



PostgreSQL on Kubernetes



Jan Mußler
@JanMussler

12-07-2017



TABLE OF CONTENTS

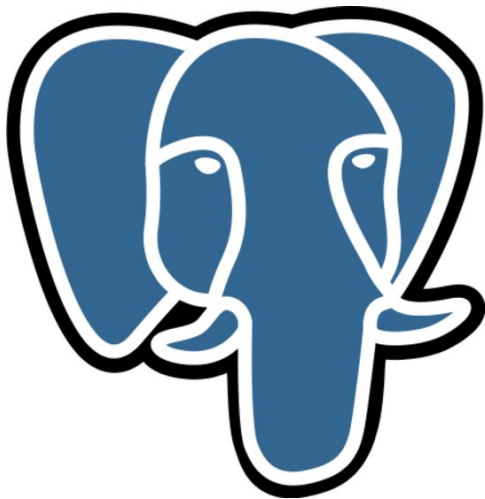


The past

The present

And the future?

PostgreSQL



Powerful, open-source, relational database

MVCC, lock free reads, great at concurrency

Transactional DDL!

Safe to modify tables in production (**yes!**)

Solid owner and role system for privileges

Highly customizable

Supports PL/pgSQL, Python, JS V8

ZALANDO AT A GLANCE

>300 databases
In data centers

>100

Cloud databases
Managed by DB team

>50

staging databases
in Kubernetes



Running PostgreSQL in two data centers

Bare metal with LXC containers

Single Git repository with all configs

Database discovery service

Script to initialize new nodes

Init from slave to lower impact

Time delayed slaves in one data center

PostgreSQL versions: 9.3+

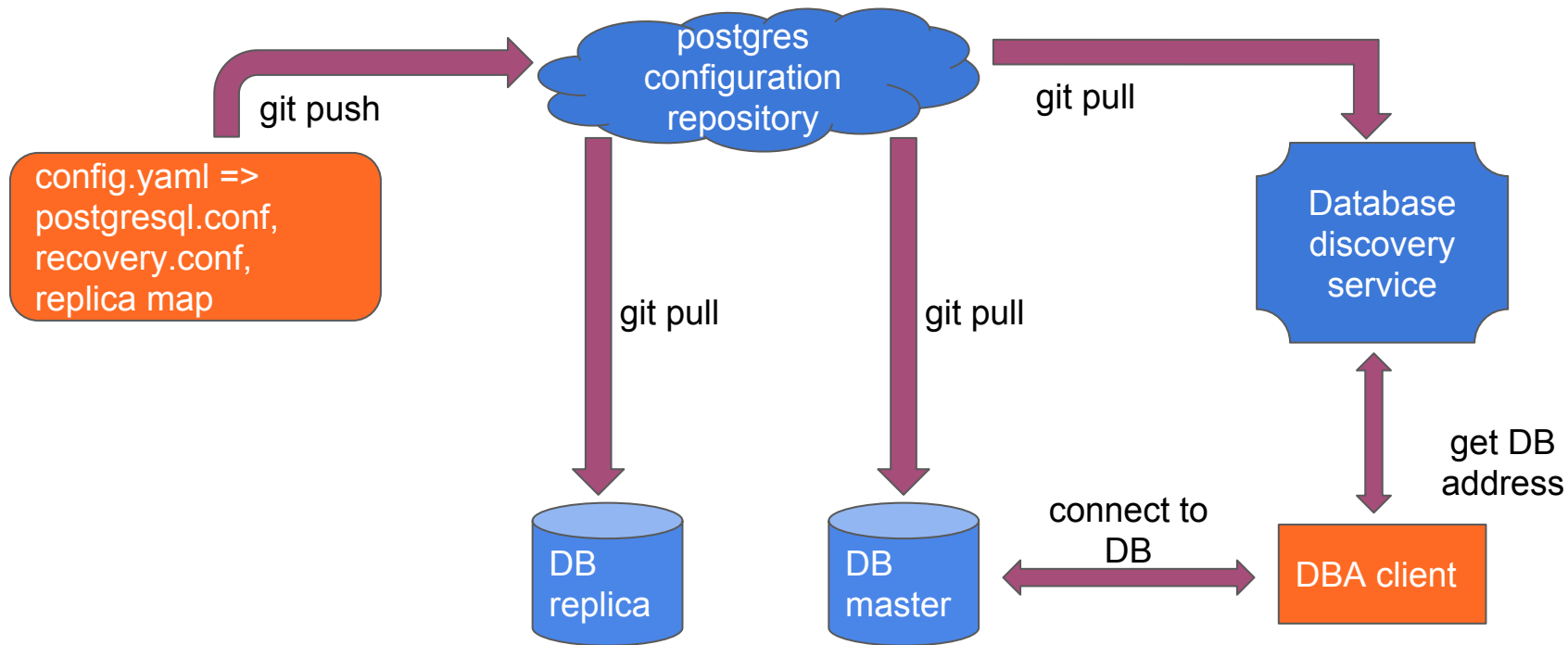
catalogdb01	Master	Slave	Slave
OK	IP: catalogdb01 Host: catalogdb01 Load: 5.23 Delay: 0 B Received: 0 B	IP: itr-catalogdb01-repl Host: catalogdb01-repl Load: 3.64 Delay: 0 B Received: 0 B	IP: itr-catalogdb01 Host: itr-catalogdb01 Load: 0.1 Delay: 1360 MB Received: ∞

