# Scalable and Fault-Tolerant Data Stream Processing on Multi-Core Architectures
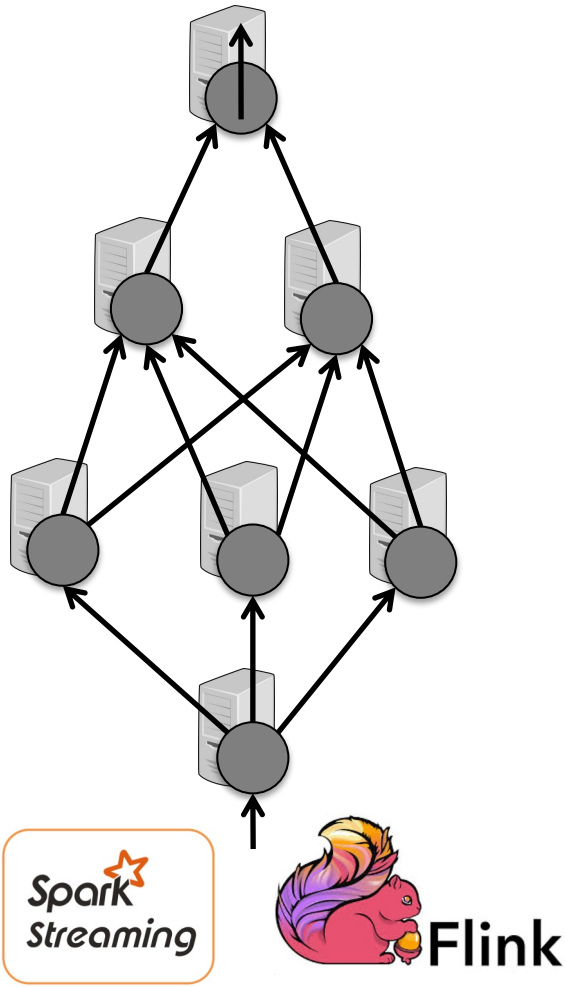
George Theodorakis
Systems Research Group, Neo4j
(was Imperial College London)

# Throughput and Result Freshness Matter



**Data-intensive system**

**Low-latency results**

**High-throughput processing**

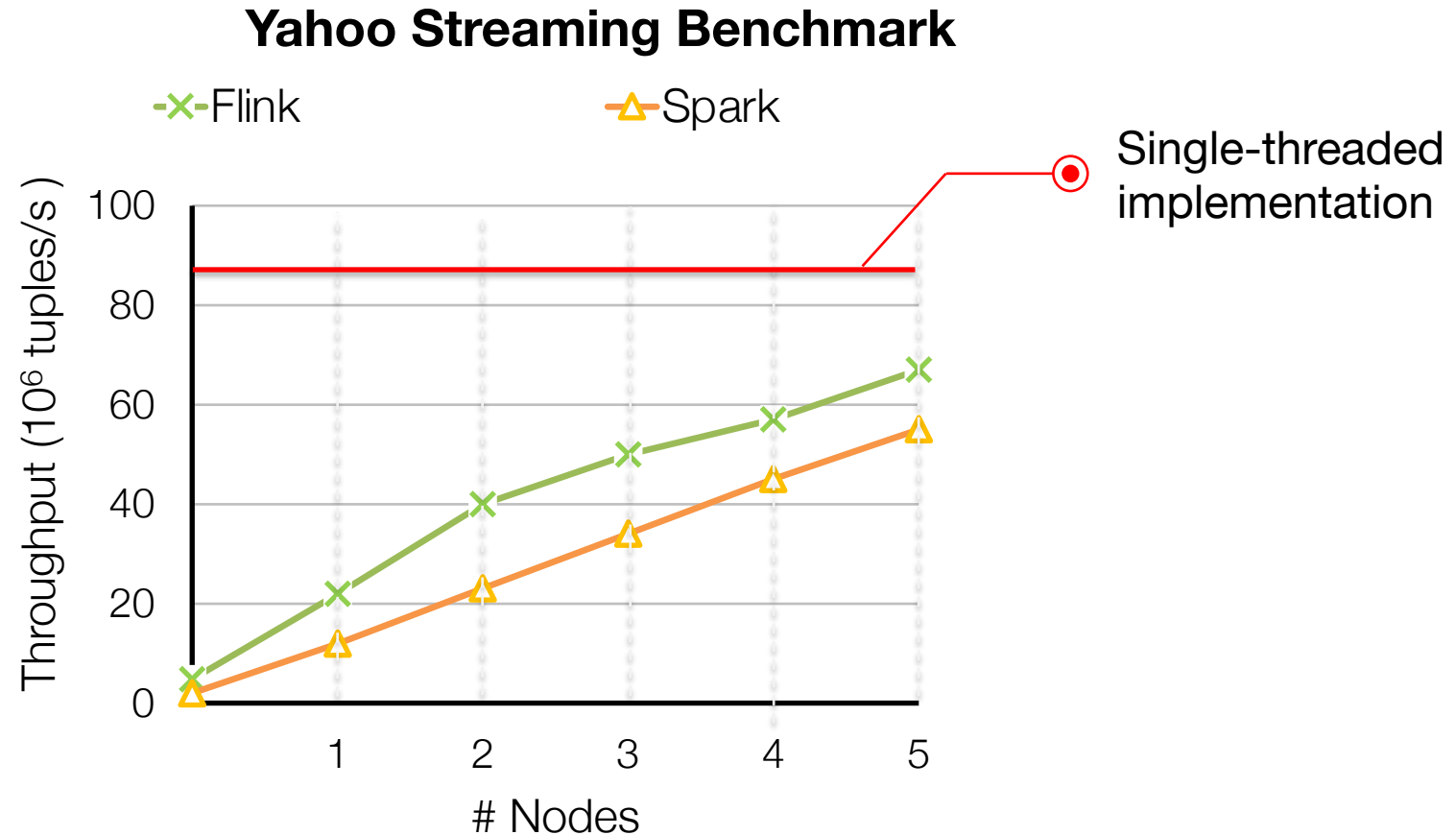| | | |
|---|---|---|
| Facebook Insights: | 12 GB generated content/s | < 10 sec latency |
| Feedzai: | 24M credit card transactions/user | < 10 ms latency |
| Uber: | PB data/day | < 1 ms latency |
| NovaSparks: | 150M trade options/s | < 1 ms latency |

