# Database reliability engineering for MySQL

**Matthias Crauwels**
**ConFoo 2022 - Online**
**Wed Feb 23rd 2022**
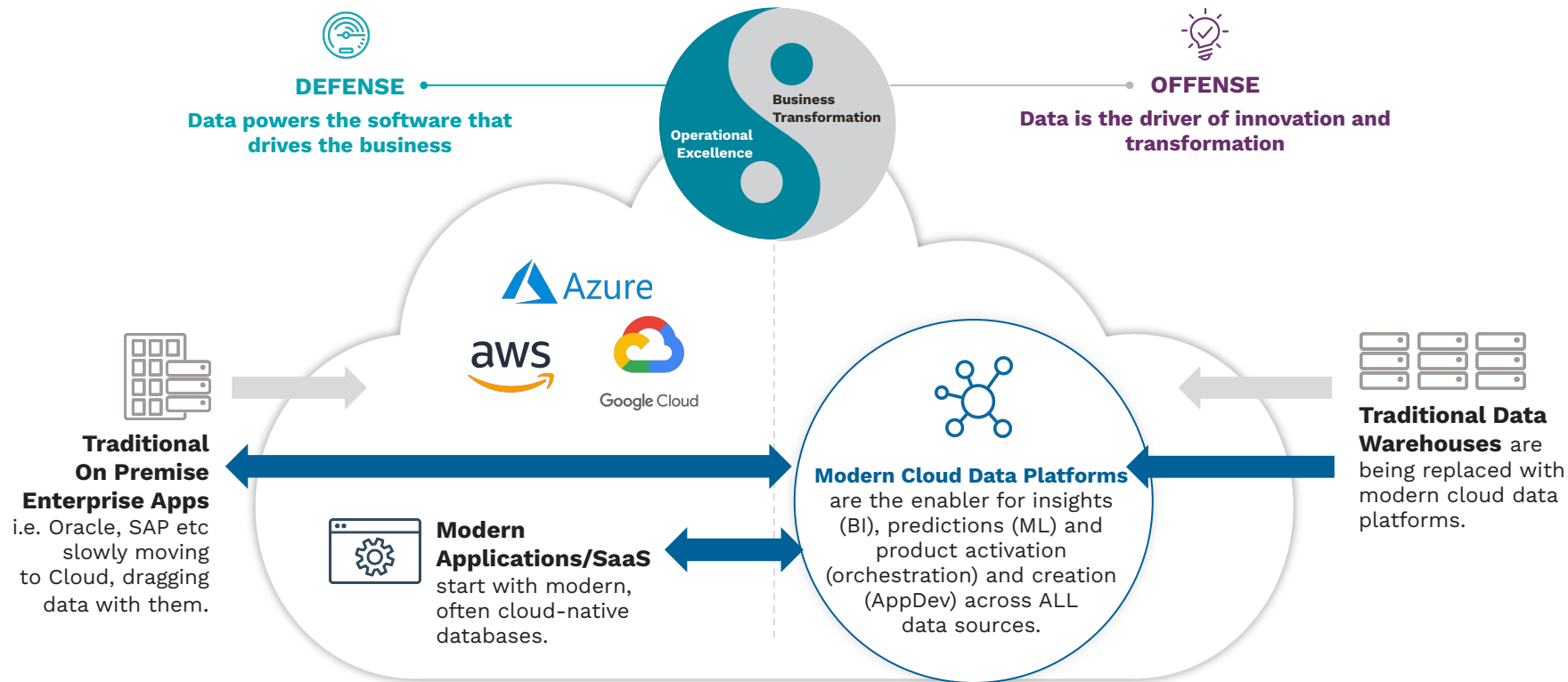
Pythian

# Speaker
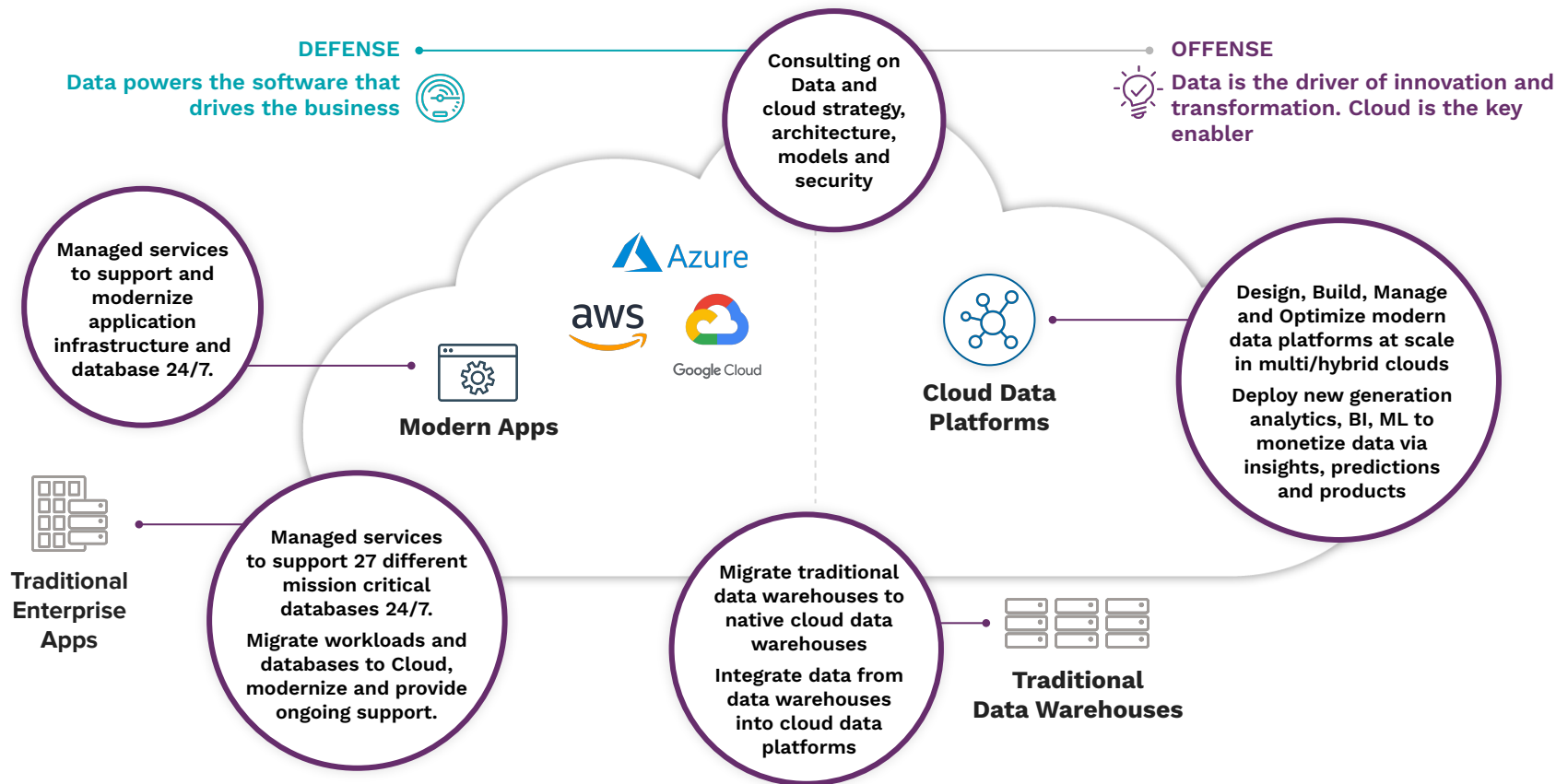


**Matthias Crauwels**
Principal Consultant
Pythian - OSDB

# How the data estate is evolving



**DEFENSE**
Data powers the software that drives the business

Business Transformation

Operational Excellence

**OFFENSE**
Data is the driver of innovation and transformation

Azure

aws

Google Cloud

**Traditional On Premise Enterprise Apps**
i.e. Oracle, SAP etc slowly moving to Cloud, dragging data with them.

**Modern Applications/SaaS**
start with modern, often cloud-native databases.

**Modern Cloud Data Platforms**
are the enabler for insights (BI), predictions (ML) and product activation (orchestration) and creation (AppDev) across ALL data sources.