catalogdb03	Master	Slave	Slave
OK	IP: catalogdb03 Host: catalogdb03 Load: 2.67 Delay: 0 B Received: 0 B	IP: itr-catalogdb03-repl Host: catalogdb03-repl Load: 3.2 Delay: 0 B Received: 0 B	IP: itr-catalogdb03 Host: itr-catalogdb03 Load: 0.08 Delay: 1477 MB Received: ∞

catalogreporting	Master	cm	Master
OK	IP: itr-pgcatalogreporting01 Host: itr-pgcatalogreporting01 Load: 0.08 Delay: 0 B Received: 0 B	OK	IP: gth-pgcmds01 Host: gth-pgcmds01 Load: 4.33 Delay: 0 B Received: 0 B

cms	Master	Slave	Slave
OK	IP: gth-pgcmds01 Host: gth-pgcmds01 Load: 0.3	IP: itr-pgcmds02 Host: itr-pgcmds02 Load: 3.18	IP: itr-pgcmds01 Host: itr-pgcmds01 Load: 0.03

Git-driven workflow in data centers



PostgreSQL on AWS

Faster database provisioning

Flexible hardware configurations

CPU, Memory, Storage, Price

Docker is enforced

Expected more node failures

Needs more automation



Patroni to the rescue

PostgreSQL management “daemon”

Platform independent using callback scripts

Implemented in Python

Master election (using ETCD, ...)

