

Особенности адаптации классических СУБД к платформе DBaaS/k8s

Андрей Лярский

Engineer B DBA/SQL

PostgreSQL DBA.

Работаю в Авито в команде SQL@DBA с 2021 года.

Стек команды: PostgreSQL, CockroachDB, Golang.

Переобувался из Linux Admin → DBA, теперь DBA → SRE.

tg: @oxumorron



Как мы адаптиров ми PostgreSQL к платформе DBaaS/k8s

PostgreSQL является одной из наиболее распространённых «классических» СУБД в Авито

PostgreSQL B Avito

2 / 2 / 4
5 / 2 / 4
6 / 4
6 / 4
6 / 4
7 / 4
8 / 4
8 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4
9 / 4</

Без учёта временных баз для тестов



Базы живут на локальных SSD



Машины разных поколений, могут быть существенные отличия по ресурсам



Инженеры тоже разных поколений

Ожидания от баз

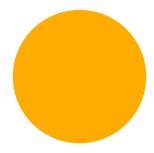
1 2 3 4



High Availability & Fault Tolerance



High Availability & Fault Tolerance

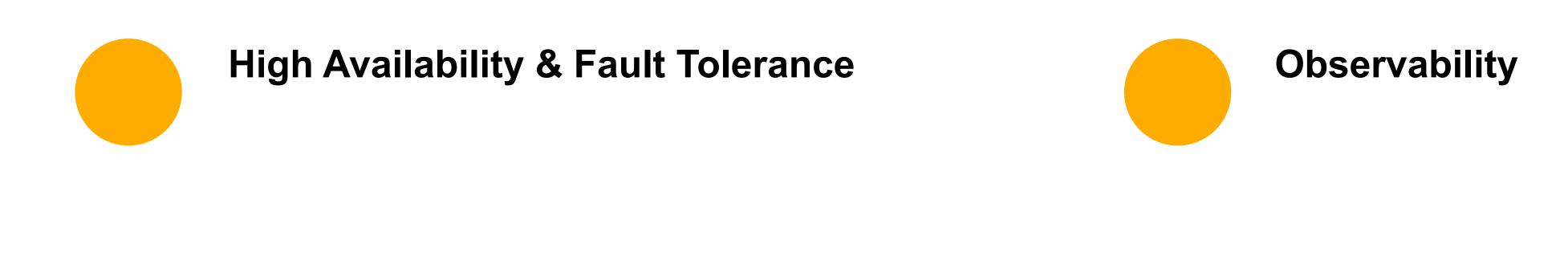


Durability



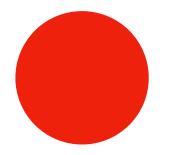
Durability







Durability



High Availability & Fault Tolerance

Native: Streaming Replication



High Availability & Fault Tolerance

Native: Streaming Replication



DurabilityNative: Backup + WAL-archive

= PITR

High Availability & Fault Tolerance

Native: Streaming Replication



DurabilityNative: Backup + WAL-archive

= PITR



Scalability Native: вертикальное (с

оговорками)

High Availability & Fault Tolerance

Native: Streaming Replication



ObservabilityNative: pg_stat_*, csvlog/jsonlog



DurabilityNative: Backup + WAL-archive

= PITR



ScalabilityNative: вертикальное (с

оговорками)



High Availability & Fault Tolerance

Native: Streaming Replication3-rd Party:

Patroni + DCS

High Availability & Fault Tolerance

Native: Streaming Replication3-rd Party:

Patroni + DCS



DurabilityNative: Backups + WAL-

archive = PITR3-rd Party: wal-g

High Availability & Fault Tolerance

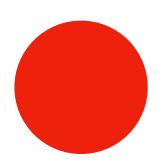
Native: Streaming Replication3-rd Party:

Patroni + DCS



DurabilityNative: Backups + WAL-

archive = PITR3-rd Party: wal-g



ScalabilityNative: вертикальное (с

оговорками)3-rd Party: spqr (c

оговорками)/app-side



High Availability & Fault Tolerance

Native: Streaming Replication3-rd Party:

Patroni + DCS

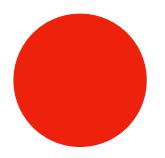


ObservabilityNative: pg_stat_*, csvlog/jsonlog3-rd Party: metrics-scraper



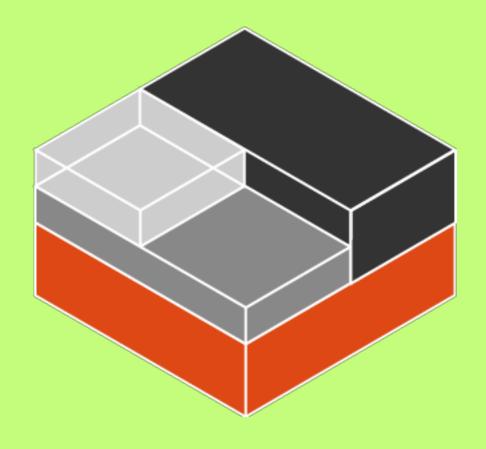
DurabilityNative: Backups + WAL-

archive = PITR3-rd Party: wal-g



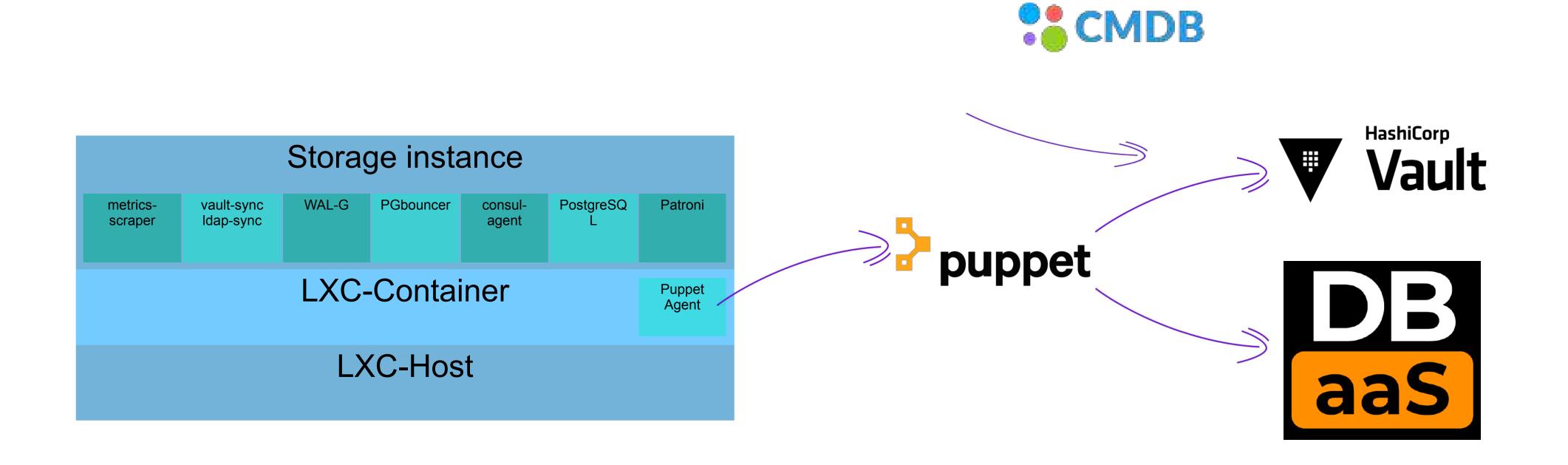
ScalabilityNative: вертикальное (с оговорками)3-rd Party: spqr (с оговорками)/app-side

