



Sailing the Database Seas: Applying SRE principles at scale

SRECon EMEA
Tuesday 29 October 2024, Dublin

Booking.com

martin.alderete@booking.com
ioannis.androulidakis@booking.com

\$ whoami

- ❑ Ioannis Androulidakis (@ioandr)
- ❑ Site Reliability Engineer at Booking.com
- ❑ Database Engineering Team, Application Data Services
- ❑ Diploma in Electrical & Computer Engineering, NTUA
- ❑ Operating Systems, Storage, Observability
- ❑ Linux enthusiast, OSS contributor



\$ whoami

- ❑ Martin Alderete (@**ma1derete**)
- ❑ Principal Site Reliability Engineer at Booking.com
- ❑ Database Engineering Team, Application Data Services
- ❑ 15+ years of experience in Distributed Systems,
System-level Programming and Site Reliability
- ❑ Passionate about OSS (also contributor), member of
different technical groups





The Journey

1. Database Reliability at Booking.com
2. SLIs and SLOs for Distributed Database Systems
3. Automating MySQL Capacity Planning
4. Postmortem Culture
5. Q/A



The Journey

1. **Database Reliability at Booking.com**
2. SLIs and SLOs for Distributed Database Systems
3. Automating MySQL Capacity Planning
4. Postmortem Culture
5. Q/A

Databases at Booking.com

MySQL is our main relational datastore

Single primary for WRITES, **tens/hundreds** of replicas for READs
In general, choose **availability** over consistency ¹



MySQL
clusters
> 250

MySQL servers
> 10K

Multi-cloud

Our fleet spans across different platforms:

- Bare-Metal (ServerDB)
- Private Cloud (Openstack)
- Public Cloud (AWS)
 - Self-managed (EC2)
 - Managed (RDS)

Multi-region

Always deploy in > 2 regions
(typically **3 regions** in Europe)

QPS

100M+

Services

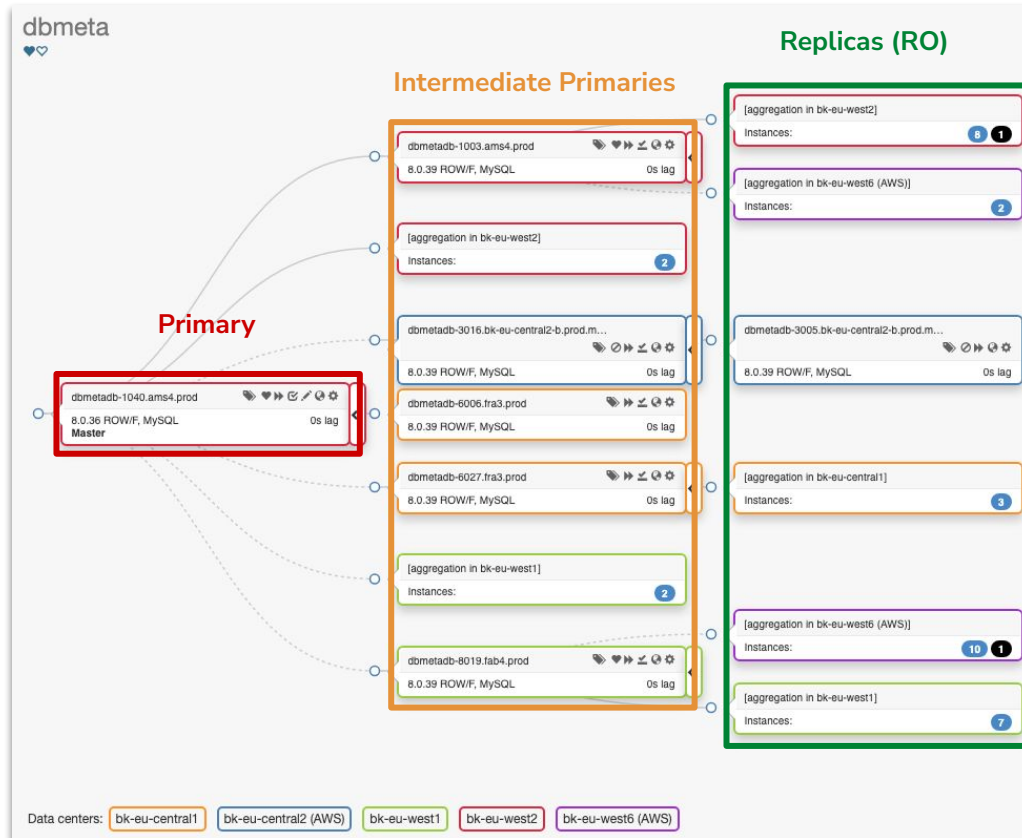
Transactions, Fintech,
Payments, Frontend,
Customer Support

Ownership

We do **not** own data
We do **not** own tables or schemas
We **do** own database infrastructure

1. https://en.wikipedia.org/wiki/CAP_theorem / https://en.wikipedia.org/wiki/PACELC_theorem

Topology of semi-sync MySQL cluster



Database Reliability Engineering (DBRE)

Set of **tools** and **practices** around database systems to ensure they are **reliable, scalable** and **compliant** with regulations.



Provide database expertise

Promote best practices for databases, automate database-related operations, provide self-service tools



Participate in on-call rotation

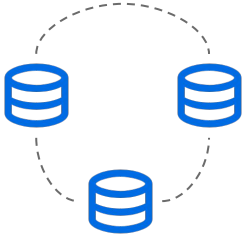
Support and debug database issues across services, join firefighting calls, cooperate with peer SREs



Deliver database observability

Implement database monitoring, alert on SLOs, measure performance

SRE Principles for Database Engineering



Monitoring Distributed Systems ¹

Database clusters are distributed across **clouds** and **locations**.
Database **primaries** and **replicas** are different by nature.



Eliminating Toil ²

Running databases at scale comes with a variety of **repetitive operational** tasks that are common across clusters. Aim for automation.



Postmortem culture ³

Outages **WILL** happen. Do not solve the same problem twice. Educate people and understand systems better.

1. <https://sre.google/sre-book/monitoring-distributed-systems>
2. <https://sre.google/sre-book/eliminating-toil>
3. <https://sre.google/sre-book/postmortem-culture>



The Journey

1. Database Reliability at Booking.com
2. **SLIs and SLOs for Distributed Database Systems**
3. Automating MySQL Capacity Planning
4. Postmortem Culture
5. Q/A

Database performance is sensitive to users

- Full table scans → long-running processes
- Non-indexed, complex queries → slow, expensive
- Database-as-a-queue → locking issues, timeouts
- READ from primary instead of replicas → increased load on primary
- Bulk UPDATE/DELETE → replication delay
- Online Schema Changes (OSC) on huge tables → contention, I/O load
- Excessive number of client connections → hit max connections limit
- ...

Which SLIs we **reject** for database SLOs

1

Query Latency

The time it takes for a MySQL server to execute a query and return result

2

Throughput

The number of queries a MySQL server can handle per second

3

Replication Delay

The amount of time a MySQL replica is behind its replication source

“If you cannot measure it (correctly) you cannot improve it”



Decide on the right depth!

Customers **don't** care about
`innodb_buffer_pool_pages_total`
or `cpu_usage_total`



Which SLIs we **use** for database SLOs

1

Query Latency

The time it takes for a MySQL server to execute a query and return result

2

Throughput

The number of queries a MySQL server can handle per second

3

Replication Delay

The amount of time a MySQL replica is behind its replication source

1

Read Availability

MySQL replicas respond to read requests from clients

2

Write Availability

MySQL primaries respond to write requests from clients

3

Replication Running

MySQL replicas receive changes from source and apply them locally

Read Availability SLO

If at least one successful read probe is achieved the data point is a **pass / OK**, otherwise a **miss / KO**.

The individual measurements are composed over time by taking the percentage of **OK** measurements with respect to the total number of measurements.

```
ten_secondly.mysql.$cluster.read_availability.$region.ok  
ten_secondly.mysql.$cluster.read_availability.$region.ko
```

Target

99.99% ("four nines")

Schedule

The SLI is calculated **24/7 every 10 sec**

Who

MySQL Availability Reporter (Golang). Runs with multiple instances across all regions

How

Execute a read probe (SELECT) against **two random replicas** of every MySQL cluster

Timeouts

Connect: **2 sec**, Read: **5 sec**

Write Availability SLO

If at least one successful write probe is achieved the data point is a **pass / OK**, otherwise a **miss / KO**.

The individual measurements are composed over time by taking the percentage of **OK** measurements with respect to the total number of measurements.

```
secondly.mysql.$cluster.write_availability.$region.ok  
secondly.mysql.$cluster.write_availability.$region.ko
```

Target

99.97% ("three nines seven")

Schedule

The SLI is measured **24/7 every 1 sec**

Who

MySQL Availability Reporter (Golang). Runs with multiple instances across all regions

