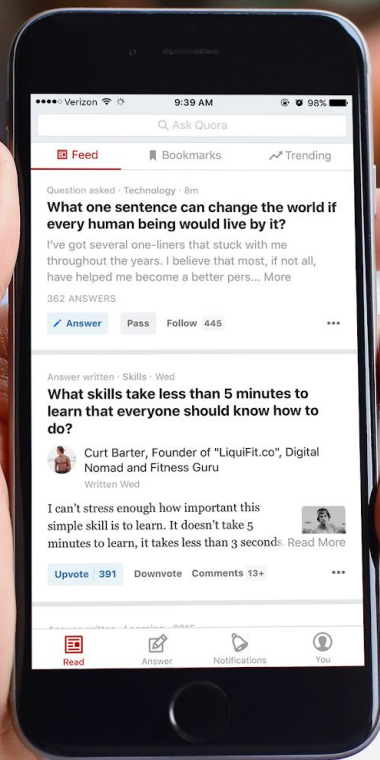


Trino at Quora: Speed, Cost, Reliability Challenges and Tips

Yifan Pan - Software Engineer @ Quora (epan@quora.com)

Reviewed by Gabriel Fernandes de Oliveira - Software Engineer @ Quora

The Quora logo, featuring the word "Quora" in a bold, red, serif font.



Verizon 9:39 AM 98%

Ask Quora

Feed

Bookmarks

Trending

Question asked · Technology · 8m

What one sentence can change the world if every human being would live by it?

I've got several one-liners that stuck with me throughout the years. I believe that most, if not all, have helped me become a better pers... More

362 ANSWERS

Answer

Pass

Follow 445

...

Answer written · Skills · Wed

What skills take less than 5 minutes to learn that everyone should know how to do?



Curt Barter, Founder of "LiquiFit.co", Digital Nomad and Fitness Guru

Written Wed

I can't stress enough how important this simple skill is to learn. It doesn't take 5 minutes to learn, it takes less than 3 seconds. Read More



Upvote 391

Downvote

Comments 13+

...



Read



Answer



Notifications

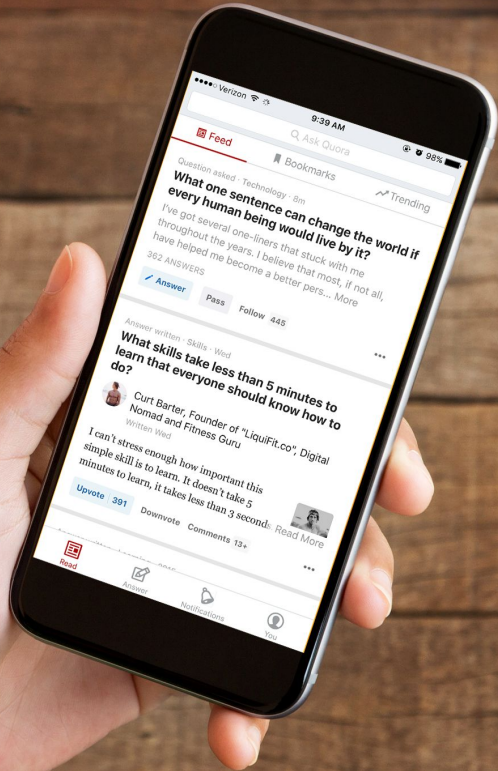


You

Quora

What is Quora?

Our mission is to share and grow the world's knowledge.