**Traditional Data Warehouses** are being replaced with modern cloud data platforms.

# Pythian's Services Across the Data Estate

**DEFENSE**
Data powers the software that drives the business

**OFFENSE**
Data is the driver of innovation and transformation. Cloud is the key enabler

Consulting on Data and cloud strategy, architecture, models and security

Managed services to support and modernize application infrastructure and database 24/7.

**Modern Apps**

Azure

aws

Google Cloud

**Cloud Data Platforms**

Design, Build, Manage and Optimize modern data platforms at scale in multi/hybrid clouds

Deploy new generation analytics, BI, ML to monetize data via insights, predictions and products

**Traditional Enterprise Apps**

Managed services to support 27 different mission critical databases 24/7.

Migrate workloads and databases to Cloud, modernize and provide ongoing support.

Migrate traditional data warehouses to native cloud data warehouses

Integrate data from data warehouses into cloud data platforms

**Traditional Data Warehouses**

# Other presentation

- Tomorrow at 10:30 I have another presentation

  **Getting started with InnoDB Cluster in MySQL 8**

  *Since MySQL 5.7 InnoDB Cluster saw the light. Back then it was not as popular because it was considered too new technology. As with any new product, there were numerous bug reports. With MySQL 8.0 the solution has matured a lot, many of the problems have been addressed and so it has became a full scale high availability solution for MySQL.*

# AGENDA

- (High) Availability

- Service Discovery

- Observability

- Disaster Recovery

# (High) Availability

# Availability for MySQL

- **When do we consider MySQL available?**
- **We need:**
  - **Just one server to write transactions to**
  - **At least one server to read transactions from**

# Availability vs High Availability

# Availability vs High Availability

# High Availability for MySQL

# MySQL - Single point of failure

- The "writer" is the single-point-of-failure in this topology
- No more writes can happen
- Application will likely be down or at least degraded
- We shall introduce a topology management tool to handle these failures.

# Orchestrator

Orchestrator is a High Availability and replication management tool.

It can be used for:

- Discovery of a topology
- Visualisation of a topology
- Refactoring of a topology
- Recovery of a topology

# Orchestrator: Discovery

Orchestrator can (and will) discover your entire replication topology as soon as you connect it to a single server in the topology.

It will use regular DBA commands such as: `SHOW SLAVE HOSTS`, `SHOW PROCESSLIST`, `SHOW SLAVE STATUS` to try and connect to the other servers in the topology.

Requirement: the `orchestrator_topology_user` needs to be created on every server in the cluster so it can connect.

# Orchestrator: Visualization

Orchestrator comes with a web interface that visualizes the servers in the topology.

# Orchestrator: Refactoring

Orchestrator can be used to refactor the topology.

This can be done from the **command line tool**, via the **API** or even via the web interface by **dragging and dropping**.

You can do things like

- Repoint a replica to a new master
- Promote a server to a (co-)master
- Start / Stop replica
- …

# Orchestrator: Recovery

All of these features are nice, but they still require a human to execute them. This doesn't help you much when your master goes down at 3AM and you get paged to resolve this.

Orchestrator can be configured to automatically recover your topology from an outage.

# Orchestrator: How recovery works?

To be able to perform a recovery, Orchestrator first needs to detect a failure.

Typical monitoring tools (think nagios) will probe the master, but what to do on failure? Take immediate action? Retry? How many retries?

As indicated before Orchestrator connects to every server in the topology and gathers information from each of the instances.