How

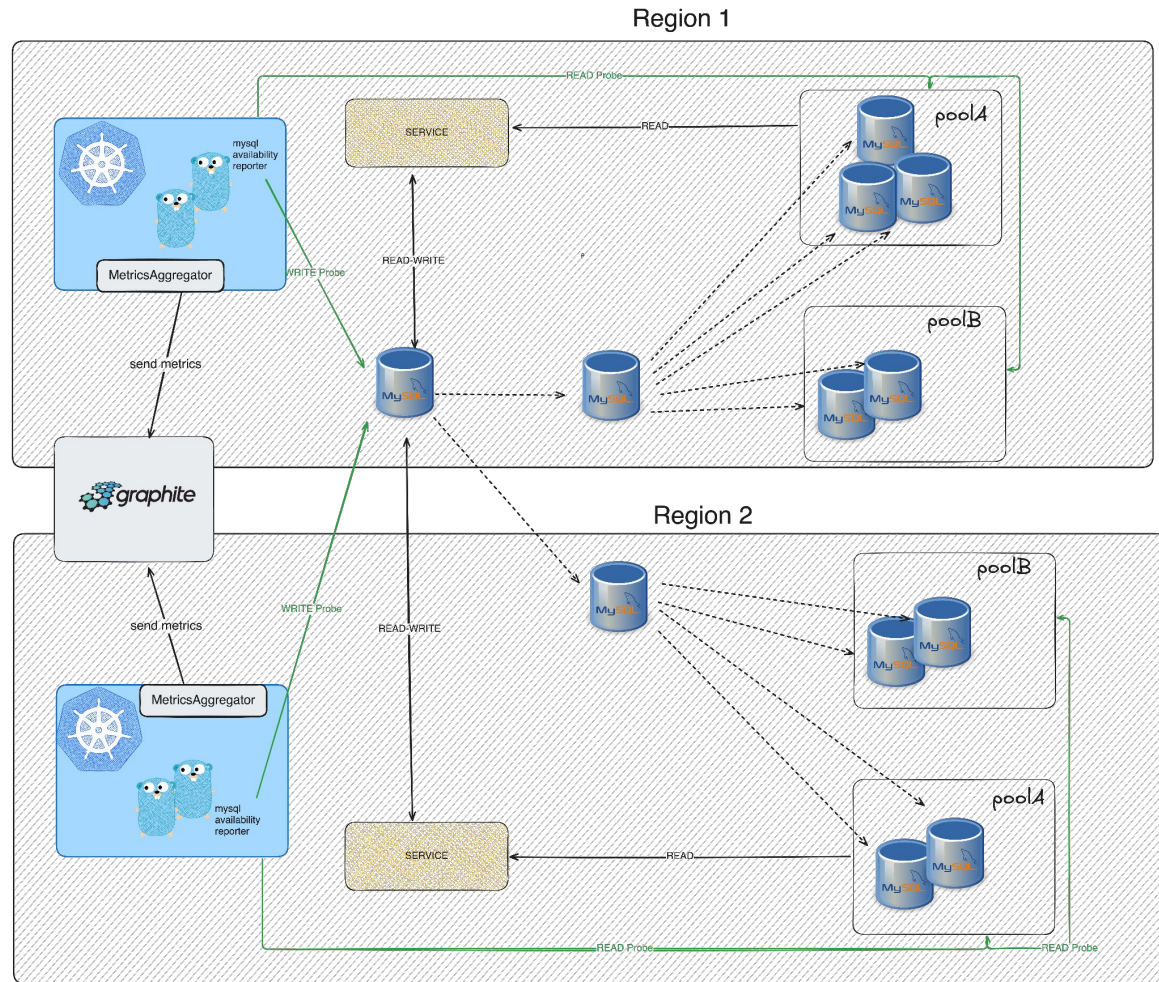
Execute a write probe (INSERT) against **the (single) primary** of every MySQL cluster

Timeouts

Connect: **2 sec**, Write: **8 sec**

MySQL Availability Reporter

- In-house **Golang** project
- Runs on Kubernetes and works **across platforms**
- **Imitates** MySQL's consumers (apps)
- **Knows** about Service Discovery
- **Monitors** READ and WRITE availability with **probes**
- Stores (aggregated) metrics in remote **Graphite** instance



Replication Running SLO

If replication is running between the replica and the primary (source) the data point is considered **pass / OK**, otherwise a **miss / KO**.

The individual measurements are composed over time by taking the percentage of **OK** measurements with respect to the total number of measurements.

```
ten_secondly.mysql.$cluster.$region.replication_running.ok  
ten_secondly.mysql.$cluster.$region.replication_running.ko
```

1. <https://dev.mysql.com/doc/refman/8.4/en/performance-schema.html>
2. <https://dev.mysql.com/doc/refman/8.4/en/show-replica-status.html>

Target

99.95% ("three nines five")

Schedule

The SLI is measured **24/7 every 10 sec**

Who

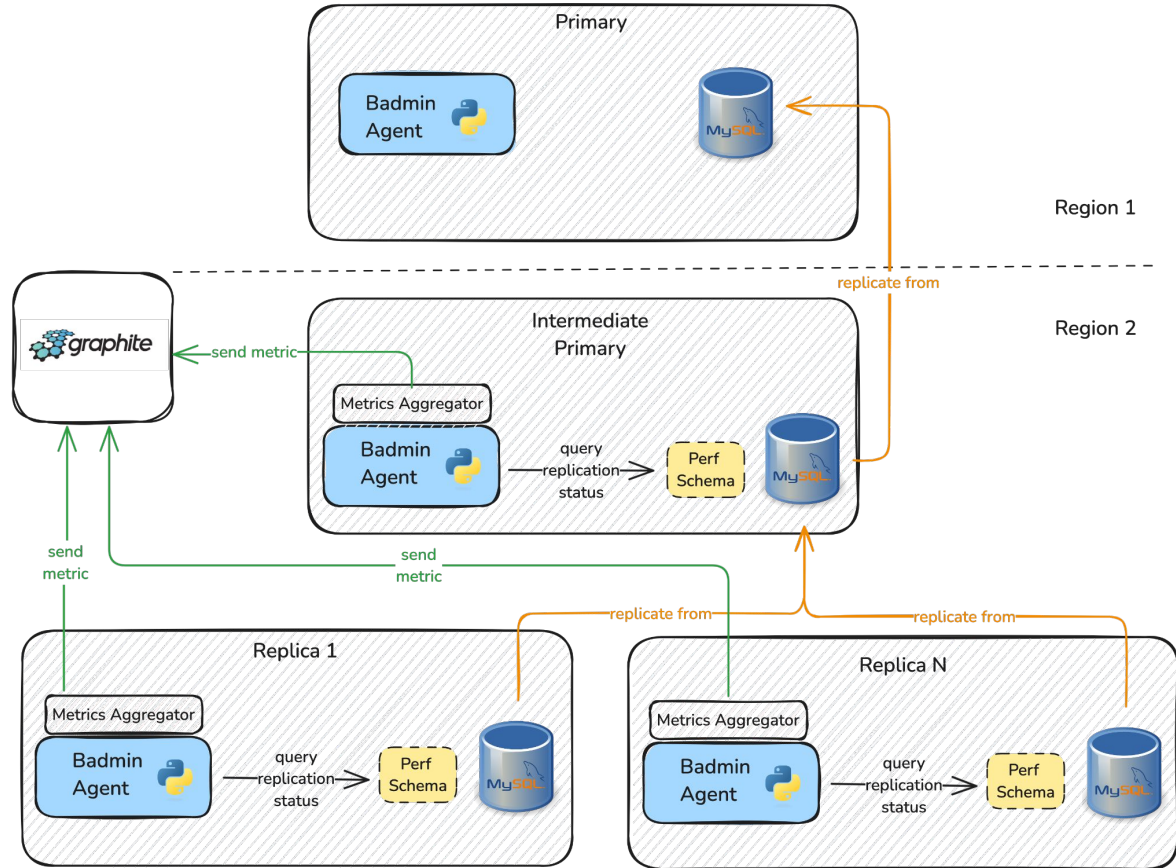
Badmin Agent (Python). Runs with a single instance inside every MySQL server

How

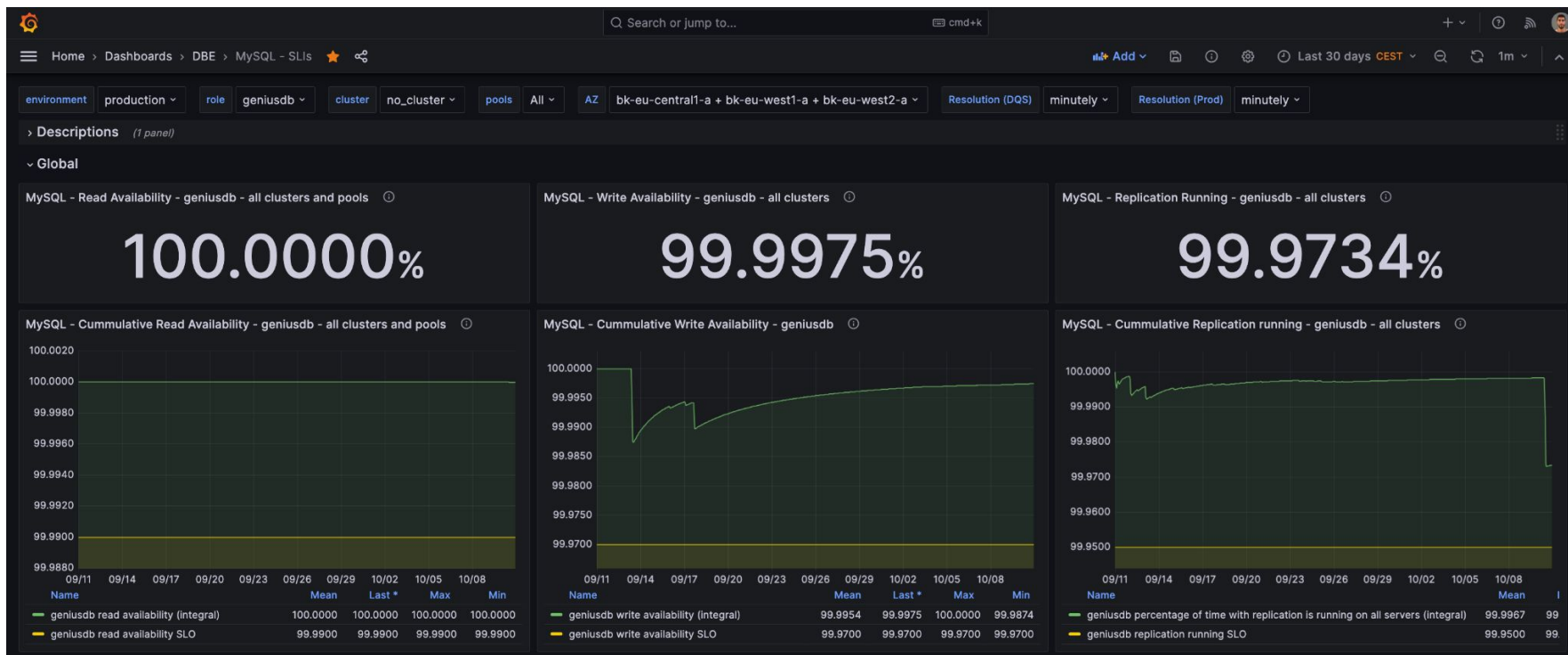
Query the performance schema¹ for replication status² of **every replica** of every MySQL cluster