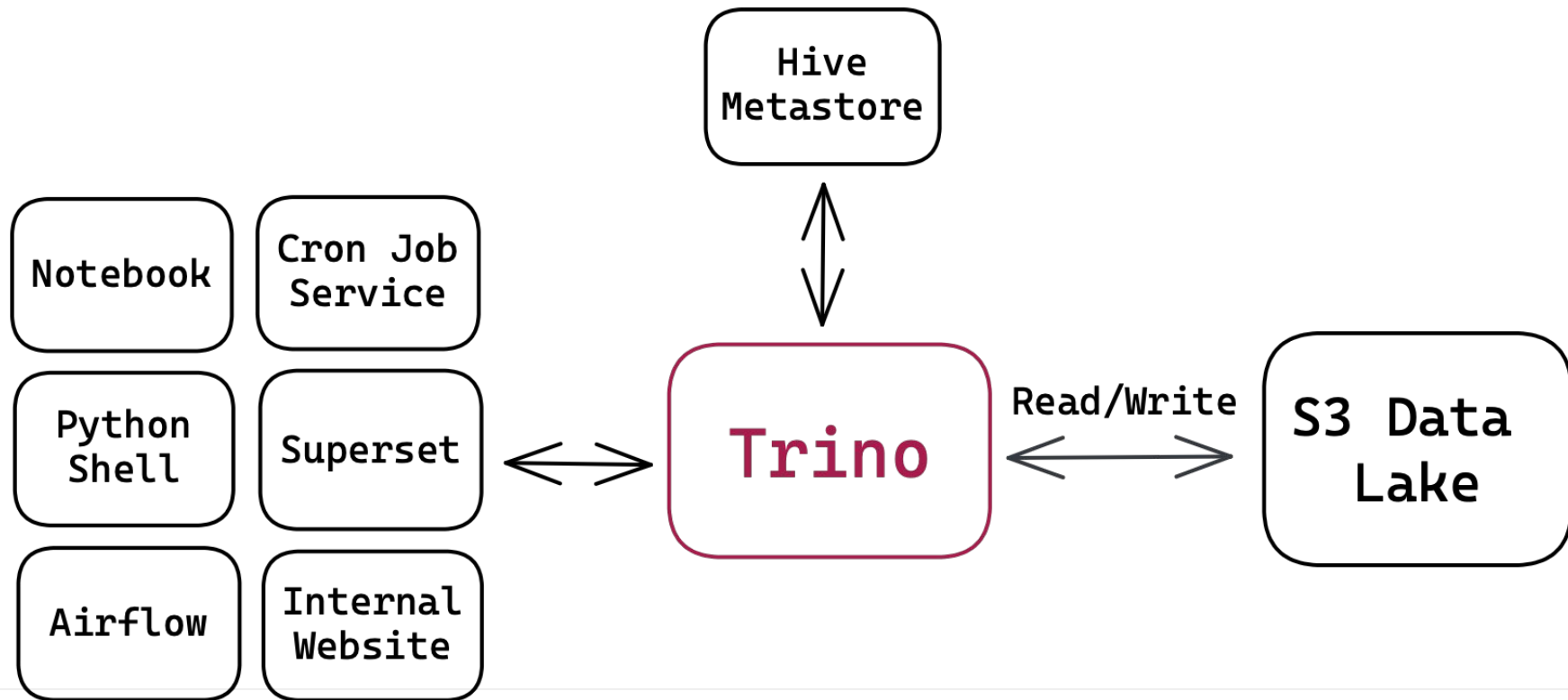


Overview

- How is Trino used at Quora?
- Cost Challenges
- Performance Challenges
- Reliability Challenges
- Summary

How is Trino used at Quora?

How is Trino used at Quora?



How is Trino used at Quora?

Main use cases of Trino at Quora:

ETL

Ad-hoc

A/B
Testing

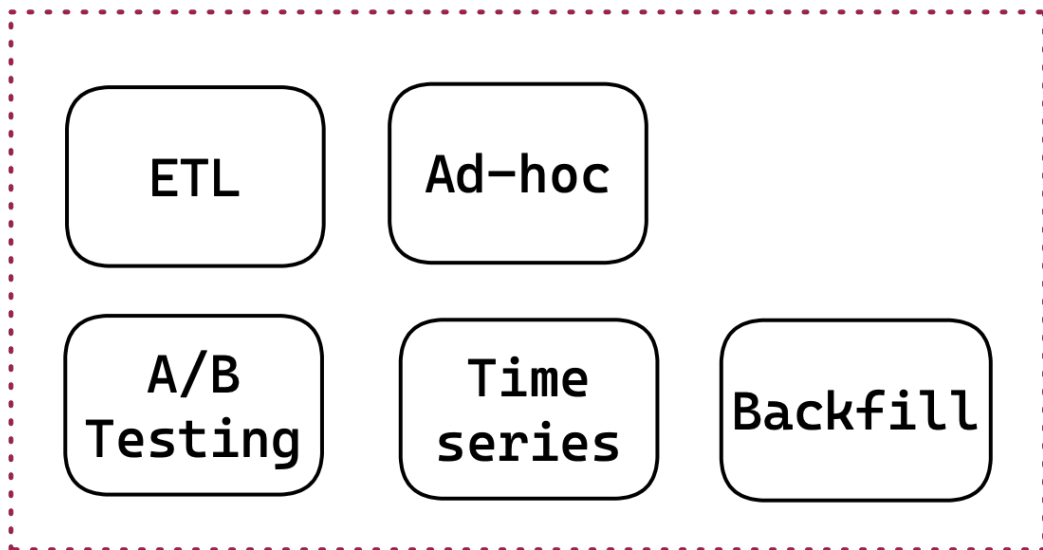
Time
series

Backfill

How is Trino used at Quora?

Quora maintains many Trino clusters, one dedicated to each use case.

