# Organize the migration of a hundred databases to the cloud

Percona Live ONLINE
May 12th 2021

"Hello

Meet
# Maxime Fouilleul
Engineering Manager for DBRE
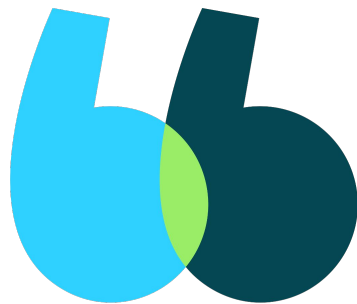
"Make the database not a problem."

## 🎨 BUILD

> Package and support the database catalog for BlaBlaCar application services.

## 🤝 ADVISE

> Provide expertise in software engineering to help teams choose the right database for them and to ensure they use it the right way.

# The go-to marketplace for shared road travel

BlaBlaCar is a community-based marketplace allowing members to book seats in individual cars and buses alike. From carpool to buses, we have one common moto: #ZeroEmptySeats.

| | Carpool | Bus Marketplace | BBC Branded Buses |
|---|---|---|---|
| **Geographies** | Global (22 countries) | Russia, Ukraine, Poland, Brazil | France, Germany |
| **Position** | Leader in all our markets | Leader in Eastern Europe Early stage in Brazil | Leader in France |

# BlaBlaCar

## "The go-to marketplace for shared travel"

**90 million** members

**25 million** travelers
per quarter

**22** countries

BlaBlaCar

BlaBlaCar Bus

BlaBlaCar Daily

# 90 million members

Google Cloud Platform

# Our production database infra in 2019

**MySQL**
30 Production Clusters

**Elasticsearch**
6 Production Clusters

**Cassandra**
7 Production Clusters

**PostgreSQL**
5 Production Clusters

**RabbitMQ**
13 Production Clusters

**Kafka**
8 Production Clusters

**Redis**
19 Production Clusters

**Couchbase**
5 Production Clusters

# The mission

**2019**
*sign cloud provider*

**end of 2020**
*close on-premise*

## Consolidate the DBRE team
*staffing plan includes 4 SRE database enthusiastic*

## Migrate 100+ databases
*Package reliable systems, accompany the migration and decommissioning*

# The dream team

Database Reliability Engineering (DBRE)

**Engineering Manager**

**Product Owner**

**SRE**
Distributed Databases

**SRE**
Cloud & Kubernetes

**SRE**
Kafka

**SRE**
Databases

"Fly me to the cloud"

# The DBRE vision

Google Compute
Engine

Google Kubernetes
Engine

Google Managed
Services

GCP Marketplace

# The DBRE vision

**Google Compute Engine**

Try to avoid ✋

**Google Kubernetes Engine**

**Google Managed Services**

**GCP Marketplace**

# The DBRE vision



**Google Compute Engine**

Try to avoid ✋

**Google Kubernetes Engine**

Prefer 👌

**Google Managed Services**

**GCP Marketplace**

# The DBRE vision

**Google Compute Engine**

Try to avoid ✋

**Google Kubernetes Engine**

Prefer 👌

**Google Managed Services**

Do 🤌

**GCP Marketplace**

# The DBRE vision

**Google Compute Engine**

Try to avoid ✋

**Google Kubernetes Engine**

Prefer 👌

**Google Managed Services**

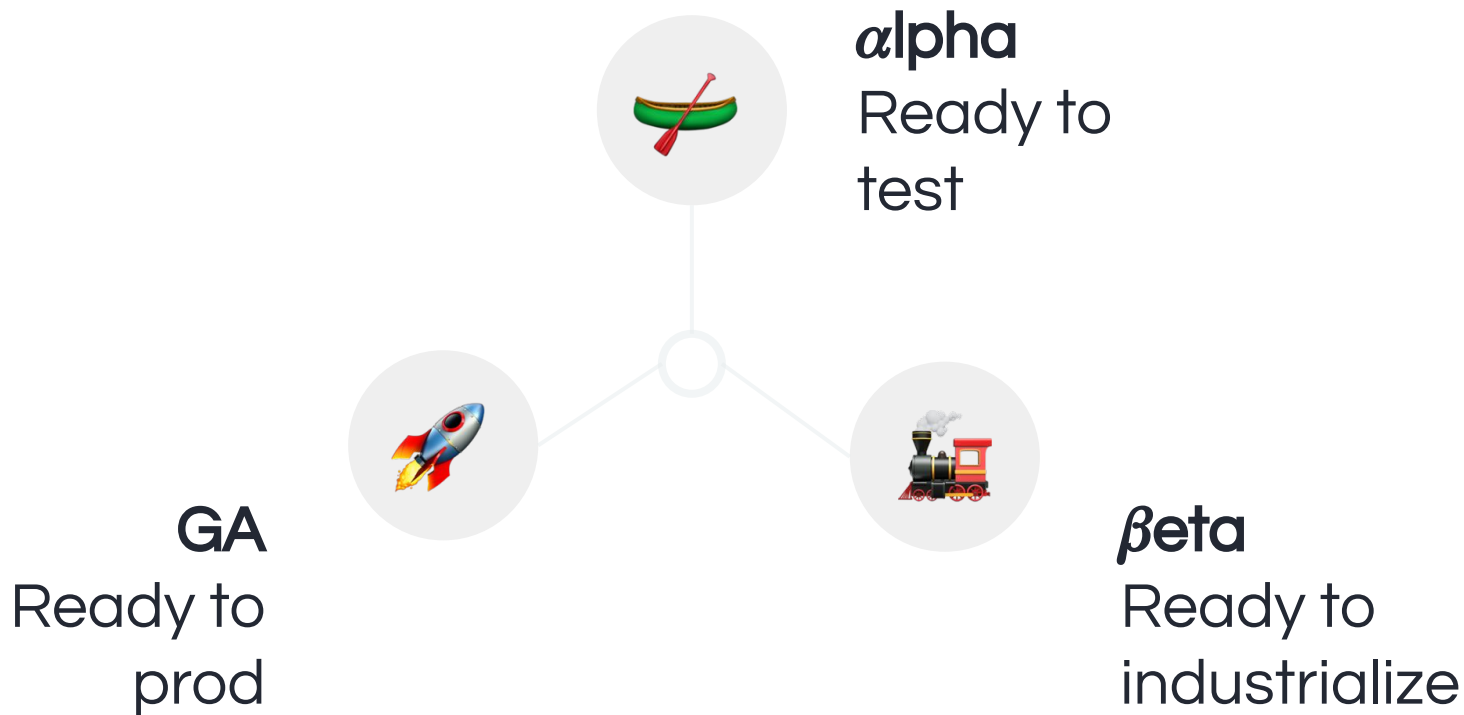Do 🤏

**GCP Marketplace**

Don't ✋

"Be transparent to ensure buy-in"

# Be clear on iterations

🛶 **𝛼lpha**
Ready to test

🚀 **GA**
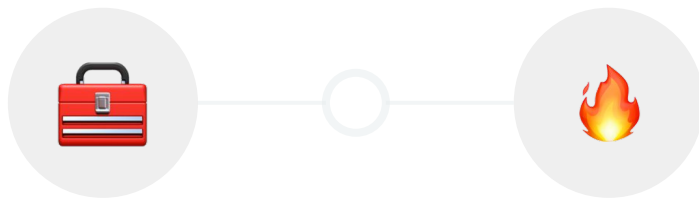Ready to prod

🚂 **𝛽eta**
Ready to industrialize

# Documentation as log

- 2020-07-20 - MariaDB - SLI/SLO - Specs and first implementation
- 2020-01-29 - MariaDB - What version to use in GCP?
- 2020-01-29 - MariaDB - Rely on disk snapshots for backups
- 2020-01-29 - MariaDB - using MaxScale as a layer 7 proxy
- 2020-01-29 - CloudSQL - New HA design supported
- 2019-12-06 - CloudSQL user management
- 2019-12-05 - CloudSQL limitations - Can't purge binary logs, let it grow, let it grow!
- 2019-12-03 - MariaDB - New features for the chart would be documented in release notes
- 2019-11-13 - MariaDB - Reboot the packaging of MariaDB in GKE
- 2019-10-25 - MariaDB - Using deported Prometheus exporter to monitor CloudSQL...
- 2019-10-03 - CloudSQL limitations - Does the lack of triggers can impact us
- 2019-10-02 - CloudSQL limitations - Replicate from MariaDB to CloudSQL is not possible.
- 2019-10-02 - CloudSQL network optimizations (bye-bye CloudSQL Proxy)
- 2019-10-01 - CloudSQL provisioning, why a terraform module?
- 2019-09-06 - CloudSQL everywhere?
- 2019-07-31 - MariaDB - Adding safe_to_bootstrap override capability to avoid getting stuck after full crash
- 2019-04-04 - MariaDB - Graceful restart - Kubernetes Probes + PDB
- 2019-01-19 - Make or Buy Study, why not CloudSQL?
- 2018-10-03 - MariaDB - Docker image + Helm chart (Alpha = stable release for testing)
- 2018-09-28 - MariaDB - Performance Benchmarks, BBC Baremetal VS GCP

# Runbooks

How do I? 🧰 ⭕ 🔥 What to do when?