Badmin Agent

- In-house **Python** project (10+ years old)
- **Sits** next to every MySQL server
- **Provides** functionality and automation for DBA tasks
- **Queries** MySQL performance schema for replication status
- Stores (aggregated) metrics in remote **Graphite** instance

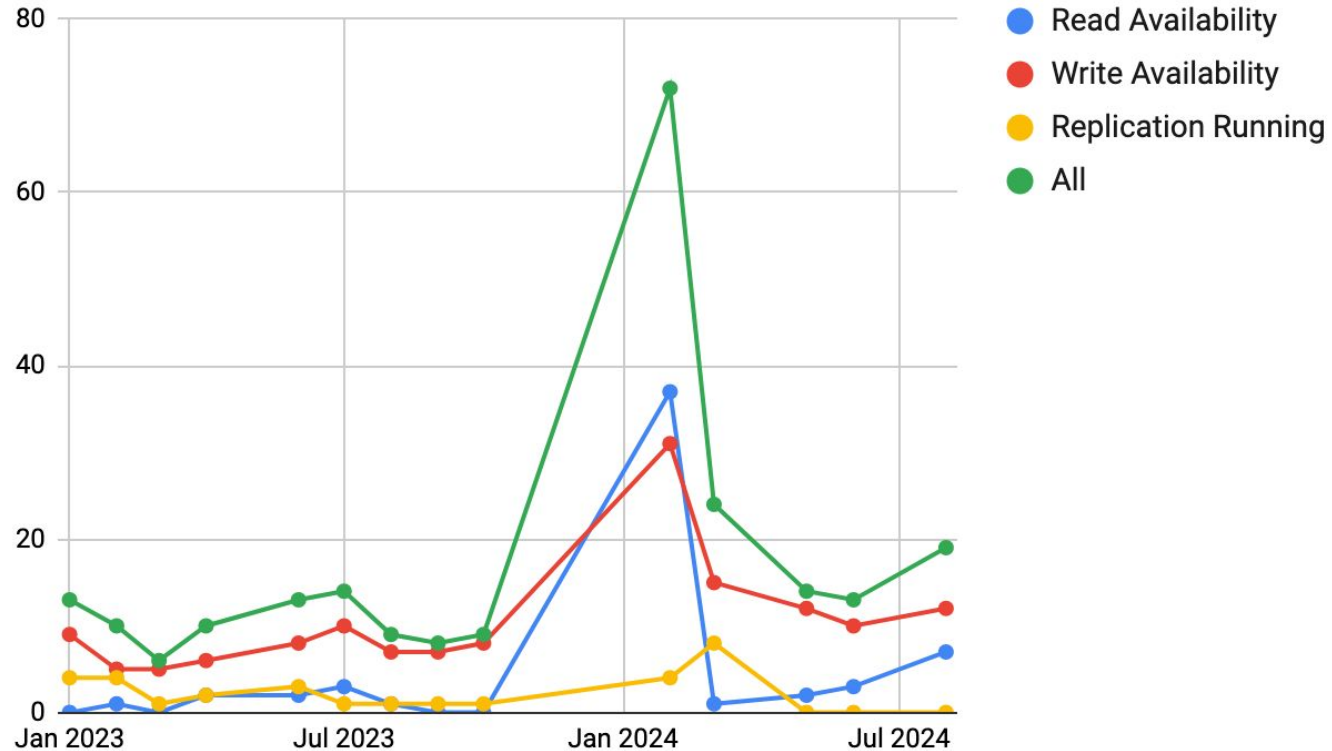


MySQL SLIs/SLOs in action - Genius DB

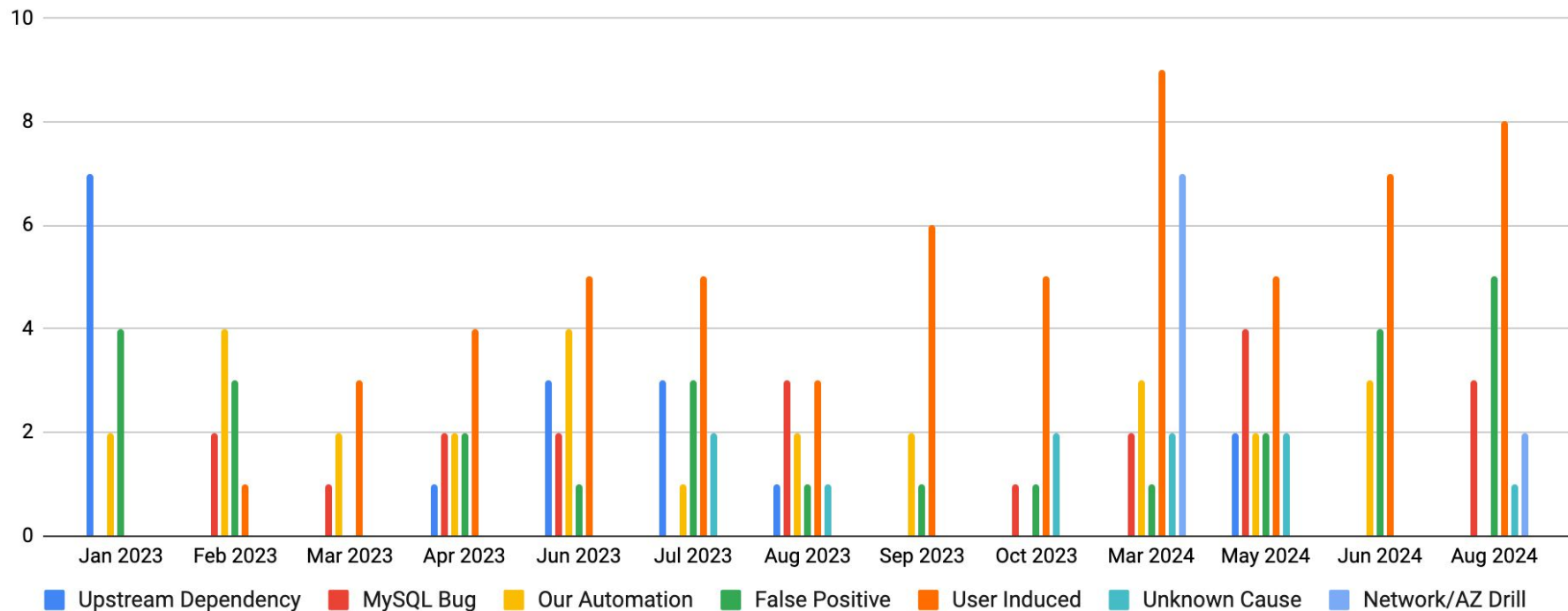


MySQL SLO Violations

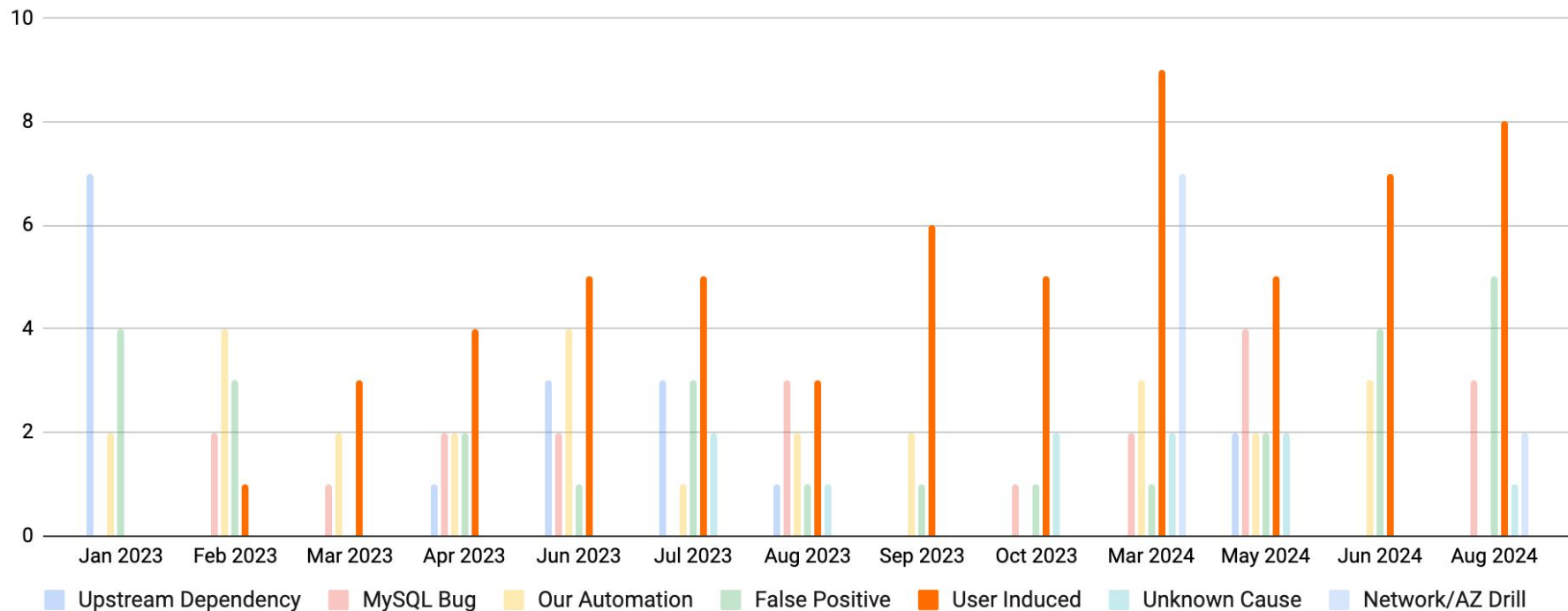
- Write Availability SLO is violated more often
- The number SLO violations is small given the number of MySQL clusters
- Spikes can happen



Who violates our MySQL SLOs?



Who violates our MySQL SLOs?