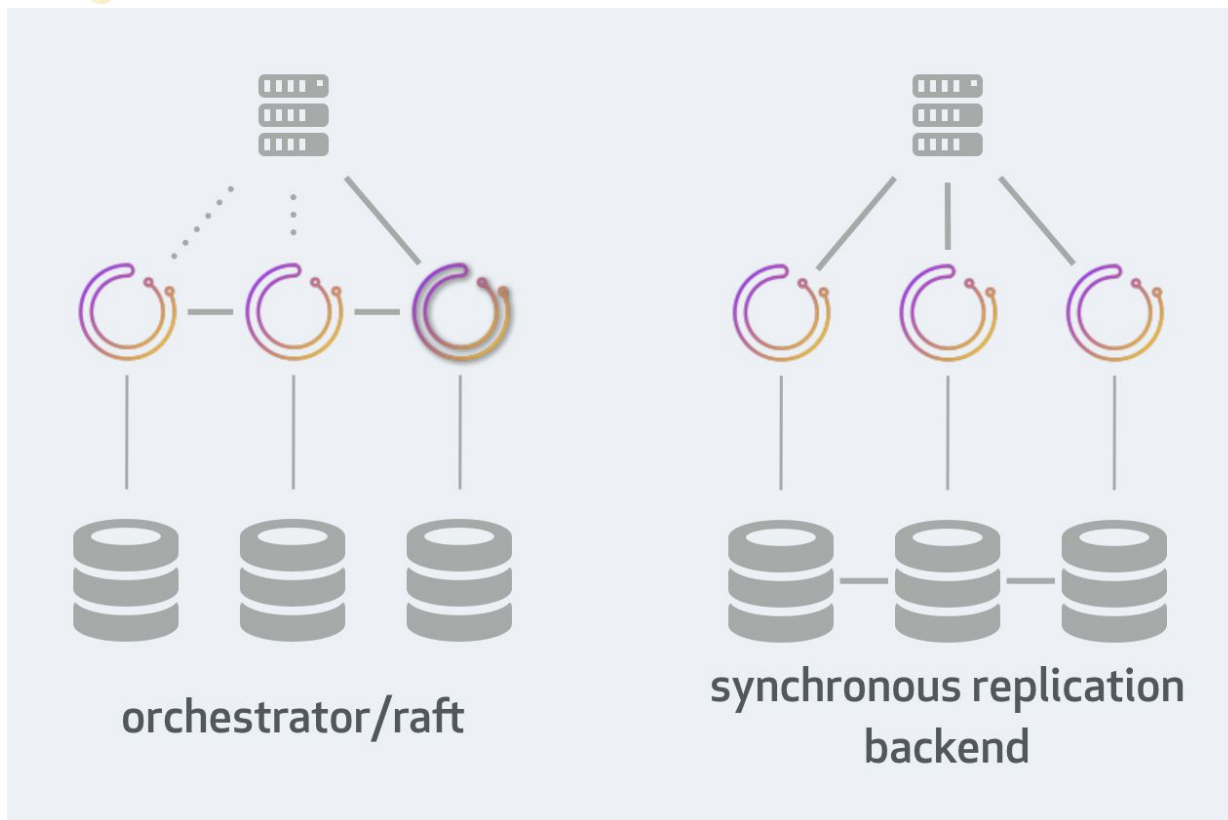Orchestrator uses this information to make decisions on the best action to take. They call this the holistic approach.

# MySQL HA – back to our example

# Orchestrator/raft

- Since Orchestrator 3.x
- Orchestrator backend db becomes standalone
  - requirement for MySQL as backend db was dropped (but still available)
  - introduced sqlite as possible backend db
- Implements RAFT consensus protocol, responsible for
  - leader election
  - distribution of data
- Minimal 3 nodes for a HA setup (50% + 1 node quorum)

# Orchestrator shared db vs raft



orchestrator/raft

synchronous replication backend

# Service Discovery

# What is service discovery?

- Know where the writer server is
- Know where the reader server is
- Several options
  - Application deploy?
  - DNS change?
  - Floating Virtual IP?
  - Proxy-server?

# ProxySQL: What?

ProxySQL is a high performance layer 7 proxy application for MySQL.