# Distributed Stream Processing Engines Face Challenges
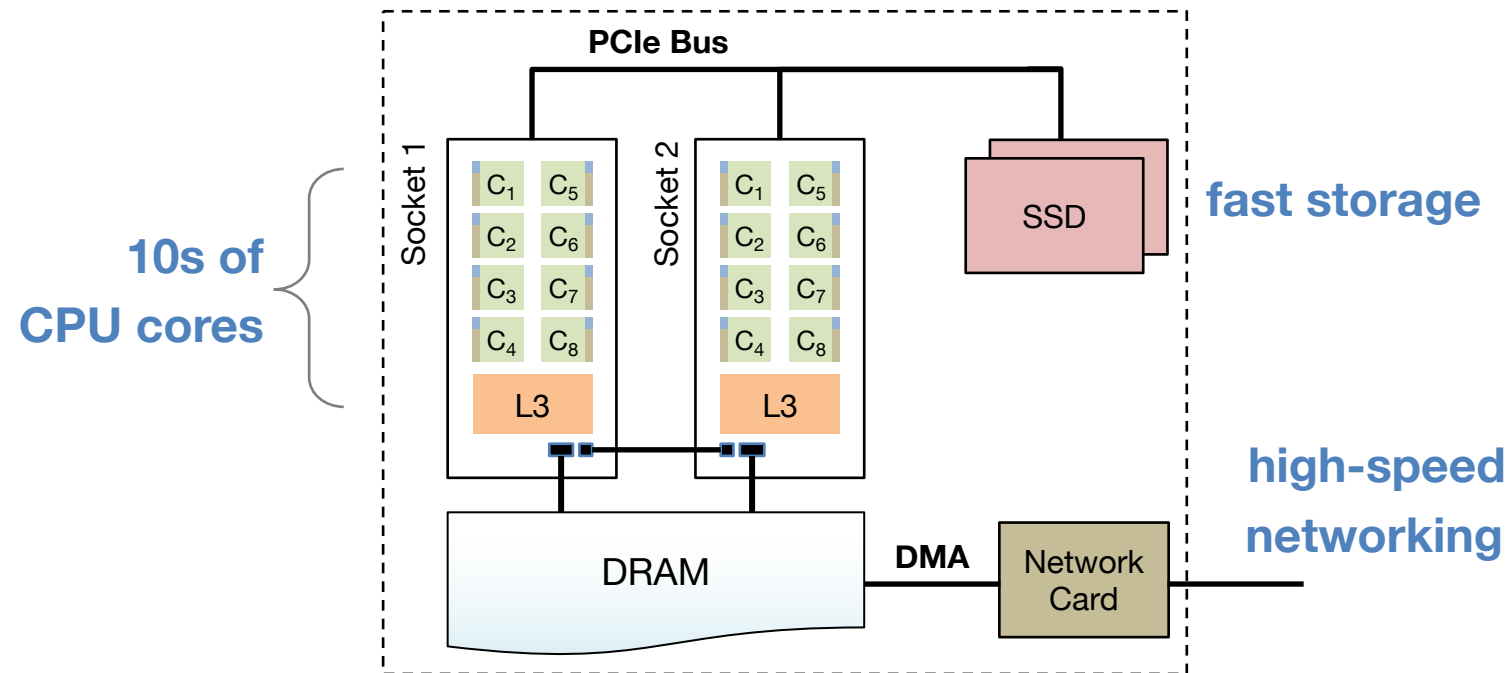


## Pros

> Complex analytics scalability
> Fault-tolerance

## Cons

> Cross-process and network overheads
> Unpredictable latency guarantees
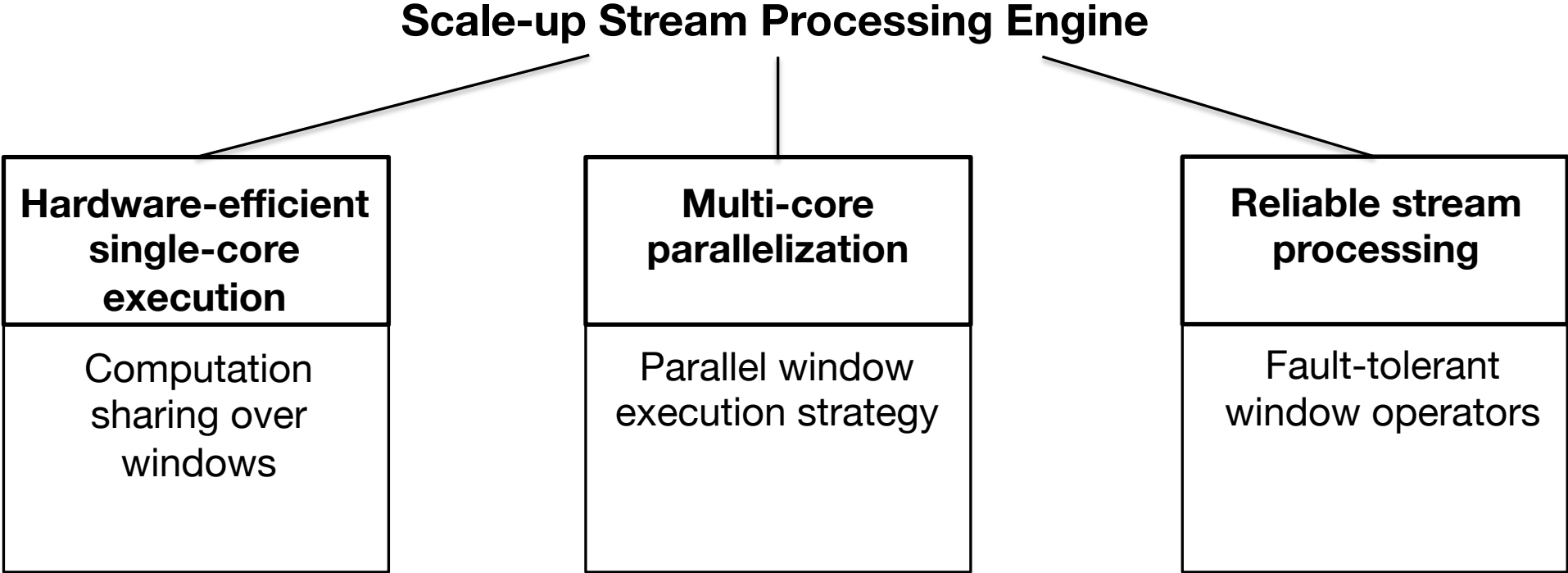> Inefficient execution strategies

# COST of Distributed Execution

**Yahoo Streaming Benchmark**

# Highly-Parallel Scale-up Architectures in Data Centers



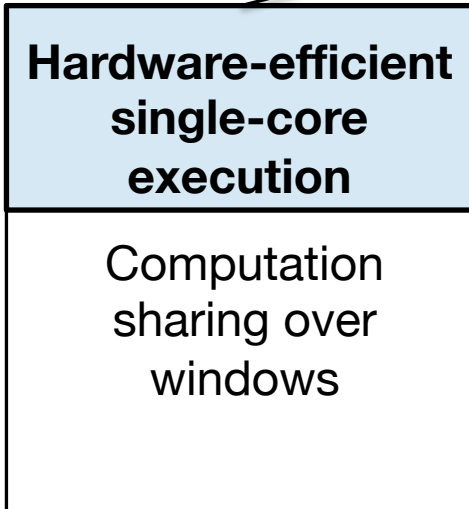👉 **Are scale-up systems a practical alternative for scalability and fault tolerance?**

