# Franck Pachot

**Developer Advocate on YugabyteDB**
(PostgreSQL-compatible distributed database)

Past:

    20 years in databases, dev and ops

    Oracle ACE Director, AWS Data Hero

    Oracle Certified Master, AWS Database Specialty

    ...

fpachot@yugabyte.com

dev.to/FranckPachot

@FranckPachot

# Databases in Oracle Cloud



What about applications build for PostgreSQL?

What about HA with shared nothing?

What about cloud-native scale-out SQL databases?

# Why Kubernetes?

High Availability
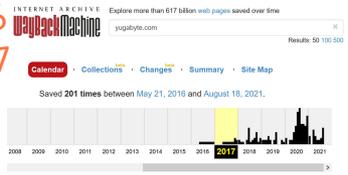
Scalability

Disaster Recovery

+ declarative

yugabyte**DB**

# And Databases?

K8s StatefulSets
- Alpha    Jul. 2016
- Beta     Dec 2016
- Stable   Dec 2017

Persistent

Shared

Replicated

yugabyteDB

# Distributed SQL Database

**yugabyteDB**

A cluster is called 'universe'

Private or public (internet) network
On-premises, cloud, hybrid, open-source
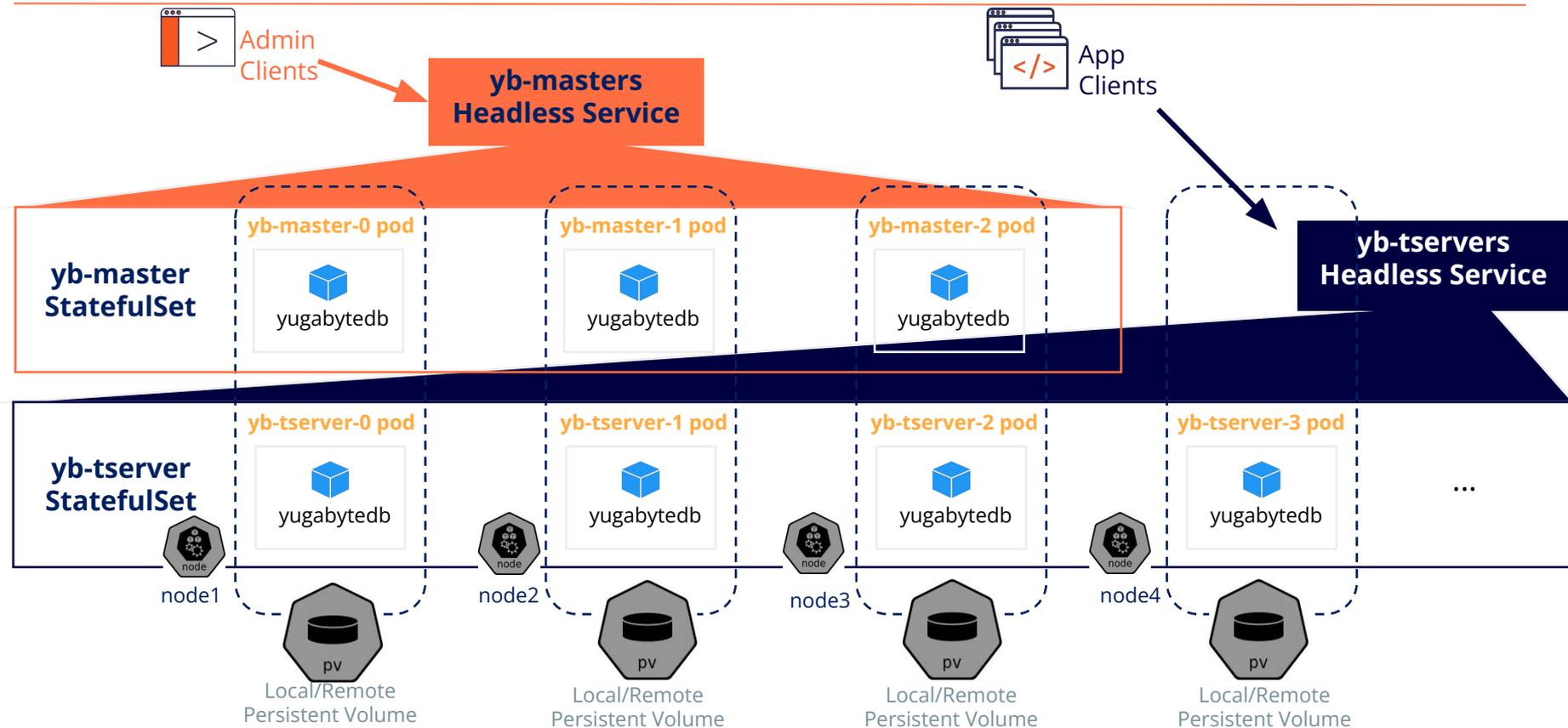
PostgreSQL API