- It provides 'intelligent' load balancing of application requests onto multiple databases
- It understands the MySQL traffic that passes through it, and can split reads from writes.
- It understands the underlying database topology, whether the instances are up or down
- It shields applications from the complexity of the underlying database topology, as well as any changes to it
- …

# ProxySQL: Terminology

- Hostgroup

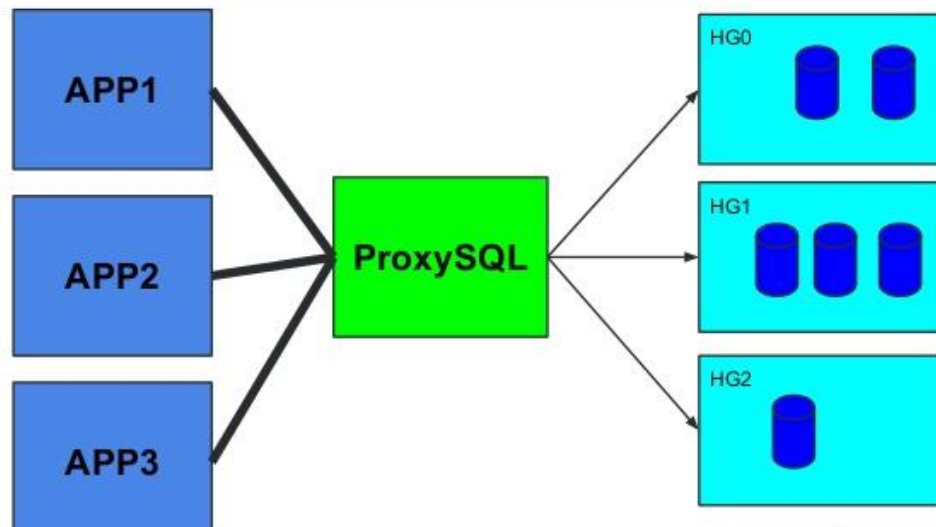  All backend MySQL servers are grouped into hostgroups. These "hostgroups" will be used for query routing.

- Query rules

  Query rules are used for routing, mirroring, rewriting or blocking queries. They are at the heart of ProxySQL's functionalities
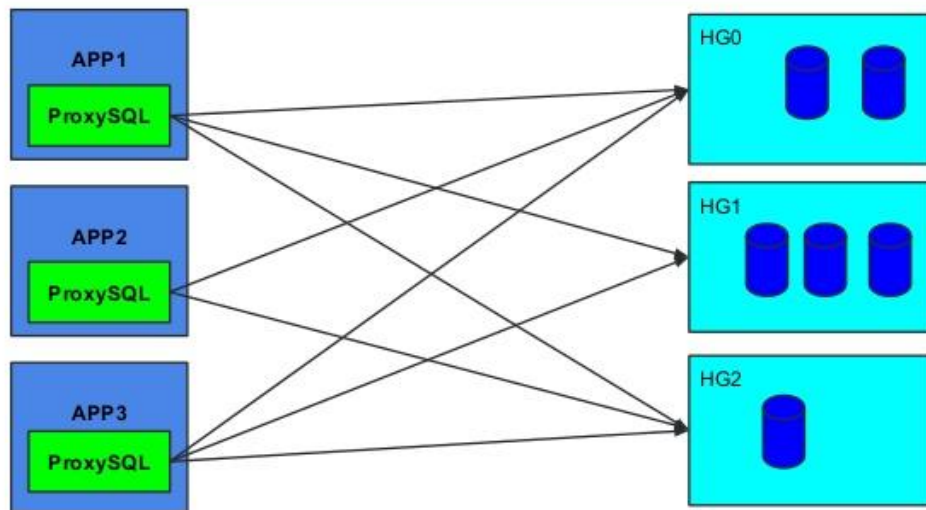
- MySQL users and servers