# Gamify the knowledge sharing process

## Level 1
Actions are basic tasks that should be mastered by each team member.

| Level | Action codename |
|---|---|
| 🏅 Level 1 | BasicHealthcheck |
| 🏅 Level 1 | RolloutMinorChanges |
| 🏅 Level 1 | Connect&Read |
| 🏅 Level 1 | ManageAccess |
| 🏅 Level 1 | PrepareClusterBoostrap |
| 🎖️ Owner | ActiveOwnership |
| 🎖️ Owner | ModifyDataset |
| 🎖️ Owner | RecoverDataloss |
| 🎖️ Owner | Advisory |

## Level "Owner"
Actions allow the component to be actively supported, they should be mastered by at least 2 members.

# Implementations

and migration paths

Timeline: 2018 — Oct — 2019 — Apr — Jul — Oct — 2020 — Apr — Jul — Oct — 2021 — 2021

- MariaDB — Oct
- Cassandra — Apr
- Elasticsearch — Feb
- PostgreSQL — May
- Kafka — Dec
- Redis — Jun
- RabbitMQ — Oct
- Elasticsearch — Oct
- RabbitMQ — Nov
- MariaDB — May
- Kafka — Mar
- Cassandra — Jul
- CloudSQL — Jun
- MemoryStore — Aug
- Cassandra — Jan
- RabbitMQ — Jan
- Redis — May
- Kafka — Oct
- Elasticsearch — Feb
- PostgreSQL — Sept
- MariaDB — Jan

Legend:
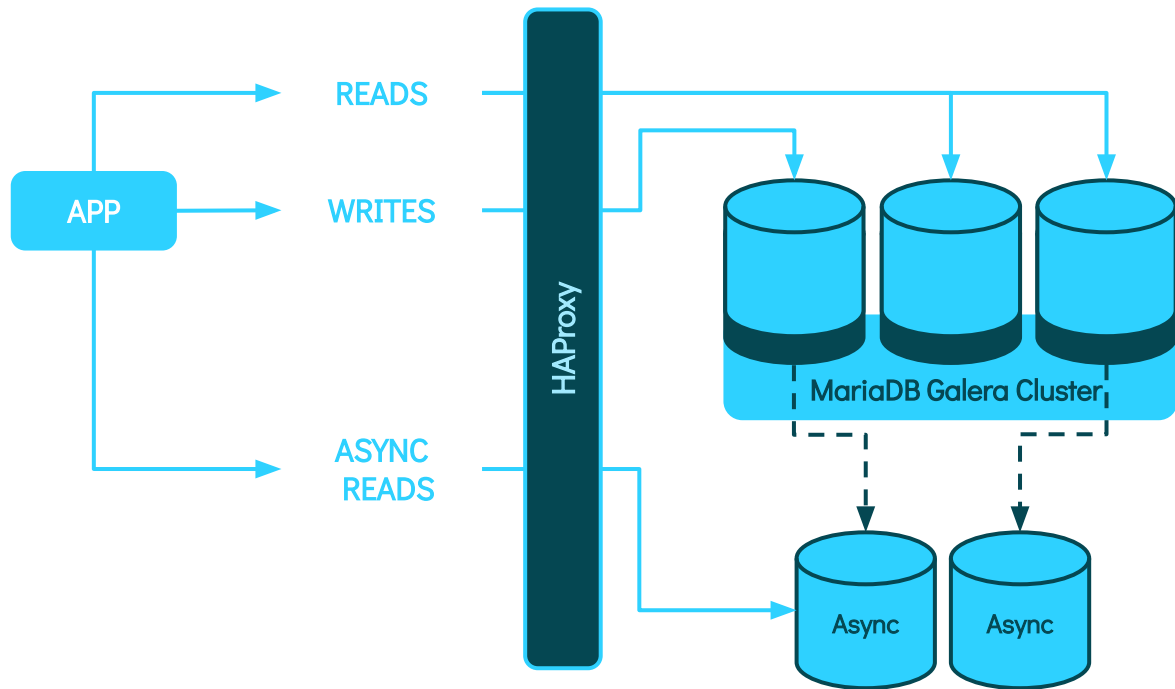- Kickoff of the $\alpha$lpha release
- First production cluster in GCP
- Last production cluster decomissioned from on-premise

# Key elements for the implementations

- Leverage **Kubernetes Statefulset**

- master **affinity** in the PodSpec

- prefer **Google Persistent disks** over Local SSD

- leverage Persistent disk **Snapshots**

- promote **distributed ownership**

- use **Terraform** for Google managed services

# A production MySQL service in 2019

# Gather requierements

- Will your database be migrated or abandoned?

- What is the tolerated downtime for the migration?

- Can we migrate the reads separately from the writes?

**Google CloudSQL**

1. The application can tolerate several minutes of downtime
2. Writes can be stopped during the dataset migration
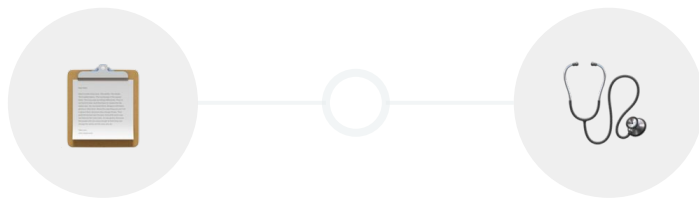3. The need is fairly lightweight

**MariaDB in Kubernetes**

# What DBRE is packaging for CloudSQL?

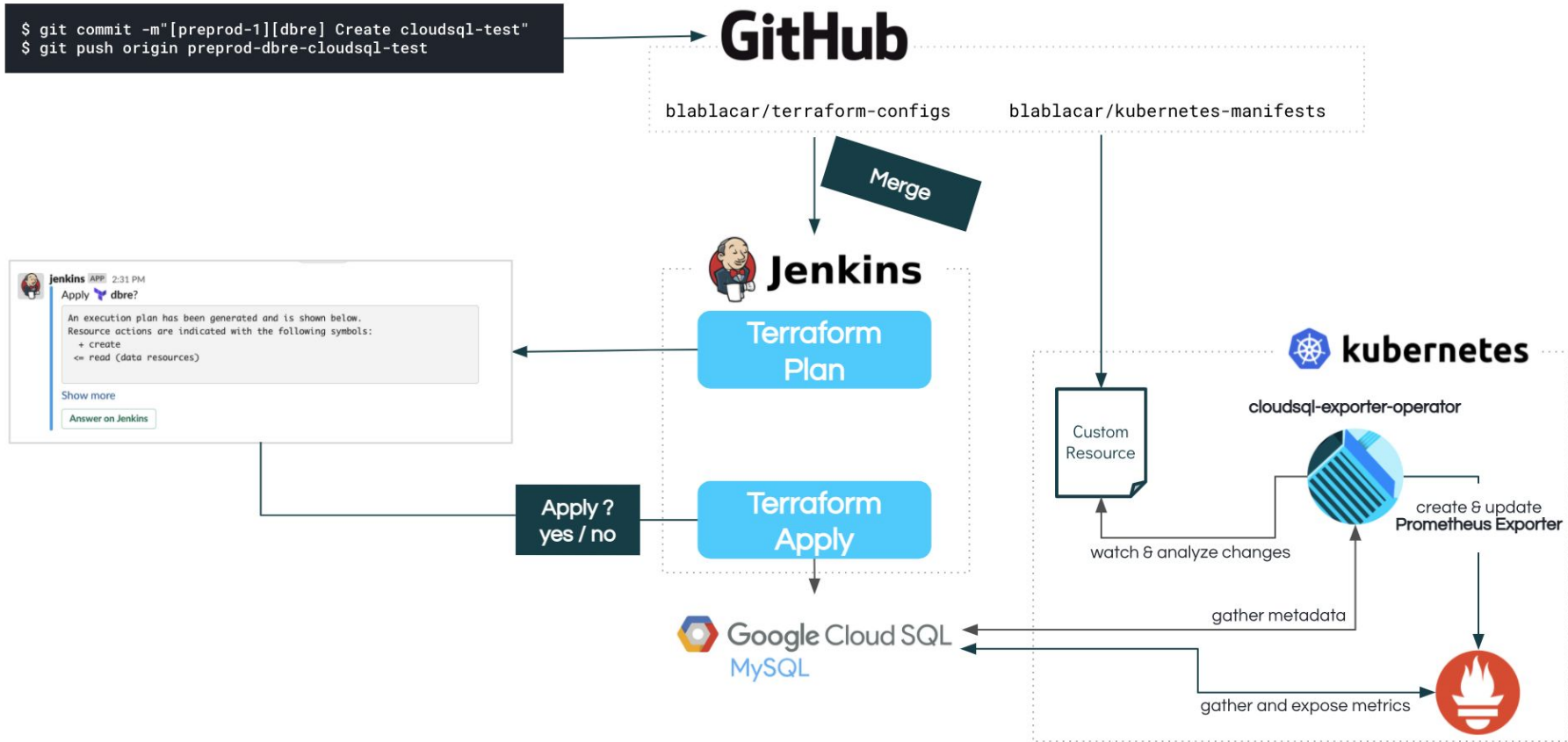**Terraform Module**
To ease and standardize the usages

**Kubernetes Operator**
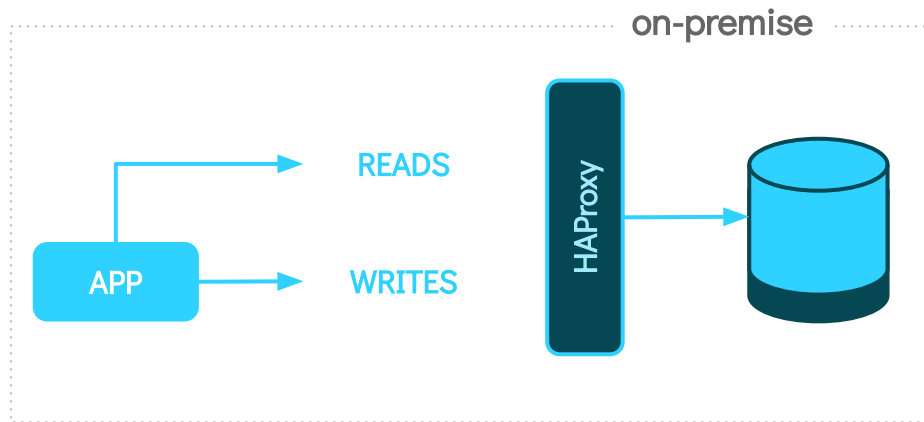To setup a Prometheus exporter

# What DBRE is packaging for CloudSQL?

```
$ git commit -m"[preprod-1][dbre] Create cloudsql-test"
$ git push origin preprod-dbre-cloudsql-test
```

## GitHub

blablacar/terraform-configs          blablacar/kubernetes-manifests

Merge

## Jenkins

jenkins APP 2:31 PM
Apply 🔱 dbre?

