

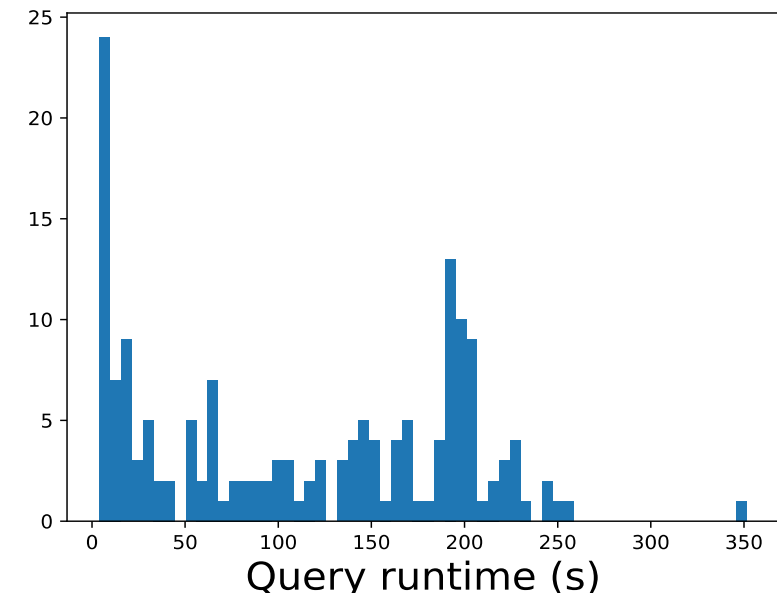
# Improving DBMS Scheduling Decisions with Accurate Performance Prediction on Concurrent Queries

---

Ziniu Wu, Markos Markakis, Chunwei Liu, Peter Baile Chen,  
Balakrishnan Narayanaswamy, Tim Kraska, Samuel Madden

# Concurrency impacts latency heavily

- Deciding **when** to execute a query can be crucial
- Bad scheduling decisions can have a significant impact on query latency.
- Microexperiment:
  - 4 clients issuing closed-loop queries against IMDB dataset on Postgres.
  - Look at runtime of the same query under different concurrent settings.
- Scheduling can have a large impact.

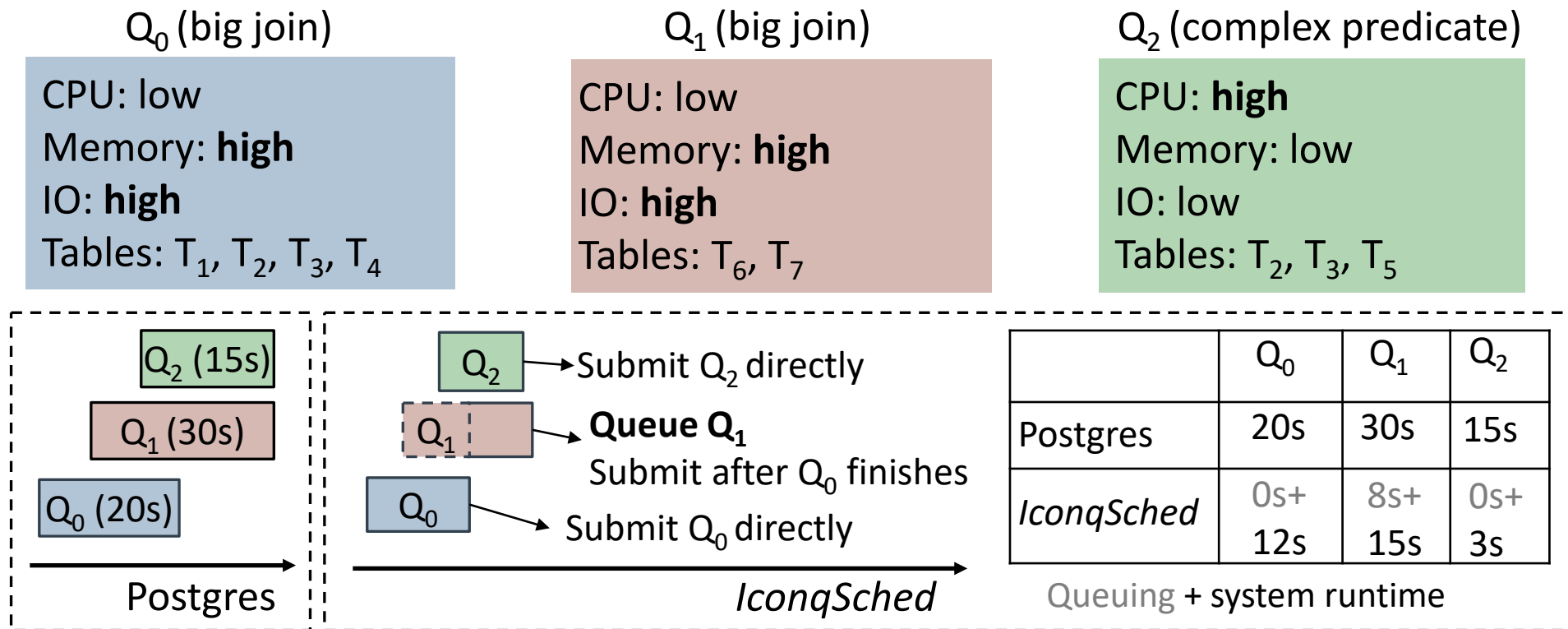


# Can we build a reusable scheduler?

---

- We explore the design of a **non-intrusive** scheduler:
  - Coarser granularity (e.g. cannot schedule individual operators).
  - But can be applied to any underlying engine.
- Goal: improve query **end-to-end time** (queueing time + system runtime)