# High-Performance Streaming and Fault-Tolerance is Hard

**Scale-up Stream Processing Engine**

| **Hardware-efficient single-core execution** | **Multi-core parallelization** | **Reliable stream processing** |
|---|---|---|
| Computation sharing over windows | Parallel window execution strategy | Fault-tolerant window operators |

# Scale-up Stream Processing Engine

| Hardware-efficient single-core execution |
| :---: |
| Computation sharing over windows |

# Scale-up Stream Processing Engine

| Hardware-efficient single-core execution |
| :---: |
| Computation sharing over windows |
| - - - - - - - - |
| **HammerSlide**<br><br>**SlideSide** |

| Multi-core parallelization |
| :---: |
| Parallel window execution strategy |

# Scaling Window Operators on Multi-Core Processors

**Performance** = **Incremental Computation** + **Parallelism**

# Tension Between Parallelism & Incremental Computation



**Tumbling Windows**

Nothing to optimize

**Sliding Windows**

Parallel Execution

Incremental Execution

+ Parallel

- Work Efficient

- Sequential

+ Work Efficient

**Performance** = **Incremental Computation** + **Parallelism**

# Existing System Implement Ad-Hoc Solutions

**Conflicting Objectives**

**Parallelism**

**Incremental Computation**

**Performance**

# Let's Double the Window Slide!

# Let's Double the Window Slide!

# Two Sides of the Same Coin

Partial Aggregates



Sashes

Panes

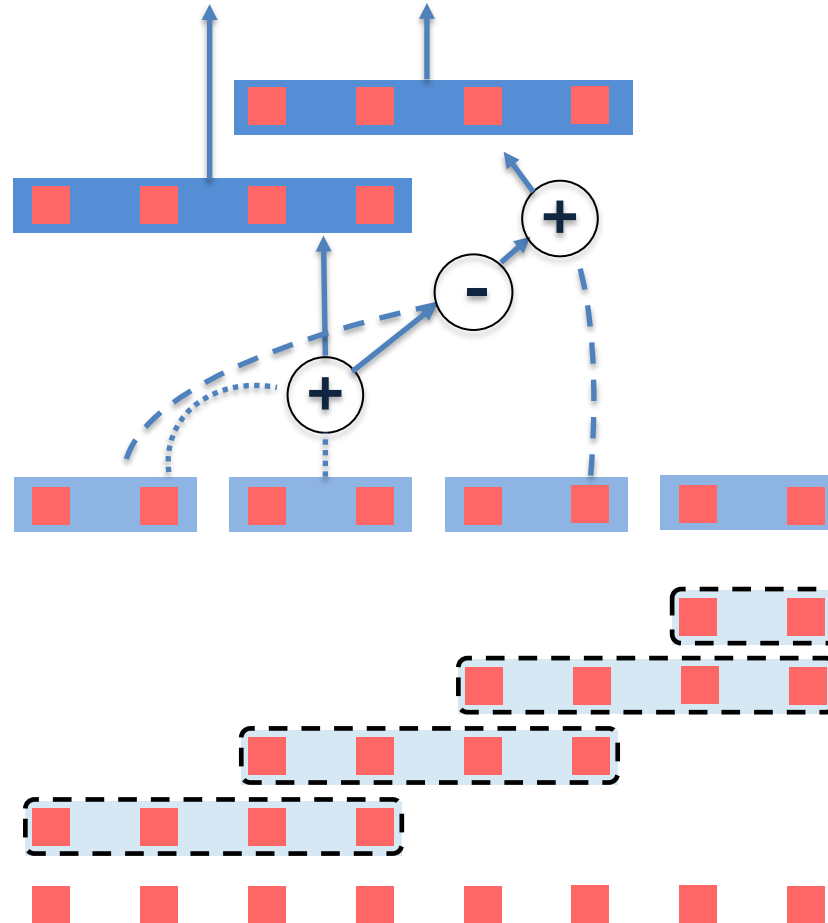Incremental Computation

Data-dependent computation

No data dependencies

Parallel Computation

