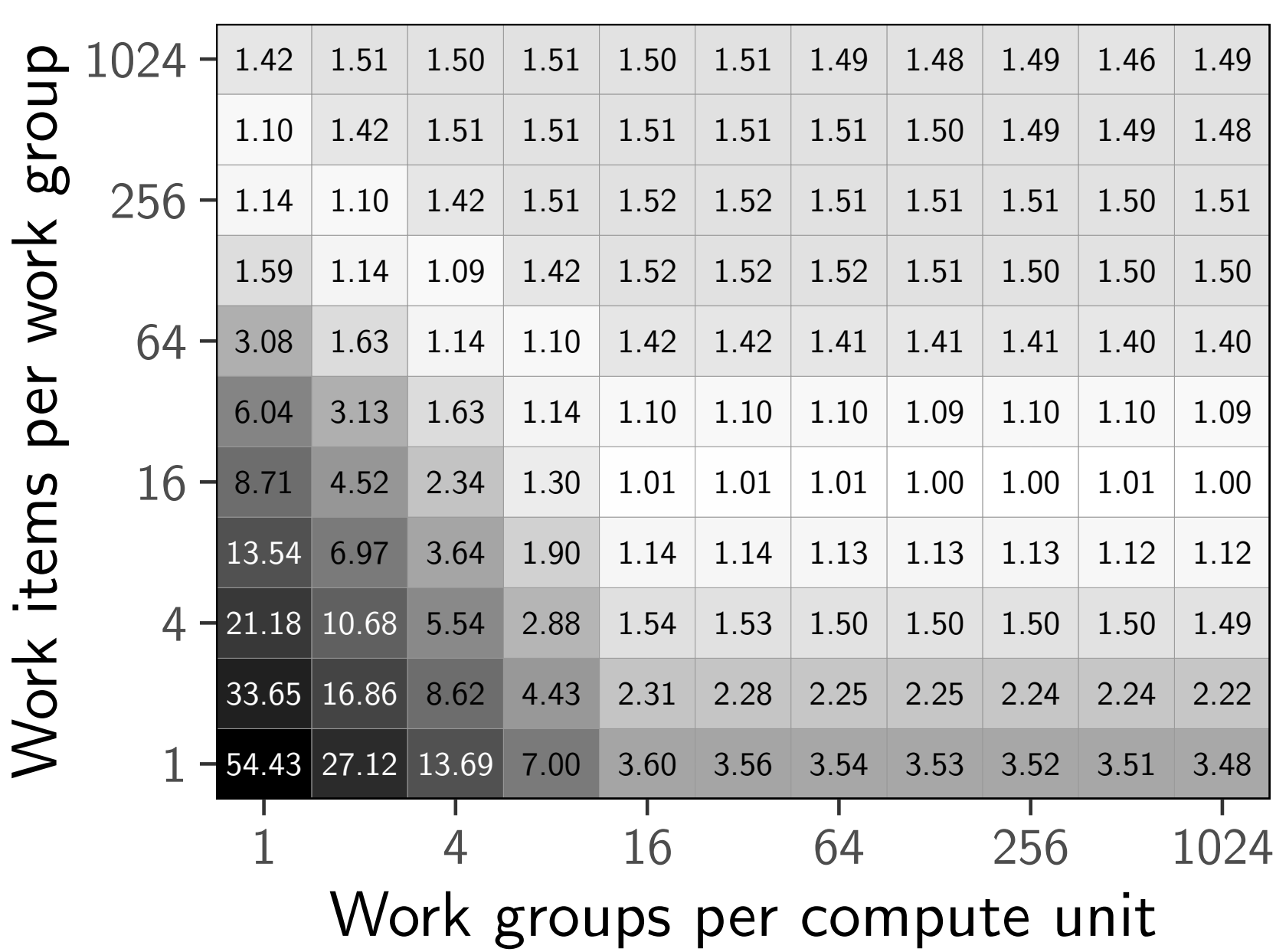
- Performance analysis on six diverse AMD and NVIDIA GPUs.
- Automatic tuning of execution parameters at runtime.

Parallelization strategies



- INDEPENDENT aggregation is not competitive on newer GPUs which implement fast atomics on local memory.
- Aggregation throughput limited by global GPU memory latency (and not transfer bandwidth) when hash table exceeds L2 cache.