An execution plan has been generated and is shown below.
Resource actions are indicated with the following symbols:
  + create
  <= read (data resources)

Show more

Answer on Jenkins

Terraform Plan

Apply ?
yes / no

Terraform Apply

## kubernetes

Custom Resource

cloudsql-exporter-operator

watch & analyze changes

create & update
**Prometheus Exporter**

gather metadata

Google Cloud SQL
MySQL

gather and expose metrics
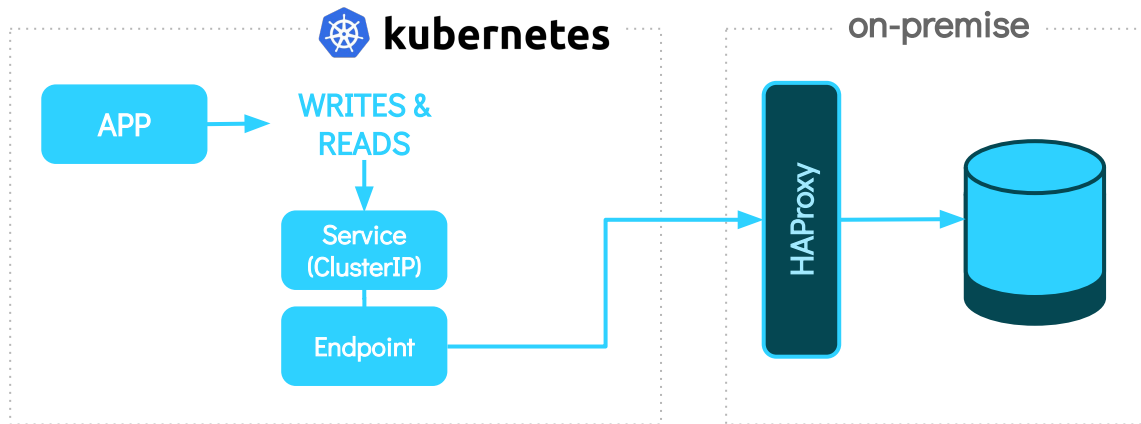
# CloudSQL Migration Path
## Initial stage

# CloudSQL Migration Path
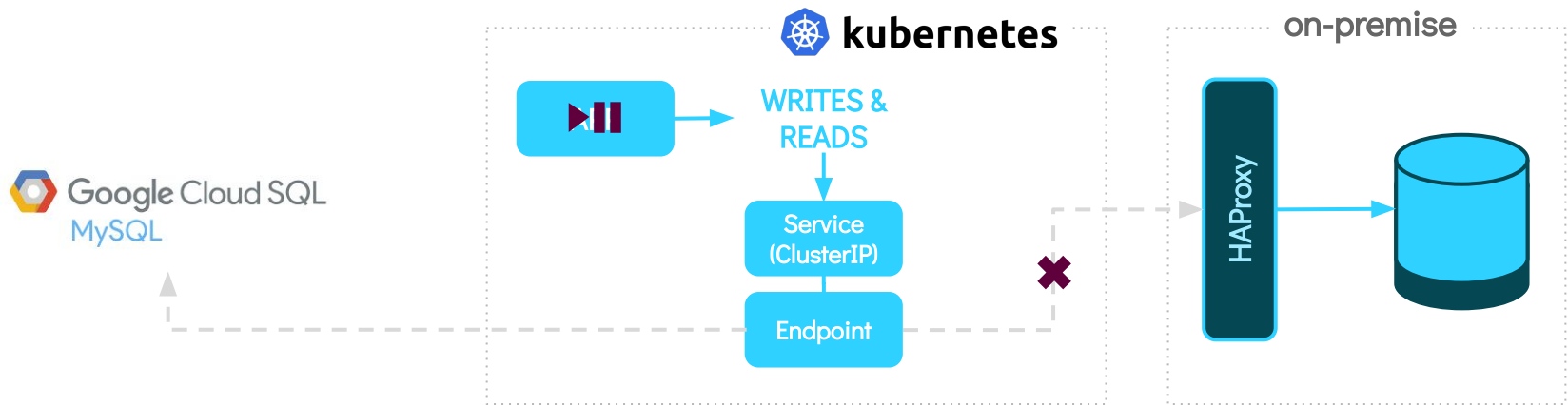## Move the application

```yaml
1    ---
2    apiVersion: v1
3    kind: Service
4    metadata:
5      name: cloudsql-demo
6      namespace: demo
7    spec:
8      type: ClusterIP
9      clusterIP: None
10     ports:
11       - protocol: TCP
12         port: 3306
13         targetPort: 3306
14         name: mysql
15   ---
16   kind: Endpoints
17   apiVersion: v1
18   metadata:
19     name: cloudsql-demo
20     namespace: demo
21   subsets:
22     - addresses:
23         - ip: <database-ip-address>
24       ports:
25         - port: 3306
```
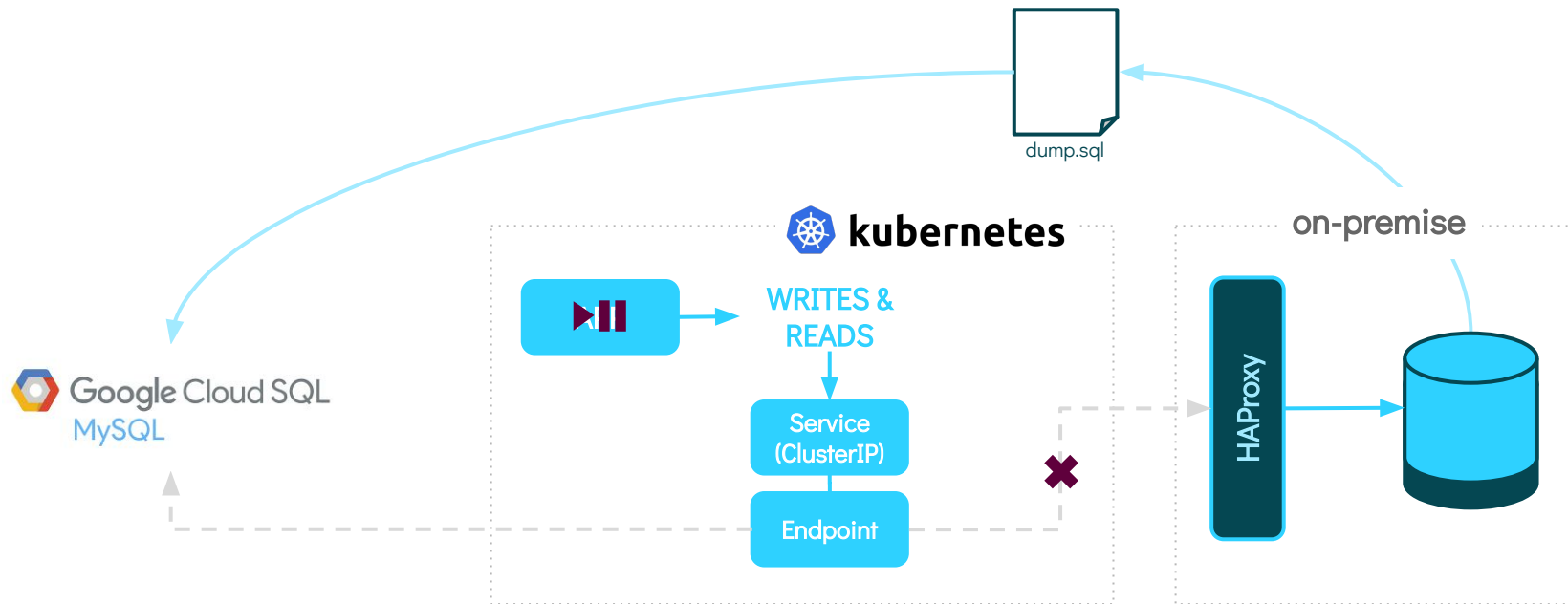


kubernetes

on-premise

APP → WRITES & READS → Service (ClusterIP) → Endpoint → HAProxy → (database)

# CloudSQL Migration Path
## Stop the application and switch endpoint

# CloudSQL Migration Path
## Copy dataset



dump.sql

kubernetes

WRITES & READS

Google Cloud SQL
MySQL

Service (ClusterIP)