☛ **How to partition streams into intermediate steps?**

# Create a Model That Splits Aggregation Into Steps



Sequential

Parallel

Incremental

Design choice

# LightSaber: Combine Parallelism With Incremental Execution

# Parallel Aggregation Tree: Multi-level Window Aggregation



Window Result

Sash Merge

Window Fragments

Sash Merge

Sash Merge

Sash Merge

Sash Merge

Pane Merge

Pane Merge

Tuple

Tuple

**Multi-core parallel**

**Incremental**

**SIMD-parallel**

# How to Generate Efficient Code for Incremental Execution

Incremental Algorithms

Pane Merge

SoE   SlideSide
Two-Stacks
SlickDeque

**op** open proceedings

**SlideSide: A fast Incremental Stream Processing Algorithm for Multiple Queries**

Georgios Theodorakis        Peter Pietzuch        Holger Pirk

**?**

No best approach

Generate workload- & query- specific code

# General Aggregation Graph: Capture Low-Level Dependencies

> Aggregation functions
> Window Types



```cpp
int leafIter = 0;
for (auto &t: input) {
  if (leafIter == WINDOW_SIZE) {
    for (int i = 0;i<WINDOW_SIZE;++i)
      s[i+1] = min(ss[i],
                l[WINDOW_SIZE-1-i]);
    leafIter = 0; ps = INT_MAX;
  }
  ps = min(ps, t);
  emit_result(min(ps,
            ss[WINDOW_SIZE-leafIter]));
  leafIter++;
}
```

# Efficient Multi-core Execution



5x to more than one order of magnitude better throughput

Throughput ($10^6$ tuples/s)

Legend:
- Flink
- Scotty
- SABER
- LightSaber

Values above bars:
- 28552 MB/s
- 28042 MB/s
- 10971 MB/s
- 1927 MB/s
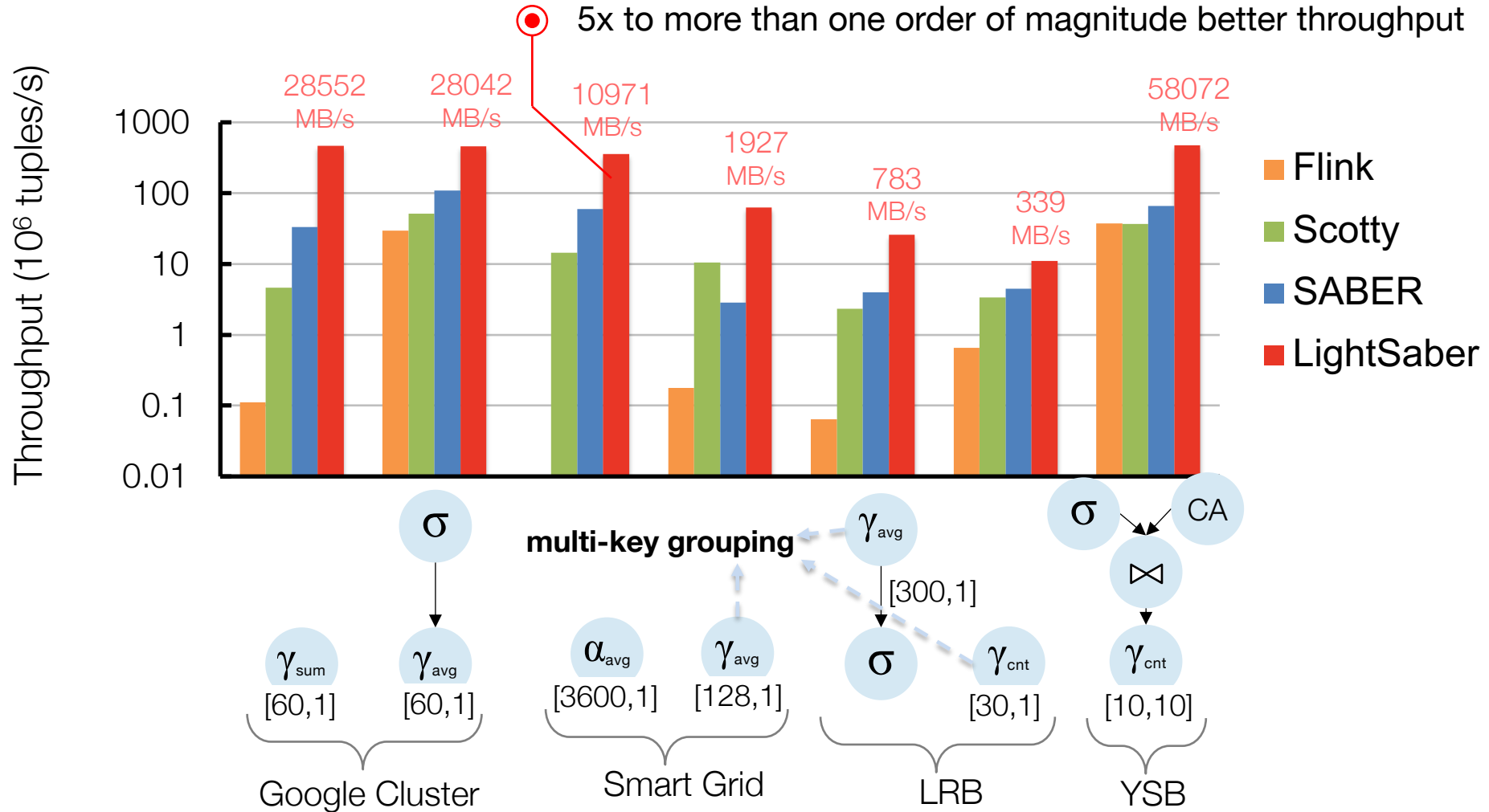- 783 MB/s
- 339 MB/s
- 58072 MB/s