# We need to capture resource contention

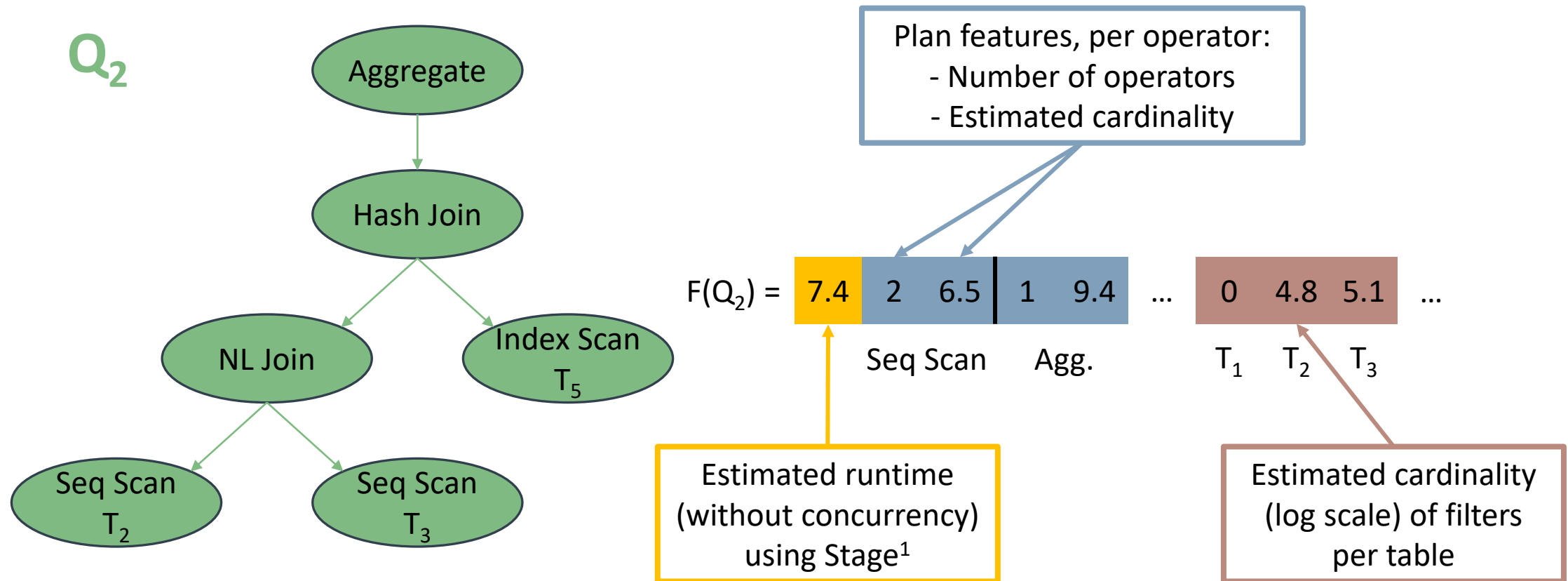


# Key challenge: splitting past and future

---

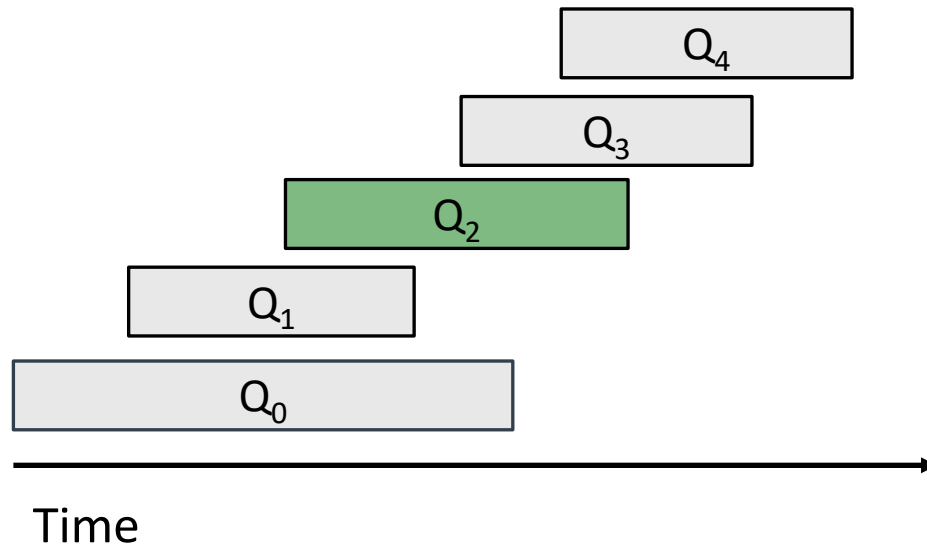
- The latency of each query is affected by both:
  - Queries that were submitted before it and are still running.
  - Queries that are submitted while it is still executing.
- The latter are **not available at scheduling time**.
- How can we **update estimates for running queries as new ones arrive?**
- Solution preview: **Bidirectional LSTM**

