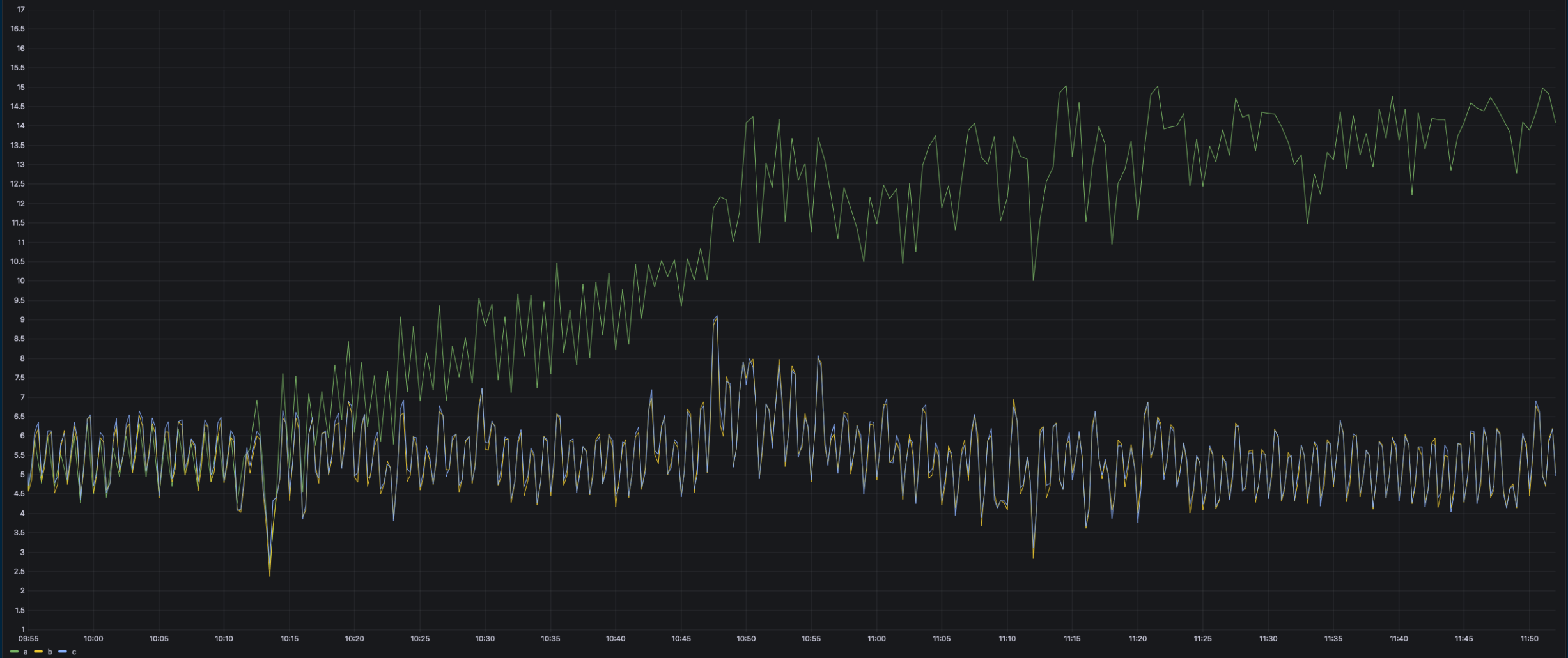


How a single API endpoint saved 3000 CPU

Maersk Observability Platform

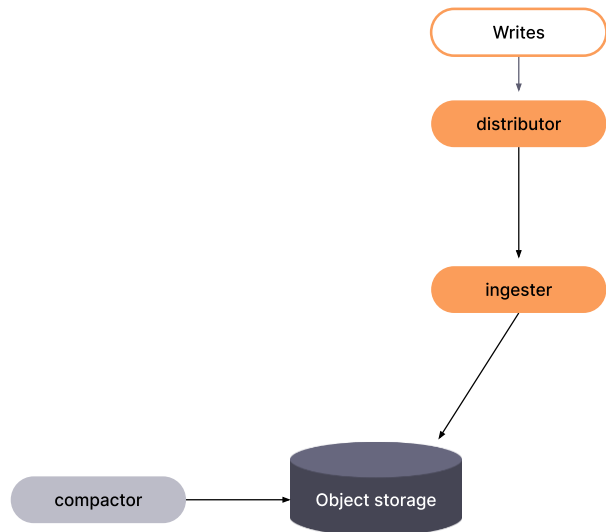
Avg CPU per zone



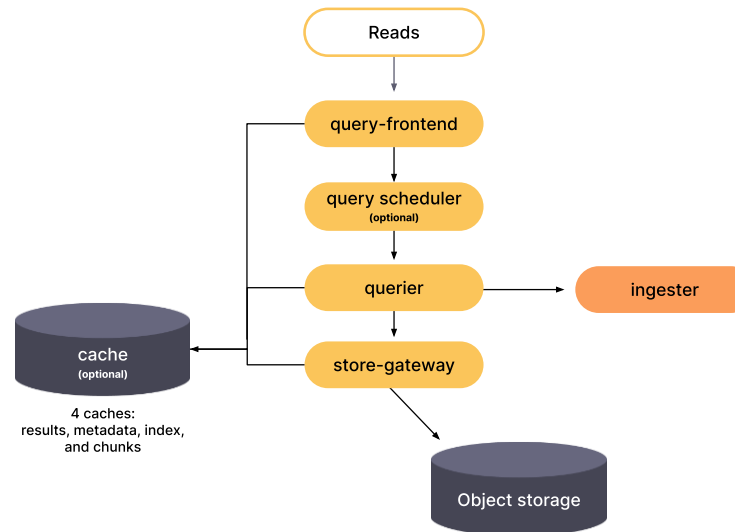
What is Mimir? What is an ingester?