Operator diagrams:

$\gamma_{sum}$ [60,1], $\gamma_{avg}$ [60,1] — Google Cluster (with $\sigma$ above $\gamma_{avg}$)

$\alpha_{avg}$ [3600,1], $\gamma_{avg}$ [128,1] — Smart Grid

multi-key grouping, $\gamma_{avg}$ [300,1], $\sigma$, $\gamma_{cnt}$ [30,1] — LRB

$\sigma$, CA, ⋈, $\gamma_{cnt}$ [10,10] — YSB

# Scale-up Stream Processing Engine

| **Hardware-efficient single-core execution** |
|---|
| Computation sharing over windows |
| |
| **HammerSlide** |
| **SlideSide** |

| **Multi-core parallelization** |
|---|
| Parallel window execution strategy |
| |
| LIGHTSABER |

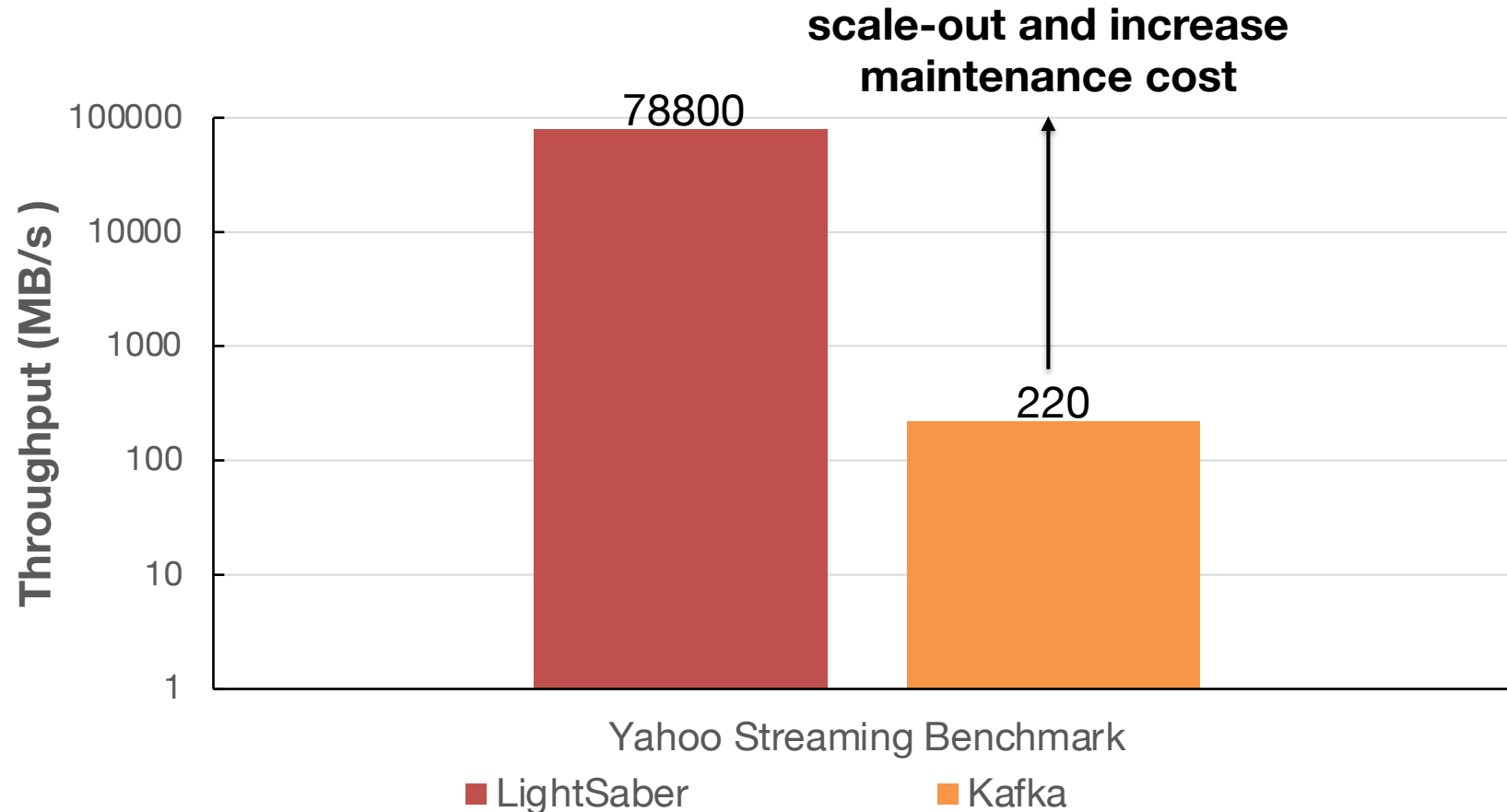| **Reliable stream processing** |
|---|
| Fault-tolerant window operators |

