

GPU Databases – The New Modality of Database Analytics

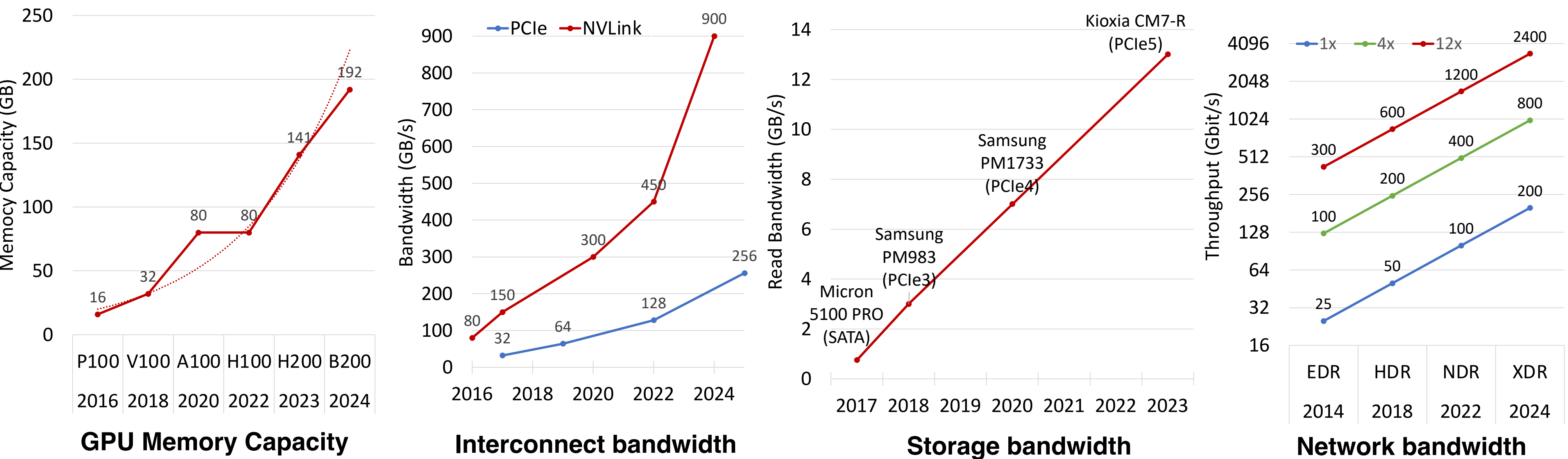
Bobbi Yogatama, Yifei Yang, Kevin Kristensen, Devesh Sarda, Abigale Kim, Xiangyao Yu (yxy@cs.wisc.edu)



Why GPU Databases?

1. GPUs are becoming the new center of computing

- Massive growth in computation power, memory capacity, and bandwidth in the past decade.
- With Grace superchips, **GPU memory capacity is no longer a bottleneck.**
- As storage and network performance grows, **CPU computation is the new bottleneck.**



2. Accessible and Affordable GPU hardware.

- GPU prices drop and availability increases over time.
- GPUs have comparable costs to CPUs in the cloud.

CPU	On-Demand Rate (AWS)
4-192 Cores	\$0.18-\$11.12 / hour

GPU	On-Demand Rate (Lambda)
1 x H100	\$2.5/hour — \$8/hour in 2023
1 x GH200	\$3.2/hour — \$4.3/hour in 2024

SiriusDB Roadmap

Completed Research Works

Single-node in-memory GPU DBMS

1. Crystal Library (**SIGMOD 2020**^[1])
2. Data Compression (**SIGMOD 2022**^[2])
3. Hybrid CPU-GPU DBMS (**VLDB 2022**^[3])
4. GPU-accelerated UDF (**DaMoN 2023**^[4])
5. Multi-GPU DBMS (**VLDB 2024**^[5])