# We first featurize a query using its plan

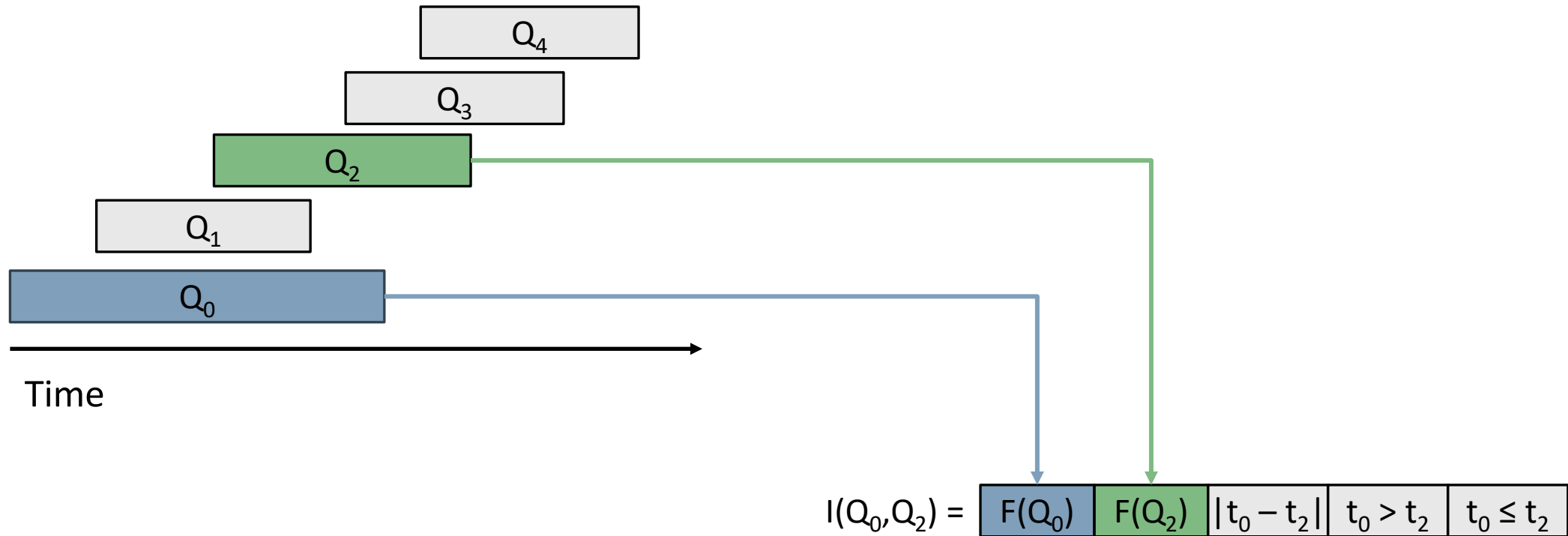


# We next add query interaction details

---

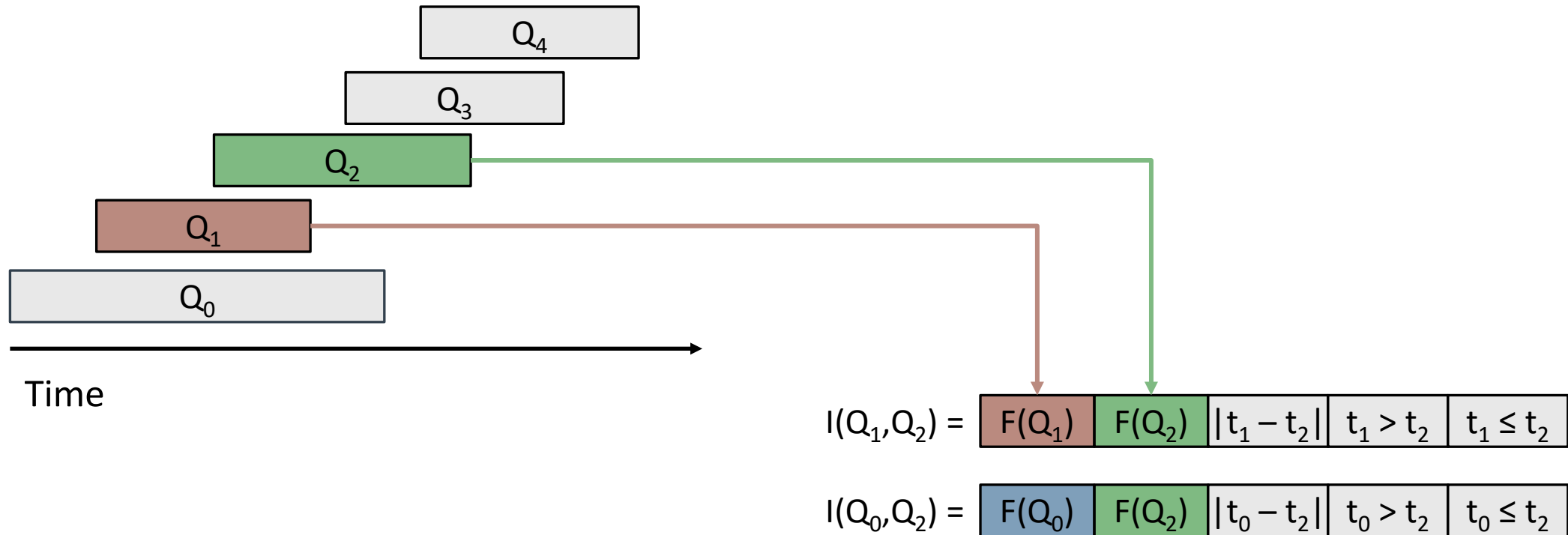


# We next add query interaction details

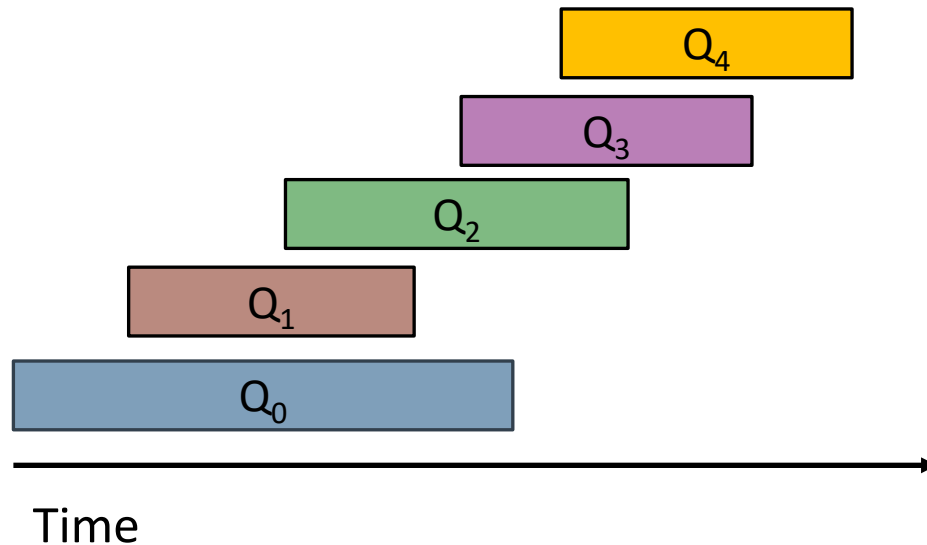




# We next add query interaction details



# We next add query interaction details



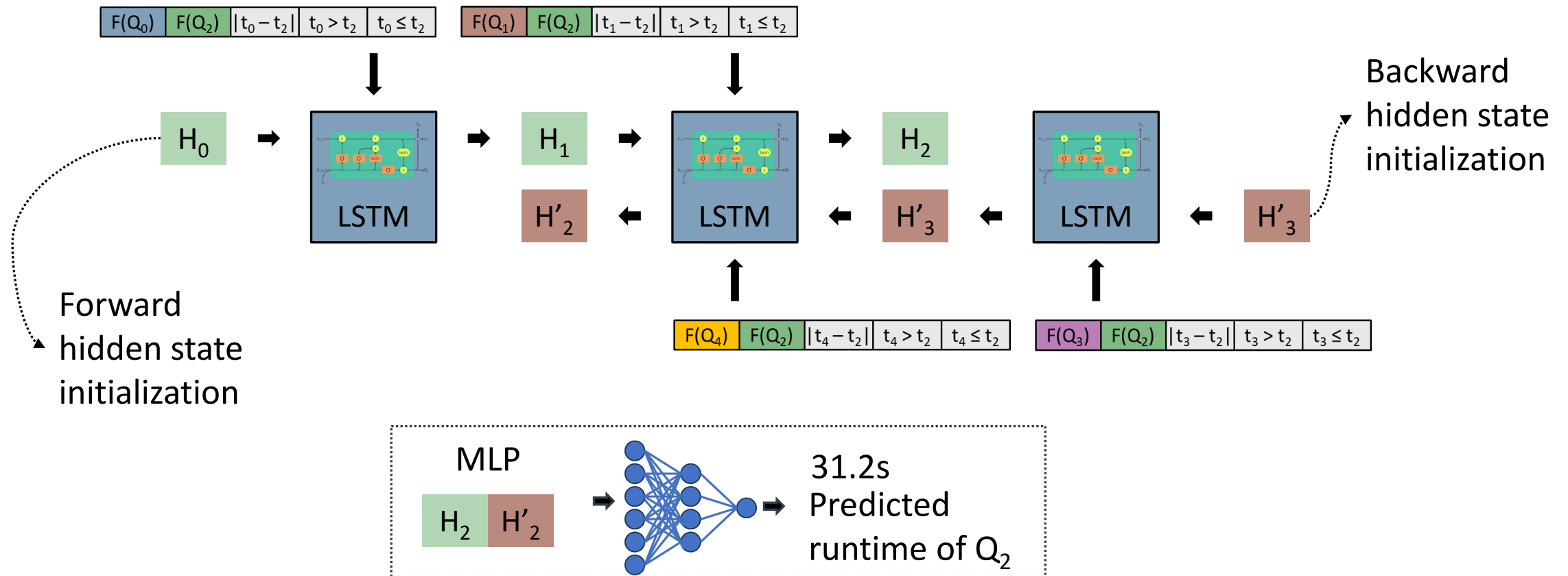
$$I(Q_1, Q_2) = \begin{array}{|c|c|c|c|c|} \hline F(Q_4) & F(Q_2) & |t_4 - t_2| & t_4 > t_2 & t_4 \leq t_2 \\ \hline \end{array}$$

$$I(Q_3, Q_2) = \begin{array}{|c|c|c|c|c|} \hline F(Q_3) & F(Q_2) & |t_3 - t_2| & t_3 > t_2 & t_3 \leq t_2 \\ \hline \end{array}$$

$$I(Q_1, Q_2) = \begin{array}{|c|c|c|c|c|} \hline F(Q_1) & F(Q_2) & |t_1 - t_2| & t_1 > t_2 & t_1 \leq t_2 \\ \hline \end{array}$$