Zalando's first open-source repo reaching 1000 ★



Spilo: Packaging Patroni in Docker

Configuration via ENV variables

From self contained for AWS deployment

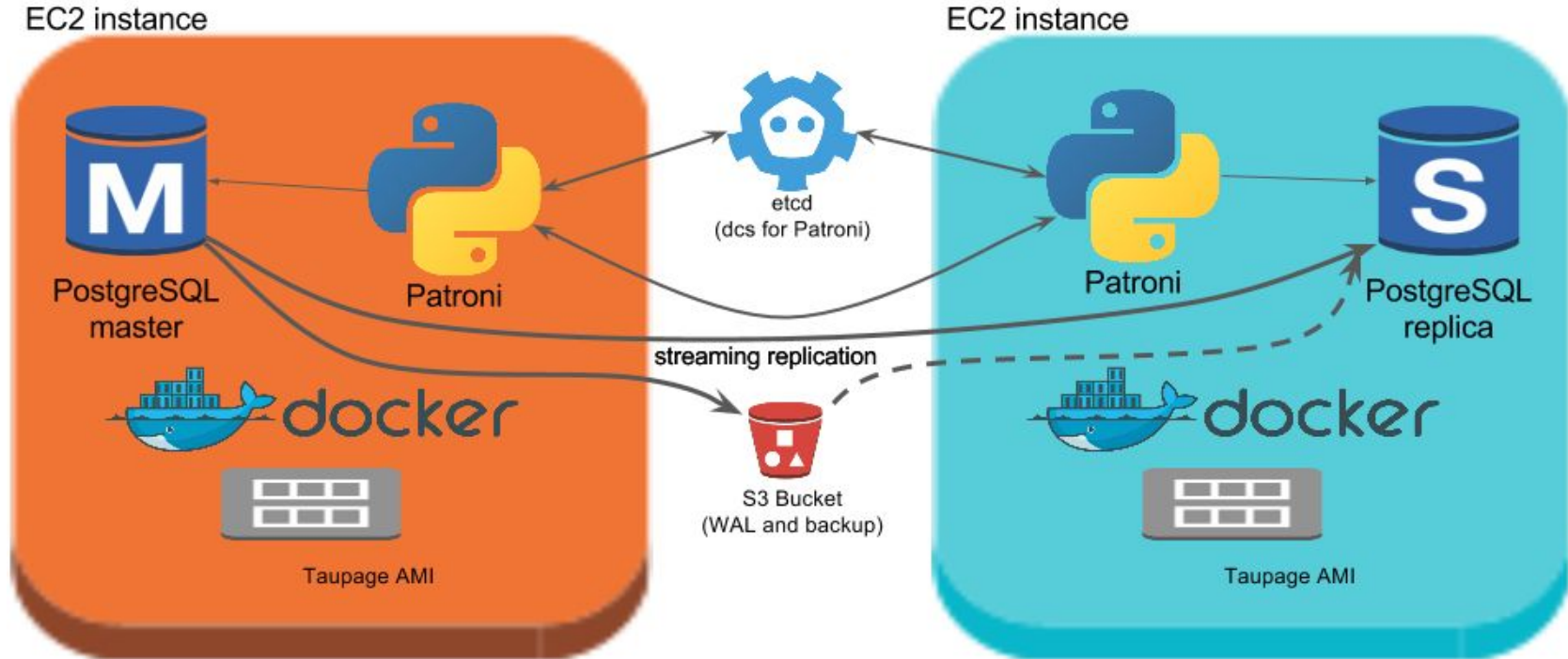
Including all PG versions since 9.3

```
115 # Install patroni and WAL-e
116 ENV PATRONIVERSION=1.2.5
117 ENV WALE_VERSION=1.0.3
118 RUN export DEBIAN_FRONTEND=noninteractive \
119     export BUILD_PACKAGES="python3-pip" \
120     && apt-get update \
121     && apt-get install -y \
122         # Required for wal-e
123         daemontools lzop \
124         # Required for /usr/local/bin/patroni
125         python3 python3-setuptools python3-pystache python3-prettytable python3-six \
126         ${BUILD_PACKAGES} \
127
128     && pip3 install pip --upgrade \
129     && pip3 install --upgrade patroni==${PATRONIVERSION} \
130         gccloud boto wal-e==${WALE_VERSION} \
131
132     # https://github.com/wal-e/wal-e/issues/318
133     && sed -i 's/^(    for i in range(0,\) num_retries):.*$/\1 100):/g' /usr/local/lib/python3.*/site-packages/boto/lib/urllib3/packages/urllib3/packages/six.py
134
135     # Clean up
136     && apt-get purge -y ${BUILD_PACKAGES} \
137     && apt-get autoremove -y \
138     && apt-get clean \
139     && rm -rf /var/lib/apt/lists/* /root/.cache
140
141
142     # install etcdctl
143     ENV ETCDVERSION 2.3.8
```