Endpoint

on-premise

HAProxy
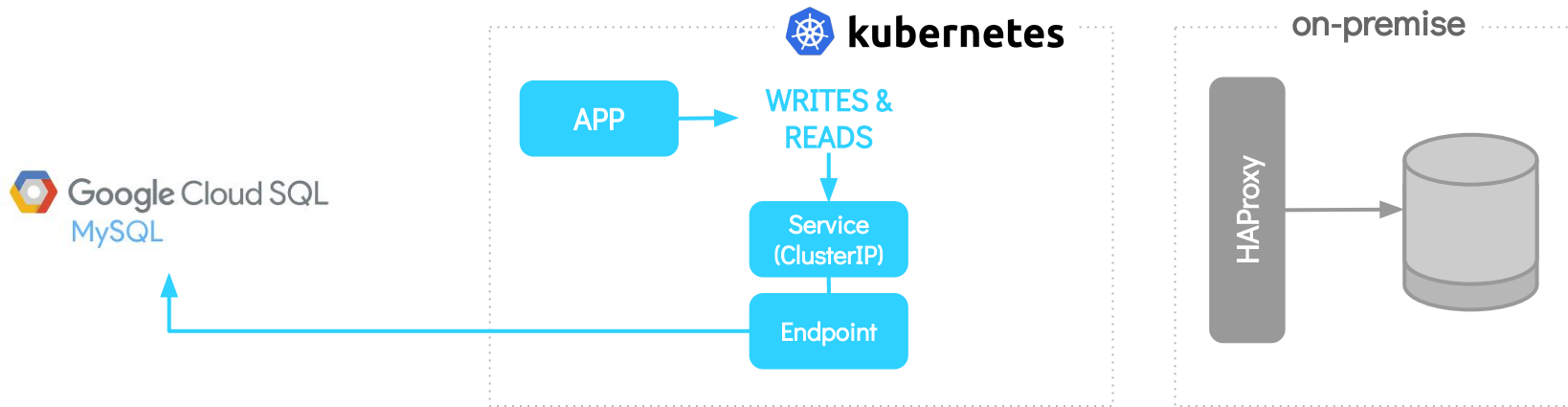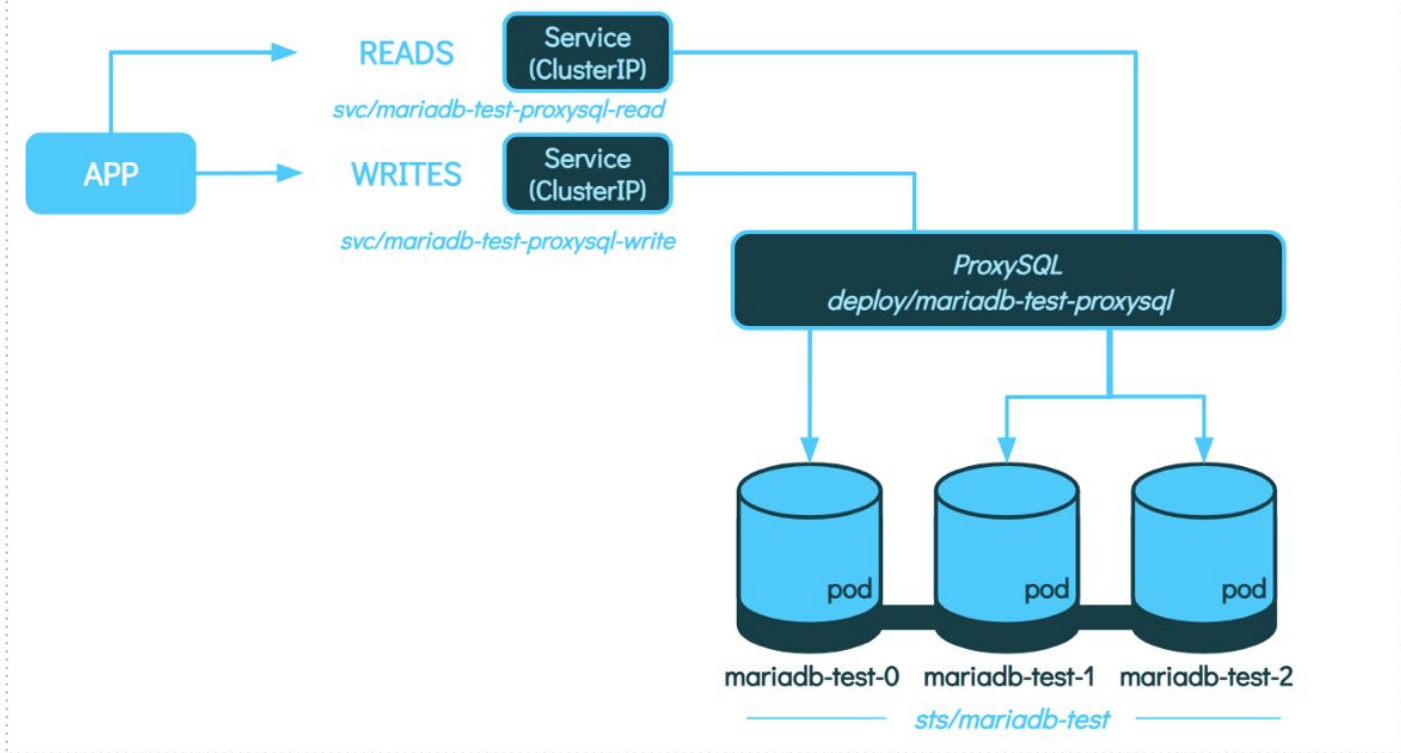
# CloudSQL Migration Path
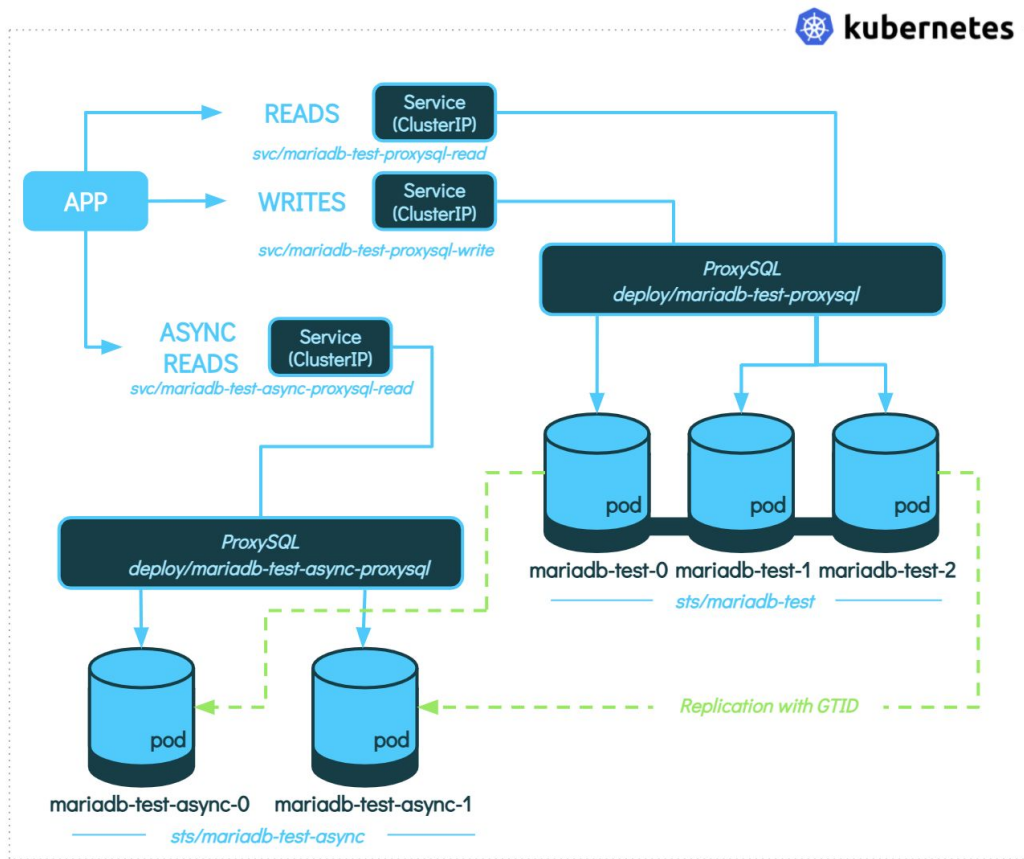## Restart the application

# MariaDB in-house packaging via a Helm Chart

- A **StatefulSet** with Galera enabled...or not

- A **Deployment** running **ProxySQL**

- **Prometheus** exporter sidecars to export metrics

- A bunch of **Jobs** that manipulate disk snapshots

- A **Deployment** running an **SLI Prober**

- Services, RBAC, and PDB...

# A production MySQL service in 2021

# With asynchronous replicas

# MariaDB packaging tips

# Dynamically find Galera seeds
## wsrep_cluster_address

```
init_galera_config.sh: |-
  #!/bin/bash
  set -ex

  # 1. Get seeds
  API_ENDPOINT="https://kubernetes.default.svc.cluster.local/api/v1/namespaces/{{ .Release.Namespace }}"
  TOKEN=$(cat /var/run/secrets/kubernetes.io/serviceaccount/token)

  EP_JSON=$(curl -sSk \
      -X GET \
      -H "Authorization: Bearer ${TOKEN}" \
      ${API_ENDPOINT}/endpoints/{{ .Release.Name }}-headless \
  )
  if [ "$(echo $EP_JSON | jq -r .kind)" == "Endpoints" ]; then
    if [ "$(echo $EP_JSON | jq -r .subsets)" != "null" ]; then
      # Endpoint = joining cluster
      SEEDS="$(echo $EP_JSON | jq -r '.subsets[0].addresses[].ip' | paste -sd, -)"
    else
      # No endpoint = create a cluster
      SEEDS=""
    fi
  else
    exit 1
  fi
```