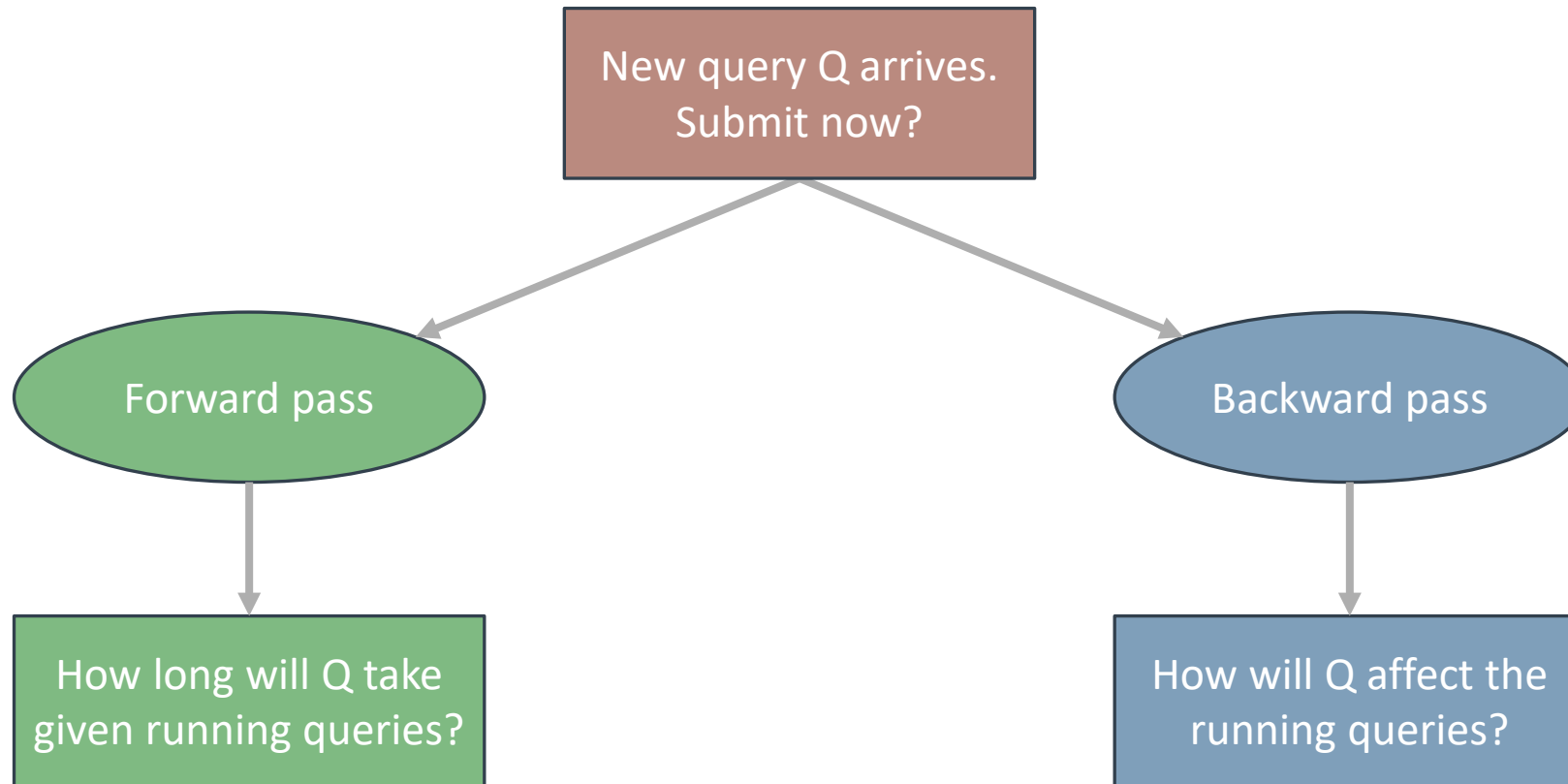
$$I(Q_0, Q_2) = \begin{array}{|c|c|c|c|c|} \hline F(Q_0) & F(Q_2) & |t_0 - t_2| & t_0 > t_2 & t_0 \leq t_2 \\ \hline \end{array}$$