Open-source time-series database for metrics

The write path

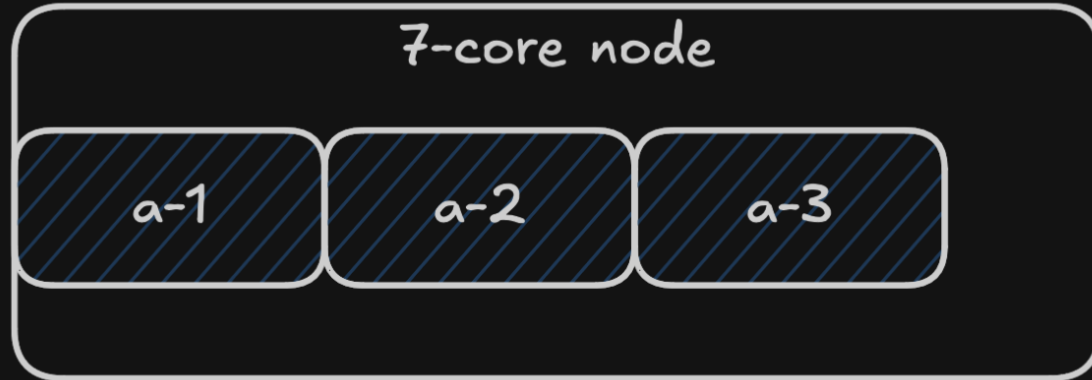


The read path