Takeaways



Choose SLIs that are 100% under your control to define your database SLOs

This helps you avoid “false positives” and stay within your SLO budget. Reduce “alert fatigue”



Decide where it makes sense to measure your SLI

Depending on the metric, you might need to measure close to the database server or mimic “how the application sees it”



Don't pursue much reliability than what's strictly necessary

Aiming at a 100% SLO is a bad decision. Use your SLO budget and always educate developers





The Journey

1. Database Reliability at Booking.com
2. SLIs and SLOs for Distributed Database Systems
3. **Automating MySQL Capacity Planning**
4. Postmortem Culture
5. Q/A

Database tasks introduce TOIL

- Server provisioning, configuration, maintenance
- Schema splits, changes
- Password rotation, grants
- Switchover/failover of primaries
- Daily backups and restores of volumes
- MySQL upgrades
- Fleet sizing
- Connection management
- Heal broken replication
- ...





Our approach: Relentless Automation

Reliability is our top priority: alerts alone are not **enough**

Choose **proactive** over reactive

Our fleet **grows faster** than our database engineering teams



Why we automated MySQL Capacity Planning

Capacity planning is critical for **Keeping The Lights On** (KTLO)

- **Risk mitigation:** undersized MySQL clusters are a risk for business (peak load)
- **Cost control:** oversized MySQL clusters increase our bill
- **Scalability:** time-consuming, repetitive task common for all MySQL clusters
- **Continuous Learning:** examine new hardware profiles and workloads, know your tools and infrastructure better
- **Analytics:** vertical vs horizontal scaling, forecasting, reporting

MySQL Capacity Planning: Design Decisions

01

Schedule at peak hour daily

We choose to run capacity tests on MySQL servers at the time of peak traffic, every day

02

Run against our **production fleet**

We choose to run capacity tests on MySQL servers serving production traffic

03

Ability to stop capacity test at any time

In case of emergency (e.g., outage) we want a killswitch to stop a capacity test immediately

04

Prioritize **business continuity over cost efficiency**

Introduce safety measures to ensure that MySQL pools always have sufficient capacity to serve traffic

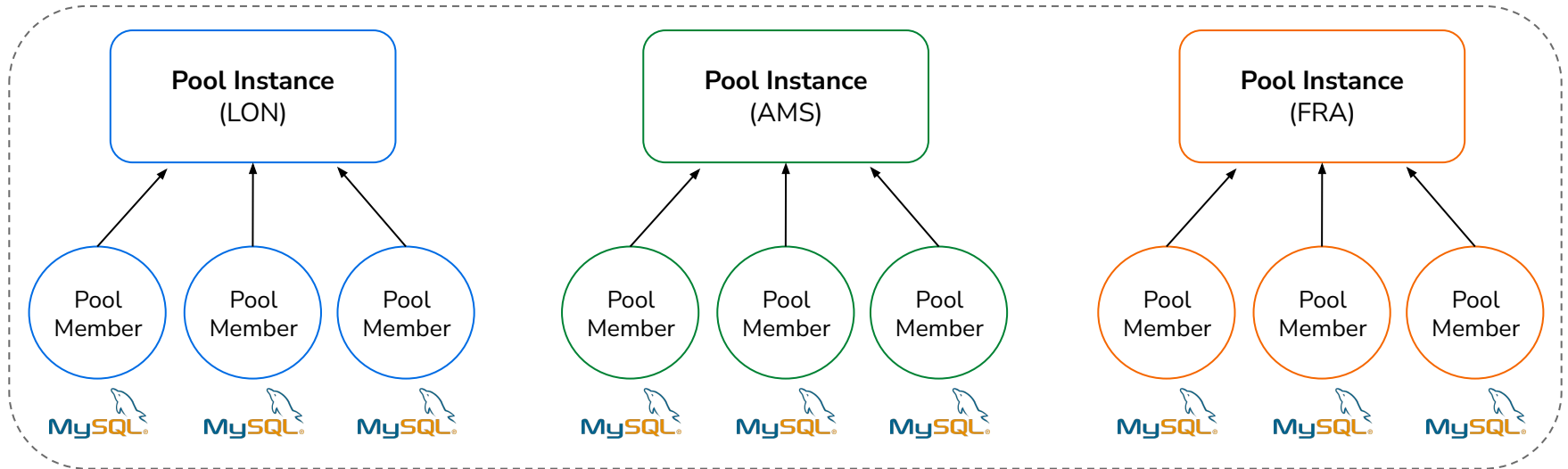
The building block for capacity planning: MySQL pools

Logical definition of **READ** query workload coming from specific services

→ Isolation, separation of concerns, solve “noisy neighbour” problem

Pools have multiple **instances** in different **regions**

Pool



MySQL Capacity Metrics



Retrieved from Prometheus `mysqld_exporter` ¹

- `mysql.capacity_planning.$cluster.$pool.$region.Com_select`

`irate()` [2m] of the number of executed SELECT queries on the MySQL server

- `mysql.capacity_planning.$cluster.$pool.$region.Threads_connected`

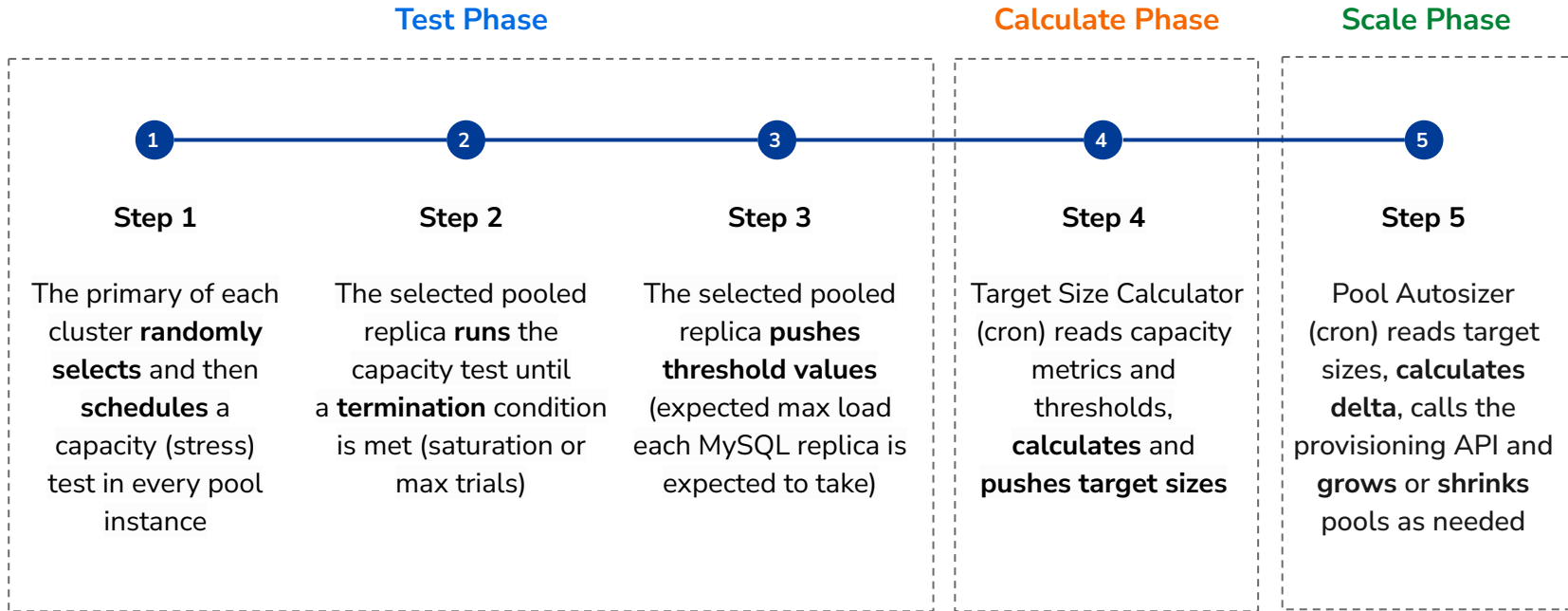
The number of threads (clients) connected to the MySQL server

- `mysql.capacity_planning.$cluster.$pool.$region.Threads_running`

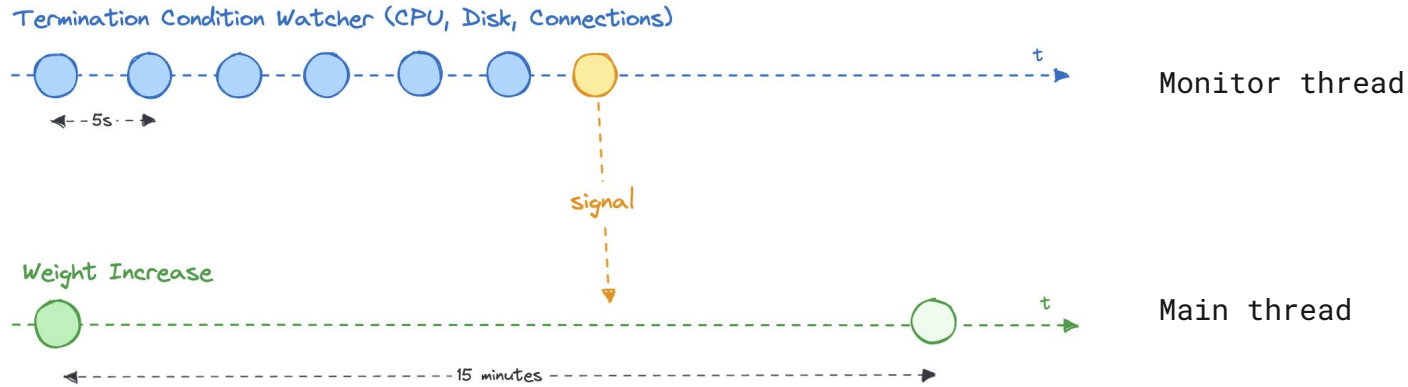
The number of database connections with an active query

1. https://github.com/prometheus/mysqld_exporter (Our fork, sorry!)