# We then train Iconq: a bidirectional LSTM

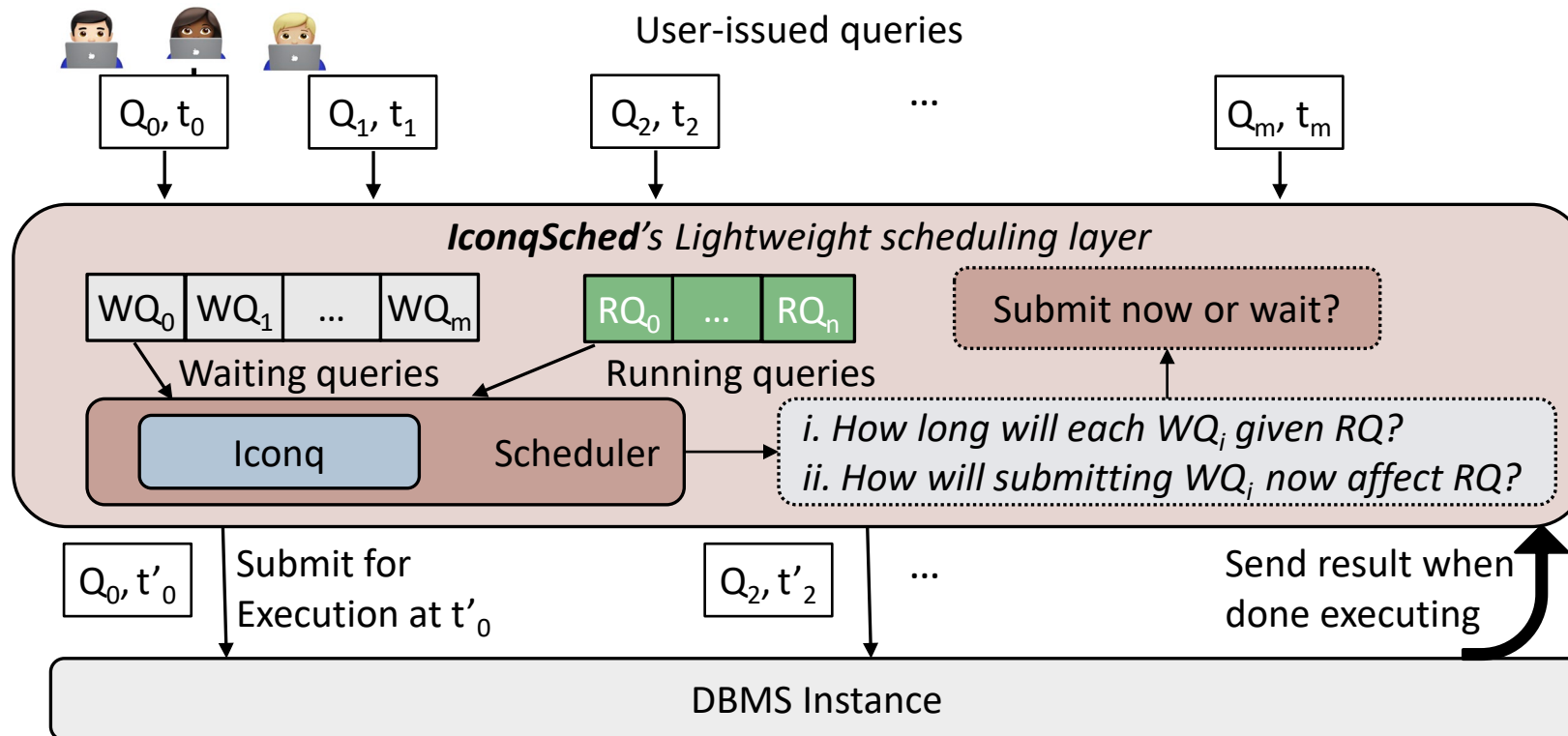


# At inference time, separate the passes

---



# Iconq is at the center of our scheduler



\* For short queries (Stage<sup>1</sup> estimate < 5s), always submit, since Iconq requires around 100ms for inference

# Setting up the evaluation

---

- Two different DBMS engines:
  - PostgreSQL
  - Amazon Redshift
- Two different open-source benchmarks:
  - CAB<sup>1</sup>: Cloud Analytics Benchmark
  - BRAD<sup>2</sup>: IMDB data + Snowset query arrival timestamps
- Two different state-of-the-art baselines:
  - PGM Scheduler<sup>3</sup>: Estimate memory and keep total bounded.
  - Qshuffler<sup>4</sup>: Represent system state as counts per query cluster.

[1] Alexander van Renen and Viktor Leis. 2023. Cloud Analytics Benchmark.

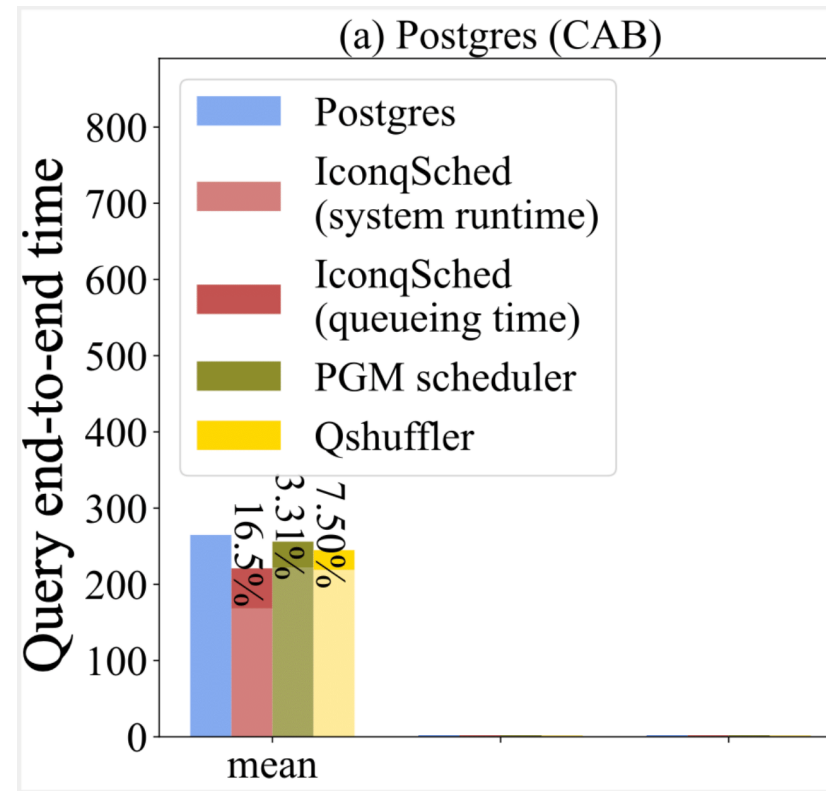
[2] Geoffrey X. Yu, Ziniu Wu, Ferdi Kossmann, Tianyu Li, Markos Markakis, Amadou Ngom, Samuel Madden, and Tim Kraska. 2024. Blueprinting the Cloud: Unifying and Automatically Optimizing Cloud Data Infrastructures with BRAD.

