# Performance Analysis and Automatic Tuning of Hash Aggregation on GPUs

**Viktor Rosenfeld**, Sebastian Breß, Steffen Zeuch, Tilmann Rabl, Volker Markl
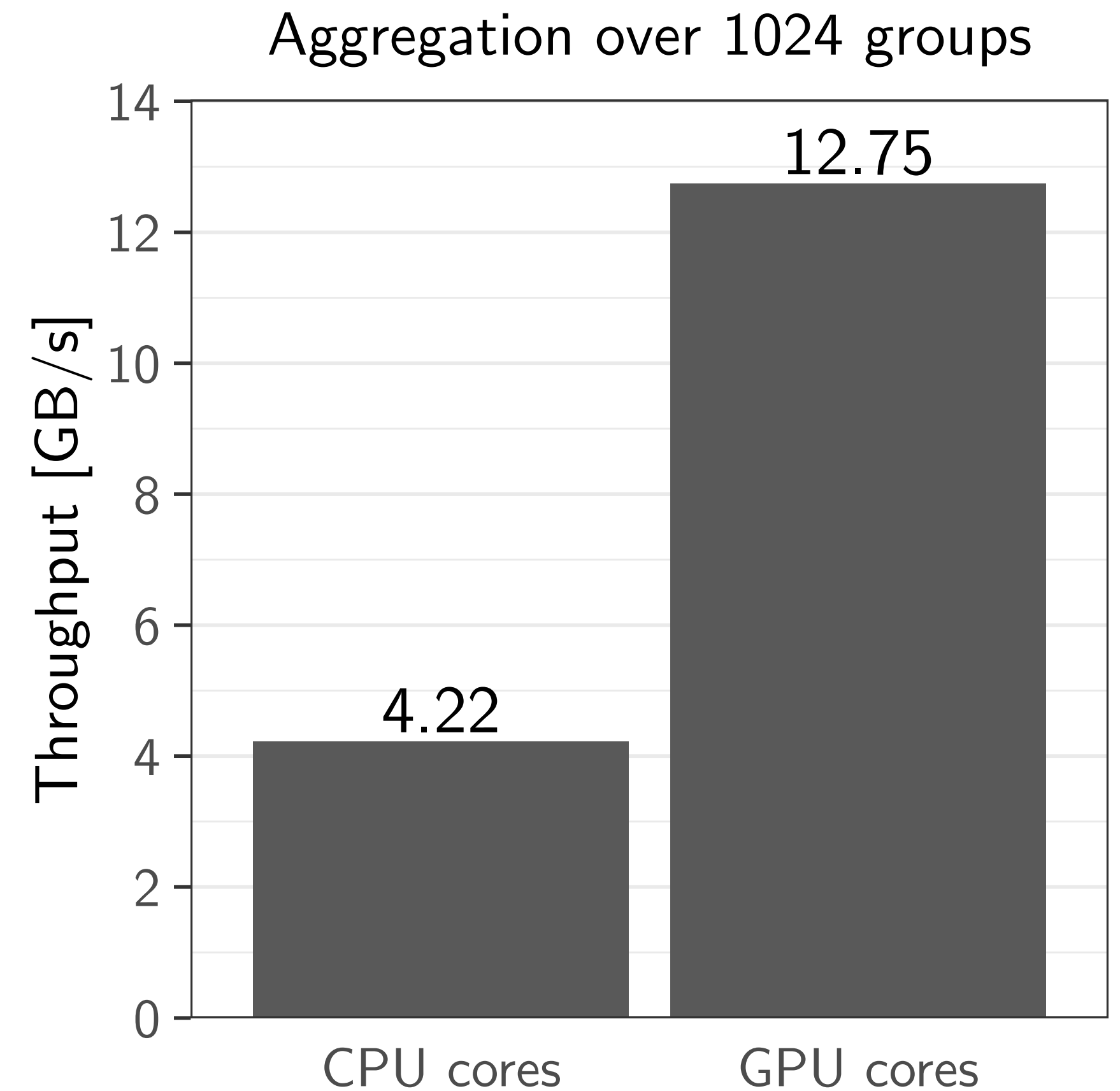
# Hash aggregation on GPUs

**Hash aggregation:**

- Used to implement `GROUP BY` and `DISTINCT`.

- Can be significantly accelerated on the GPU.

**Example:**

- Query:
`SELECT G, sum(A) FROM R GROUP BY G`

- Processor: AMD A10-7850K APU.

  - CPU and GPU cores integrated on the same die.

➡ **Aggregation on GPU cores 1.6× – 4.8× faster across different group cardinalities.**

Aggregation over 1024 groups

Throughput [GB/s]

12.75

4.22

CPU cores    GPU cores

# Previous work

- Hash aggregation extensively studied on CPUs.