  These are configuration items which the proxy uses to operate

# ProxySQL: Basic design (1)

# ProxySQL: Basic design (2)

# Observability

# What is observability

- Know if your systems have an issue (alerting)
  - Pager
  - IM
  - Email
- Troubleshoot performance and/or diagnose issue (graphs)
  - dashboards
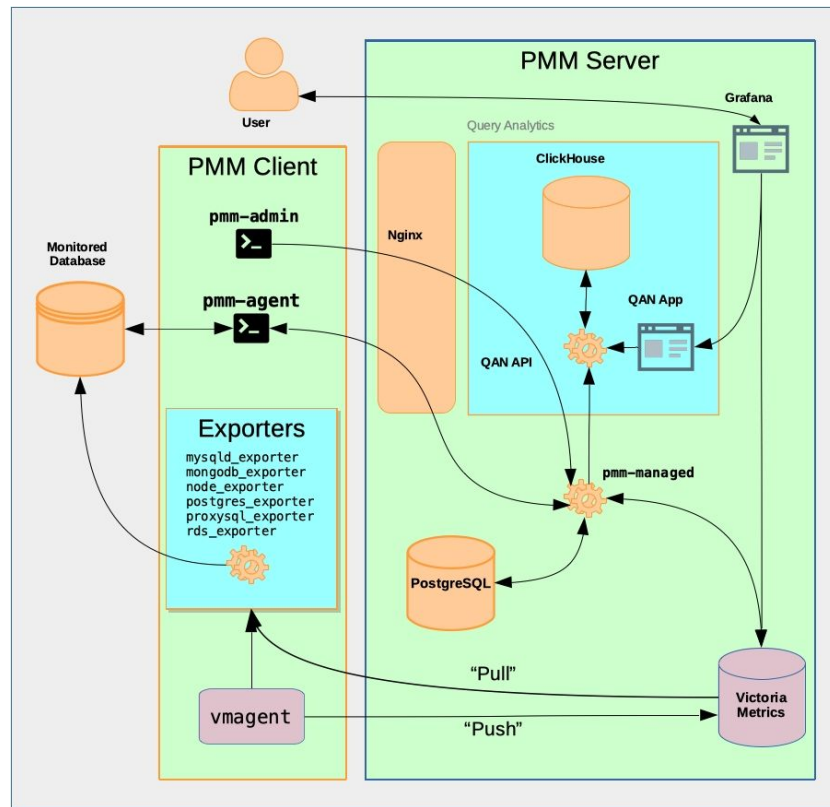  - comparison to a past situation
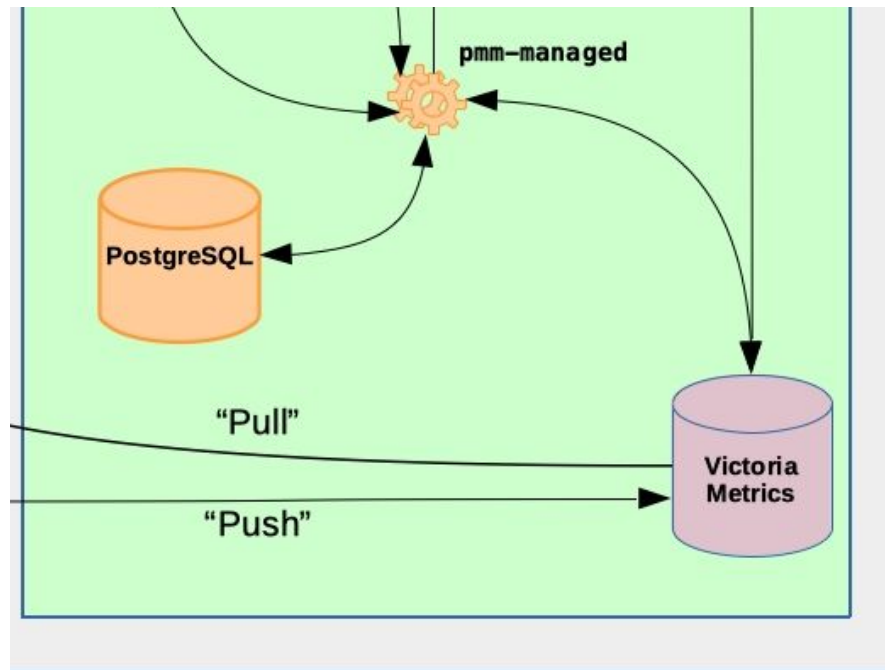
# Percona Monitoring and Management (PMM)