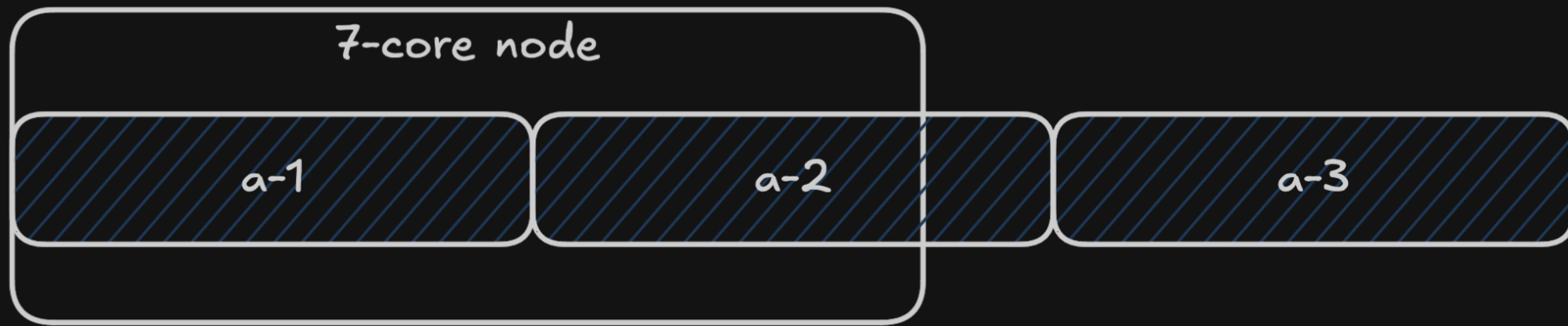


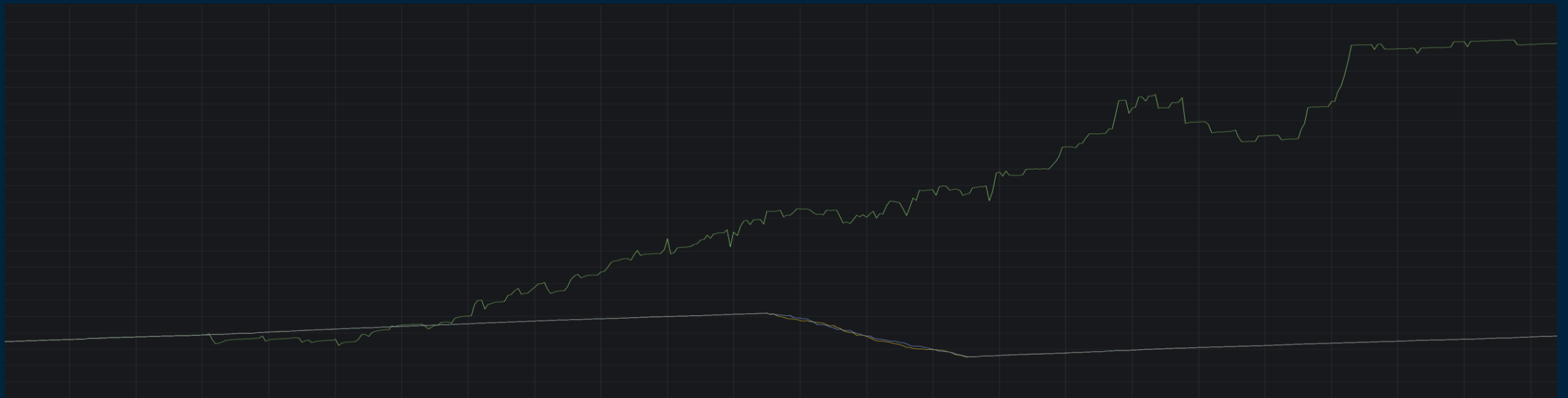
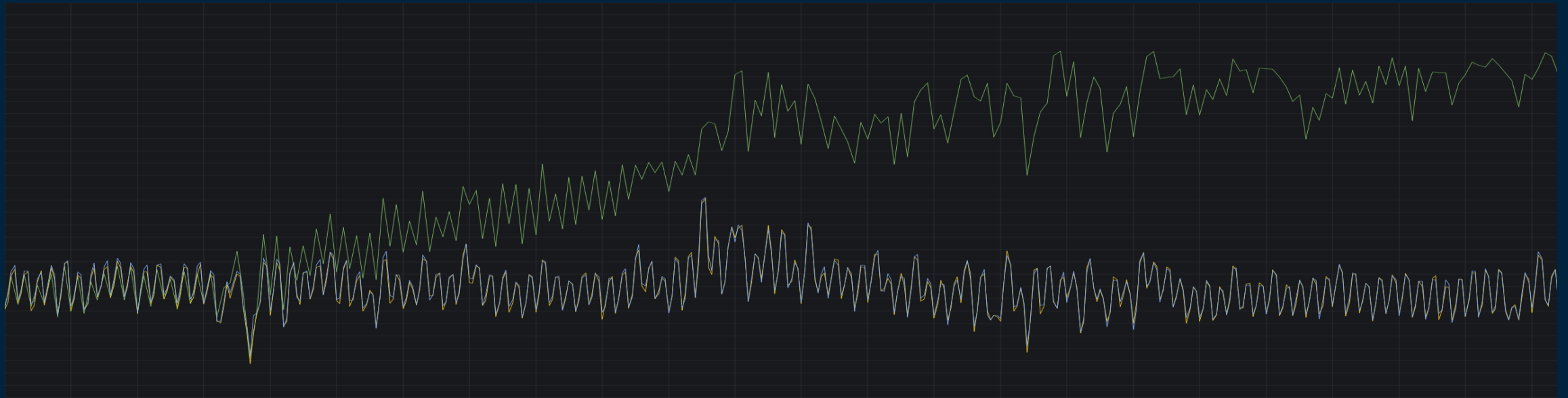
Why are CPU spikes an issue?