- Only a single in-depth study on GPUs:
  Karnagel et al., *Optimizing GPU-accelerated Group-By and Aggregation*, ADMS@VLDB, 2015

  - Evaluated influence of **parallelization strategies** and **thread configuration** based on group cardinality.

  - Formulated heuristics based on a **single NVIDIA Kepler GPU.**

**Do these heuristics yield good performance on other GPUs?**

**Part 1**

# Performance analysis of hash aggregation on various GPUs

# Tested GPUs

| GPU | Microarchitecture | Integration |
|-----|-------------------|-------------|
| Tesla K40m | Kepler | PCIe 3.0 |
| GeForce GTX 980 | Maxwell | PCIe 3.0 |
| GeForce GTX 1080 | Pascal | PCIe 3.0 |
| Tesla V100 | Volta | NVLink 2.0 |
| A10-7850K | Graphics Core Next 2nd Gen. | on die |
| Radeon R9 Fury | Graphics Core Next 3rd Gen. | PCIe 3.0 |

# Parallelization strategies

SHARED



INDEPENDENT



WORKGROUPLOCAL



- Concurrent accesses to same hash bucket resolved with atomics.

- SHARED and INDEPENDENT also commonly used on CPUs.

- WORKGROUPLOCAL uses fast local GPU memory.

- Fastest strategy is data and query dependent (amount of contention and cache efficiency).

**How do these parallelization strategies perform on different GPUs?**

# Parallelization strategies

Parallelization strategy — SHARED — INDEPENDENT — WORKGROUPLOCAL

Tesla K40m (Kepler) — Tesla V100 (Volta)

INDEPENDENT is best

LOCAL is best

Native atomics on local memory (since Maxwell)

Different region

PCIe 3.0 — NVLink 2.0

INDEPENDENT aggregation not competitive on newer GPUs that implement fast atomics on local memory.

# Thread configurations



Work items per work group (y-axis): 1024, 512, 256, 128, 64, 32, 16, 8, 4, 2, 1

Work groups per compute unit (x-axis): 1, 2, 4, 8, 16, 32, 64, 128, 256, 512, 1024

Faster ← → Slower

- Number of threads determined by two variables (OpenCL terminology).

  - *Work groups per compute unit*

  - *Work items per work group*

- Performance influenced by hardware and kernel properties.

  - Warp size, register file, TLB cache, ...

  - Number of registers used in kernel, local memory usage, ...

**Can we find optimal thread configurations across GPUs?**

# Performance penalty

- When a thread configuration optimized for a specific GPU (rows) is executed on another GPU (columns).



Normalized runtime (different scales)