We use a Kubernetes **Headless Service** to get available (ready) endpoints

If we find endpoints we join a cluster

If we don't find any endpoint we bootstrap a cluster

# Having accurate **Liveness** and **Readiness**

Simple ping to report the MySQL is live or not

Prevent killing a node doing an SST (Galera full resync)

```
liveness_probe.sh: |-
  #!/bin/bash

  # If mysql ping or SST in progress
  # Use TCP instead of Unix socket to be usable from side cars.
  mysqladmin -h 127.0.0.1 -u monitoring ping 2> /dev/null || [ -d {{ .Values.config.mysqld.datadir }}/.sst ]
readiness_probe.sh: |-
  #!/bin/bash
  # wsrep_local_state vs wsrep_local_state_comment
  # 1 = Joining
  # 2 = Donor/Desynced # Ready as we use non-blocking SST (xtrabackup/mariabackup).
  # 3 = Joined
  # 4 = Synced
  wsrep_local_state=$(mysql -u monitoring -BN -e "SHOW GLOBAL STATUS LIKE 'wsrep_local_state'" | awk '{ print $2 }')
  if [ -z $wsrep_local_state ]; then
    exit 1
  fi
  if [ $wsrep_local_state == 1 ] || [ $wsrep_local_state == 3 ]; then
    exit 1
  fi
```
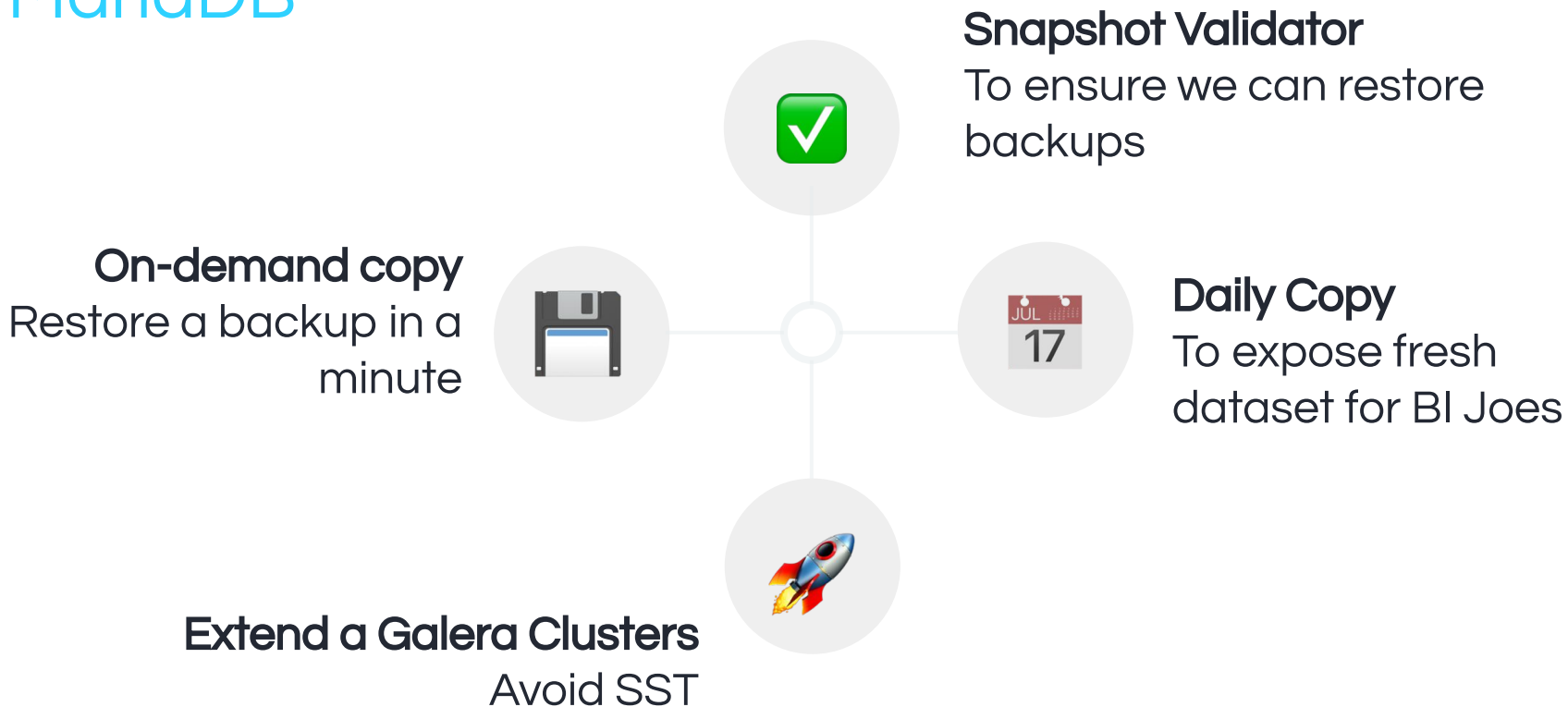
Only nodes **Synced** and **Donor** are considered "ready"

# Having fun with the **Persistent disk Snapshots**
## MariaDB

**Snapshot Validator**
To ensure we can restore backups

**On-demand copy**
Restore a backup in a minute

**Daily Copy**
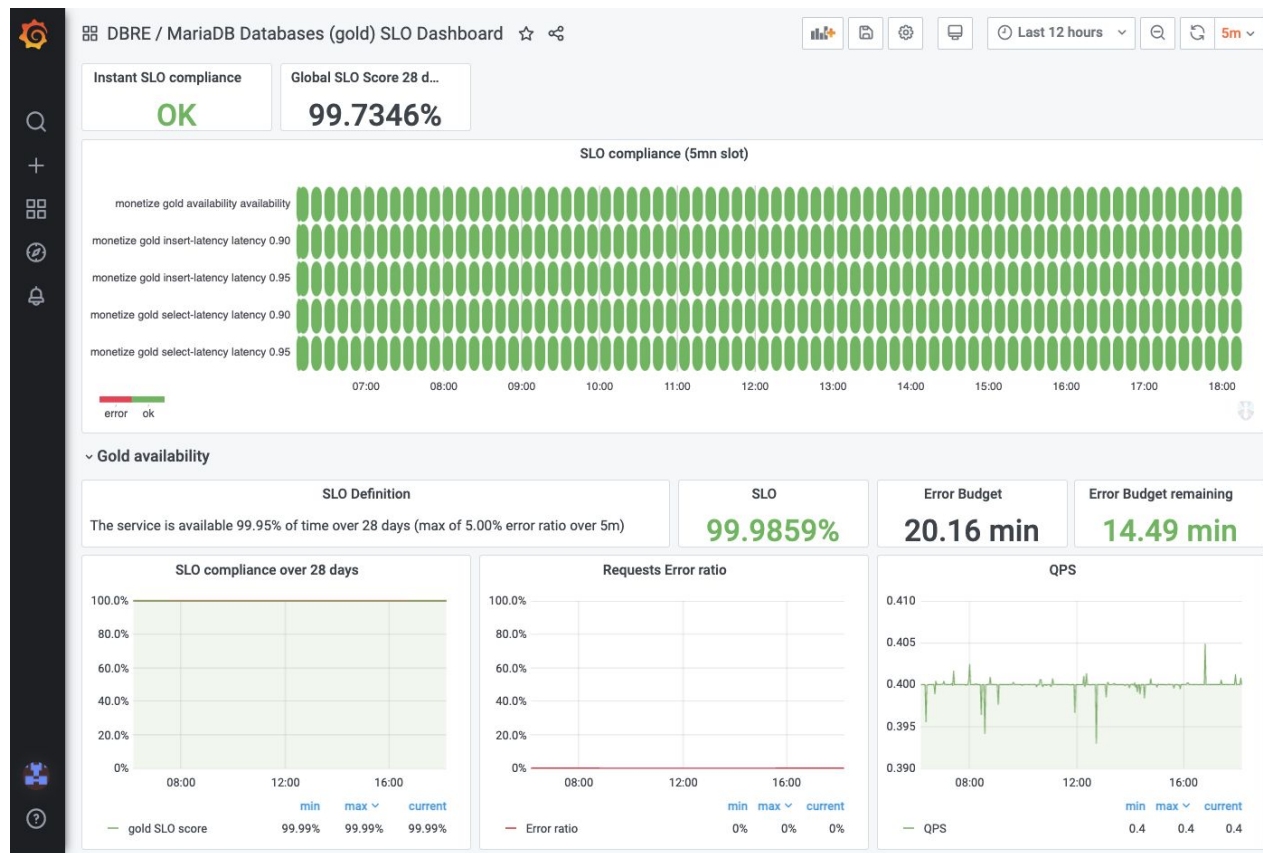To expose fresh dataset for BI Joes

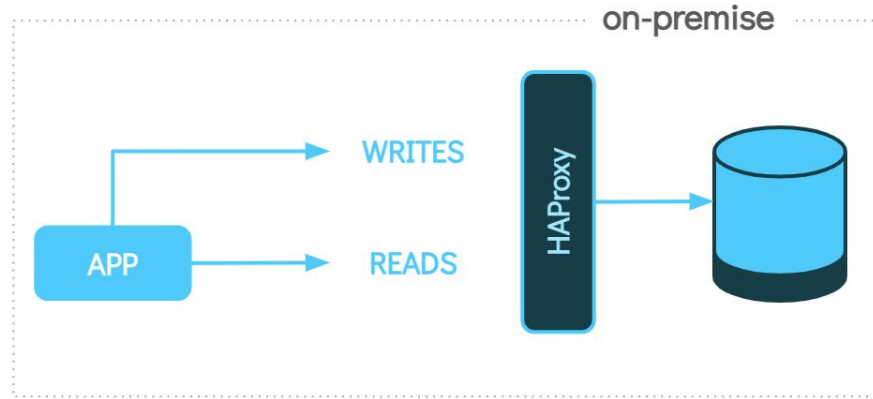**Extend a Galera Clusters**
Avoid SST

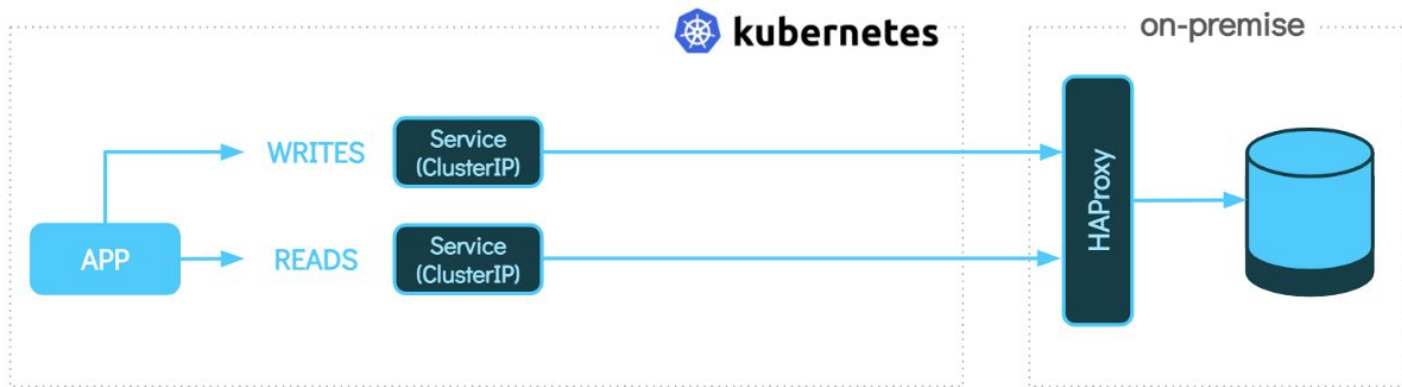# Having a SLI Prober to implement **SLO**