Under normal operations



During rollout of zone 4





What is a metric series?

response_time{}



response_time{status_code="200"}
response_time{status_code="400"}
response_time{status_code="500"}



response_time{status_code="200", path="/read", continent="eu"}
response_time{status_code="400", path="/read", continent="eu"}
response_time{status_code="500", path="/read", continent="eu"}
response_time{status_code="200", path="/write", continent="eu"}
response_time{status_code="400", path="/write", continent="eu"}
response_time{status_code="500", path="/write", continent="eu"}
response_time{status_code="200", path="/read", continent="na"}
response_time{status_code="400", path="/read", continent="na"}
response_time{status_code="500", path="/read", continent="na"}
response_time{status_code="200", path="/write", continent="na"}
response_time{status_code="400", path="/write", continent="na"}
response_time{status_code="500", path="/write", continent="na"}
response_time{status_code="200", path="/read", continent="oc"}
response_time{status_code="400", path="/read", continent="oc"}
response_time{status_code="500", path="/read", continent="oc"}
response_time{status_code="200", path="/write", continent="oc"}
response_time{status_code="400", path="/write", continent="oc"}
response_time{status_code="500", path="/write", continent="oc"}



response_time{status_code="200", path="/read"}
response_time{status_code="400", path="/read"}
response_time{status_code="500", path="/read"}
response_time{status_code="200", path="/write"}
response_time{status_code="400", path="/write"}
response_time{status_code="500", path="/write"}

Each unique set of label values is one metric series

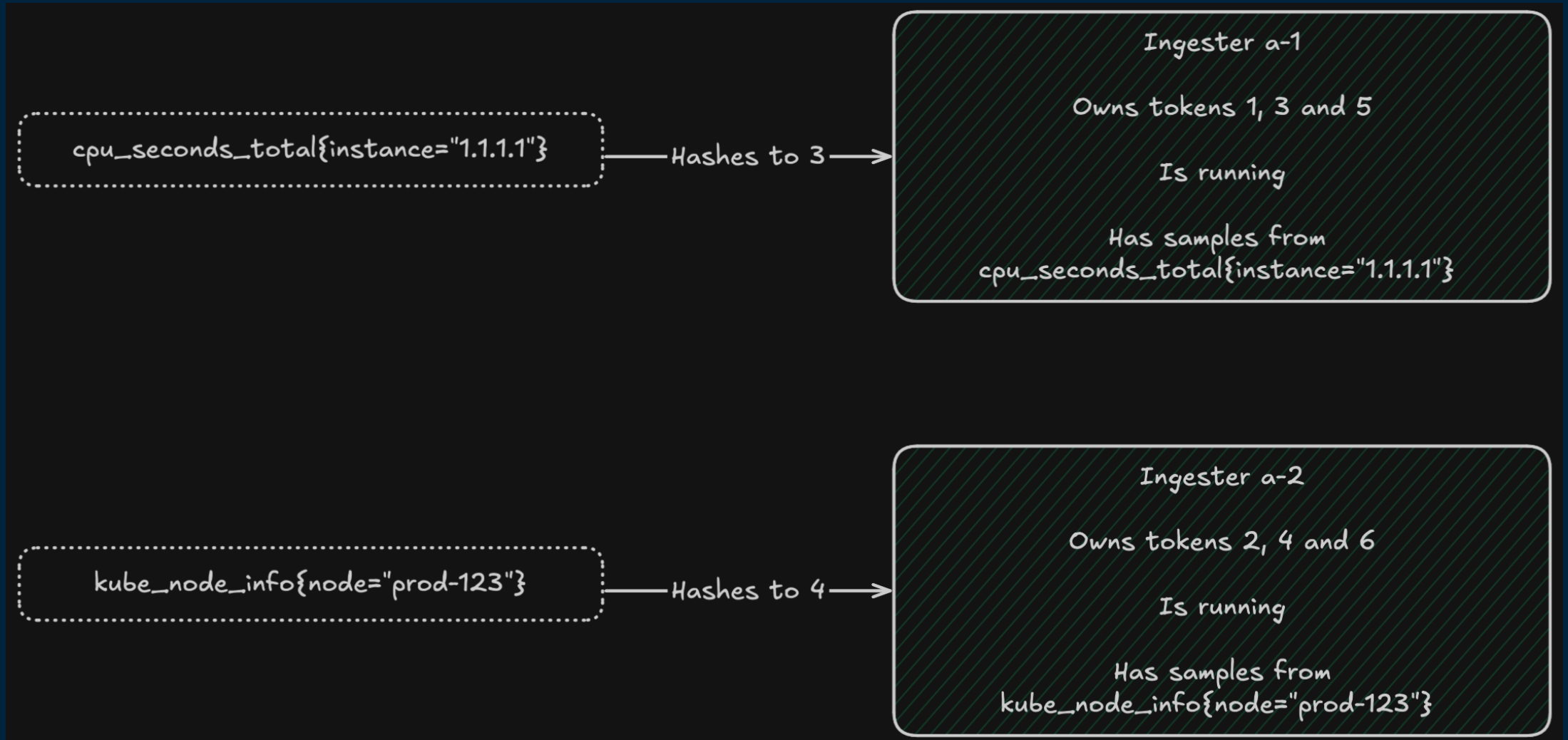
What makes up an ingester's CPU usage?