Trino Clusters



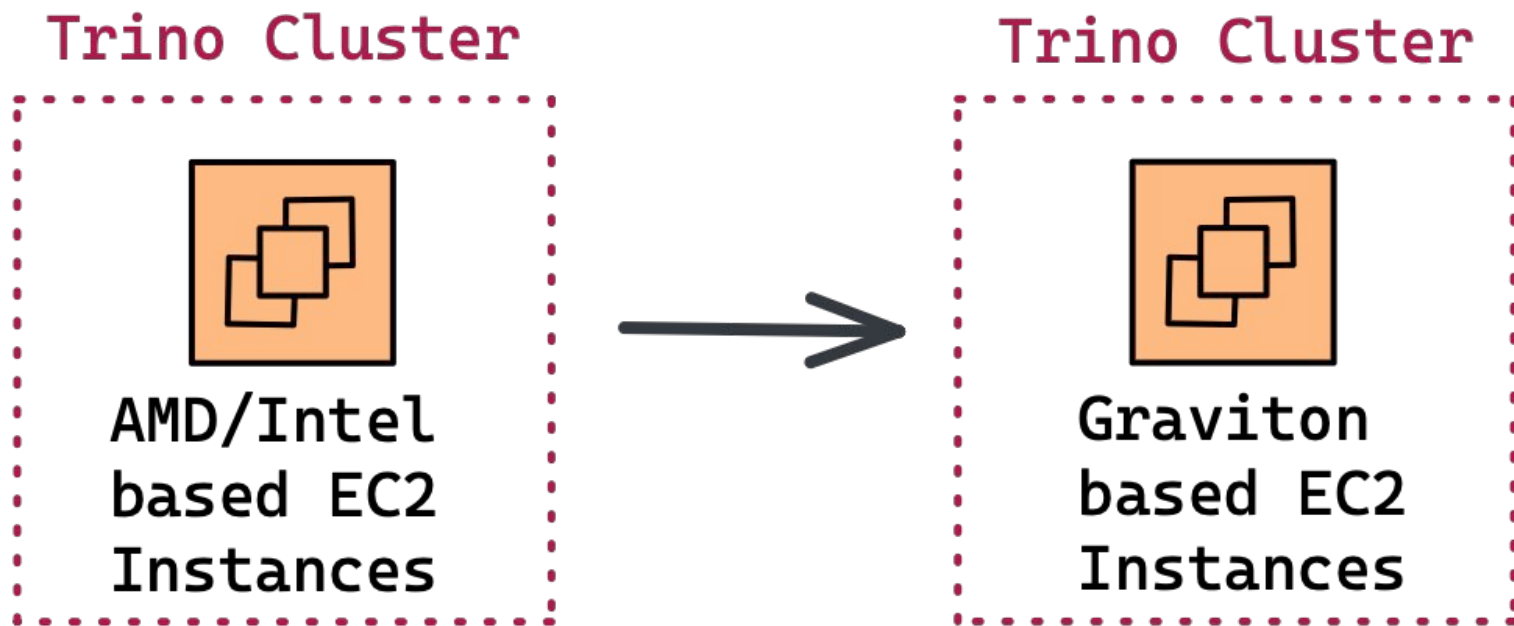
Cost Challenges

Strategies to Reduce Infrastructure Cost

- Use Graviton instances
- Auto-scale
- Optimize ETL Query Efficiency

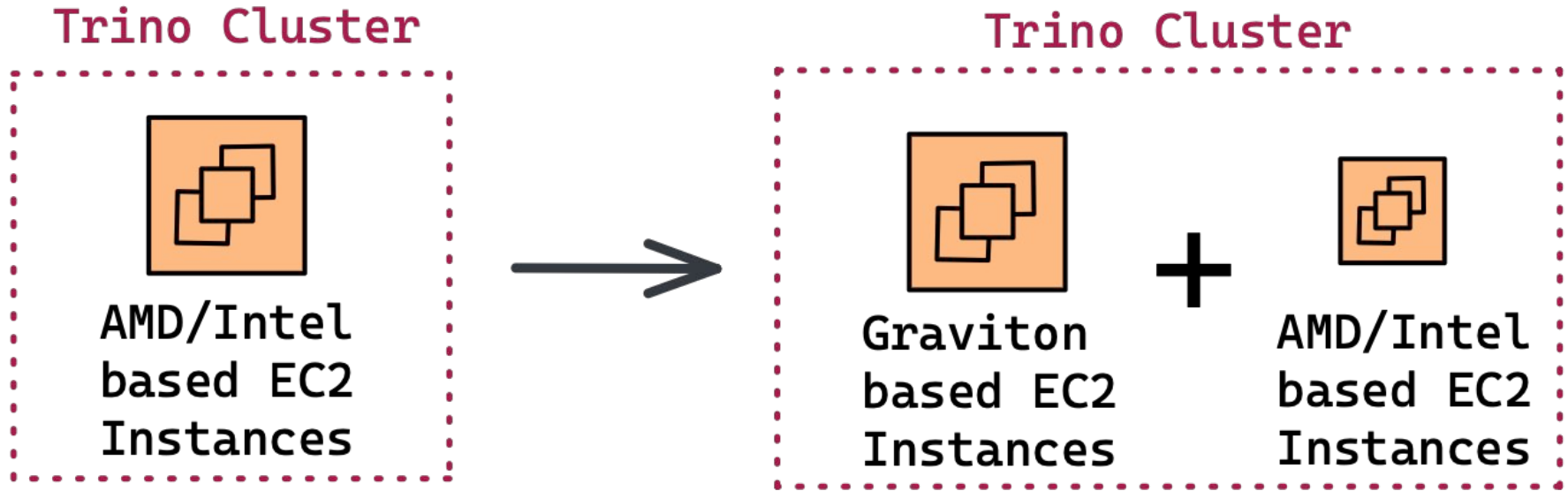
1) Use Graviton instances

In August 2020, we moved all of Trino clusters to Graviton EC2 Instances.