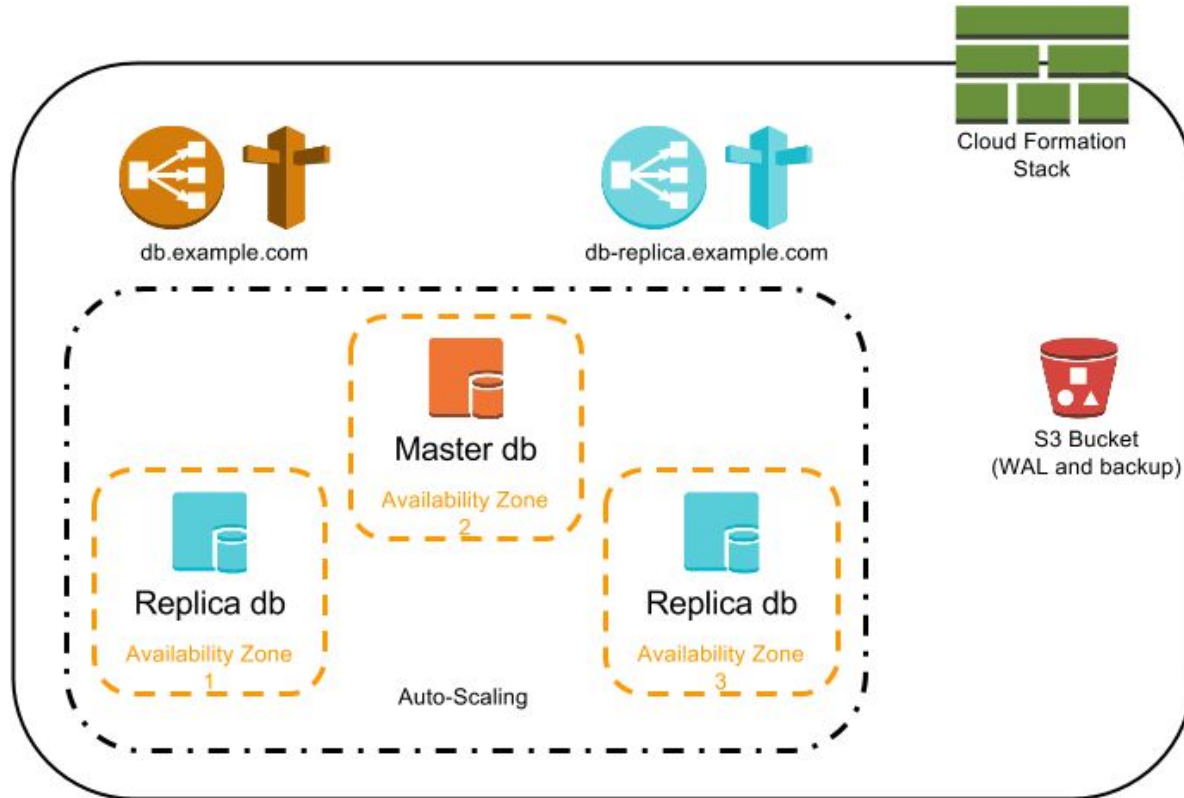
Try it locally!

```
docker run -p 8080:8080 registry.opensource.zalan.do/acid/spilo-9.6:1.2-p27
```

Spilo deployment on AWS



Spilo-based deployments in the cloud



Why not AWS RDS or Aurora PostgreSQL



Not an easy answer :)

Full control

- Independent of cloud provider
- Real super user available
- Custom extensions
- Streaming/WAL replication in and out
- Local storage not supported on RDS

Costs? Cost of development? ...