Thread configuration search space

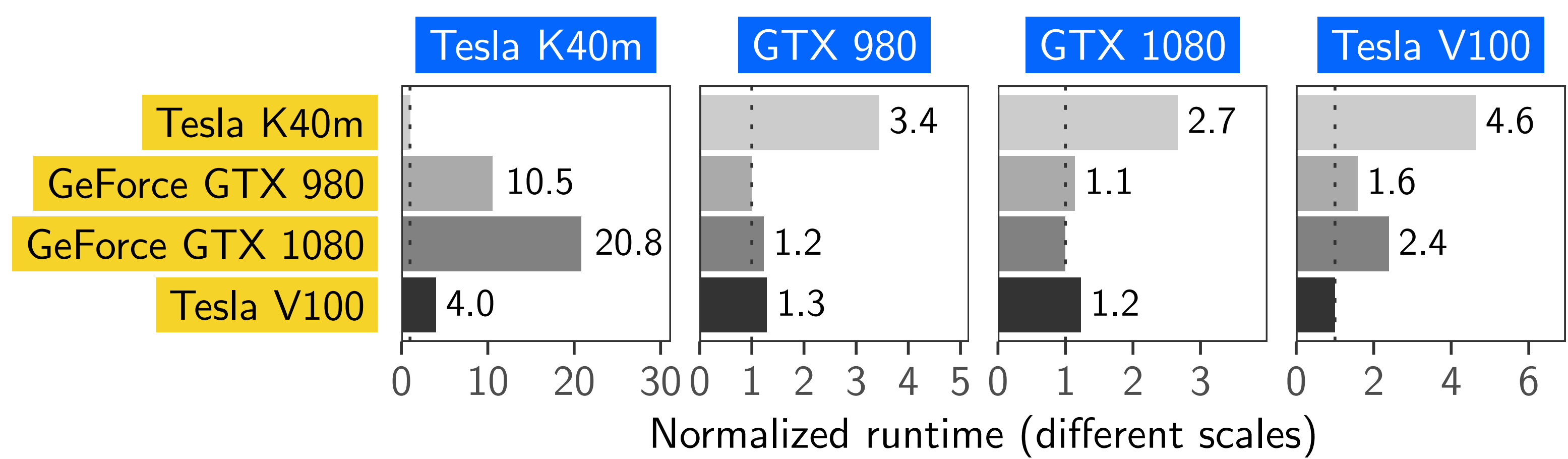


- Number of threads determines execution runtime.
- Performance plateaus:** Adjacent thread configurations have similar runtimes.
- AMD GPUs exhibit high degree of runtime variation.

- Single local minimum when we account for runtime variation.