Future/Ongoing Works

Terabyte-scale GPU Analytics

1. Query Exec with GPU Direct Storage.
2. Distributed GPU DBMS.
3. Supporting Complex SQL Operators.

GPU-accelerated RAG Database

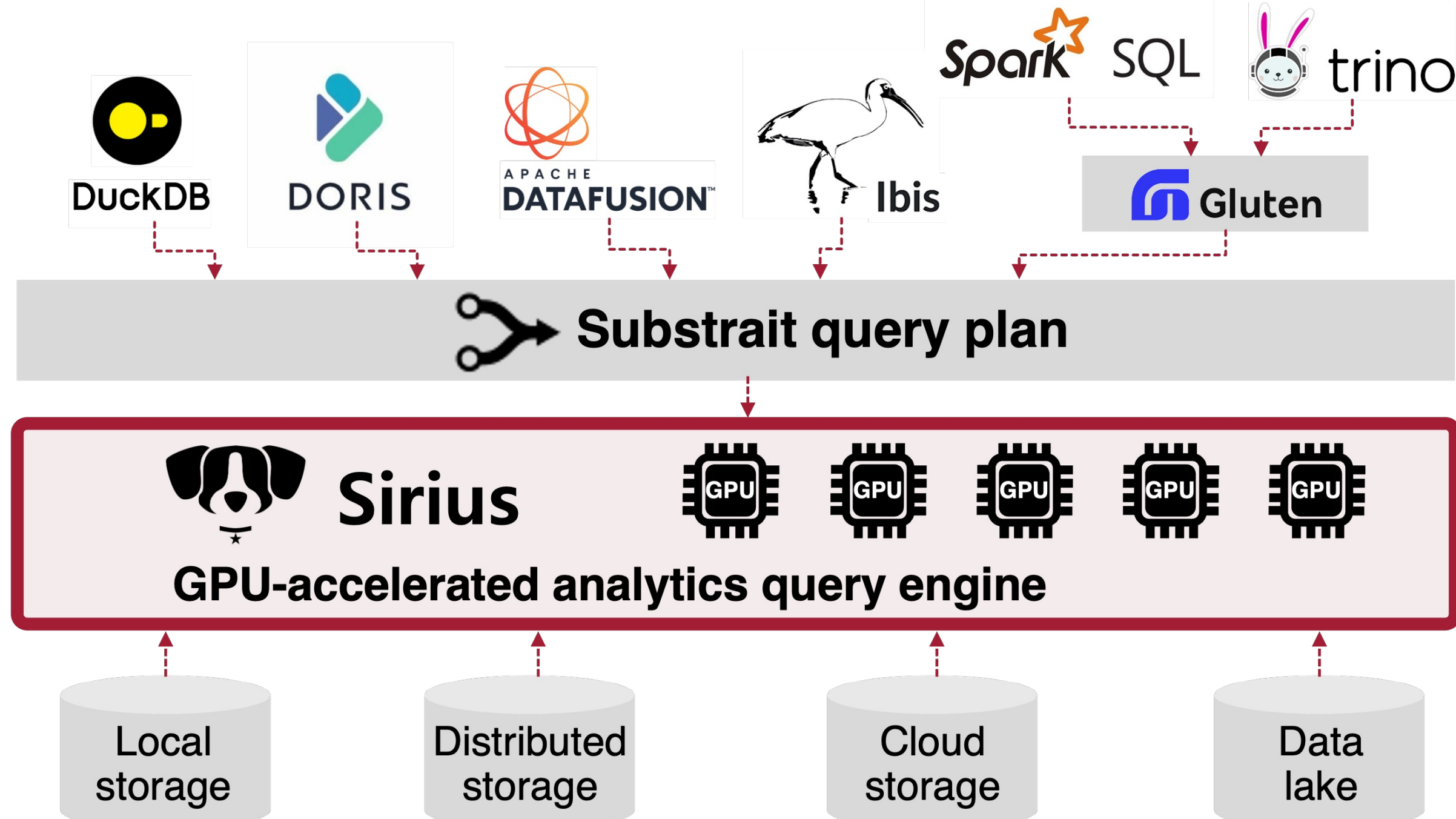
1. Integrated SQL and Vector Search on GPUs.
2. Larger than GPU memory Vector Search.