Costs: RDS vs Spilo

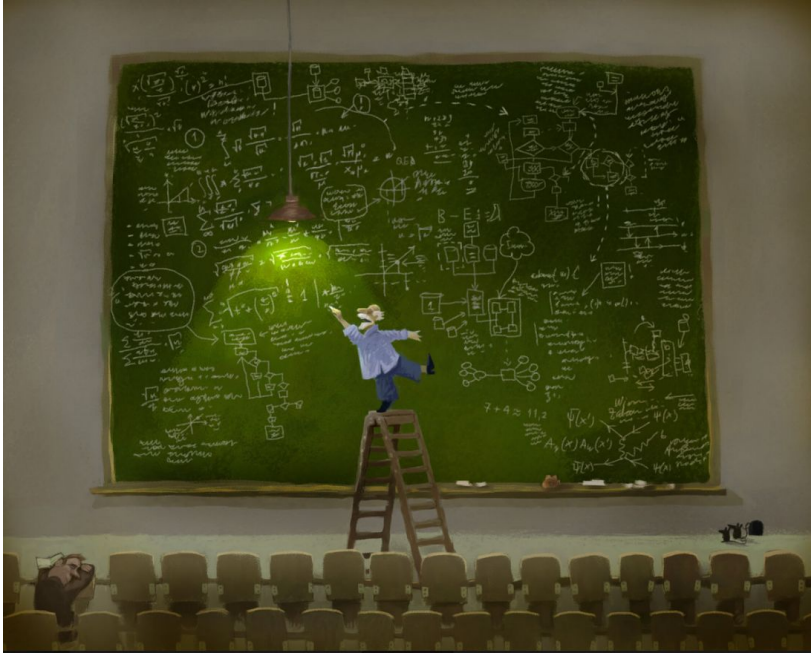
	RDS Multi AZ + 100GB	Spilo (2x EC2 + 100GB gp2)
m4.large	335€	219€
m4.2xlarge	1300€	750€
m4.4xlarge	2600€	1460€

+ ELB Traffic costs for Spilo (RDS can be VPC internal)

PostgreSQL as a Service



Goals



Automation

- Get rid of Jira and users waiting
- Convenient way to get new cluster
- Enable users to modify cluster setup

Integration

- Works with deployment pipeline
- User and Application user provisioning
- Monitoring out of the box

Kubernetes



Kelsey Hightower ✓

@kelseyhightower

Following



Kelsey's guide to running traditional databases on Kubernetes. Strongly consider using a managed service.

6:56 PM - 20 Jan 2017

92 Retweets 175 Likes



The Stateful Set

```
apiVersion: apps/v1beta1
kind: StatefulSet
metadata:
  labels:
    application: spilo
    name: zmon-nc-10-3
    namespace: default
spec:
  replicas: 2
  selector:
    matchLabels:
      application: spilo
  serviceName: zmon-nc-10-3
  template:
    metadata:
    spec:
      containers:
      - env:
        - name: ETCD_HOST
          value: etcd.o.z.d
        image: r.o.z.d/acid/spiloprivate-10:1.3-p2
```

Building block for stateful services

- Stable host or pod names
- Stable assignment of volumes
- Scale factor
- Ordering during upstart and scaling

Useful for: Cassandra, ETCD, PostgreSQL, ...

Storage on Kubernetes

```
spec:
  containers:
  - env:
    [ .. ]
    volumeMounts:
    - mountPath: /home/postgres/pgdata
      name: pgdata

[..]

volumeClaimTemplates:
- metadata:
  annotations:
    volume.alpha.kubernetes.io/storage-class: default
  name: pgdata
  spec:
    accessModes:
    - ReadWriteOnce
    resources:
      requests:
        storage: 10Gi
```

Volumes and Volume claims

Auto provisioning of volumes in the Cloud

EBS volumes in our case

No resize support in K8S itself

Volumes per failure domain (AZ)

Local storage in 1.7

Third party resources

Store user-defined data in K8S

Create your own “types”

Big YAML “blob”