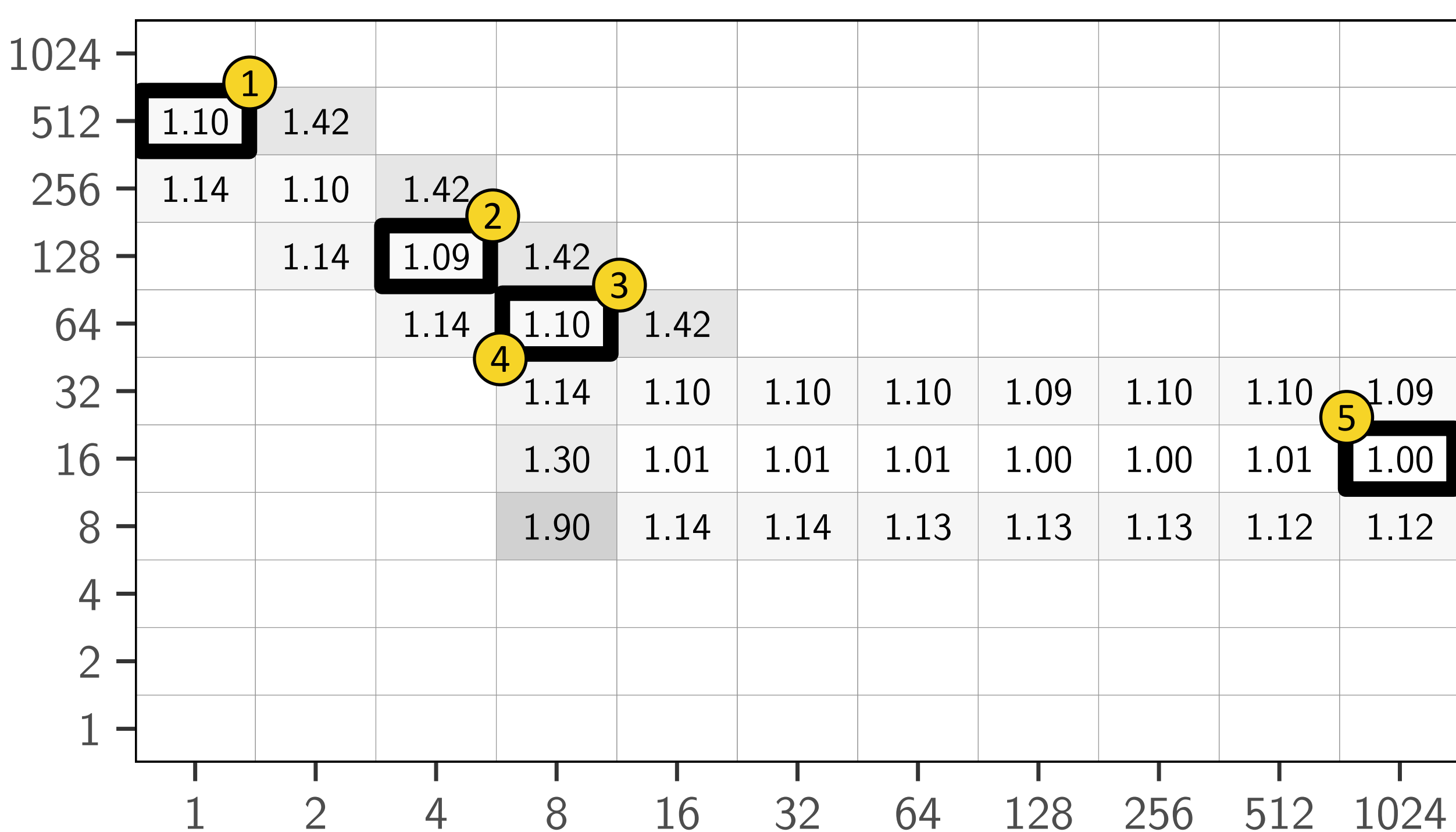
Performance penalty

A configuration optimized for a specific GPU is executed on another GPU.



- The optimal thread configuration is highly GPU-dependent.
- Suboptimal configurations lead to large performance penalties.