SiriusDB: The Next-Generation GPU-Accelerated Query Engine

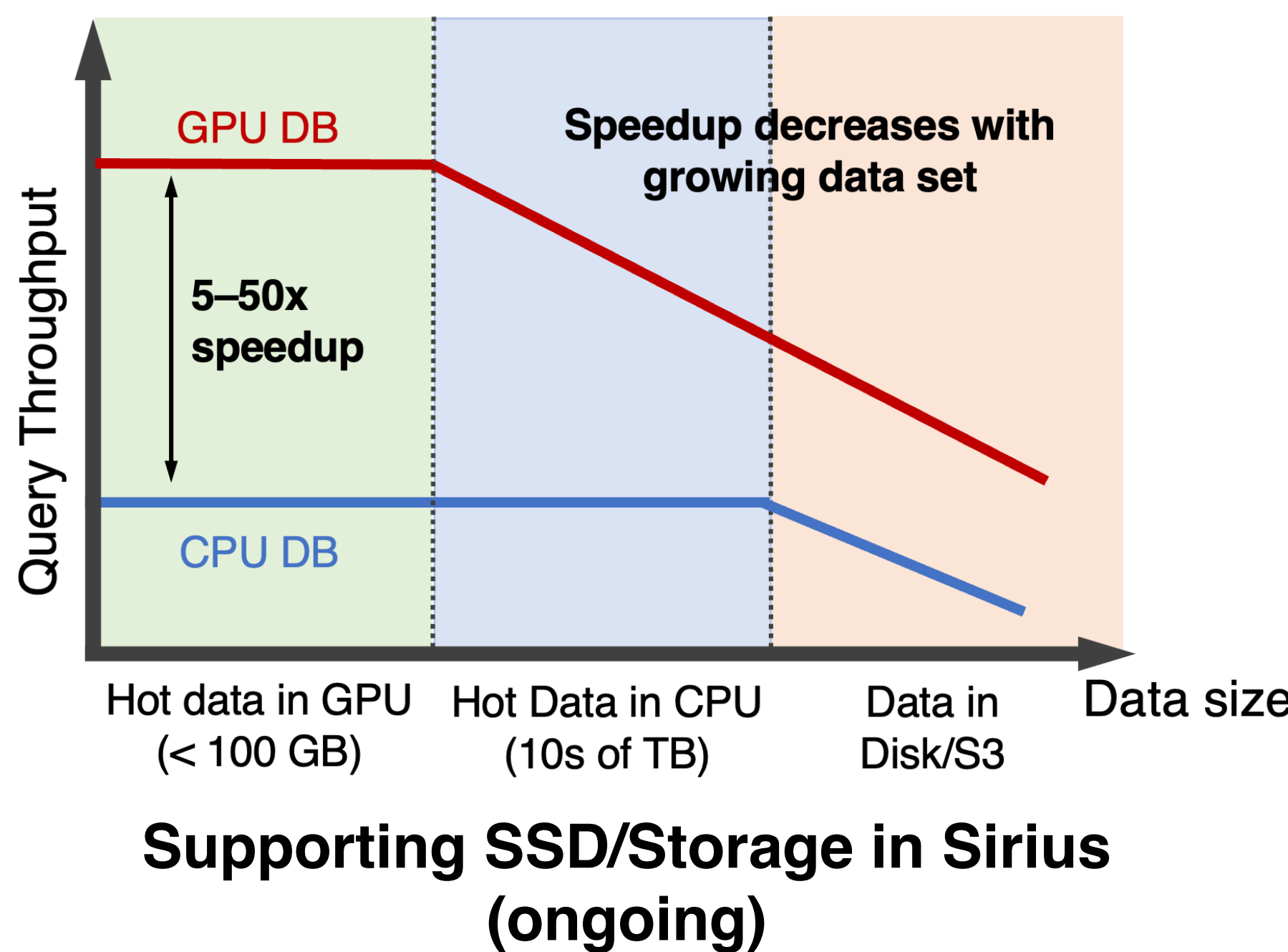
Riding the wave of GPU hardware improvement.