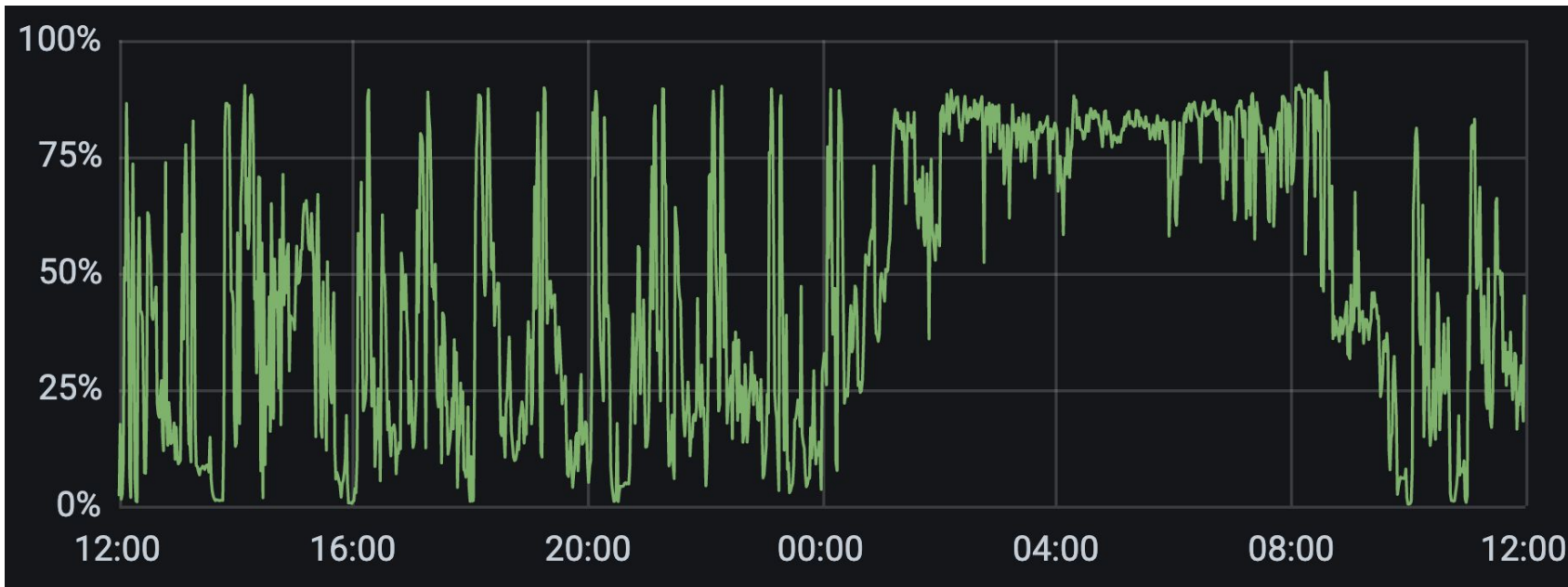


1) Use Graviton instances

Challenge: Instance availability issue.



2) Auto-scale



Average CPU utilization of our ETL cluster ***without*** auto-scaling in a single day.

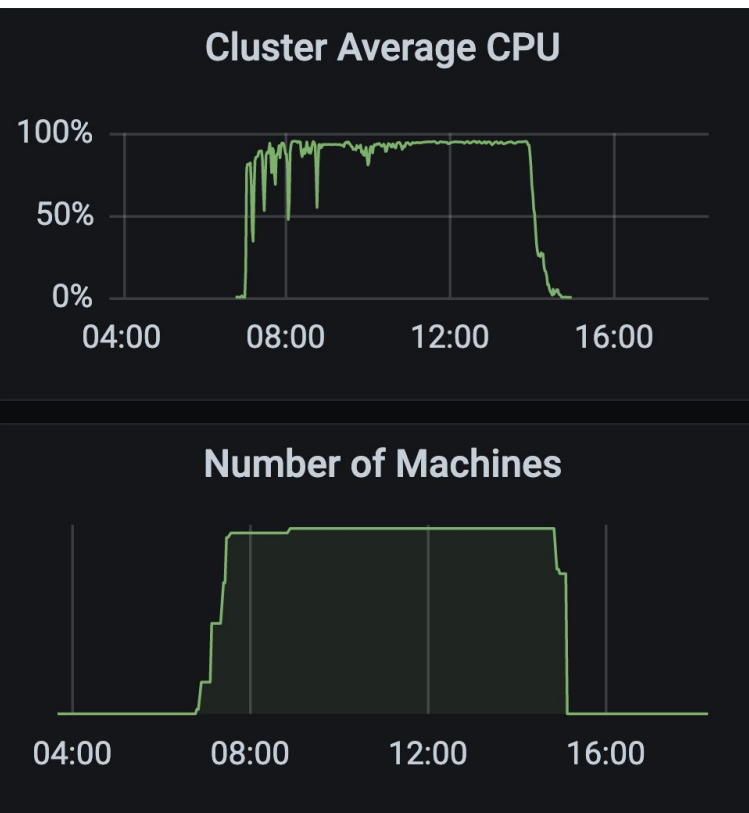
2) Auto-scale

The scaling strategies vary based on the use cases:

Trino Cluster	Auto Scaling Strategy
ETL	Scale based on CPU utilization.
Ad-hoc	Scale up during the day; Scale down during the night and weekends.
Backfill	Automatically scale up/down when users submit a backfill job.