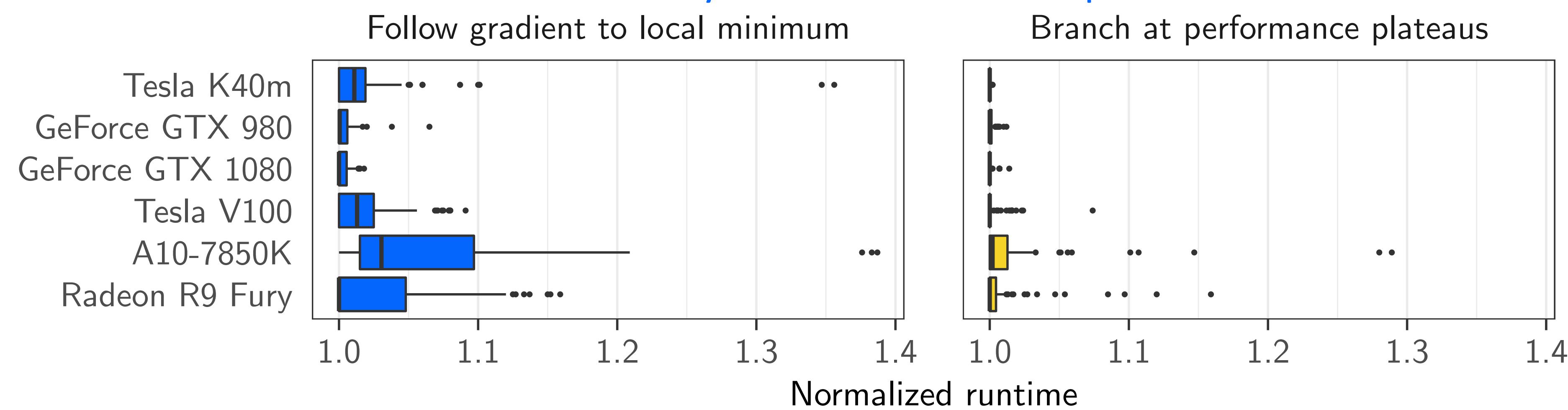
Optimization algorithm



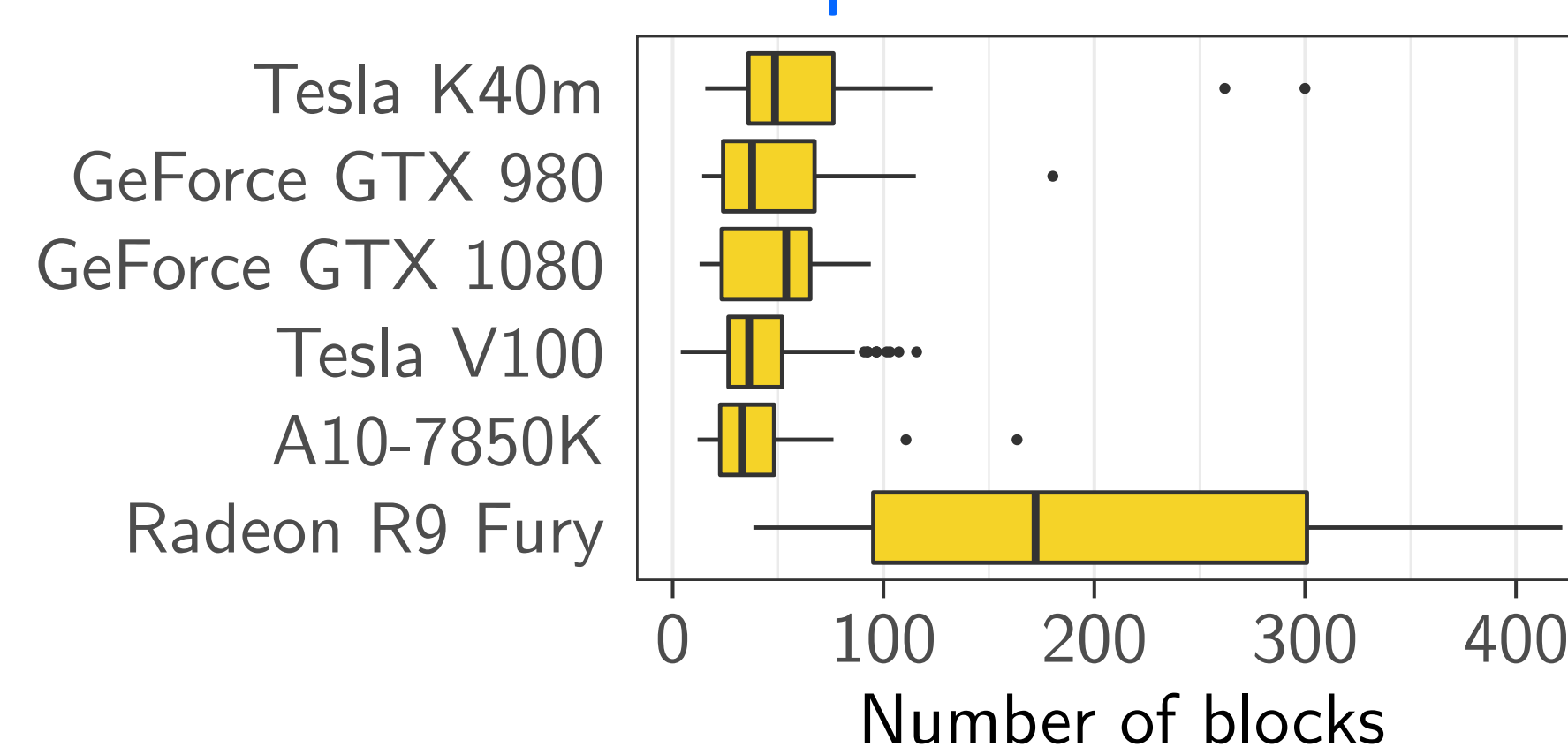
- Start with initial thread configuration.
- Follow gradient in search space to local minimum.
- Branch search path at performance plateaus.
- Prune slow branches.
- Stop at minimum when there are no more branches.