PG compatible with read/write on all nodes

No shared storage, only network

Write on all nodes (Raft protocol)

# YugabyteDB on Kubernetes

© 2020 All Rights Reserved

Demo

Oracle Cloud

OKE

YugabyteDB

# Kubernetes on Oracle Cloud

# Kubernetes on Oracle Cloud

# Helm charts

```
helm repo add yugabytedb \
  https://charts.yugabyte.com
helm repo update

kubectl create namespace yb-demo

helm install yb-demo  \
  yugabytedb/yugabyte  \
  --namespace yb-demo
  --set replicas.{master,tserver}=3
```

allows 1 node failure

# StatefulSets

```
kind: StatefulSet
metadata:
  name: yb-tserver
  namespace: yb-demo
spec:
  replicas: 3
  serviceName: yb-tservers
  podManagementPolicy: parallel
  updateStrategy:
   type: RollingUpdate
  containers:
  - name: yb-tserver
    image: 'yugabytedb/yugabyte:2.7.2.0-b216'
    command:
      - exec /home/yugabyte/bin/yb-tserver \
        --replication_factor=3 --enable_ysql=true
```

ordered or parallel
for faster scale-up

no downtime upgrade

yugabyteDB

# Services

```
kind: Service
metadata:
  name: yb-tserver-service
  namespace: yb-demo
spec:
  ports:
  - name: tcp-ysql-port
    protocol: TCP
    port: 5433
    targetPort: 5433
    nodePort: 31874
  - name: tcp-yql-port
    port: 9042
  selector:
    app: yb-tserver
  clusterIP: 10.244.147.170
  type: LoadBalancer
```

# Storage:

## Ephemeral storage:

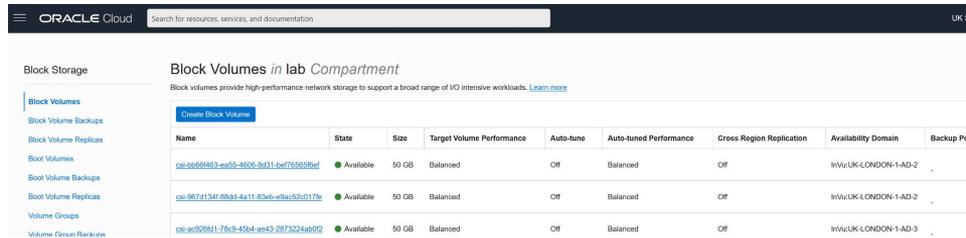if a container dies, data must be read from the other nodes:
- quorum must still be there or data is lost (RPO>0)
- Full availability is up when all data has been transferred (RTO>0)

## Shared storage

A shared remote storage (NFS)
is **not** necessary in a distributed DB



## Local or Persistent Volume

Can be pre-provisioned in the worker node or outside (cloud block storage)

# Persistent Volumes:

✔ provisioned dynamically by K8s from block storage

✔ resilient to node failures (without reconstruction)

```
"volumeClaimTemplates": [
 {
  "kind": "PersistentVolumeClaim",
  "name": "datadir0"
 },
 "spec": {
  "accessModes": [ "ReadWriteOnce" ],
  "resources": { "requests": { "storage": "10Gi" }  },
  "storageClassName": "bv",
  "volumeMode": "Filesystem"
```

mounted once
= not shared

# Anti-affinity:

✔ One pod per node (privileges durability over HA)

preferred or required

```
spec:
 affinity:
  preferredDuringSchedulingIgnoredDuringExecution:
  - weight: 100
   podAffinityTerm:
    labelSelector:
     matchExpressions:
     - key: app
      operator: In
      values:
      - yb-tserver
    topologyKey: kubernetes.io/hostname
```