$(\text{metric_series_count} * \text{reads_per_second}) + \text{writes_per_second} + \text{misc} = \text{CPU usage}$

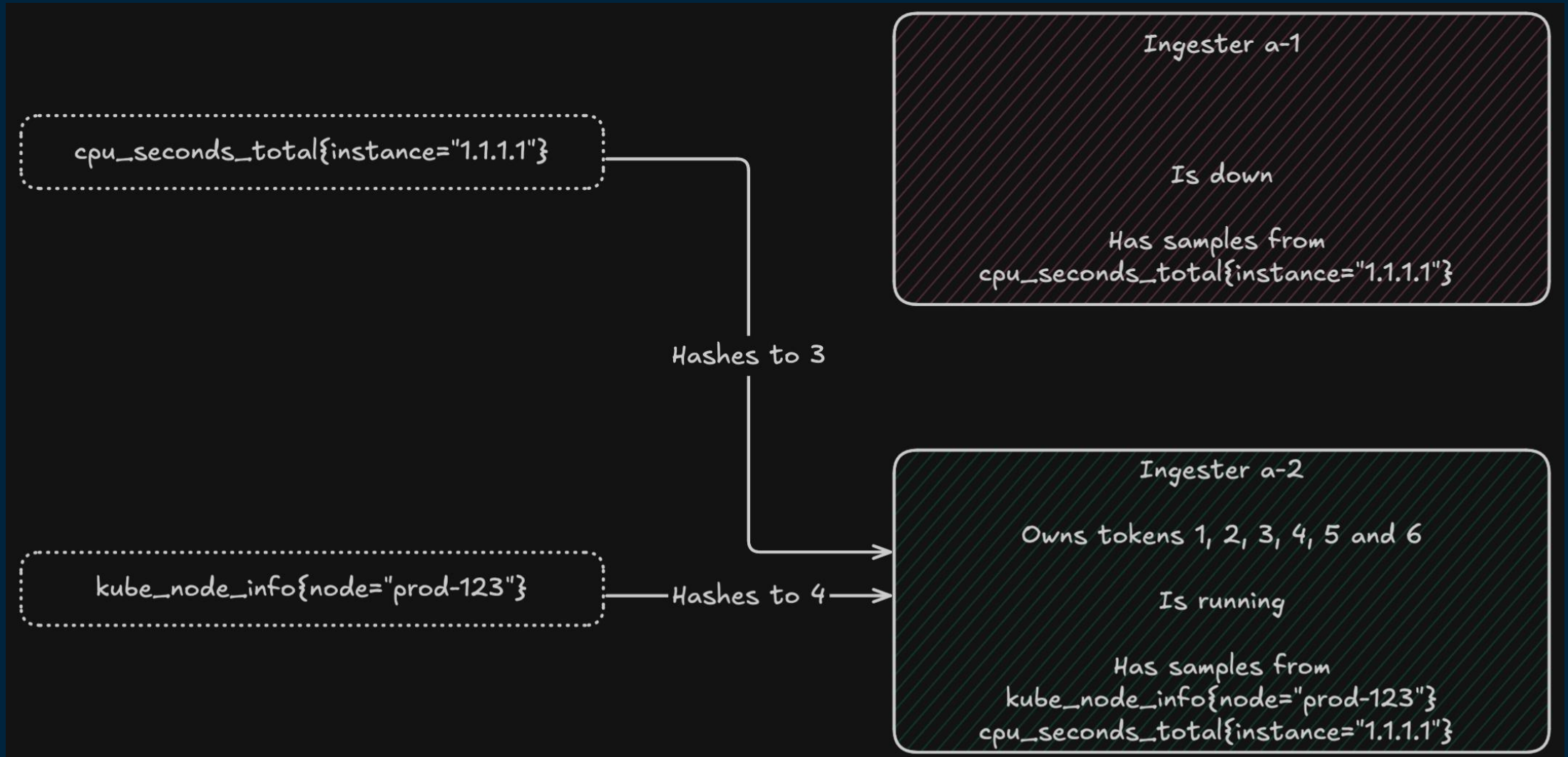
	Few reads per second	Many reads per second
Low metric series count	Low CPU usage	Medium CPU usage
High metric series count	Medium CPU usage	High CPU usage