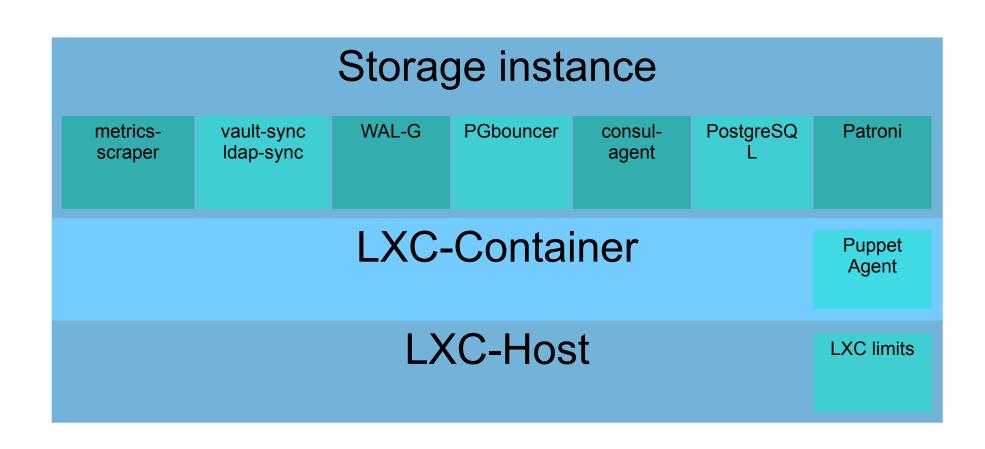
PostgreSQL на LXC Отправная точка



1 2 3 4

PostgreSQL Ha LXClaC: Puppet

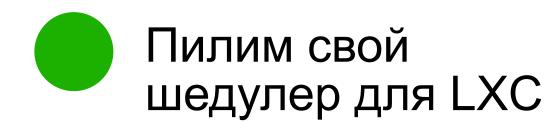


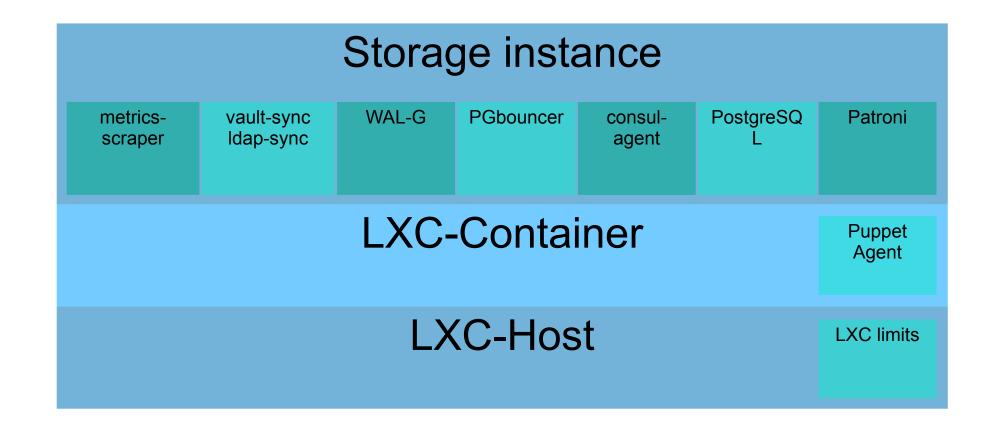


Проблема

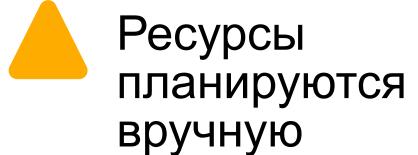
Ресурсы планируются вручную

Решени





Проблема



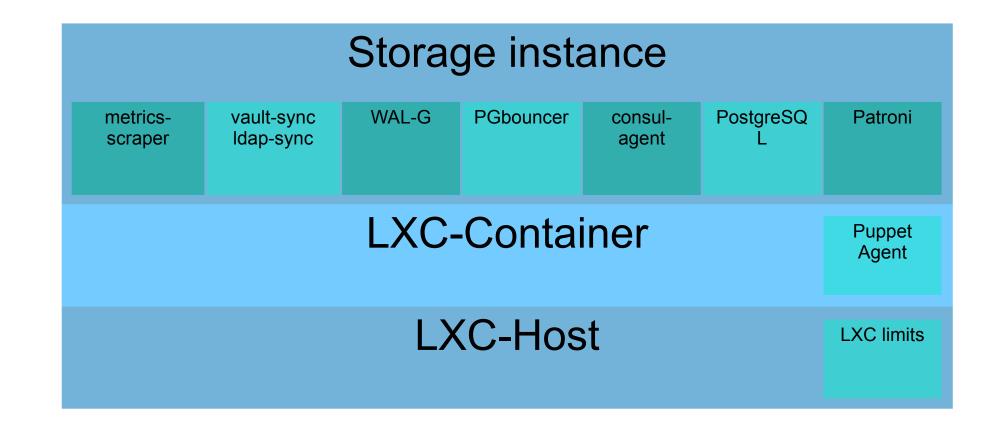


Базы создаются вручнуюСрок — 3-7 дней

Решени

Пилим свой шедулер для LXC





Проблема

Ресурсы планируются вручную

Базы создаются вручнуюСрок — 3-7 дней