Phased Approach: Test, Calculate, Scale



Increase traffic, monitor saturation



Every 15 min we increase the **weight** of the pool member until saturation point or **MAX_TRIALS**

Detect saturation: **CPU utilization, disk read latency, MySQL max connections**¹

1. https://dev.mysql.com/doc/refman/8.4/en/server-system-variables.html#sysvar_max_user_connections

Target sizes for MySQL pools

Regional Target Size

`ceil(metric_value/threshold)`

metric value = **p90** value in the last **24h**

(emitted by mysqld exporter)

threshold = **median** over the last **7 days**

(emitted by the capacity test)

- 💡 Keep the **maximum** target size found after checking all 3 capacity metrics
- 💡 Respect empirical upper and lower **thresholds** to prevent drastic changes in pool sizes

Global Target Size

Sum of regional target sizes for all **N** regions divided by **N-1** (over provisioning for redundancy)

Extra adjustments:

- 1) **Minimum** target size adjustment (business critical clusters)
- 2) **Additive** adjustment (budget for maintenance, hardware failures, etc.)

Example

Let's assume a pool with instances in 3 regions and regional target sizes of 5, 5 and 2, respectively.

Without over-provisioning we would get

$12 / 3 = 4$ replicas per region

With over-provisioning we get

$12 / 2 = 6$ replicas per region

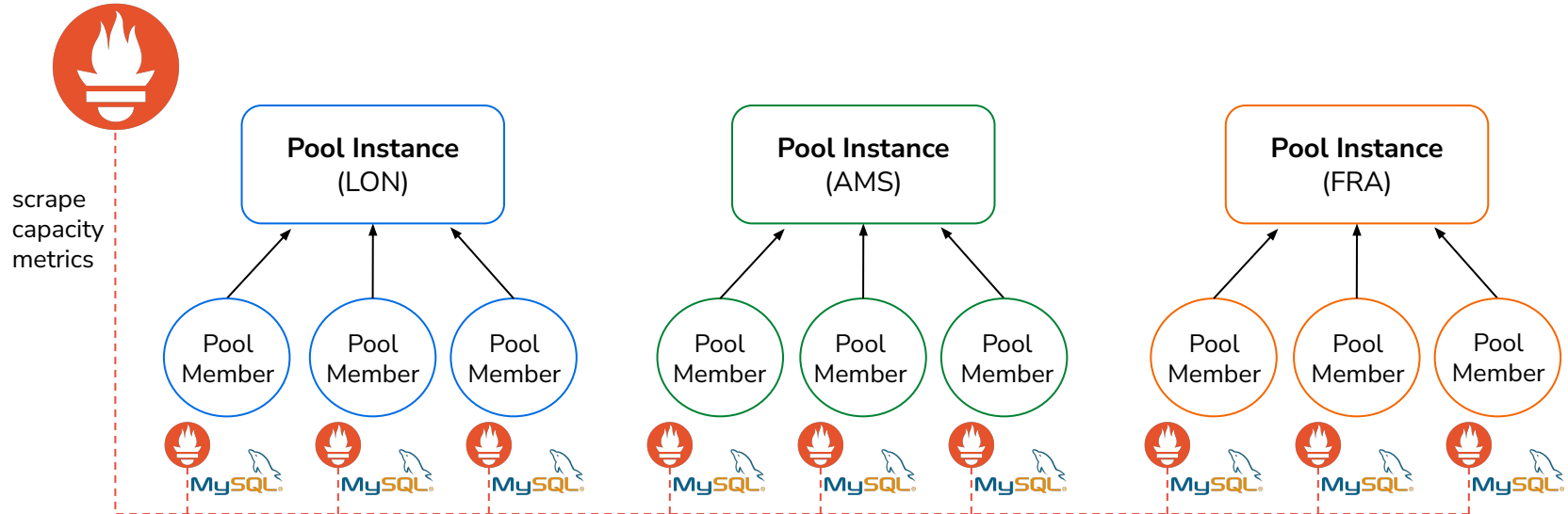


Assuming a regional failure, each pool instance should be able to handle **1.5x times** the traffic that it normally handles