[3] Abhay Mehta, Chetan Gupta, and Umeshwar Dayal. 2008.

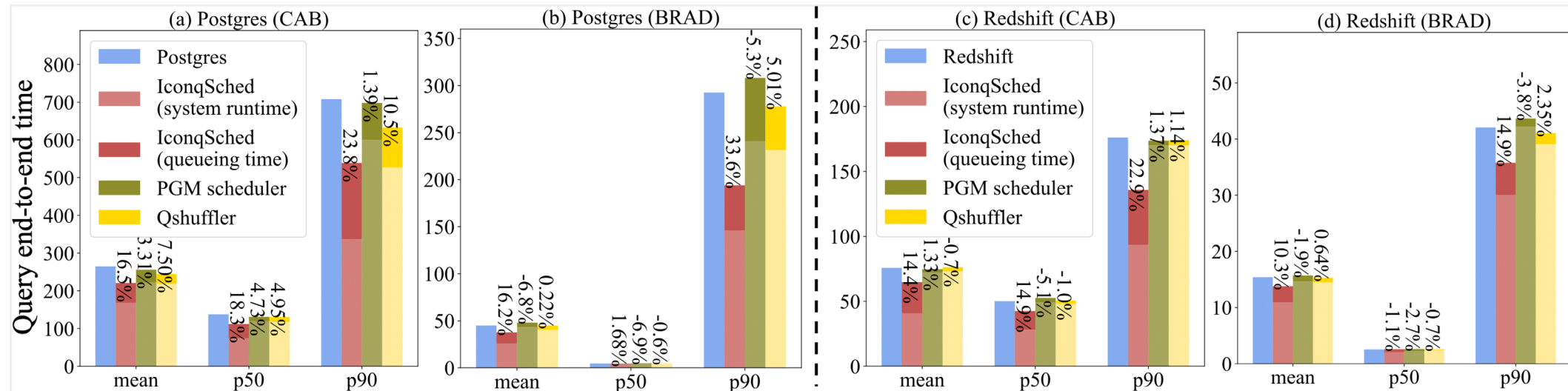
BI batch manager: a system for managing batch workloads on enterprise data-warehouses.

[4] Mumtaz Ahmad, Ashraf Aboulnaga, Shvath Babu, and Kamesh Munagala. 2011. Interaction-aware scheduling of report-generation workloads.