# MariaDB migration path

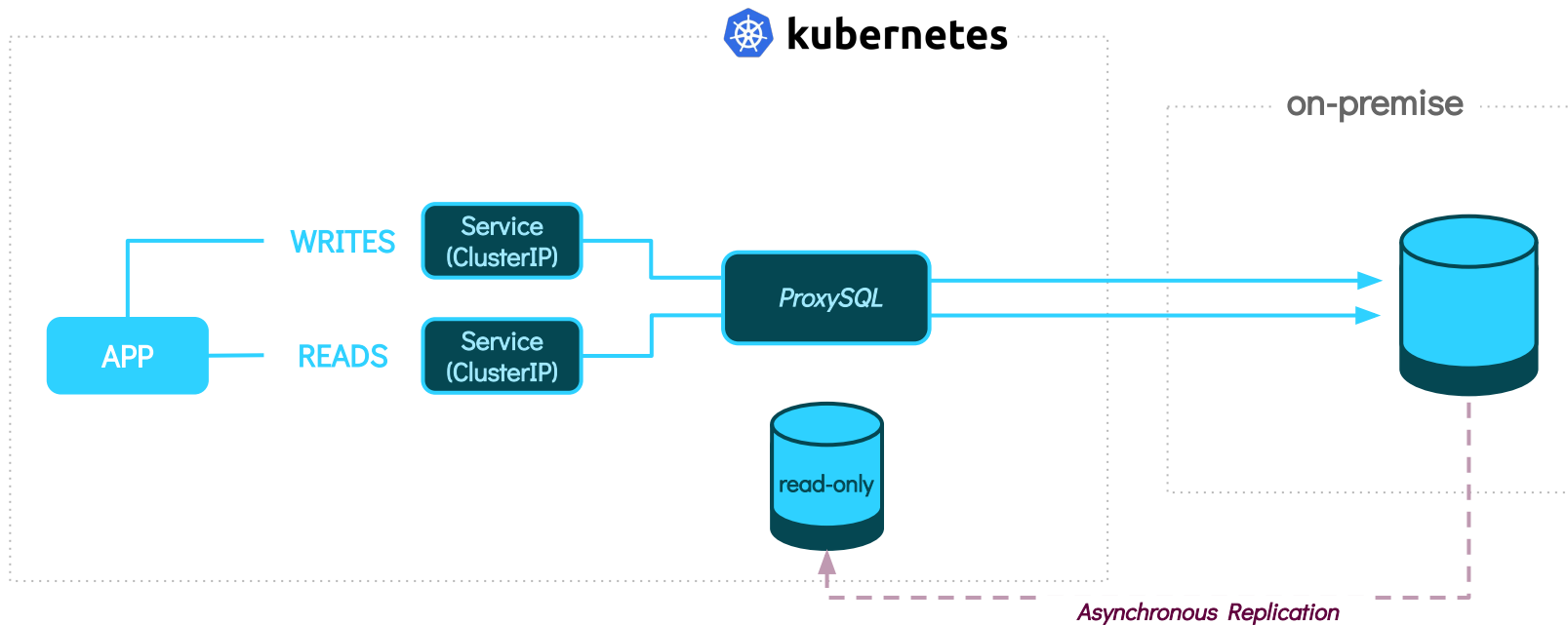# MariaDB Migration Path
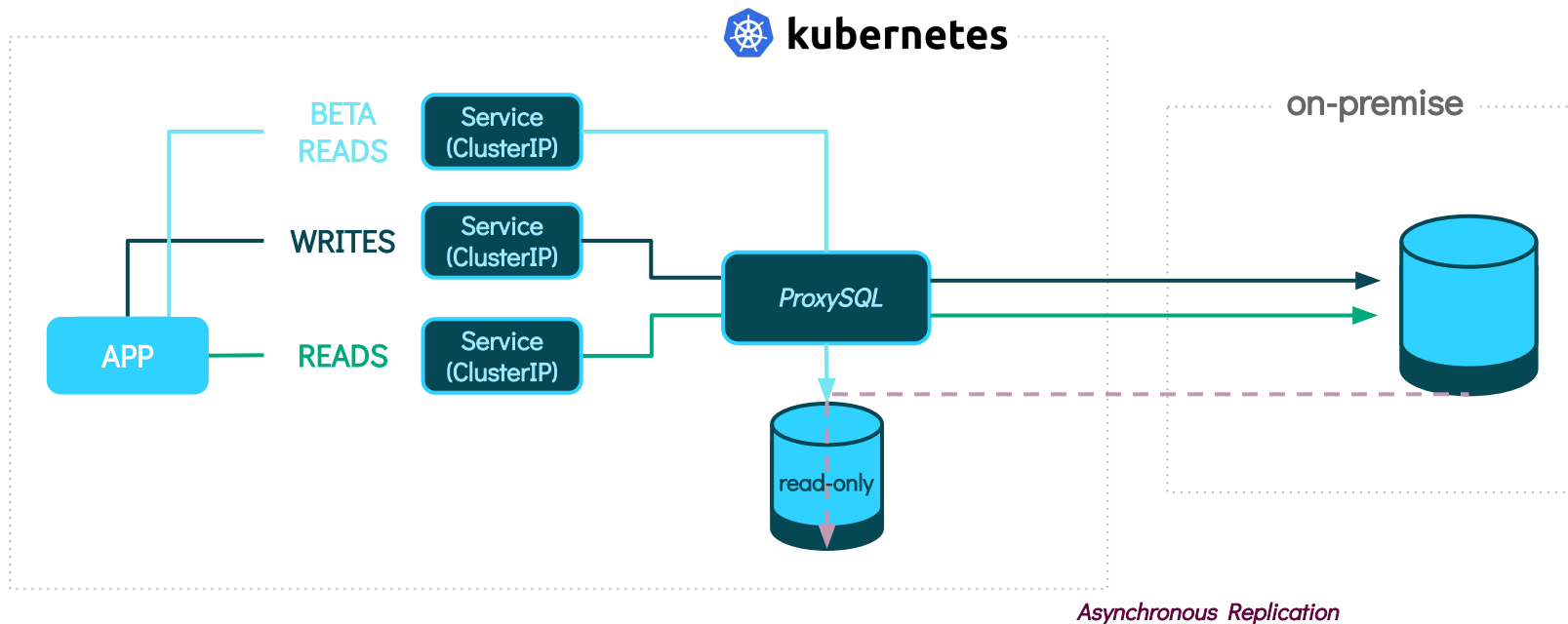## Initial stage

