Cloud Data Processing with Cost-Efficient Latency SLOs using Probabilistic Query Performance Predictions

Data Systems Group @ MIT

Markos Markakis May 10, 2025



Need to Know

Reservations can be made one day in advance, starting at 10:00 am.

10:00:00 C 10:00:02 09:59:59 C



10:00:05

Dinner 4:00 PM 10:00 PM A Notify Dining Room Dining Room



Total Covers						
Month	Days of service	Total covers	Reserved covers	Walkin covers	Waitlist covers	No show covers parties
March	12	84	51	33	0	4 / 2
Average Daily Covers (vs previous month)						
	, , ,	,				
Service	Avg covers	Avg reserved covers	Avg walkin covers	Avg waitlist covers	Avg no show covers	Avg no show res rate
Service Breakfast	Avg covers	Avg reserved covers	Avg walkin covers	Avg waitlist covers	Avg no show covers O(O -)	Avg no show res rate 33.33%(33.33% ▲)
Service Breakfast Lunch	Avg covers 2(-14) 2(-1)	Avg reserved covers 1(-3 •) 2(-1 •)	Avg walkin covers 1(-11 ▼) 0(0 −)	Avg waitlist covers 0(0 -) 0(0 -)	Avg no show covers 0(0 -) 1(1 -)	Avg no show res rate 33.33%(33.33% ▲) 16.67%(16.67% ▲)

Change in seated diners by week, 2025 vs. 2024

This graph measures the weekly change in seated diners from online reservations for 2025 vs. 2024. Hover over any given date to see how 2025 compares to the respective week in 2024. For example, in the US on the week ending on January 6, 2025, seated diners were up 25% compared to the respective week of the year in 2024.









How can we make them all "fast enough", reliably and cheaply?





Cloud Data Processing

with Cost-Efficient Latency SLOs using Probabilistic Query Performance Predictions

Databases are increasingly cloud-hosted

Look under the hood only when you create a database

Pay only when you run queries

Run queries when you want

Look under the hood only when you create a database ...but more abstraction causes latency uncertainty













2s @p95

WHERE date =

'2025-05-10'







Run queries when you want

...but risk latency-impacting interference

Look under the hood only when you create a database ...but more abstraction causes latency uncertainty Pay only when you run queries ...but big gaps can hurt both latency and cost Run queries when you want ...but risk latency-impacting interference

Cloud Data Processing with Cost-Efficient Latency SLOs using Probabilistic Query Performance Predictions



Why not just isolate the workloads?









Why not just isolate the workloads? Costly









Cloud Data Processing with Cost-Efficient Latency SLOs using Probabilistic Query Performance Predictions

Queries fight for resources like CPU and Memory



Queries fight for resources like CPU and Memory ...but the latency impact depends on relative timing



We can use arrival time differences as features



	\wedge
Q1	2 5 0
Operators: 2x01, 5x02	1 1 0 0 0 0
Tables accessed T1, T2	L ₁ 0
Q2	5 0 1
Operators: 5x01, 1x03	0 0 1 1 1 0
Tables accessed T3, T4, T5	L_2 D_2
Q3	0 2 3
Operators: 2x02, 3x03	0 1 0 0 1
Tables accessed T2, T6	L_3 D_3

Time

We can use arrival time differences as features ...and predict latency distributions



Time

The predicted distributions capture the tail okay ...but some important questions remain



 1. High concentration near the mean =
predicted distributions are very wide, what better type of distribution to assume?

2. Median inference time is 40.13ms but the maximum can be 70.25ms, how can we bring this down?



Thank you!



Ziniu Wu

Tim Kraska