Watchable by applications via K8S API

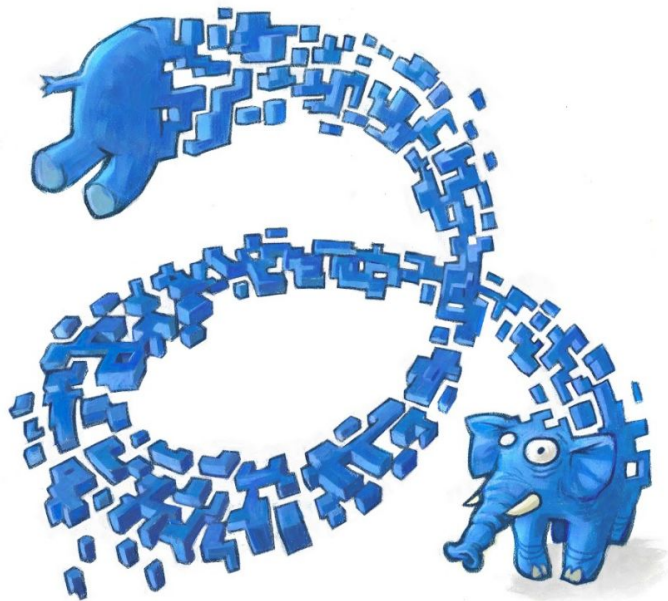
Changes in 1.7 to Custom Resource Definition

Cluster manifest with PostgreSQL configuration

```
apiVersion: acid.zalan.do/v1
kind: Postgresql
metadata:
  labels:
    team: acid
  name: acid-newdatabase-01
  namespace: default
spec:
  allowedSourceRanges:
  - 192.168.1.0/24
  numberOfInstances: 2
  postgresql:
    version: "9.6"
    parameters:
      shared_buffers: 512MB
      max_connections: 128
  teamId: acid
  volume:
    size: 10Gi
```

payload

The “postgres-operator”