Unlock GPU acceleration by modifying a single line of configuration (same data, same API).

Currently supports DuckDB, more analytical processing engine will be supported in the future.



Sirius System Architecture

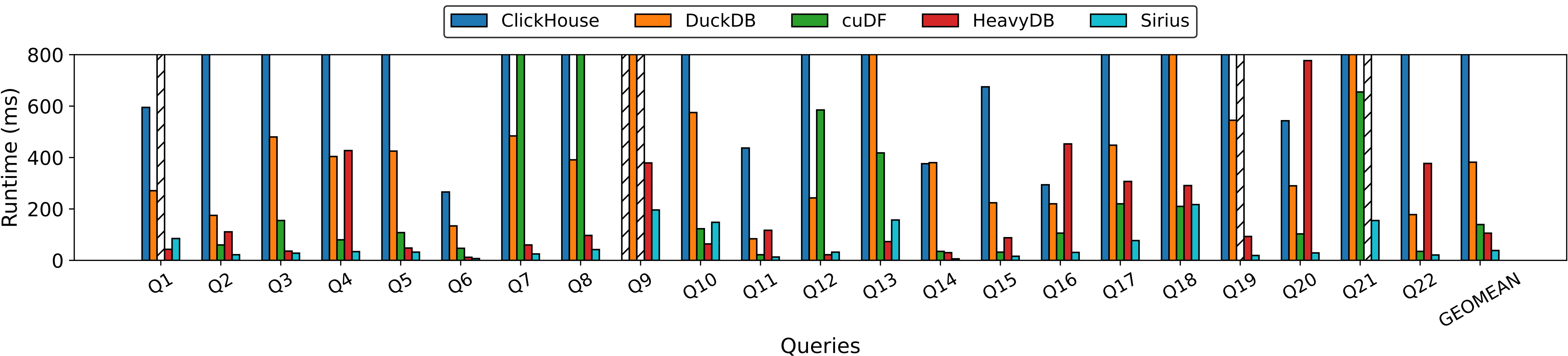


Evaluation — 100 GB TPC-H Benchmark

GPU Instance: GH200 (**\$3.2/hour** before discount), CPU Instance: m7i.16xlarge (**\$3.2/hour**)

- ClickHouse (~10 years)
- DuckDB (~6 years)
- cuDF: A GPU-accelerated data processing library by NVIDIA (~6 years)
- HeavyDB: A Commercial GPU DB for Real-time Analytics (~10 years)
- **Sirius*: our engine (~7 months)**

*Sirius has multiple pending optimizations



Sirius is **60x** faster than ClickHouse, **10x** faster than DuckDB, **3.6x** faster than cuDF, and **2.8x** faster than HeavyDB on a **cost-normalized hardware**

[1] Anil Shanbhag, Sam Madden, Xiangyao Yu. A Study of the Fundamental Performance Characteristics of GPUs and CPUs for Database Analytics, SIGMOD 2020
[2] Bobbi Yogatama*, Anil Shanbhag*, Xiangyao Yu, and Samuel Madden. Tile-based Lightweight Integer Compression in GPU, SIGMOD 2022
[3] Bobbi Yogatama, Weiwei Gong, Xiangyao Yu. Orchestrating Data placement and Query Execution in Heterogeneous CPU-GPU DBMS, VLDB 2022
[4] Bobbi Yogatama et al. Accelerating User-Defined Aggregate Functions with Block-wide Execution and JIT Compilation on GPUs, DaMoN@SIGMOD 2023
[5] Bobbi Yogatama, Weiwei Gong, Xiangyao Yu. Scaling your Hybrid CPU-GPU DBMS to Multiple GPUs, VLDB 2024