# What is PMM?

- All-in-one solution for Observability, focussed on databases
- Free and open source
- Out-of-the-box solutions for
  - MySQL
  - PostgreSQL
  - MongoDB
  - ProxySQL
- Available as
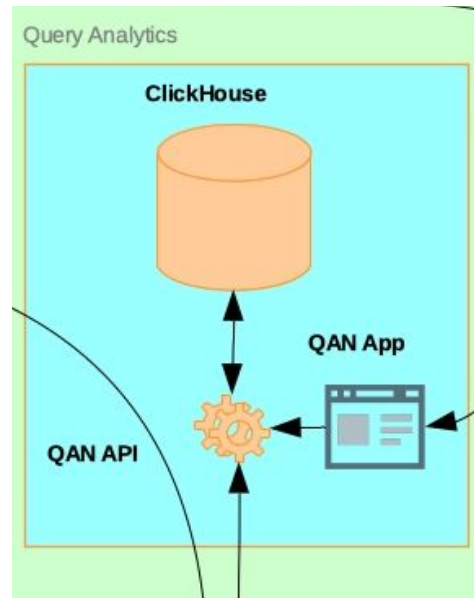  - docker container
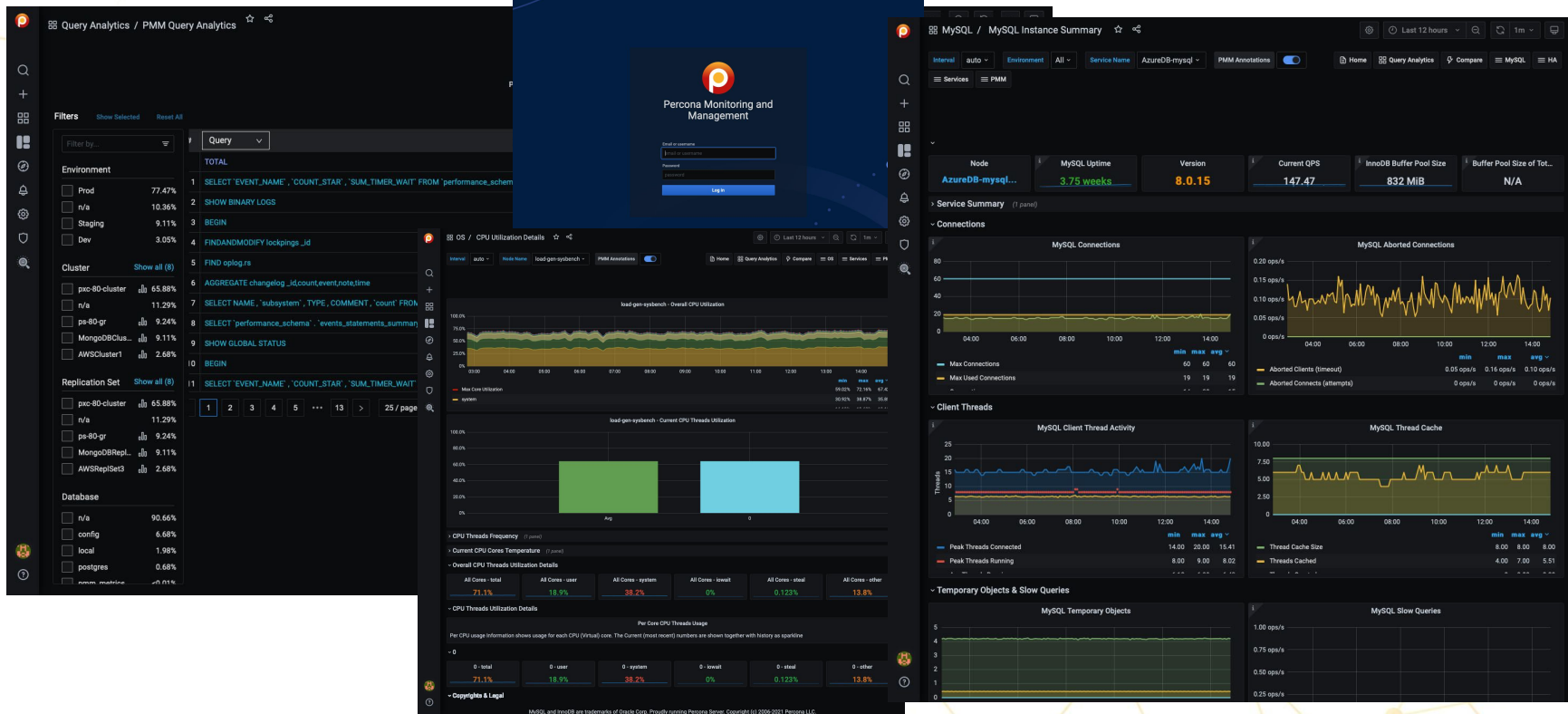  - OVF (VM image)
  - AMI (AWS)