2) Auto-scale

Trino Cluster	Auto Scaling Strategy
A/B Testing	<p>Problem:</p> <ul style="list-style-type: none">• The workload is very heavy.• Takes a couple of hours per day. <p>Solution:</p> <ul style="list-style-type: none">• Only start the cluster after dependent data is ready;• Immediately shut down the cluster after queries are finished.



3) Optimize ETL Query Efficiency

Query Optimization: Apply the "**WHERE**" clause to the partition keys.

Table A is a Hive table **partitioned by dt column*

1658448000000 is epoch in milliseconds (July 22, 2022 12:00:00 AM)

1658491200000 is epoch in milliseconds (July 22, 2022 12:00:00 PM)

Suboptimal

```
SELECT columnA FROM A
WHERE
    A.time > 1658448000000
    AND A.time < 1658491200000
```

Better

```
SELECT columnA FROM A
WHERE
    A.dt = DATE'2022-07-22'
    AND A.time > 1658448000000
    AND A.time < 1658491200000
```

3) Optimize ETL Query Efficiency


- Build a tool to automatically detect ETL queries that scan too many partition keys.
- Contact query owners.

3) Optimize ETL Query Efficiency

Use the tool to find corner cases that Predicate Pushdown couldn't handle in the *Hive Connector*.

```
WHERE
  (A.dt = DATE'2022-07-17' AND A.time ≥ B.time - 604800000000)
  OR
  (A.dt = DATE'2022-07-18' AND A.time < 1658188800000000)
  OR
  A.dt = DATE'2022-07-23'
```

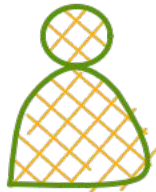
**Table A is a Hive table
partitioned by dt column*



```
WHERE
  A.dt IN (DATE'2022-07-17', DATE'2022-07-18', DATE'2022-07-23')
  AND
  (
    (A.dt = DATE'2022-07-17' AND A.time ≥ B.time - 604800000000)
    OR
    (A.dt = DATE'2022-07-18' AND A.time < 1658188800000000)
    OR
    A.dt = DATE'2022-07-23'
  )
```

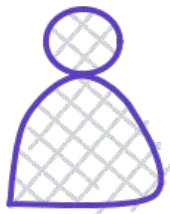
Performance Challenges

Performance Challenges



User

is Trino slow today? This is quite a heavy query but it usually takes around 15min, it's taking almost an hour now: [query](#)



User

is Trino slow for other folks today? i have a [query](#) running for close to 2 hours now, didn't expect it to take this long (but maybe it's an issue specific to my query)



Identifying and Preventing Slow Workers

Symptom:

Execution progress skewness

The elapsed time for task 2.16 was 40 minutes, while most tasks took 3 minutes; Bytes/s for task 2.16 was 17.1M while most tasks had a higher Bytes/s.