**The optimal thread configuration is highly GPU-dependent.**

**Full evaluation takes hours. How can we find fast thread configurations efficiently?**

# Finding fast thread configurations

# Thread configuration search space

Normalized runtime

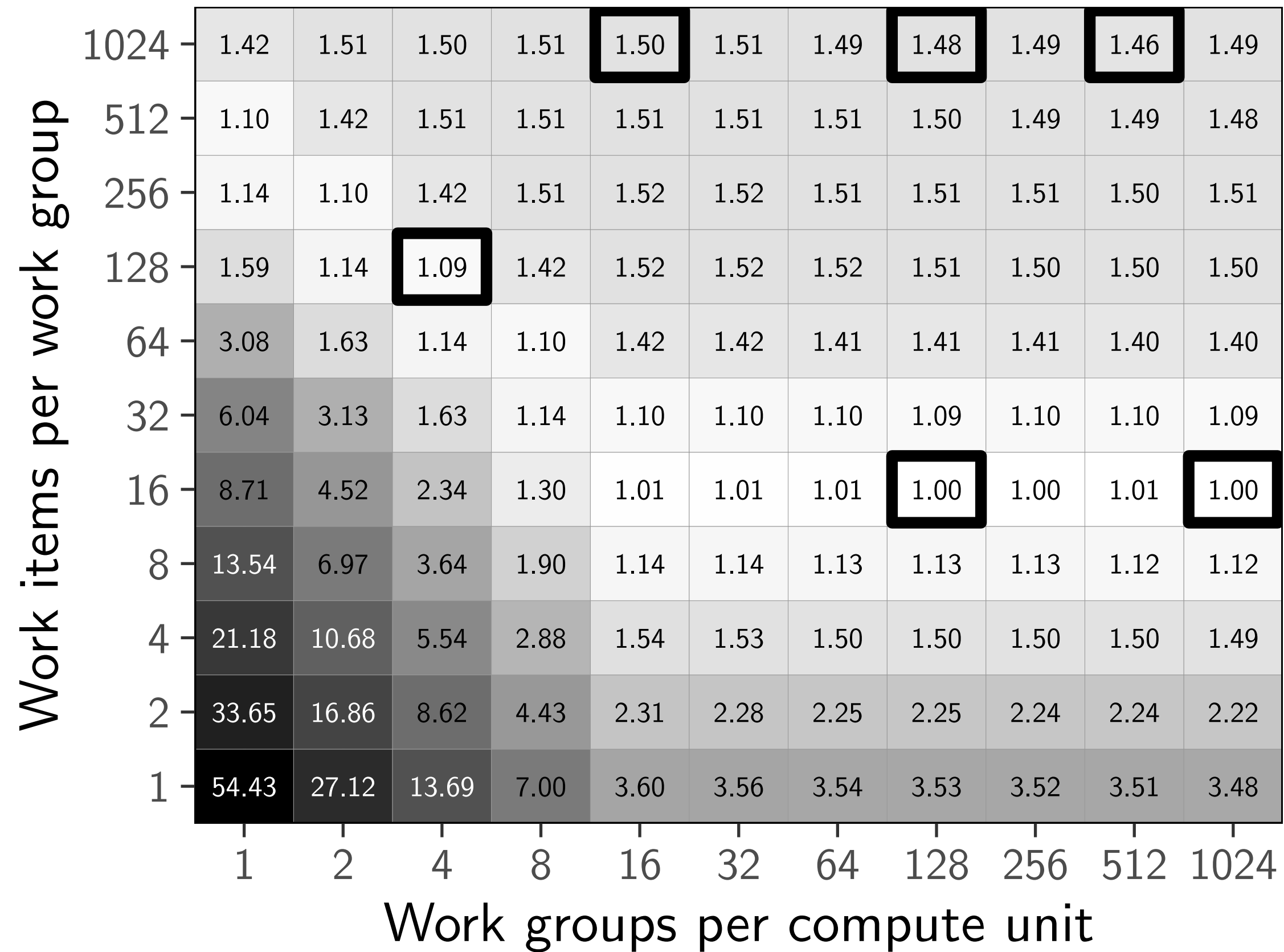| | 1 | 2 | 4 | 8 | 16 | 32 | 64 | 128 | 256 | 512 | 1024 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **1024** | 1.42 | 1.51 | 1.50 | 1.51 | **1.50** | 1.51 | 1.49 | **1.48** | 1.49 | **1.46** | 1.49 |
| **512** | 1.10 | 1.42 | 1.51 | 1.51 | 1.51 | 1.51 | 1.51 | 1.50 | 1.49 | 1.49 | 1.48 |
| **256** | 1.14 | 1.10 | 1.42 | 1.51 | 1.52 | 1.52 | 1.51 | 1.51 | 1.51 | 1.50 | 1.51 |
| **128** | 1.59 | 1.14 | **1.09** | 1.42 | 1.52 | 1.52 | 1.52 | 1.51 | 1.50 | 1.50 | 1.50 |
| **64** | 3.08 | 1.63 | 1.14 | 1.10 | 1.42 | 1.42 | 1.41 | 1.41 | 1.41 | 1.40 | 1.40 |
| **32** | 6.04 | 3.13 | 1.63 | 1.14 | 1.10 | 1.10 | 1.10 | 1.09 | 1.10 | 1.10 | 1.09 |
| **16** | 8.71 | 4.52 | 2.34 | 1.30 | 1.01 | 1.01 | 1.01 | **1.00** | 1.00 | 1.01 | **1.00** |
| **8** | 13.54 | 6.97 | 3.64 | 1.90 | 1.14 | 1.14 | 1.13 | 1.13 | 1.13 | 1.12 | 1.12 |
| **4** | 21.18 | 10.68 | 5.54 | 2.88 | 1.54 | 1.53 | 1.50 | 1.50 | 1.50 | 1.50 | 1.49 |
| **2** | 33.65 | 16.86 | 8.62 | 4.43 | 2.31 | 2.28 | 2.25 | 2.25 | 2.24 | 2.24 | 2.22 |
| **1** | 54.43 | 27.12 | 13.69 | 7.00 | 3.60 | 3.56 | 3.54 | 3.53 | 3.52 | 3.51 | 3.48 |

Work items per work group / Work groups per compute unit

Influence regions of local minima

- Search space: Tesla K40m, SHARED aggregation, ca. 8 million groups.
- Search space appears convex but has multiple local minima.