This is why we care about metric cardinality!

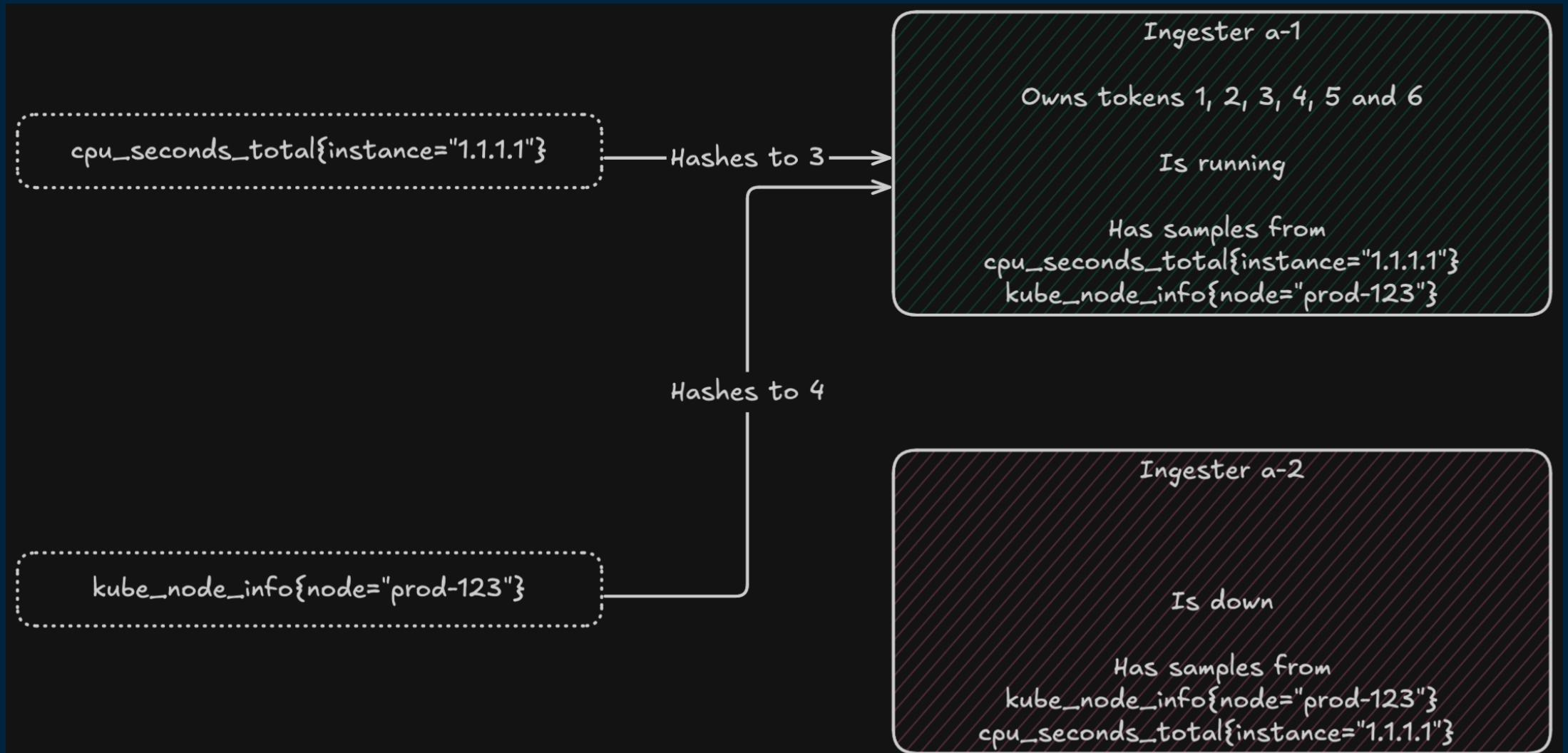
Why do metric series proliferate during rollouts?