Tasks

Show ▾

Auto-Refresh: On

ID ▴	Host	State	II	▶	🚩	✓	Rows	Rows/s	Bytes	Bytes/s	Elapsed	CPU Time	Buffered
2.0		FINISHED	0	0	0	2585	6.18B	3.64M	40.8G	24.6M	28.27m	1.70h	0
2.1		FINISHED	0	0	0	6051	9.53B	61.3M	62.9G	414M	2.59m	1.09h	0
2.2		FINISHED	0	0	0	5675	9.02B	58.3M	59.6G	394M	2.58m	1.13h	0
2.3		FINISHED	0	0	0	2712	6.08B	5.12M	40.2G	34.7M	19.79m	1.57h	0
2.4		FINISHED	0	0	0	5630	8.74B	56.5M	57.6G	381M	2.58m	1.14h	0
2.5		FINISHED	0	0	0	5555	8.64B	55.8M	57.0G	377M	2.58m	1.13h	0
2.6		FINISHED	0	0	0	2345	5.44B	3.97M	35.9G	26.8M	22.86m	1.56h	0
2.7		FINISHED	0	0	0	5688	8.93B	57.4M	58.9G	388M	2.59m	1.15h	0
2.8		FINISHED	0	0	0	5775	9.16B	59.0M	60.5G	399M	2.59m	1.11h	0
2.9		FINISHED	0	0	0	5445	8.71B	56.5M	57.5G	382M	2.57m	1.15h	0
2.10		FINISHED	0	0	0	2102	5.26B	3.31M	34.6G	22.3M	26.48m	1.78h	0
2.11		FINISHED	0	0	0	1952	5.03B	4.31M	33.2G	29.2M	19.42m	1.48h	0
2.12		FINISHED	0	0	0	5575	8.81B	57.2M	58.1G	386M	2.57m	1.15h	0
2.13		FINISHED	0	0	0	5775	9.17B	59.5M	60.5G	402M	2.57m	1.13h	0
2.14		FINISHED	0	0	0	2527	6.17B	4.04M	40.6G	27.3M	25.44m	1.13h	0
2.15		FINISHED	0	0	0	5961	9.47B	61.4M	62.4G	415M	2.57m	1.11h	0
2.16		FINISHED	0	0	0	2537	6.12B	2.54M	40.3G	17.1M	40.22m	2.33h	0
2.17		FINISHED	0	0	0	5715	8.93B	57.9M	59.0G	392M	2.57m	1.14h	0
2.18		FINISHED	0	0	0	5484	8.65B	56.1M	57.1G	379M	2.57m	1.16h	0
2.19		FINISHED	0	0	0	5801	9.04B	58.6M	59.7G	396M	2.57m	1.14h	0
2.20		FINISHED	0	0	0	5907	9.57B	62.1M	63.2G	420M	2.57m	1.09h	0
2.21		FINISHED	0	0	0	5628	8.92B	57.9M	58.9G	391M	2.57m	1.13h	0
2.22		FINISHED	0	0	0	5890	9.30B	60.1M	61.3G	406M	2.58m	1.11h	0

Identifying and Preventing Slow Workers

Other symptoms of slow workers:

- Below-average CPU utilization.
- Below-average Load.
- ...

Identifying and Preventing Slow Workers

If a Trino worker runs for a long time, it is more likely to become a “slow worker”.

Solution:

- Gracefully restart worker nodes that have been running for more than 24 hours.
- Build a detector that alerts when a worker node with a low CPU or load outlier is found.

Reliability Challenges

Be cautious when overwriting the Trino configurations

A recent example of cluster being unhealthy due to overwriting one of the configurations:



- Runnable drivers and worker parallelism drop to zero
- Long GC pauses on the coordinator
- Connection errors between workers and coordinators
- We tried killing some queries in the hope of freeing up some resources, but the cluster cannot recover from the slowness quickly

Be cautious when overwriting the Trino configurations

- After digging into the heap dump of the coordinator, we found **query history** used a lot of memory.
- Reducing **query.min-expire-age** solves the issue.

`query.min-expire-age`

- **Type:** `duration`
- **Default value:** `15m`

The minimal age of a query in the history before it is expired. An expired query is removed from the query history buffer and no longer available in the [Web UI](#).

Monitoring

We built different components to monitor the health of our Trino clusters:

- Collect information on every query sent to Trino through EventListener.
- Monitor Trino query failure rate.
- Periodically send health check queries to Trino clusters.
- Track workers' uptime, alerting if any worker runs for over 36 hours.
- Track Trino's JMX metrics.

Summary

Summary

Cost Challenges

- Using Graviton for better cost-efficiency
- Apply auto-scaling rules
- Optimize ETL query efficiency

Performance Challenges

- Gracefully restart all Trino workers in the “rolling” fashion to prevent slow workers

Reliability Challenges

- Avoid misconfiguration
- Monitoring

Thank you!
Questions?

Quora