# We improve query end-to-end time

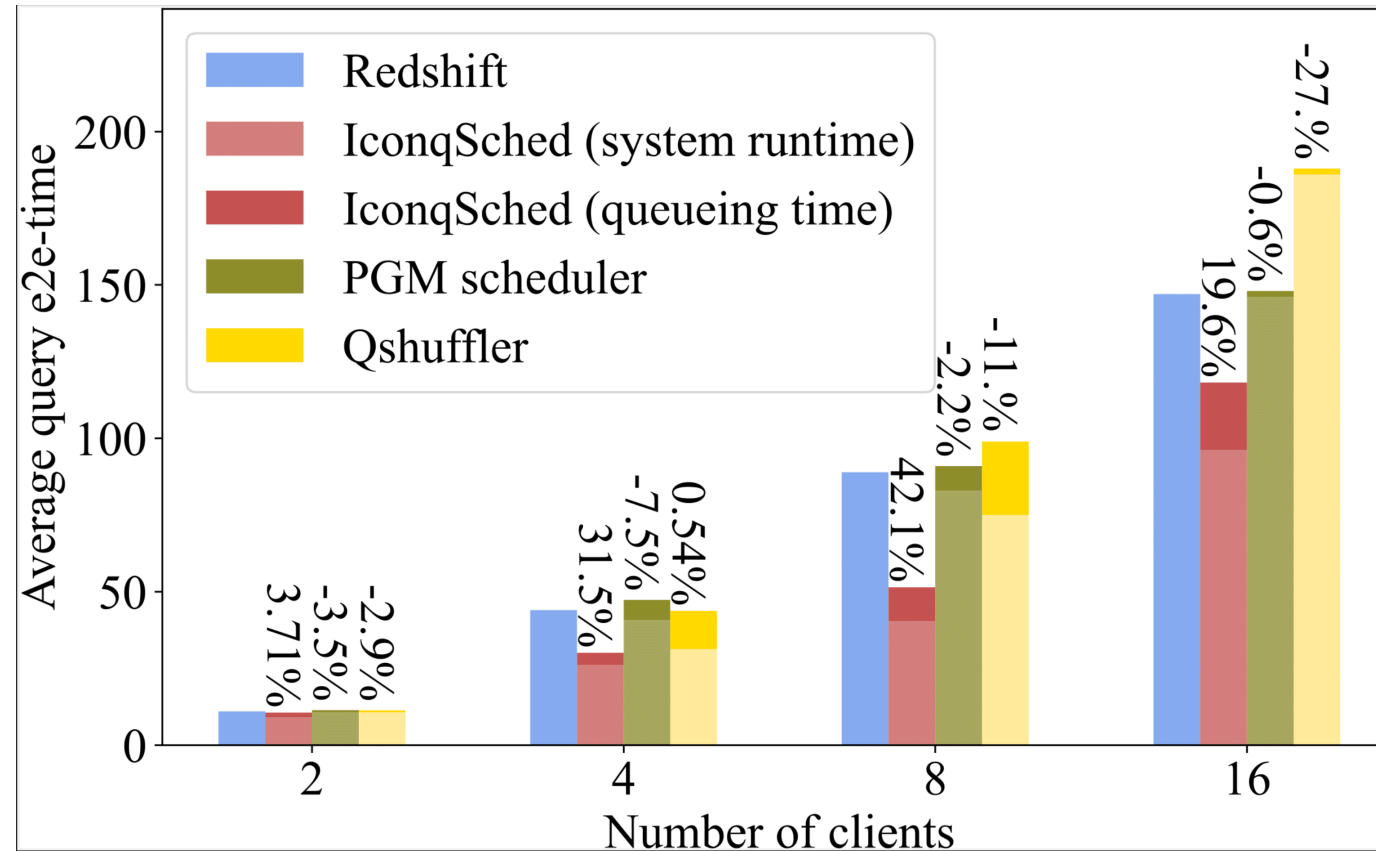


# We outperform baselines across settings





# Iconq performs under high concurrency



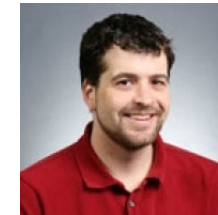
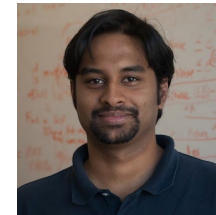
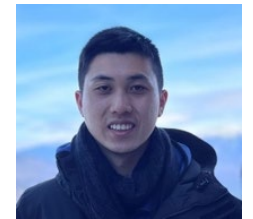
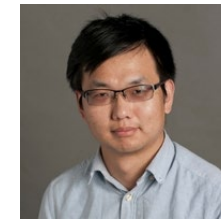
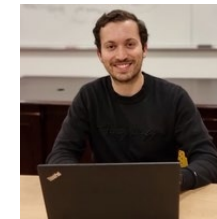
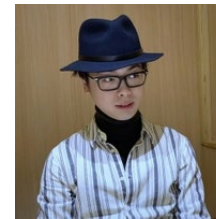
# Key Takeaway

---

IconqSched reduces  
query end-to-end time  
without accessing DBMS internals

# Thank you!

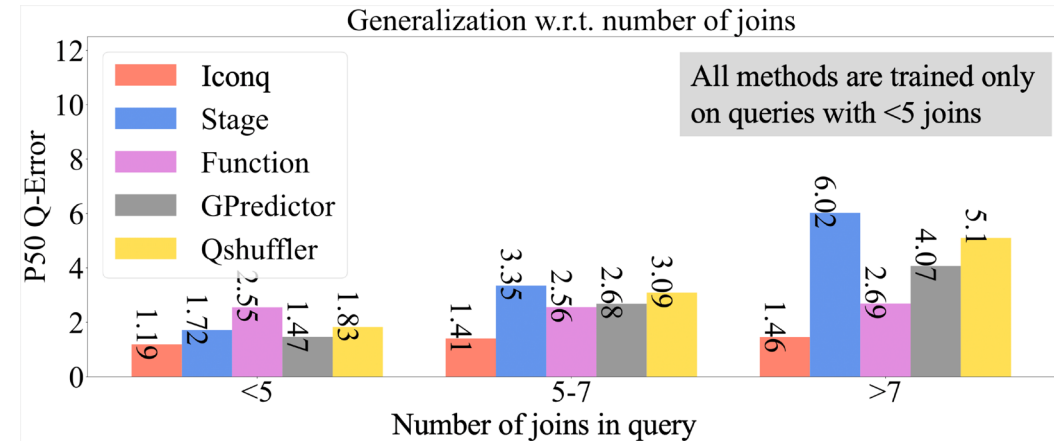
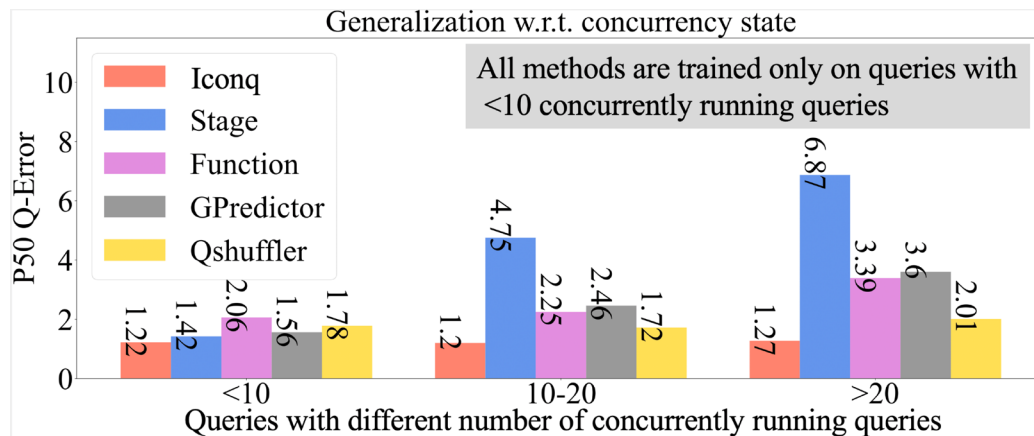
---



**Paper:** <https://www.vldb.org/pvldb/vol18/p4185-wu.pdf>

**Questions? Reach out at** [ziniuw@mit.edu](mailto:ziniuw@mit.edu)

# Iconq generalizes well to harder contexts



$$Q_{\text{error}} = \max\left(\frac{\text{true}}{\text{pred}}, \frac{\text{pred}}{\text{true}}\right)$$