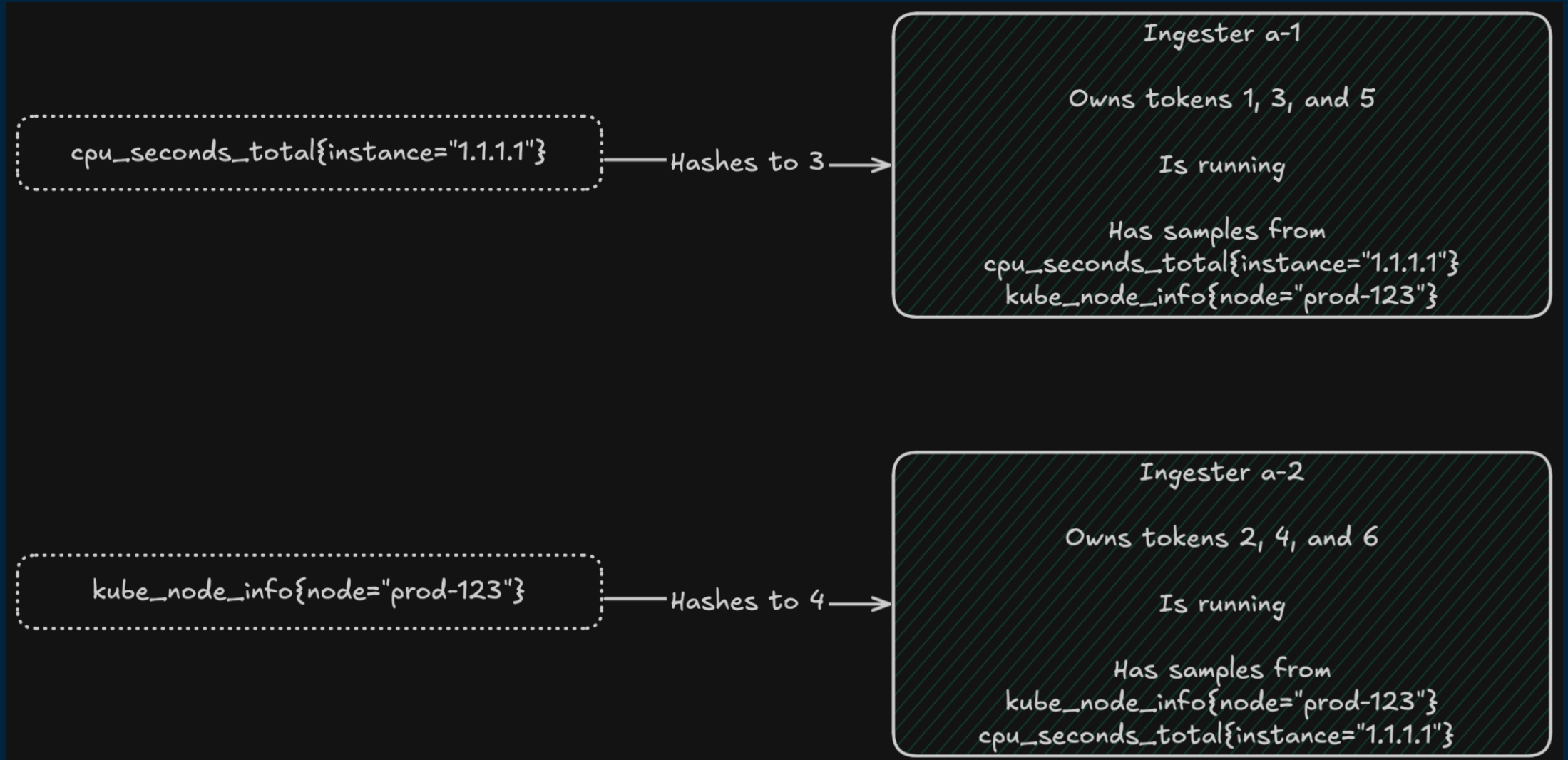
Why do metric series proliferate during rollouts?



Why do metric series proliferate during rollouts?



Why do metric series proliferate during rollouts?



Let's recap

- Read latency increases during rollouts *because*
- Nodes don't have enough CPU time to serve demand *because*
- Ingestor CPU usage surges *because*
- The amount of metric series held by each ingestor increases by ~100% *because*
- Series ownership is handed back and forth between ingestors *because*
- Ingestors are continuously terminated and leave the ring

Potential solutions

Don't do rollouts

Potential solutions

Don't leave the ring on shutdown

Pros:

- Metric series no longer proliferate during rollouts.
- Easy to implement.