# Scale-up Engines Have Limited Adoption due to Lack of Built-in Fault-Tolerance



> Fault-tolerance requires persisting data from queries
> Persistence is offloaded to external systems

# Kafka Ingestion Trails Scale-up Performance



**scale-out and increase maintenance cost**

Chart: Throughput (MB/s), log scale y-axis from 1 to 100000, Yahoo Streaming Benchmark

- LightSaber: 78800
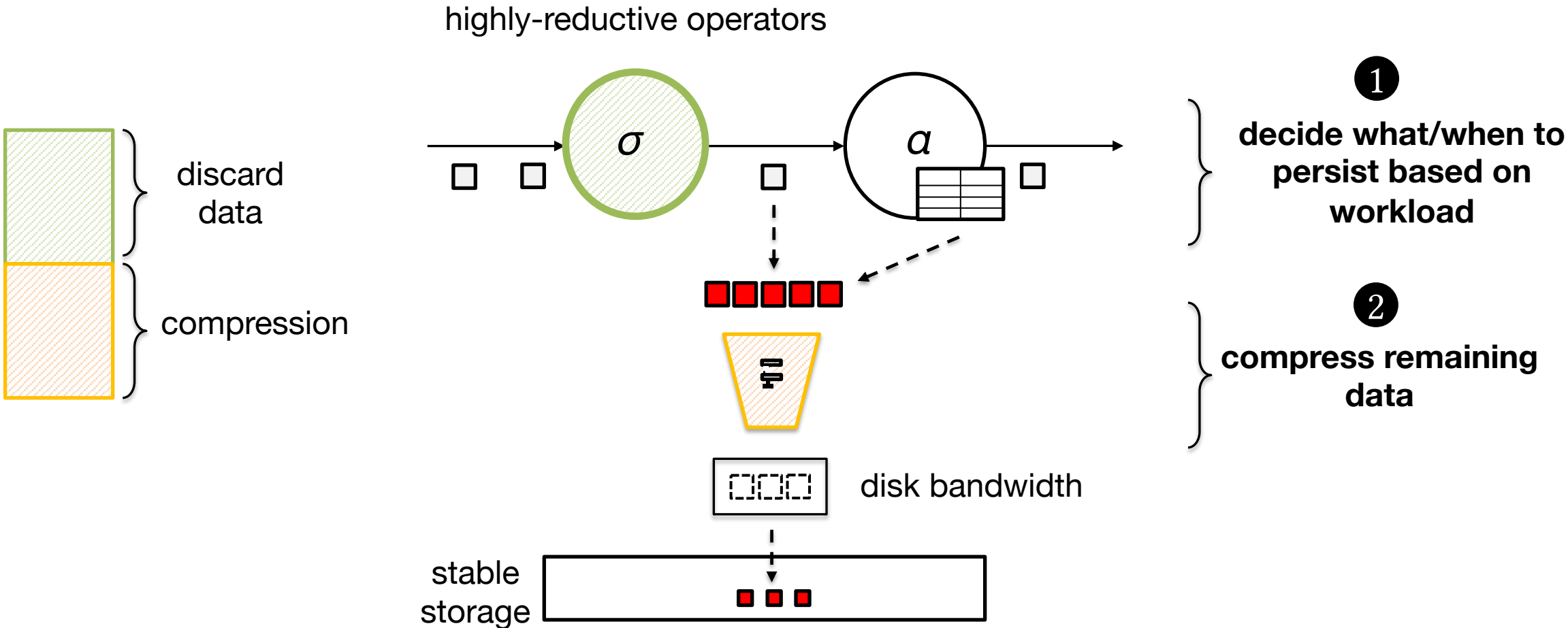- Kafka: 220

**Single-node fault-tolerance without compromising performance!**

# Key Idea: Reduce Required Disk I/O Bandwidth

# Scabbard: Reduce Required Disk I/O Bandwidth



highly-reductive operators

discard data

compression

① decide what/when to persist based on workload
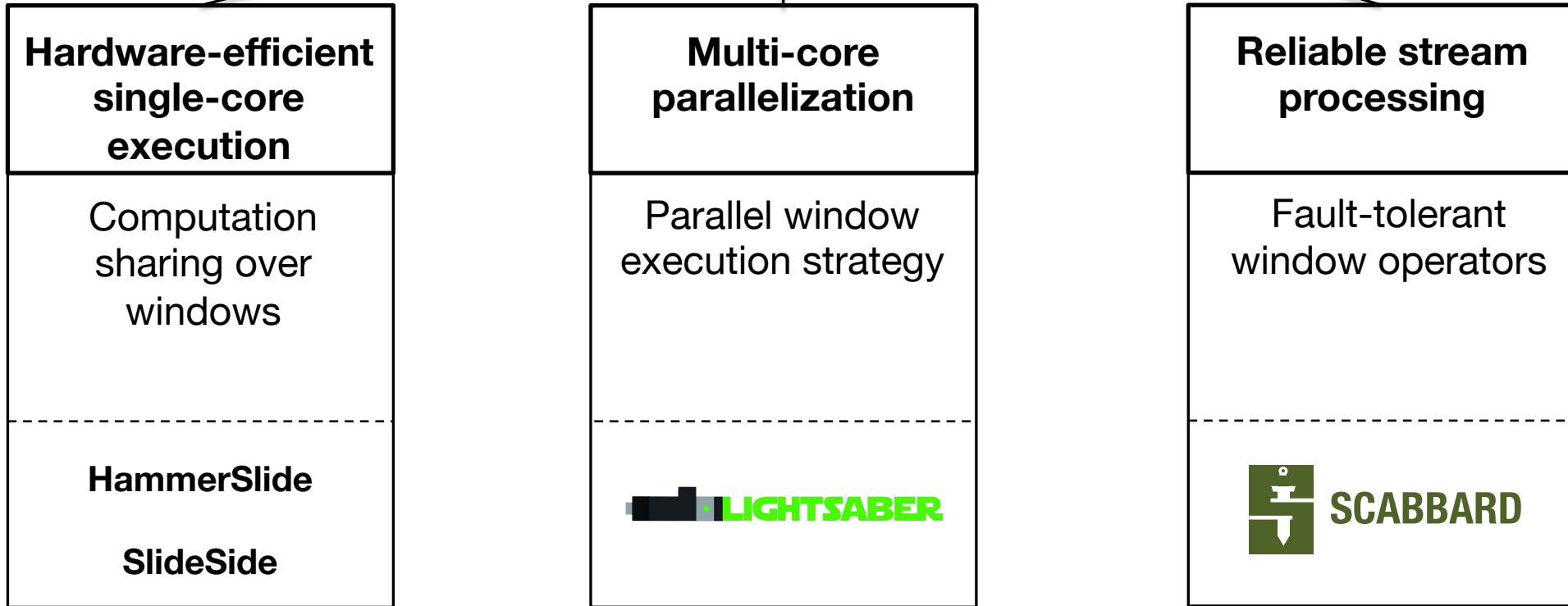
② compress remaining data

disk bandwidth

stable storage

# Single-Node Fault-Tolerant Stream Processing



> Co-optimize persistence and query execution
> JIT compile compression operators at runtime
> Use remote storage (e.g., EBS) and high-speed networking (RDMA)

# Scale-up Stream Processing Engine

| Hardware-efficient single-core execution |
| :---: |
| Computation sharing over windows |
| ----- |
| **HammerSlide**<br><br>**SlideSide** |

| Multi-core parallelization |
| :---: |
| Parallel window execution strategy |
| ----- |
| **LIGHTSABER** |

| Reliable stream processing |
| :---: |
| Fault-tolerant window operators |
| ----- |
| **SCABBARD** |

# Summary

Single-node SPEs provide a practical alternative for **scalable** and **reliable** stream processing!



https://github.com/lsds/LightSaber

**Thank you!**
**Questions?**

**George Theodorakis**

**george.theodorakis@neo4j.com**