Application to manage PostgreSQL clusters

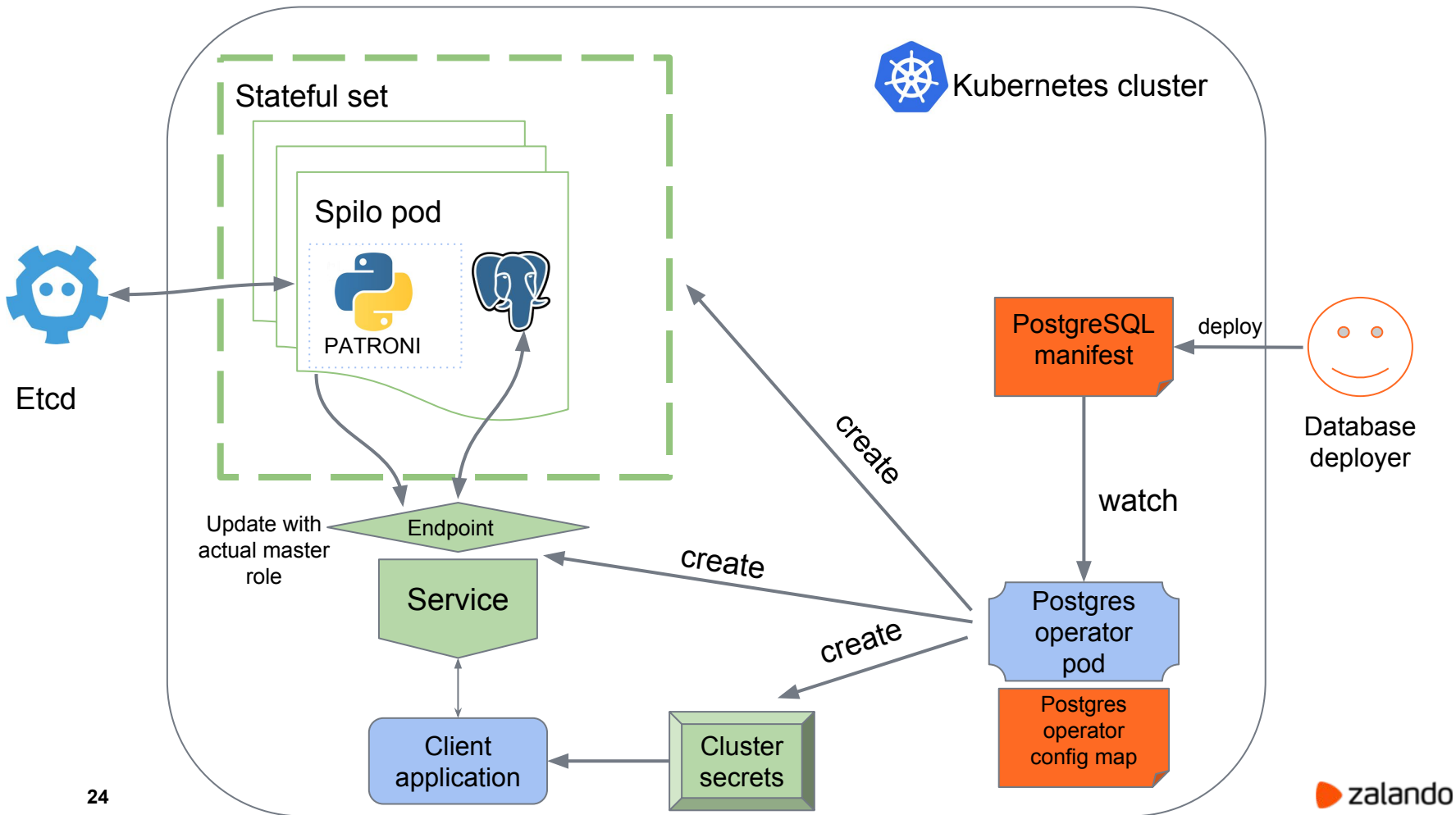
Observes 3rd party “postgres” manifests

Spawns and modifies new clusters

Syncs and provisions roles

Handles volume resize, incl. Resize2fs

Also responsible for e.g. updating Docker images



Create a new PostgreSQL cluster

Cluster YAML definition

```
kind: "Postgresql"
apiVersion: "acid.zalan.do/v1"

metadata:
  name: "guild-24x7-new-cluster"
  namespace: "default"
  labels:
    team: guild-24x7

spec:
  teamId: "guild-24x7"
  volume:
    size: "10Gi"
    numberOfInstances: 1

  postgresql:
    version: "9.6"

  allowedSourceRanges:
    # IP ranges to access your cluster go here
```

Cluster configuration

Database Name	<input type="text" value="new-cluster"/>
Owning team	<input type="text" value="guild-24x7"/>
PostgreSQL Version	<input type="text" value="9.6"/>
DNS Name:	new-cluster.acid.staging.
Number of instances	<input type="text" value="1"/>
Replica Load Balancer	<input type="checkbox"/> Enable Replica ELB
Volume Size	<input type="text" value="10"/> Gi
Office/DC IP ranges	<input type="checkbox"/> Datacenter 1 <input type="checkbox"/> Datacenter 2 <input type="checkbox"/> Berlin
AWS NAT IPs	1. <input type="text"/> / 32 2. <input type="text"/> / 32 3. <input type="text"/> / 32
Odd Host	IP <input type="text"/> / 32

[Copy Definition](#)[Create Cluster](#)

Waiting for master to become available

PostgreSQL cluster status **zmon-nc-10-3** [edit](#)

Cluster YAML definition

```
apiVersion: acid.zalan.do/v1
kind: PostgreSQL
metadata:
  creationTimestamp: '2017-06-07T16:31:33Z'
  labels:
    team: zmon
  name: zmon-nc-10-3
  namespace: default
spec:
  allowedSourceRanges:
  numberOfInstances: 1
  postgresql:
    version: '10'
    teamId: zmon
    volume:
      size: 10Gi
  status: Running
```

Checking status of Cluster



Waiting for master to become available

First PostgreSQL cluster container spawned

StatefulSet created

PostgreSQL 3rd party object created

Create request successful

Cluster create completed

PostgreSQL cluster status **zmon-nc-10-3** [edit](#)

Cluster YAML definition

```
apiVersion: acid.zalan.do/v1
kind: Postgresql
metadata:
  creationTimestamp: '2017-06-07T16:31:33Z'
  labels:
    team: zmon
    name: zmon-nc-10-3
    namespace: default
spec:
  allowedSourceRanges:
    -
  numberOfInstances: 1
  postgresql:
    version: '10'
  teamId: zmon
  volume:
    size: 10Gi
status: Running
```