Evaluation

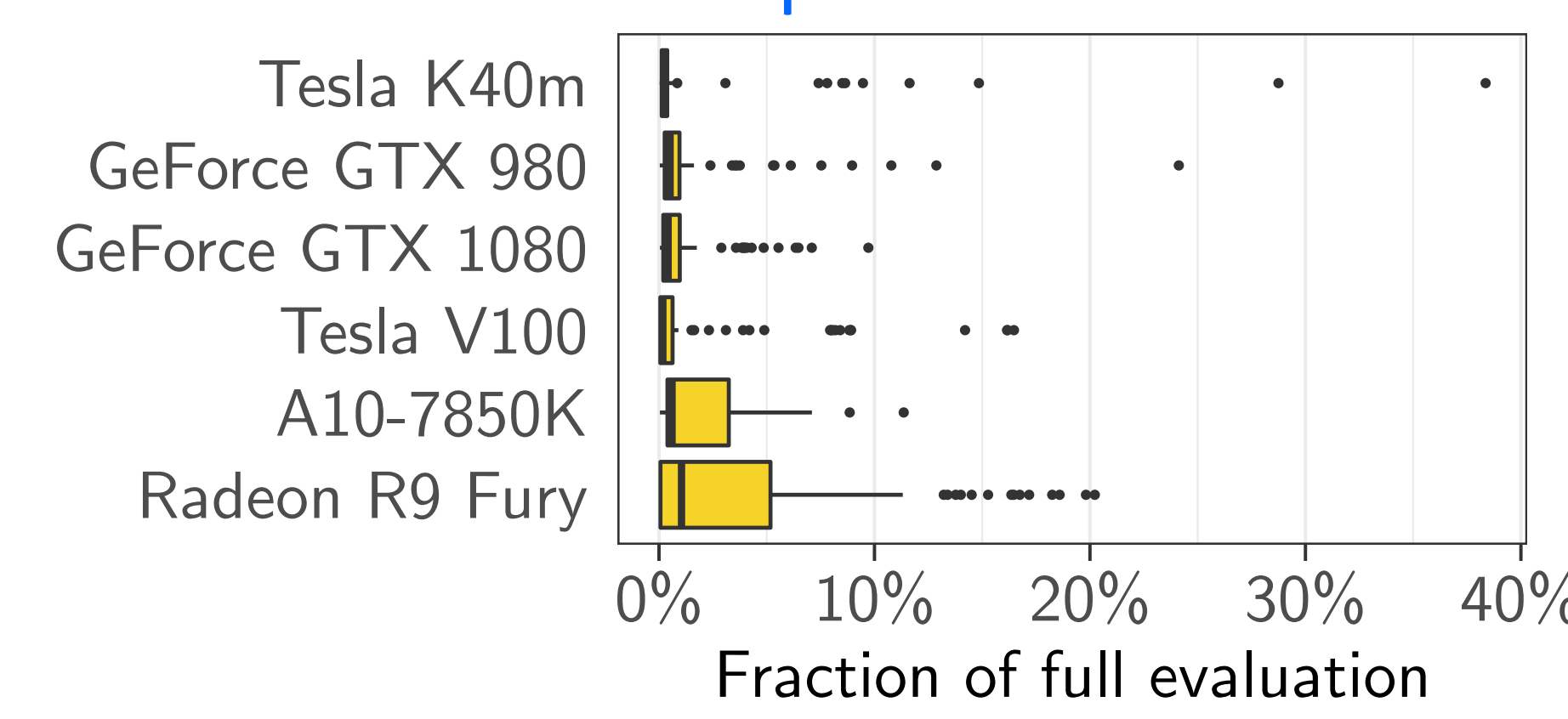
Quality of found execution parameters



Optimization costs



Optimization effort



Main takeaways

- Heuristics to select execution parameters, which are derived from analyzing a single GPU, are not generalizable to other GPUs.
- We can optimize execution parameters by following the gradient in the search space and branching the search at performance plateaus.