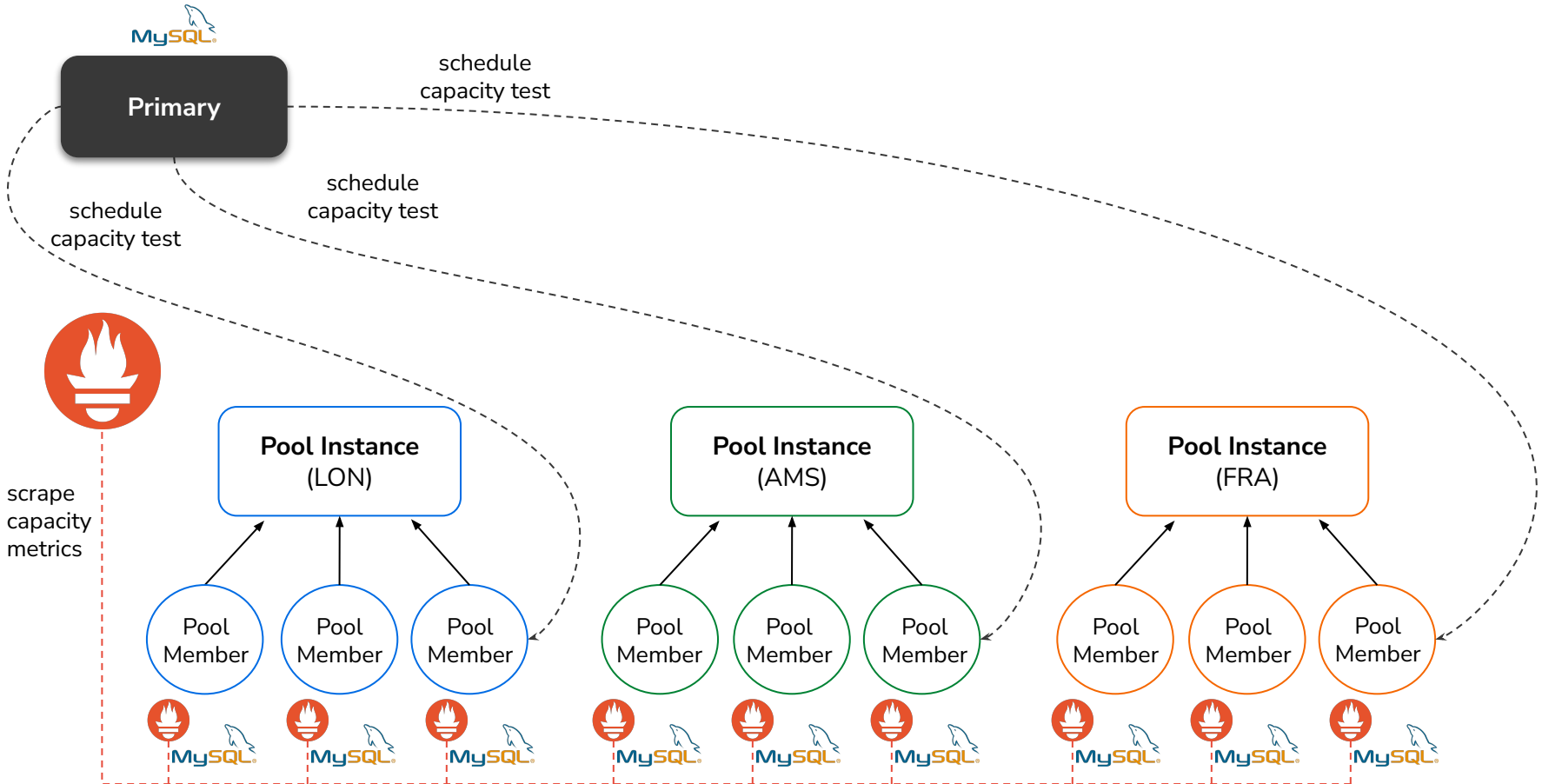
Test Phase



Primary



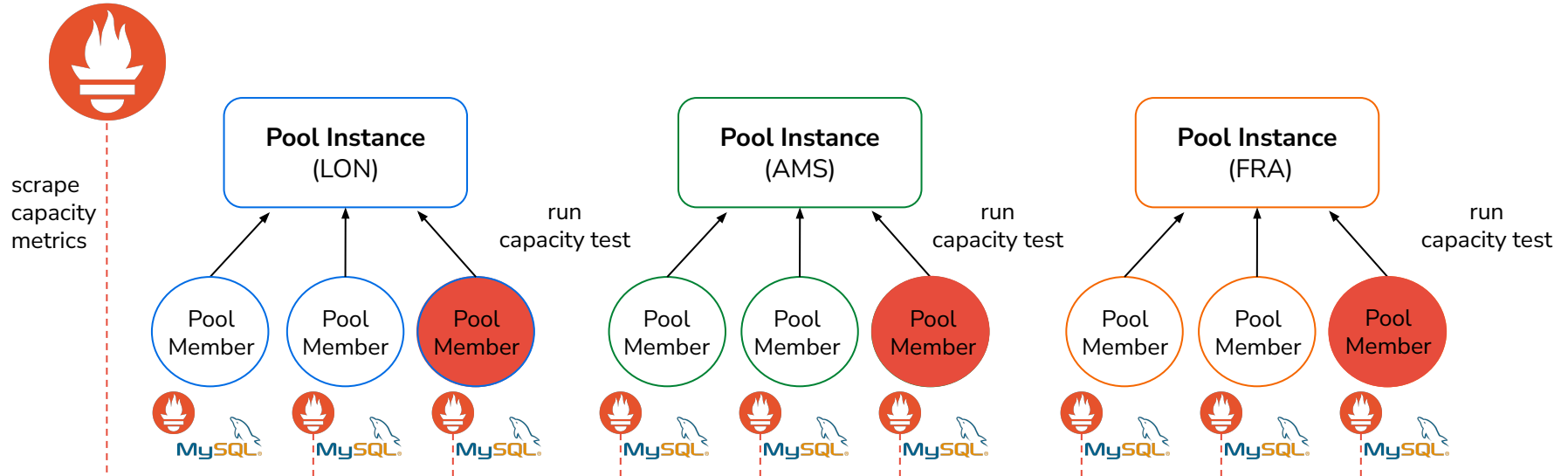
Test Phase



Test Phase



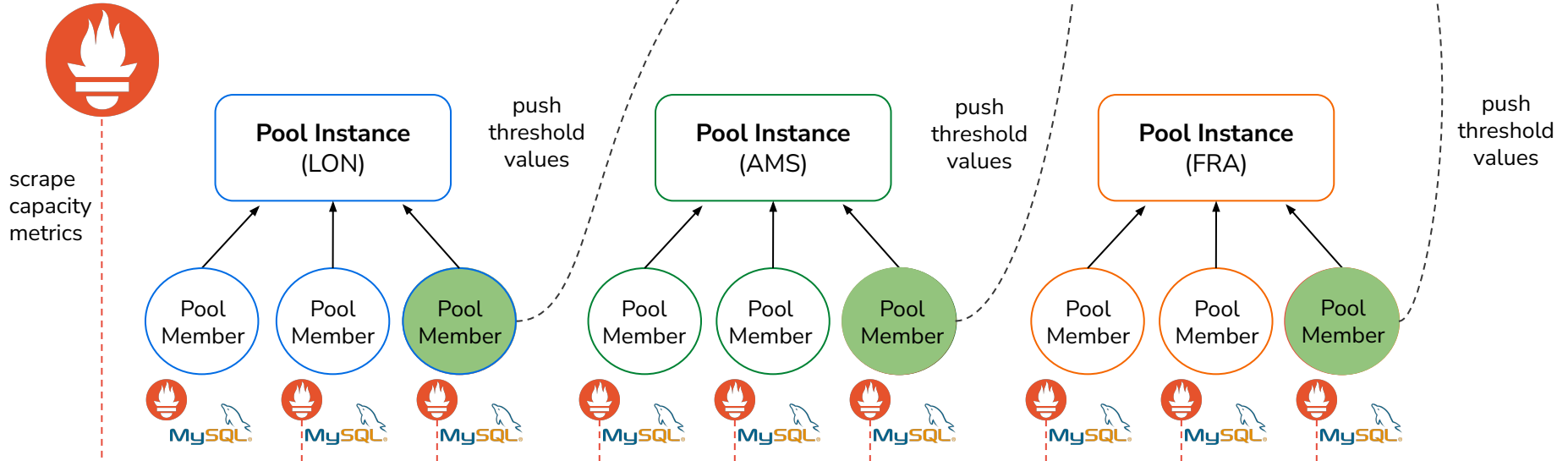
Primary



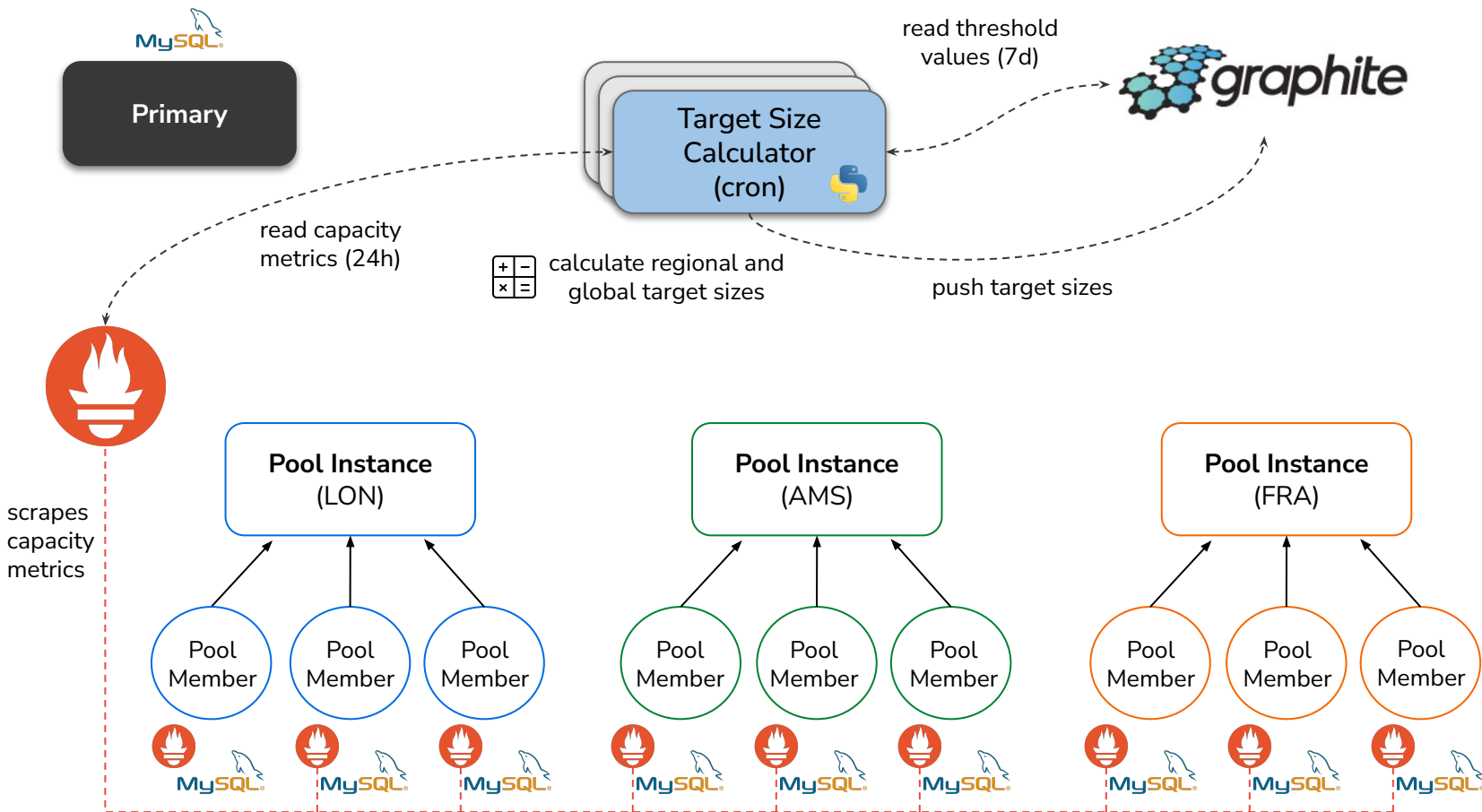
Test Phase



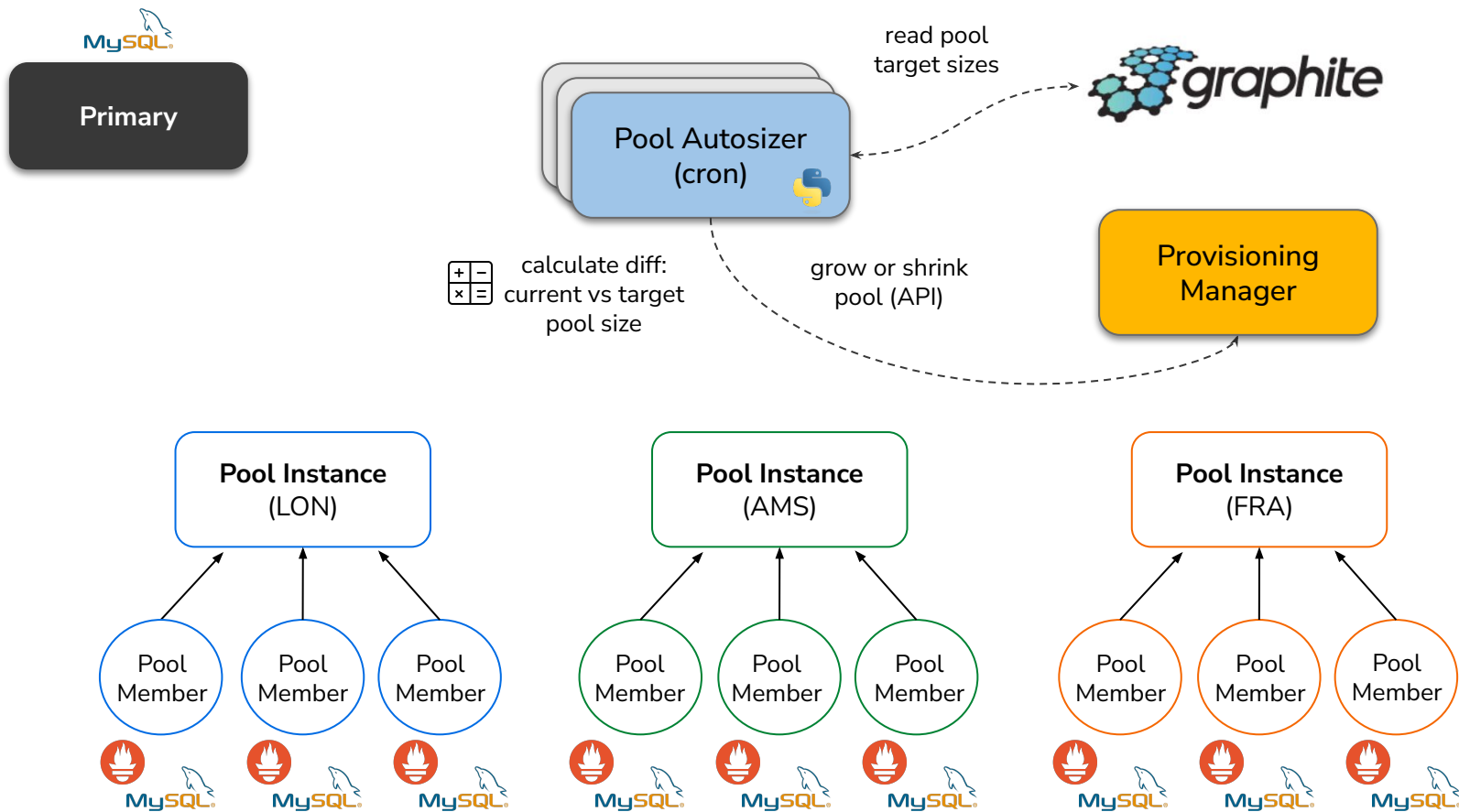
Primary



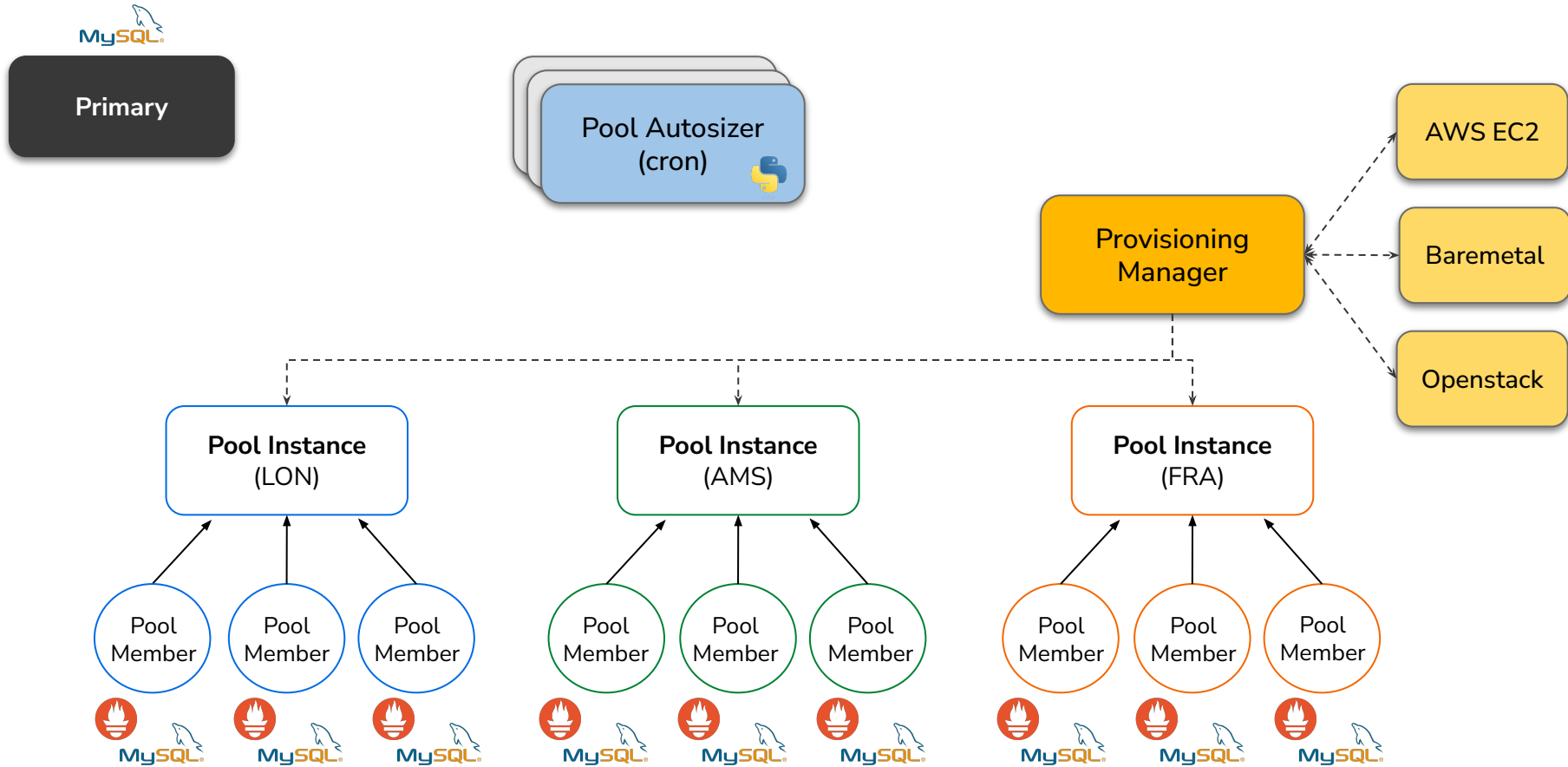
Calculate Phase



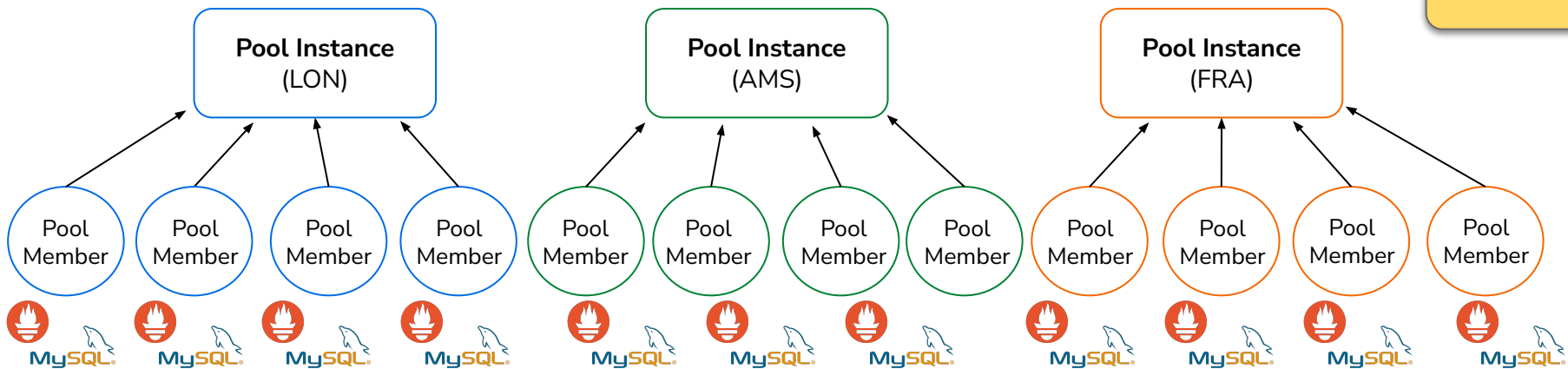
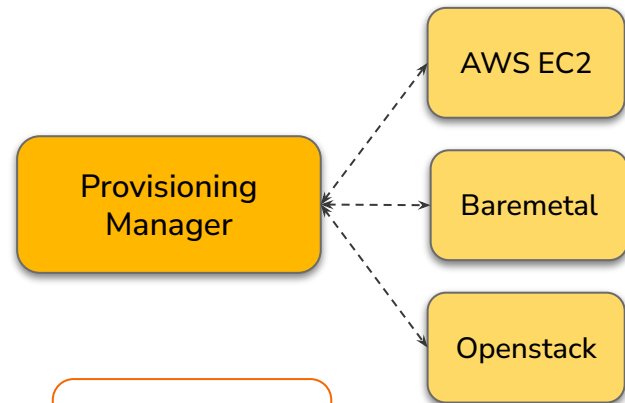