# Headless service

✔ new nodes discovered and added to DNS

✔ direct connection (placement aware smart clients)

```
apiVersion: v1
kind: Service
metadata:
  name: yb-tservers
  labels:
    app: yb-tserver
spec:
  clusterIP: None
```

# External IP



```
cloudshell$ kubectl get services -n yb-demo
```

```
NAME                TYPE          CLUSTER-IP     EXTERNAL-IP
PORT(S)

yb-master-ui        LoadBalancer  10.96.81.113   150.230.125.157
7000:30510/TCP

yb-masters          ClusterIP     None           <none>
7000/TCP,7100/TCP

yb-tserver-service  LoadBalancer  10.96.94.76    132.226.208.207
5433:31094/TCP,6379:31326/TCP,9042:32136/TCP

yb-tservers         ClusterIP     None           <none>
5433/TCP,9000/TCP,12000/TCP,11000/TCP,13000/TCP,9100/TCP,6379/TCP,9042/TCP
```

# Demo: scale-out



more nodes

more replicas

autorebalance

# Demo: node failure



Kill node

timeout

dead node detected

Other nodes still working (new leaders elected)

# Automating Day 2 Operations

**HANDLING FAILURES**

**ROLLING UPGRADES**

**BACKUP & RESTORE**

**K8s**: Pod failure is automatic
**ops**: Node failure: manually add new workers

**ops**: Local storage failure: manually add new volume
**YB**: Automatic re-sharding

**K8s**: onDelete or rollingUpdate
(pod spawned with same network id / storage)
**YB**: can run with nodes in newer version

**YB**: distributed snapshots and backup
**ops**: restore to existing or new cluster

**yugabyteDB**

# Distributed vs. Streaming Replication

## Streaming replication and sharding:

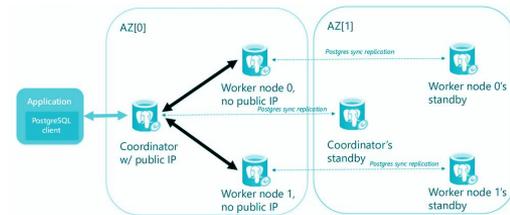The primary is still a SPOF:

      Failover may take time (RTO)

      Failover may miss transactions (RPO)

      Manual or complex automation (rolling upgrade require many failovers)

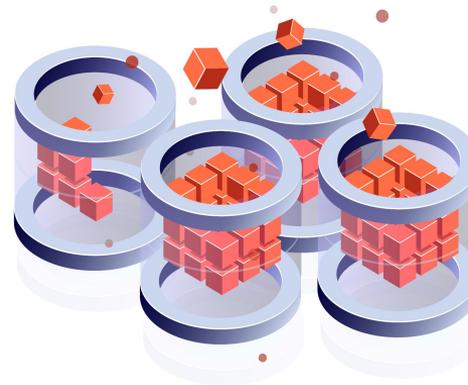*PostgreSQL on K8s at Zalando: Two years in production* [https://av.tib.eu/media/52142](https://av.tib.eu/media/52142)

## Distributed with replication factor

All nodes are equal

      Leaders are balanced over all nodes

      Followers are ready to be elected in few seconds (Raft protocol)

# Thank You

Join us on Slack:
www.yugabyte.com/slack
Star us on GitHub:
github.com/yugabyte/yugabyte-db

fpachot@yugabyte.com

dev.to/FranckPachot

@FranckPachot

yugabyteDB

**Core message:**

- A PostgreSQL database active on multiple nodes

- Operations fully automated (cloud, K8s, PaaS)

- Distributed to provide: Resilience, High Availability, Geo Distribution, Elasticity