# MariaDB Migration Path
## Move the application

# MariaDB Migration Path
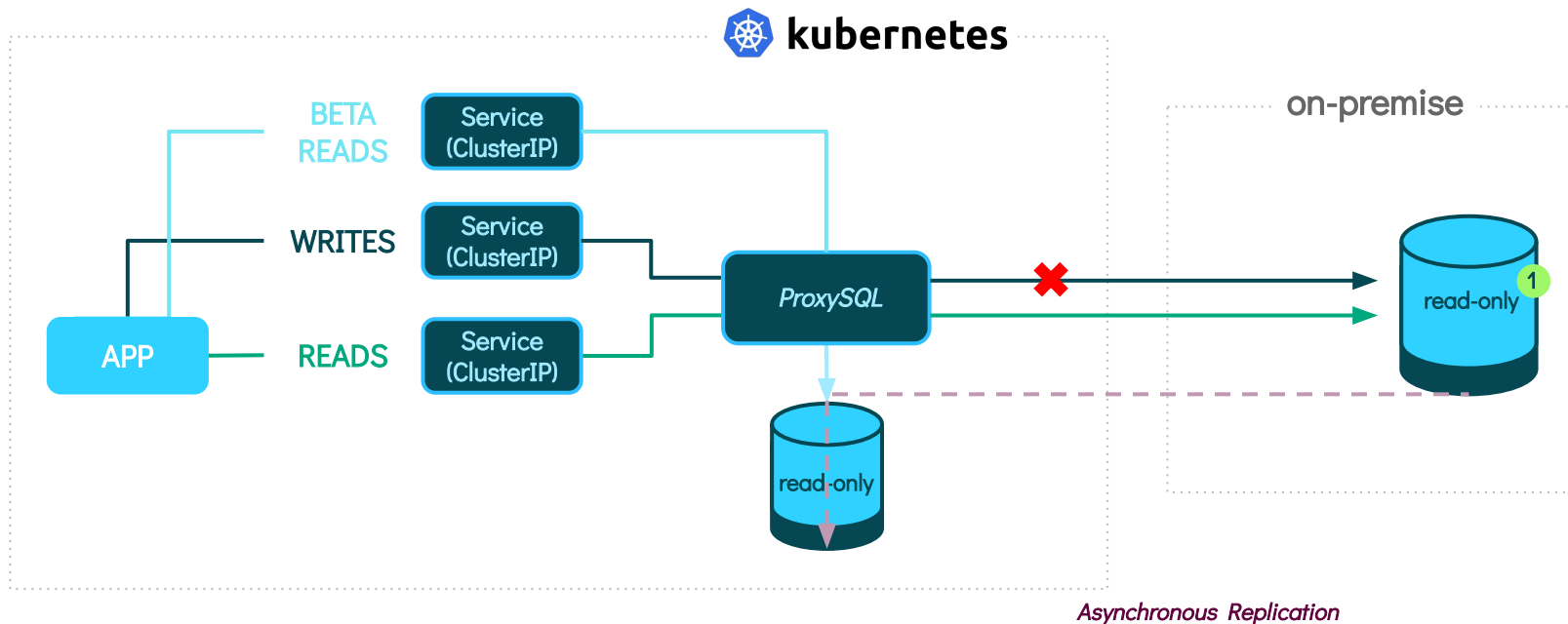## Setup the database in GCP (with replication)

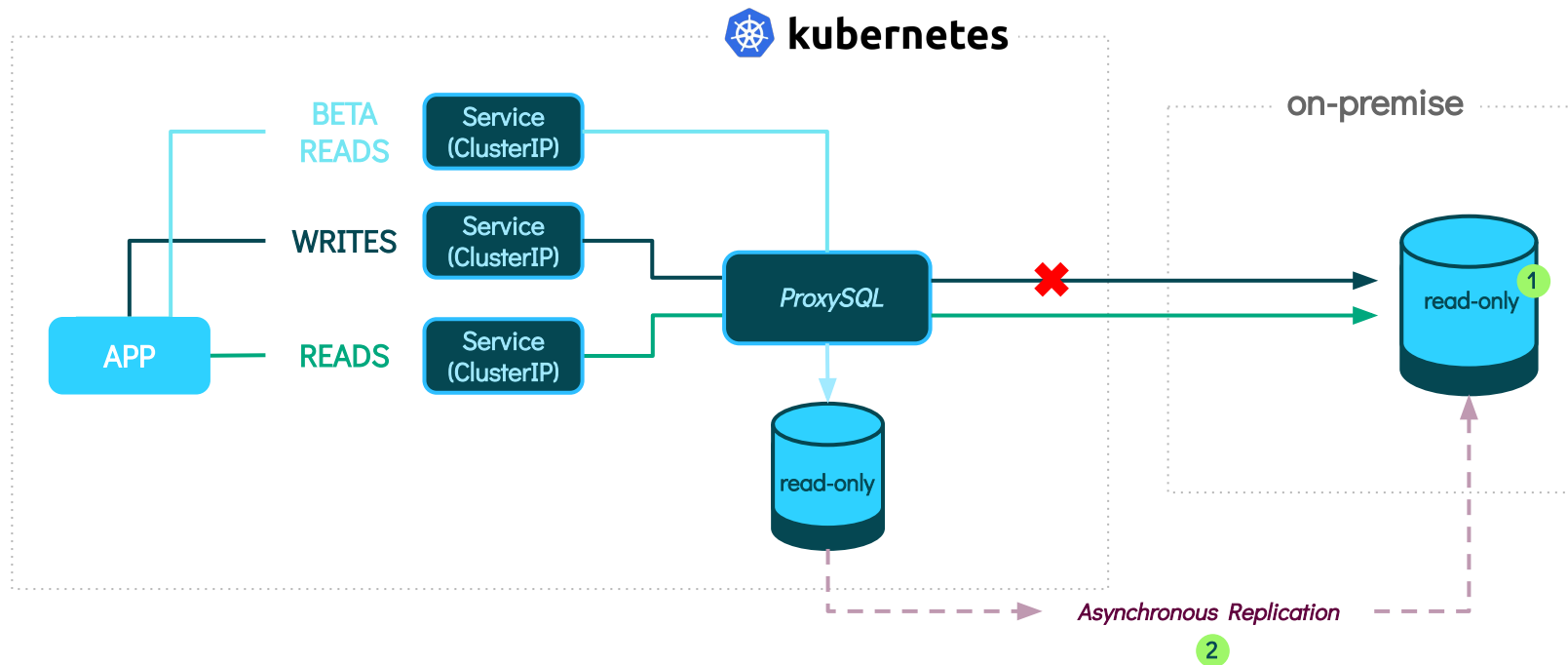# MariaDB Migration Path
## Open a Beta endpoint for reads



kubernetes

BETA READS

WRITES

APP

READS

Service (ClusterIP)

Service (ClusterIP)

Service (ClusterIP)

ProxySQL

read-only

on-premise

*Asynchronous Replication*

# MariaDB Migration Path
## D-Day: Set read-only



kubernetes

BETA
READS

Service
(ClusterIP)

WRITES

Service
(ClusterIP)

APP

READS

Service
(ClusterIP)

ProxySQL

read-only

on-premise

read-only

1

*Asynchronous Replication*

# MariaDB Migration Path
## D-Day: Reverse the replication stream
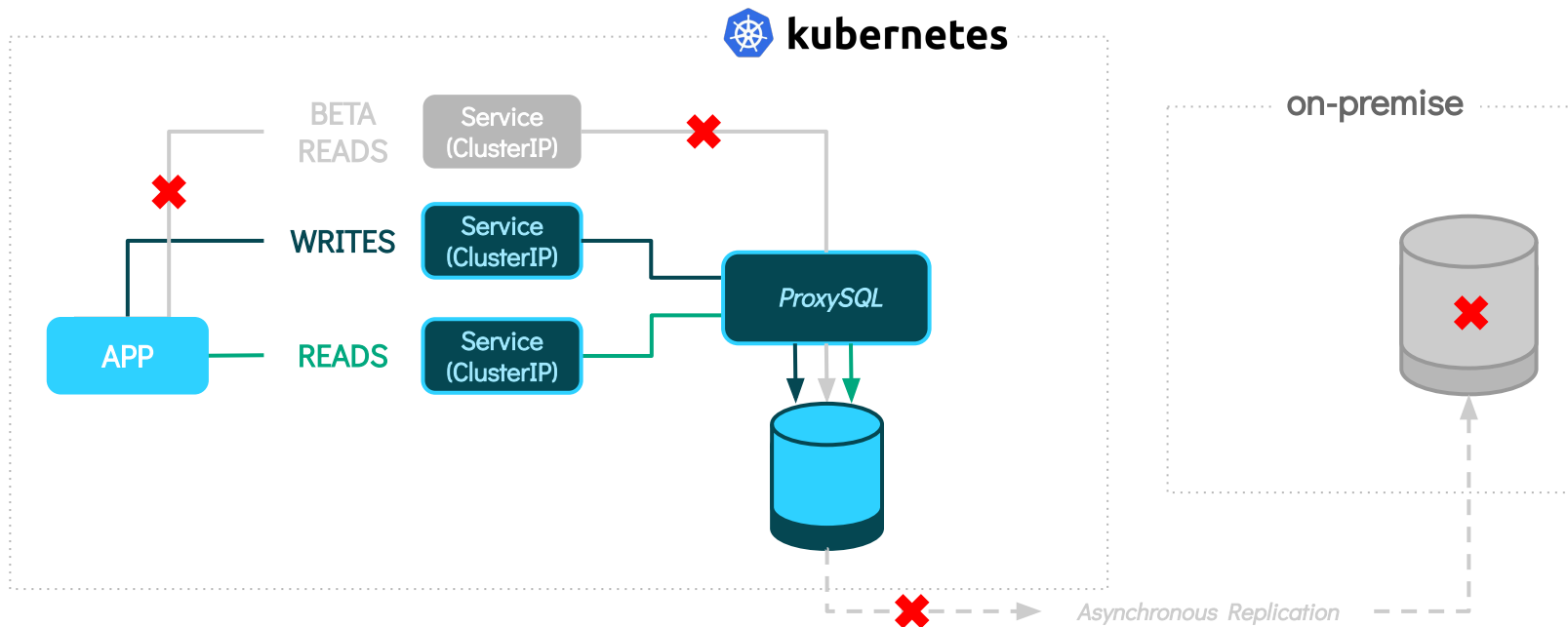
# MariaDB Migration Path
## D-Day: Enable writes in GCP