Cons:

- Multi-zone evictions cause downtime.

Zone A:

- Ingester a-1
- Ingester a-2
- Ingester a-3
- Ingester a-4

Zone B:

- Ingester b-1
- Ingester b-2
- Ingester b-3
- Ingester b-4

Zone C:

- Ingester c-1
- Ingester c-2
- Ingester c-3
- Ingester c-4

It takes only a single terminated ingester for the entire zone to be considered unhealthy!

Potential solutions

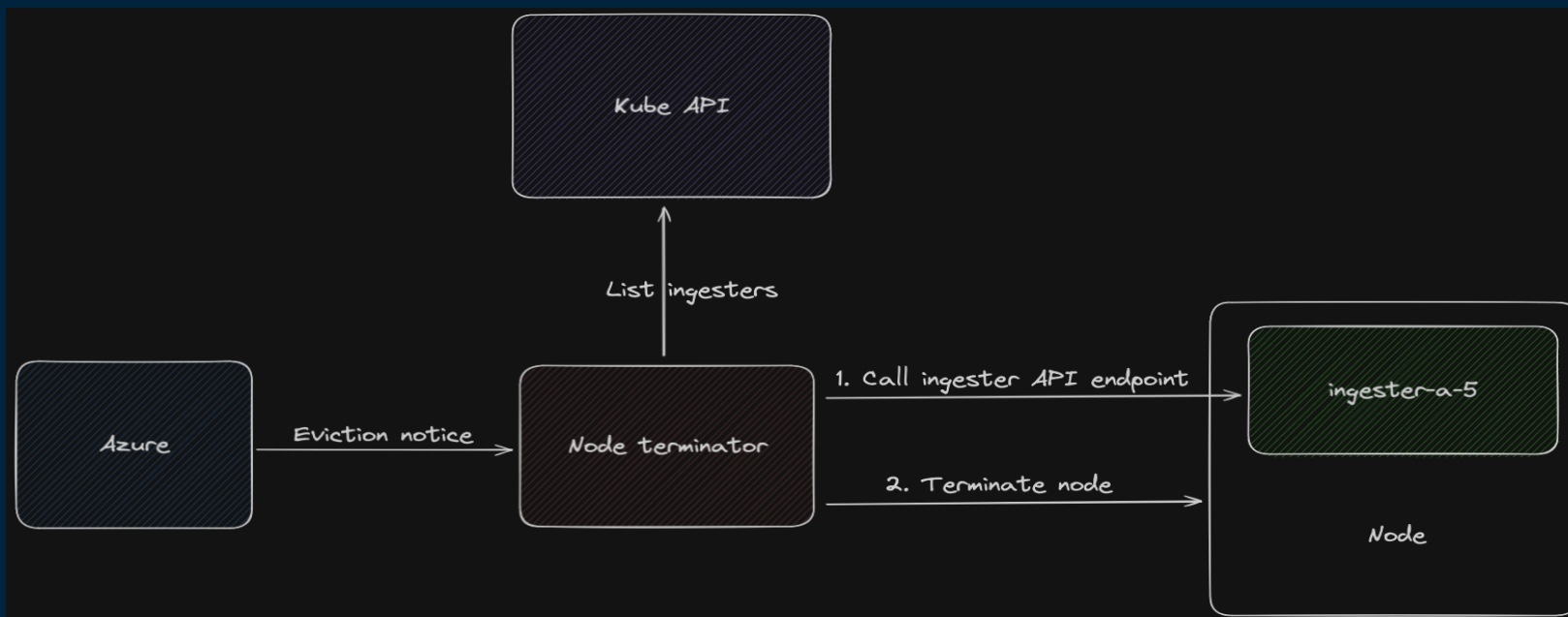
Conditionally leave the ring on shutdown

Pros:

- Metric series no longer proliferate during rollouts.
- Multi-zone evictions don't cause downtime

Cons:

- Requires buy-in from Mimir team and upstream contribution.



Gains



CPU savings are equivalent to ~\$13700 per month

Q&A

- Ingestor CPU utilisation balloons during rollouts because metric series proliferate
- Ingestors cannot always stay in the ring because of spot node evictions
- The introduction of a new API endpoint allows ingestors to conditionally leave the ring