# PMM: Components

# PMM: Components
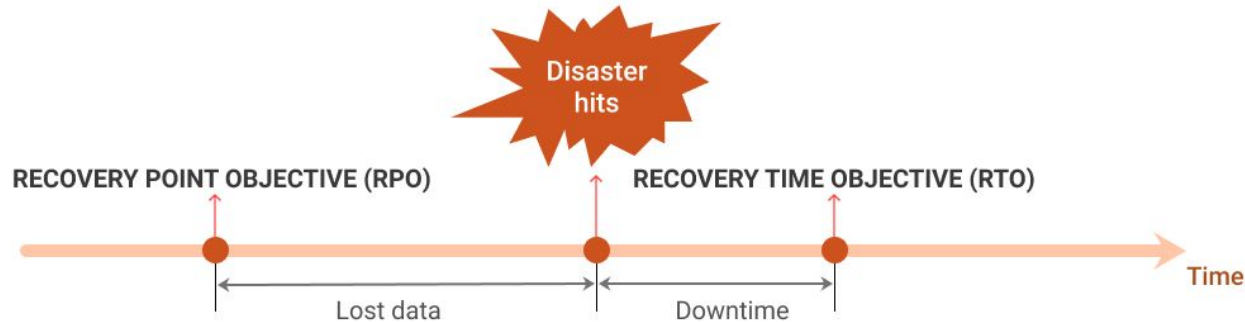
# PMM: Components

# PMM

# Disaster Recovery

# What is Disaster Recovery (DR)

- The ability to recover from a catastrophic event or a human error
- RTO vs RPO

## RPO and RTO explained

**Disaster hits**

**RECOVERY POINT OBJECTIVE (RPO)**

**RECOVERY TIME OBJECTIVE (RTO)**

Time

Lost data    Downtime

# Backups

- Fastest RTO is usually achieved with disk snapshots
- How to make consistent snapshot?
  - either stop MySQL on a replica cleanly (`innodb_fast_shutdown = 0`)
  - run a `FLUSH TABLES WITH READ LOCK` + freeze the filesystem
- Both options block normal operations, so they should be ran on a designated replica.

# Backups - Percona Xtrabackup

- supports online, non blocking backups for InnoDB
- full backups
- incremental backups
- physical backup
- make sure to match the Xtrabackup version that corresponds with your MySQL version.

# Point-in-time recovery

- use the MySQL binary logs to recover your environment to a specific point-in-time
- if done right this can reduce your RPO to 0
- stream your binary logs to a backup (DR) location
- doing PITR will increase your RTO

# Conclusion

# Conclusion

- MySQL replication is required for High Availability
- Orchestrator can automate your writer-availability
- ProxySQL will help with Service Discovery
- PMM will expose the environment metrics for observability
- Backups are important for Disaster Recovery

# Questions?

# Thank you!

crauwels@pythian.com
@mcrauwel