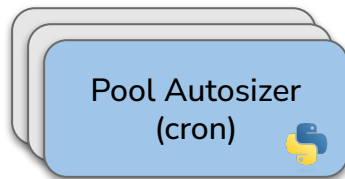
Scale Phase



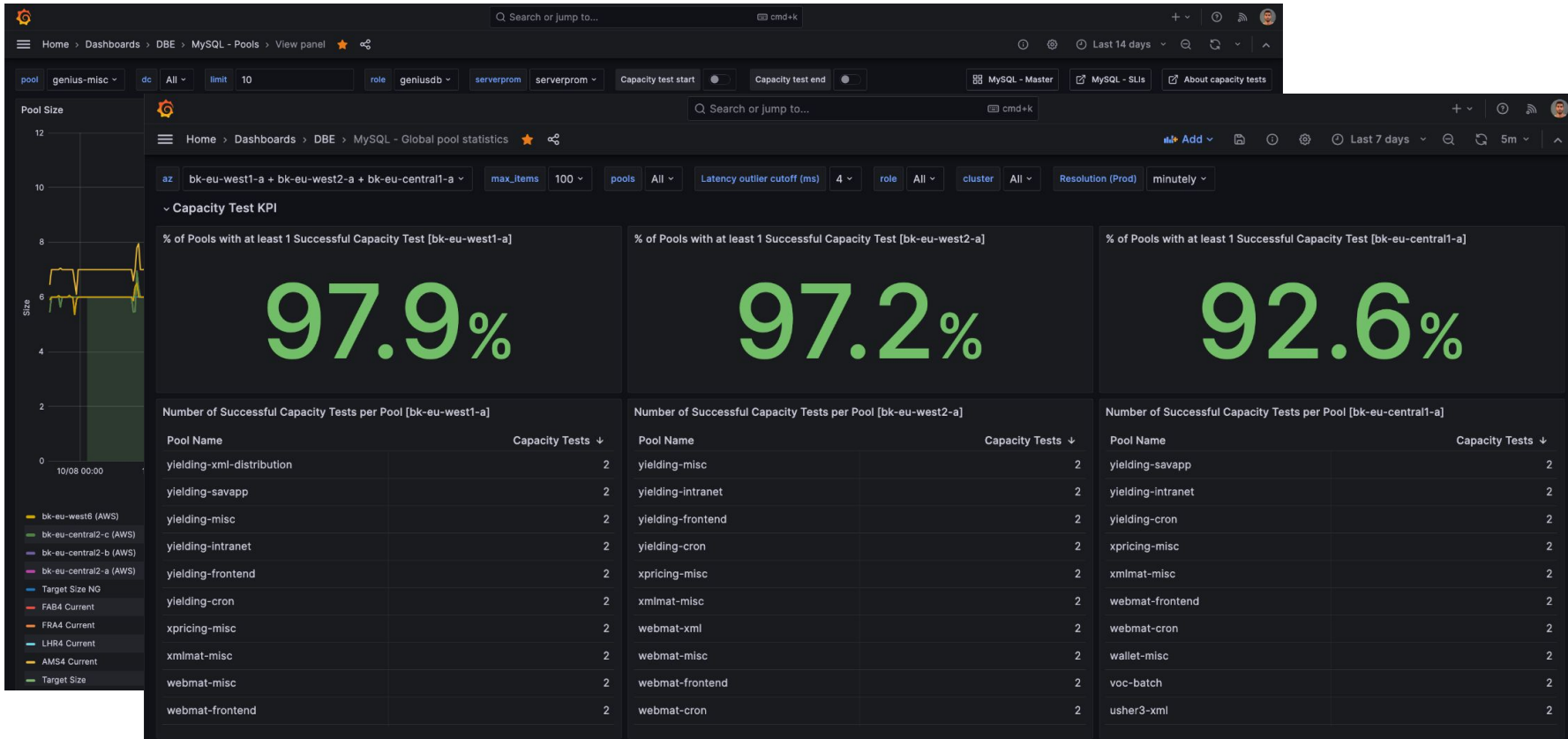
Scale Phase



Scale Phase



MySQL Pool Autosizing in action



Future Plans

- Explore **new** hardware profiles available in private and public cloud
- **Fine-tune** capacity metric thresholds towards cost efficiency
- **Minimise** errors (e.g., configure grants, database locks)
- Use a **single** sink for metrics (Prometheus vs Graphite)