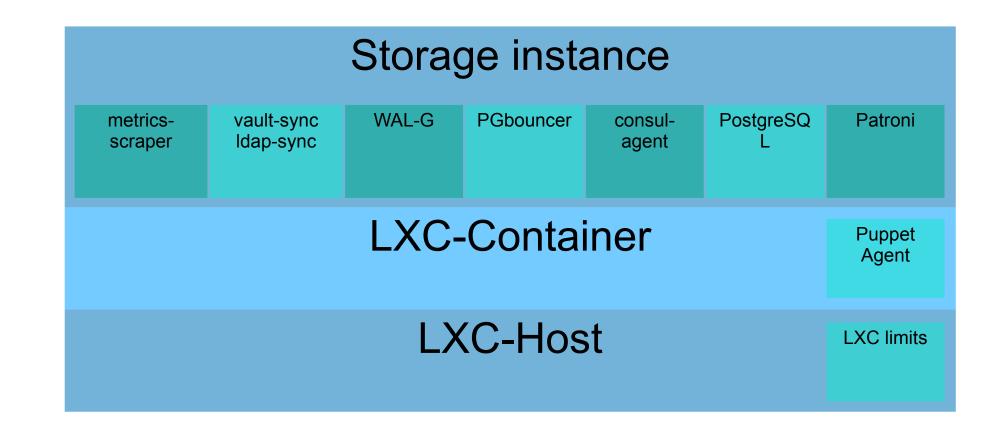
Puppet: SPOF/bottlneck

Решение

Пилим свой шедулер для LXC

Пилим свою автоматику над LXC

Immutable Infrastructure



Проблема

Ресурсы планируются вручную

Базы создаются вручнуюСрок — 3-7 дней

Puppet: SPOF/bottlneck

Consul as DCS: SPOF

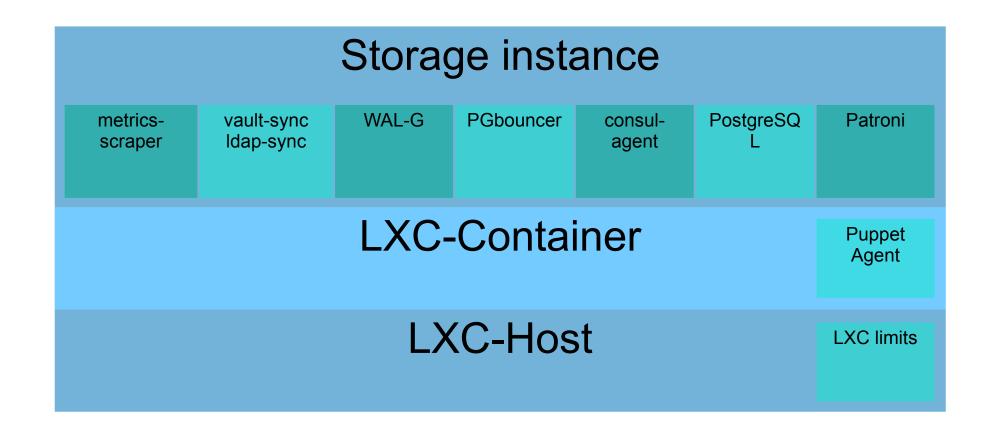
Решение

Пилим свой шедулер для LXC

Пилим свою автоматику над LXC

Immutable Infrastructure

Пилим per-storage DCS



Проблема

Ресурсы планируются вручную

Базы создаются вручнуюСрок — 3-7 дней