# MariaDB Migration Path
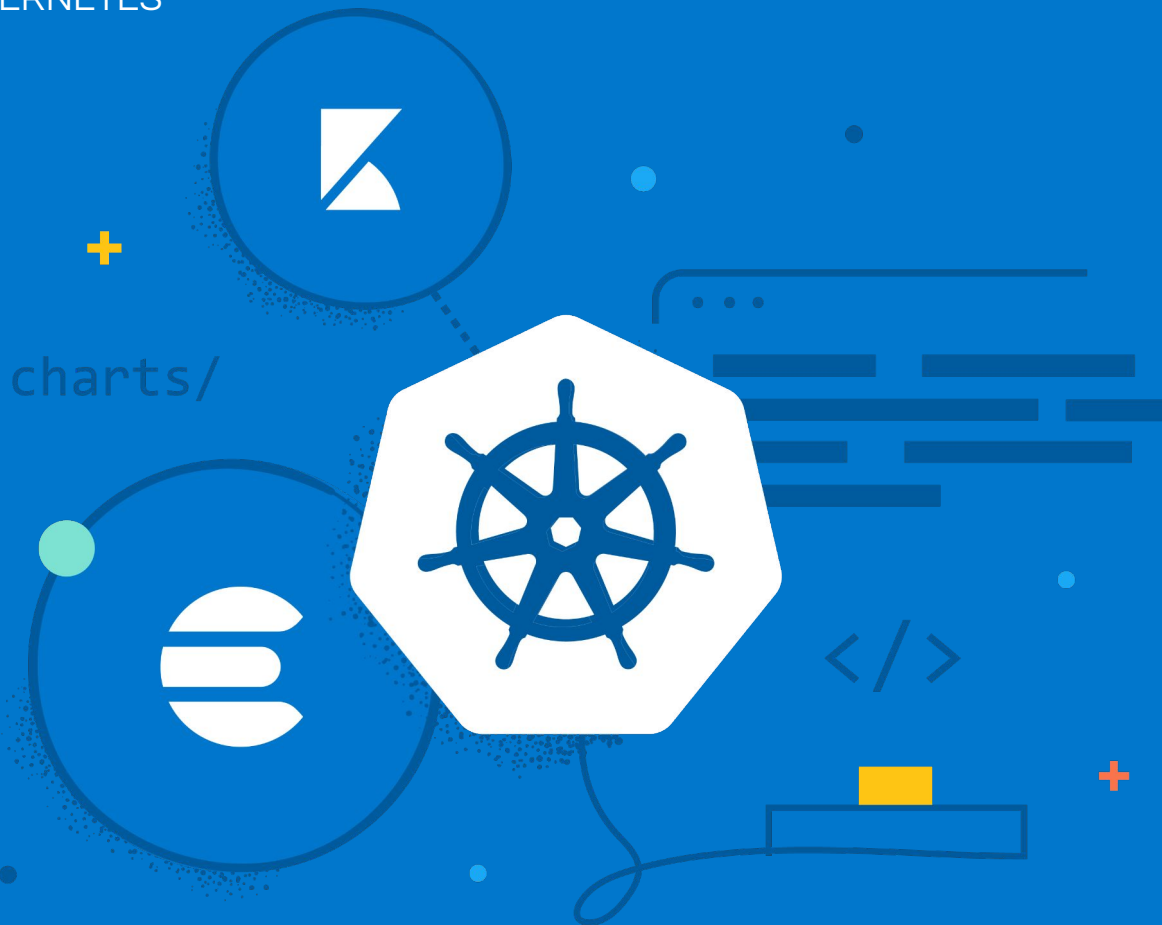## D-Day: Change the endpoints

# MariaDB Migration Path
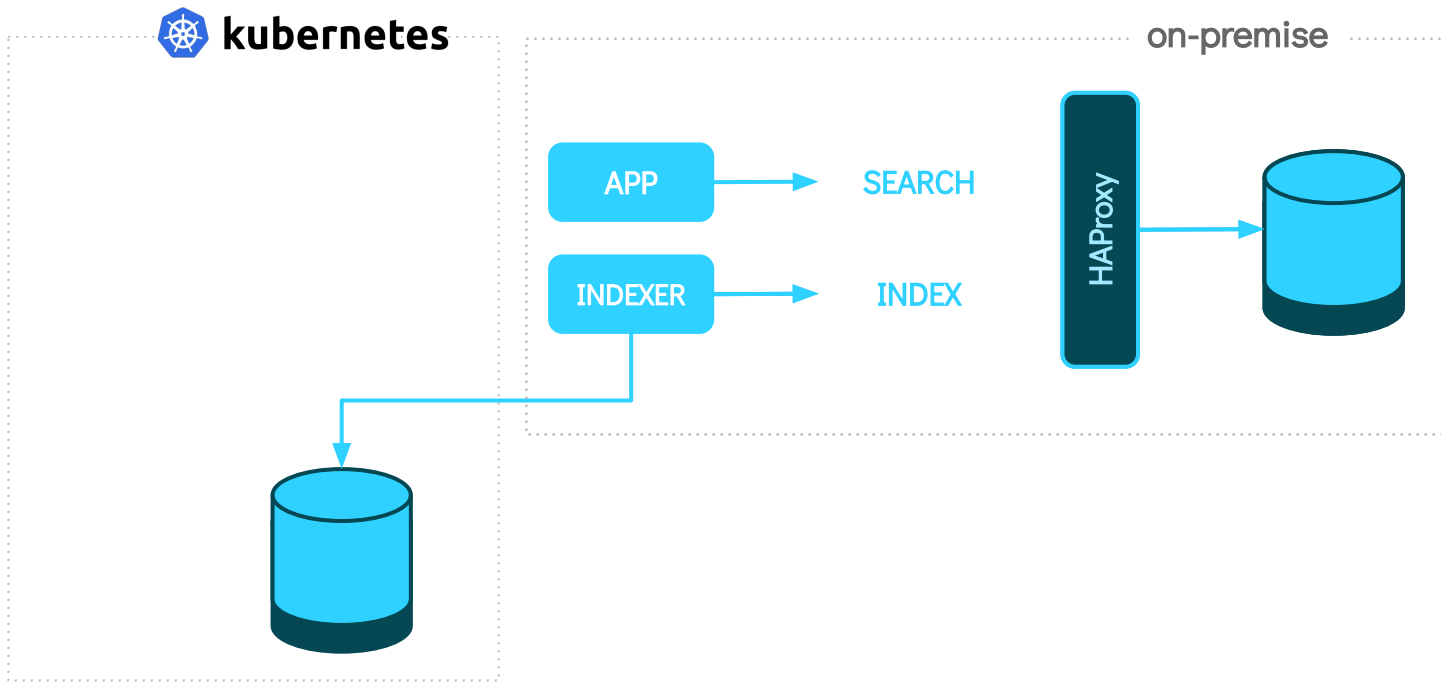## Cleaning and decommissioning
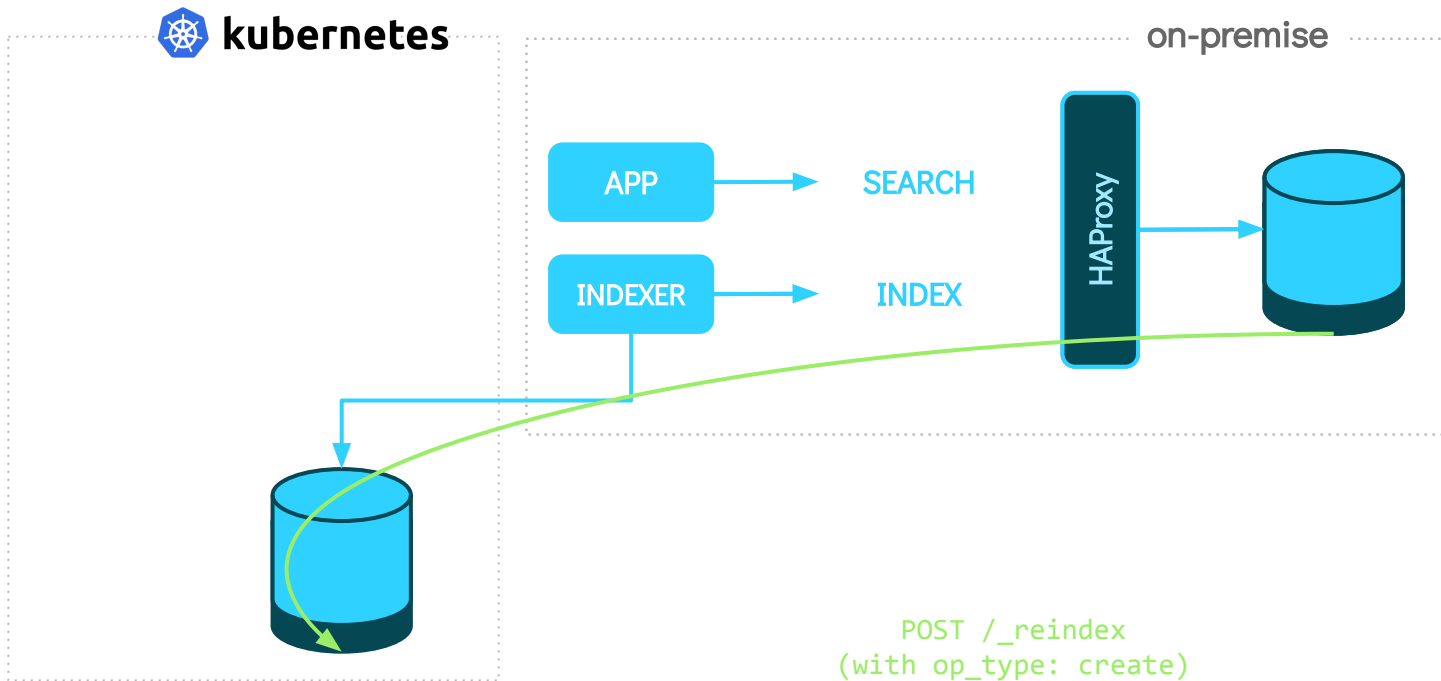
Elasticsearch

charts/

# Elasticsearch Migration Path
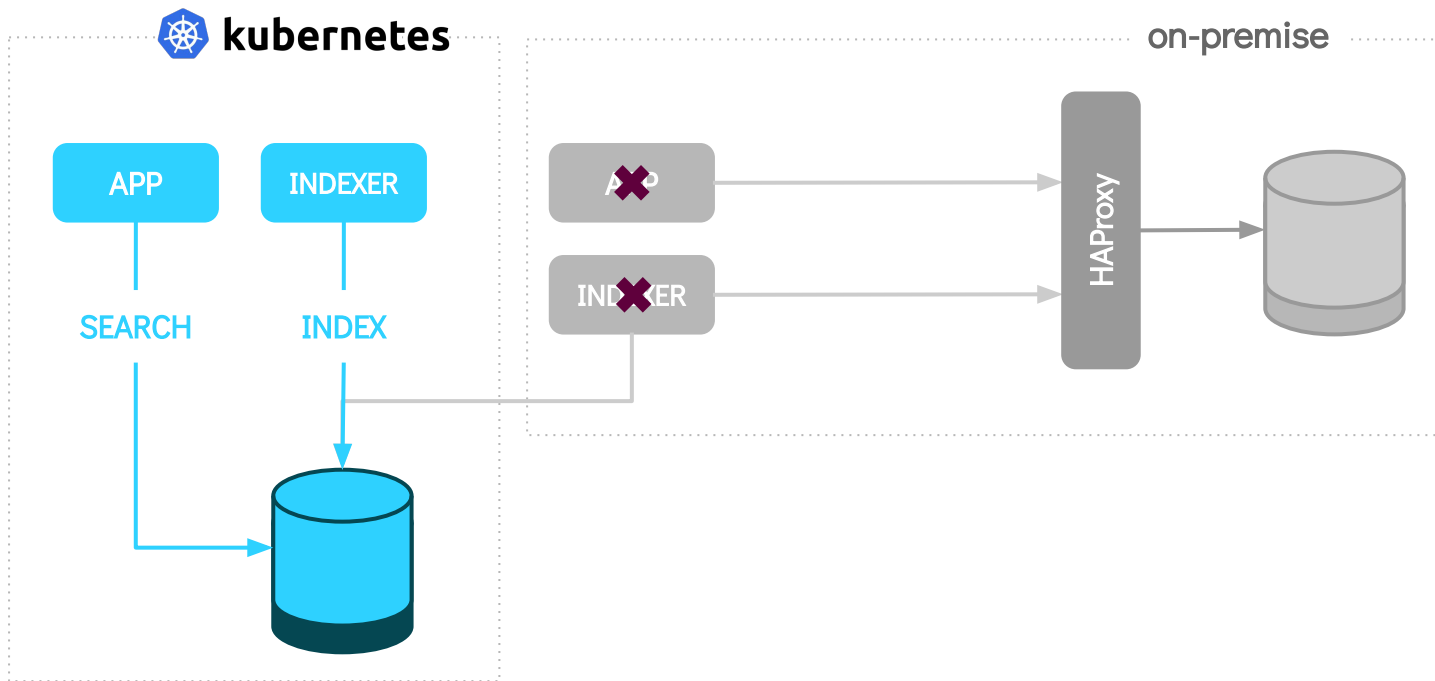## Implement double-writes in indexer

# Elasticsearch Migration Path
## Get missing data from on-premise

# Elasticsearch Migration Path
## Move the application

That's only two use cases...

# Solutions chosen in GCP

Cassandra → Helm Chart
*BlaBlaCar*

RabbitMQ → Helm Chart
*Bitnami*

Redis → MemoryStore

PostgreSQL → CloudSQL

Kafka → K8S Operator
*Strimzi*

Couchbase → MemoryStore