The Journey

1. Database Reliability at Booking.com
2. SLIs and SLOs for Distributed Database Systems
3. Automating MySQL Capacity Planning
- 4. Postmortem Culture**
5. Q/A



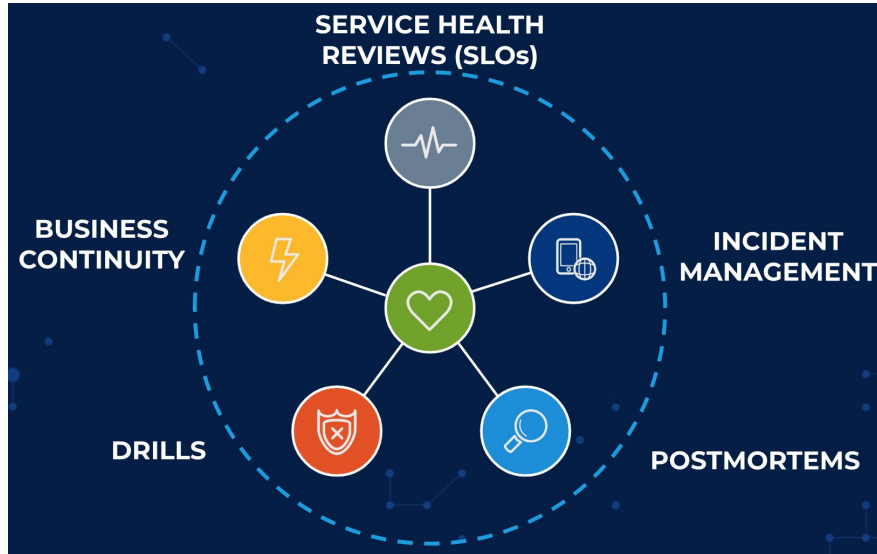
**Systems break,
that is life!**



Don't Panic!

stay calm...or try

Our Service Health Program



Service Health Review

Continuous process to reviews the statuses of services we own.

Incident Management

Standard process to enable services to deliver on SLOs

Postmortems / RFO (Reason For Outage)

Continuous learning. Incidents are a **NICE** opportunity to learn.

Network Drills

Keep our business running with the minimum impact.

Business Continuity

Automated checks of failure recovery aspects (capacity planning, failovers, redundancy, backup, restore, etc).

Postmortem - we embrace the risk together!

- Drive Service **quality** and avoid **repeated** failures
- Understanding / Eliminating the **root cause**
- **Improve** our systems and processes
- **Learning and continuous improvement**

Blameless

Blameless

Blameless



RFO (Reason For Outage) - “RFOPedia”

Enable full text search

ID	Title	Status	Category	Outage start	TTD	TTM	Owner
6394	Broken replication on multiple boxes due to lack of auth plugin	published	Resilience				
6344	P1 Drop in bookings across all platforms due to apaymentsdb database issues	published	Resilience				
6328	Expired access to Jira prevented Schema Change requests and SOX report generation	Owner		Owning Team			
		Foo Bar		MySQL DB-Engineering			
6206	Drop in accommodation bookings resulted from primary switchover for payment component	Co-owners		Primary Reviewer			
		Foo Bar 2		Foo Reviewer			
6196	Drop in accommodation resulted from puppet deployment	Primary Component pay-auth / mysql_pay_auth					
6181	Accommodations Softness due to issues on AV6 database						
6105	IAM primary database was not accepting any writes affecting the new login sessions						
6047	Badmin code bug led to MySQL pool drop in production while testing in DQS						
6006	Proxysql prevented the initiation of experiments.	PagerDuty Incidents					
5958	mktsvadb: MySQL GR primary failure caused write unavailability	Start time		Title			Priority

Follow ups

Created At	Title	Priority	Component	Owner	Completed
	Include Primary Write availability chart into MySQL - Master dashboards	High	mysql-sli / mysql-availability-reporter	Foo Bar	YES
	Improve documentation for running dba mysql primary_failover command	High	badmin / badmin-agent	Foo Bar	YES
	Enhance "badmin cname patch" to verify new Account relation of DNS record and to forbid unintended changes of such	High	dbe / rfo	Foo Bar	YES
	Update db_type and status tags in Orchestrator UI of the primary after a switchover / failover sooner	Low	orchestrator / orchestratorapp	Foo Bar	NO

on-chains.dc_master_replication_delay; started [DBA-06]	P2
ed with book rate issues due to db issues	P1
je to lack of auth plugin	P1

Some of our RFOs Follows-up

Improved the way `my_plugin_log_message()` handles `plugin_ptr` argument
#520

 Closed malderete wants to merge 1 commit into `mysql:trunk` from `malderete:my-plugin-log-message-handle-nullptr-nocrash` 

Unified behaviour when calling "`plugin->deinit`" for all plugin's types #521

 Closed malderete wants to merge 1 commit into `mysql:trunk` from `malderete:no-nullptr-for-special-plugin-types` 



 Conversation 3  Commits 1  Checks 0  Files changed 3



malderete commented on Feb 19 · edited

Handle None diff values when collecting diagnostics #20

 Closed ioandr wants to merge 1 commit into `mysql:master` from `ioandr:py-collect-diagonstics-none-diff` 

 Conversation 11  Commits 1  Checks 0  Files changed 1

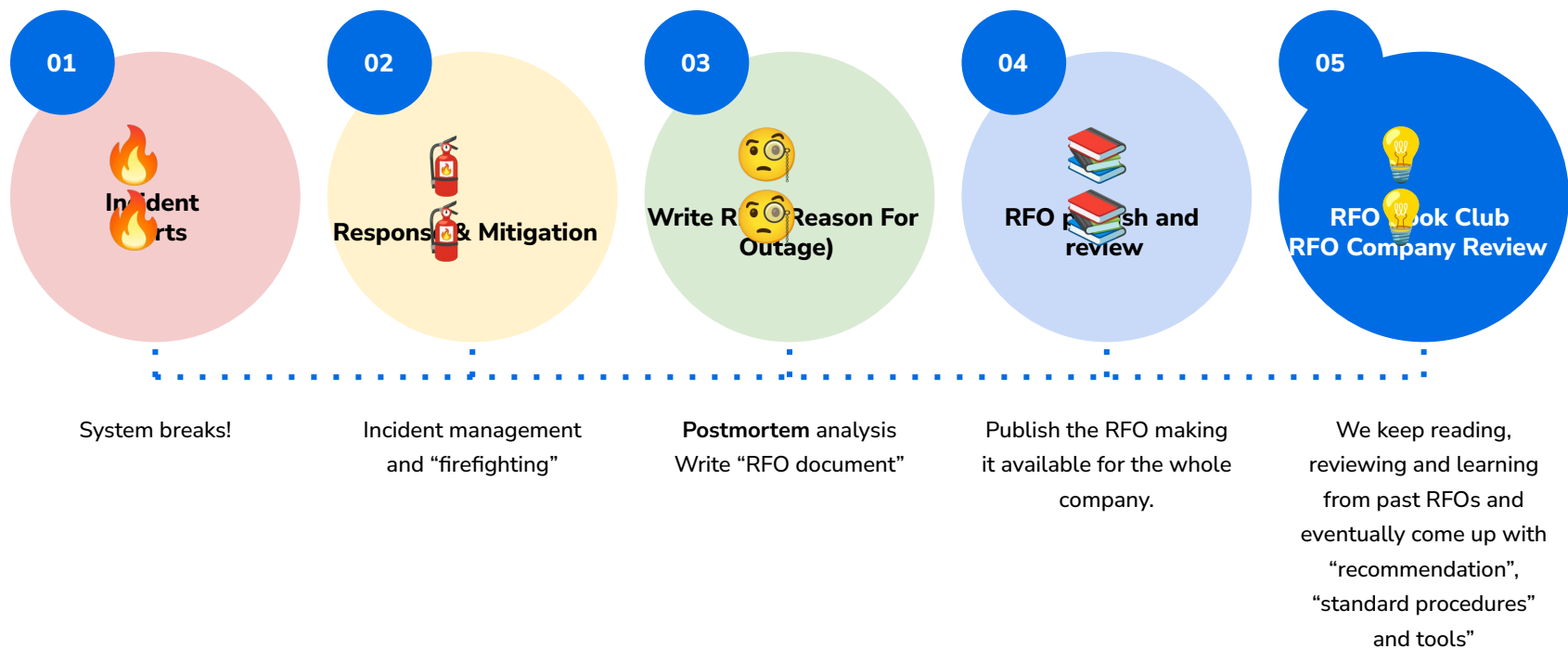


ioandr commented on Apr 19 · edited

Contributions



It isn't just a DOCUMENT but a CULTURE!





Questions ?

Thank you!

martin.alderete@booking.com

ioannis.androulidakis@booking.com

<https://jobs.booking.com/booking/jobs>

Booking .com



Other presentations

- [MySQL Keynote: Oracle OpenWorld 2018](#)
- [Booking.com: Reliable Operations and Rapid Development with MySQL](#)
- [Running MySQL on Three Different Platforms at Booking.com](#)
- and more...