# Performance plateaus
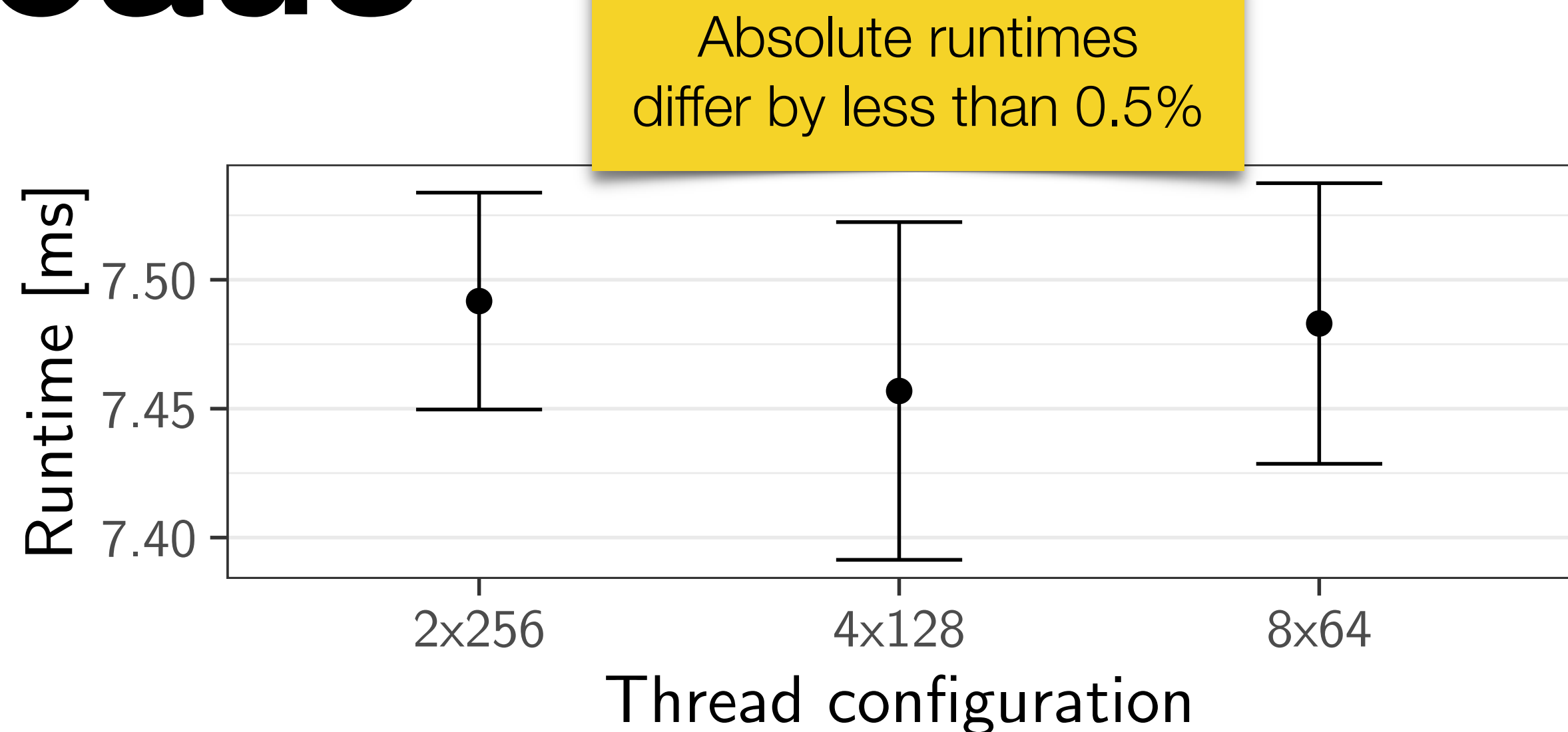
Absolute runtimes differ by less than 0.5%

- **Performance plateau**: Runtimes of two adjacent thread configurations differ by less than a small delta.
- **Nearly convex**: Single local minimum if we account for runtime variation.

**Thread configuration search spaces are nearly convex.**

# Finding fast thread configurations



① Start with initial thread configuration.

② Follow gradient in search space to local minimum.

③ Branch search path at performance plateaus.

④ Prune branches that are slower than the fastest thread configuration found so far.

⑤ Stop at minimum when there are no more branches.

**Approach: Follow gradient and branch search path at performance plateaus.**

# Runtime of found implementation

# Summary

1. INDEPENDENT aggregation is not competitive on newer GPUs that implement fast atomics on local memory. Use WORKGROUPLOCAL aggregation instead.

2. The optimal thread configuration is highly dependent on the executing GPU.

> **Heuristics derived from analyzing a single GPU
> are not generalizable to other GPUs.**

3. Thread configuration search spaces are *nearly convex*, i.e., they have a single local minimum when we account for runtime variation.

> **Follow gradient and branch at performance plateaus
> to find fast thread configurations.**

4. NVIDIA GPUs exhibit a low degree of runtime variation. AMD GPUs exhibit a high degree of runtime variation.

5. Aggregation performance is limited by global GPU memory latency (and not transfer bandwidth) when the hash table exceeds the L2 cache.