Checking status of Cluster

PostgreSQL ready: **nc-10-3.zmon** (DNS may be slow)

PostgreSQL master available, label is attached

First PostgreSQL cluster container spawned

StatefulSet created

PostgreSQL 3rd party object created

Create request successful

Deployment to production

YAML file → Git commit → Deployment → K8S apply

Employee and application roles

Employee roles get created based on owning teams

Employee login via oauth2 tokens, no credential sync, MFA

Application roles:

- Declared in manifest
- Credentials stored in kubernetes secret
- Applications can mount secret
- Credentials don't leave the system

Monitoring



















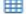









Integration with ZMON

My Team's clusters (2)

Team	Instances	Name	Monitoring
zmon	1	zmon-nc-10-3	 Node metrics  Diskspace  Tables  Indexes
zmon	2	zmon-new-cluster-replica-lb	 Node metrics  Diskspace  Tables  Indexes

All clusters (169)

search:

Team	Instances	Name	Monitoring
ACID	2	acid-testcluster	 Node metrics  Diskspace  Tables  Indexes
ACID	2	acid-testcluster-42	 Node metrics  Diskspace  Tables  Indexes
DBAAS	3	contentbrokerpicasso	 Node metrics  Diskspace  Tables  Indexes
Gerrymanders	2	prodstoreaxiom13	 Node metrics  Diskspace  Tables  Indexes
Gerrymanders	3	productstoreaxiom1	 Node metrics  Diskspace  Tables  Indexes
Gerrymanders	2	productstoreaxiom5	 Node metrics  Diskspace  Tables  Indexes
Gerrymanders	3	productstoreaxiom9	 Node metrics  Diskspace  Tables  Indexes

Monitoring with pgview.web

CPU (Cores: 4) 🌐



Memory

Total	Free	Buffers	Cached	Dirty	Commit Limit	Committed As
32.16GB	1.140GB	5.294GB	22.60GB	114.5MB	B	B

Partitions 🗄️

Device	await	read	write	Free	Total	Path	Size
data	9.37	8960	132352	8.271GB	11.49GB	/home/postgres/pgdata/pgroot/data	542.9MB
wal	9.37	8960	132352	8.271GB	11.49GB	/home/postgres/pgdata/pgroot/data/pg_xlog	201.4MB

Processes 3 / 6 of maximum 100 Play Non Backends 🗄️ 🔍

PID	Lock	Type	utime	stime	read	write	age	DB	User	Query
59		logger	0	0	0	0				
201		checkpointer	0	0	0	0				
202		writer	0	0	0	0				
203		stats collector	0	0	0	0				
312		backend	0	0	0	0		postgres	postgres	idle in transaction (aborted)
383		backend	68	31	0	108864	5	postgres	postgres	insert into test select id from generate_series(1, 10000000) id;
445		wal writer	0	0	0	0				
446		autovacuum launcher	0	0	0	0				
447		archiver	0	0	0	0				last was 000000030000000600000044
466	383	backend	0	0	0	0	1	postgres	postgres	lock table test;

EC2 Instance Metrics via ZMON



Learnings on AWS: not only credits!

EBS IOPS credits

CPU Credits (t2)

Start/Stop to change instance type
(EBS backed only!)

Prepare for instance stops:
Retirement, rolling upgrades,...

- Postgres Operator

github.com/zalando-incubator/postgres-operator

- Spilo HA Docker image

github.com/zalando/spilo

- Patroni

github.com/zalando/patroni

- Web-based monitoring via Postgres background worker

github.com/CyberDem0n/bg_mon

- PAM Oauth2 for PostgreSQL

<https://github.com/zalando-incubator/pam-oauth2>

- Postgres Operator blog post

<https://jobs.zalando.com/tech/blog/postgresql-in-a-time-of-kubernetes/>

Thank you!