Puppet: SPOF/bottleneck

Consul as DCS: SPOF

Решение

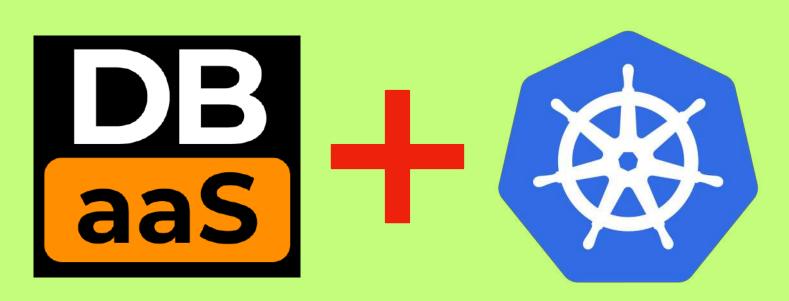
Пилим свой шедулер для LXC

Пилим свою автоматику над LXC

Immutable Infrastructure

Пилим per-storage DCS

PostgreSQL B DBaaS/k8s



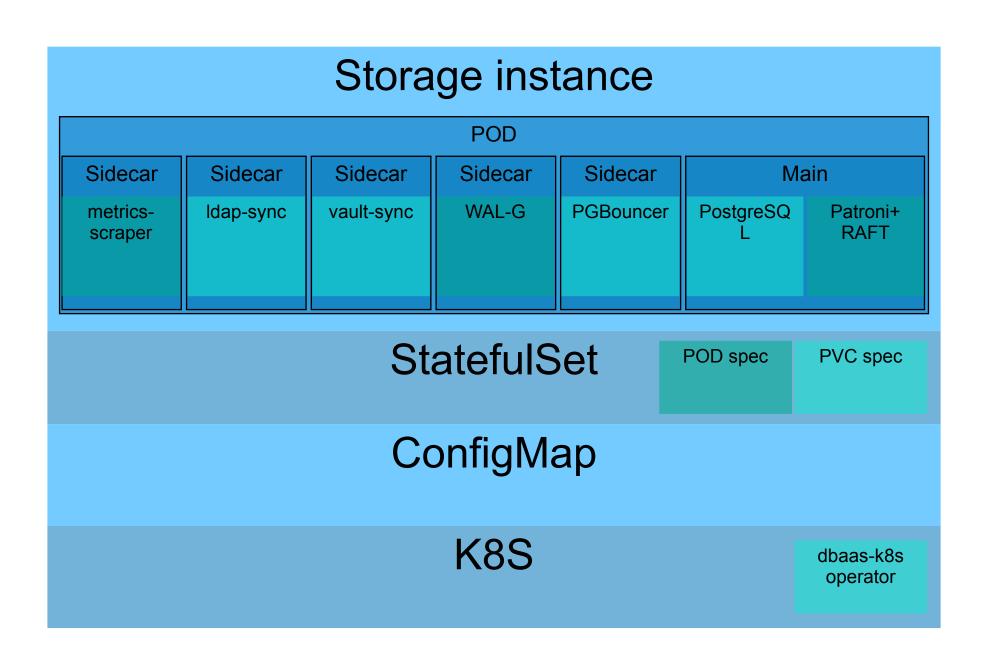
1

2

3

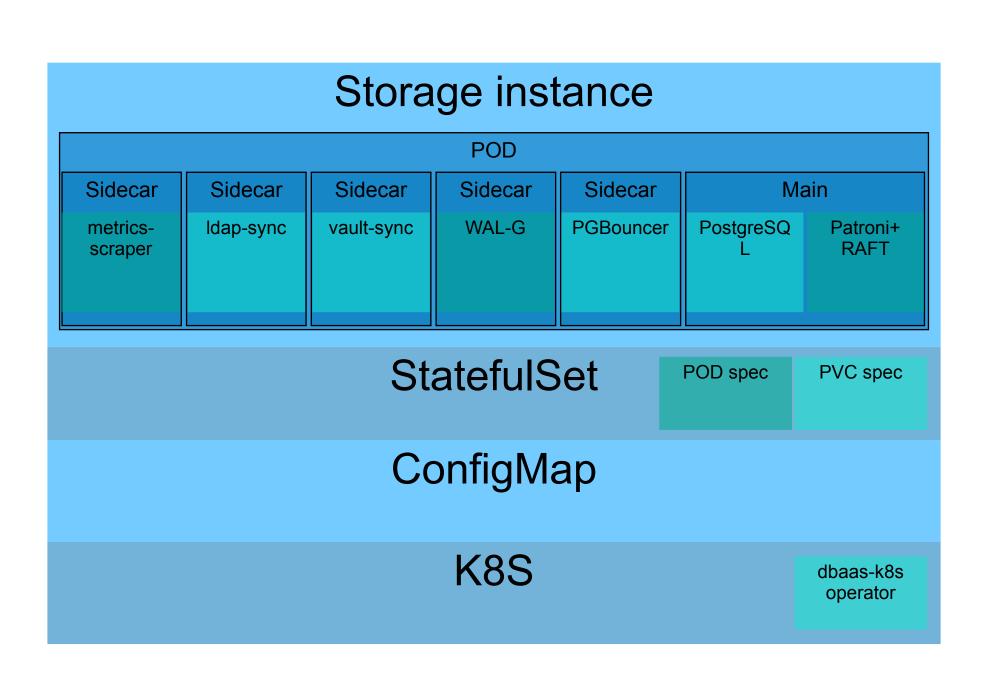
4

PostgreSQL в DBaaS/k8sНаивный подход



Отличия от LXC

- Перешли на Patroni RAFT.
- Все компоненты убрали в сайдкары.
- Получили scheduling «из коробки».
- Получили механику RollingUpdate от DBaaS.

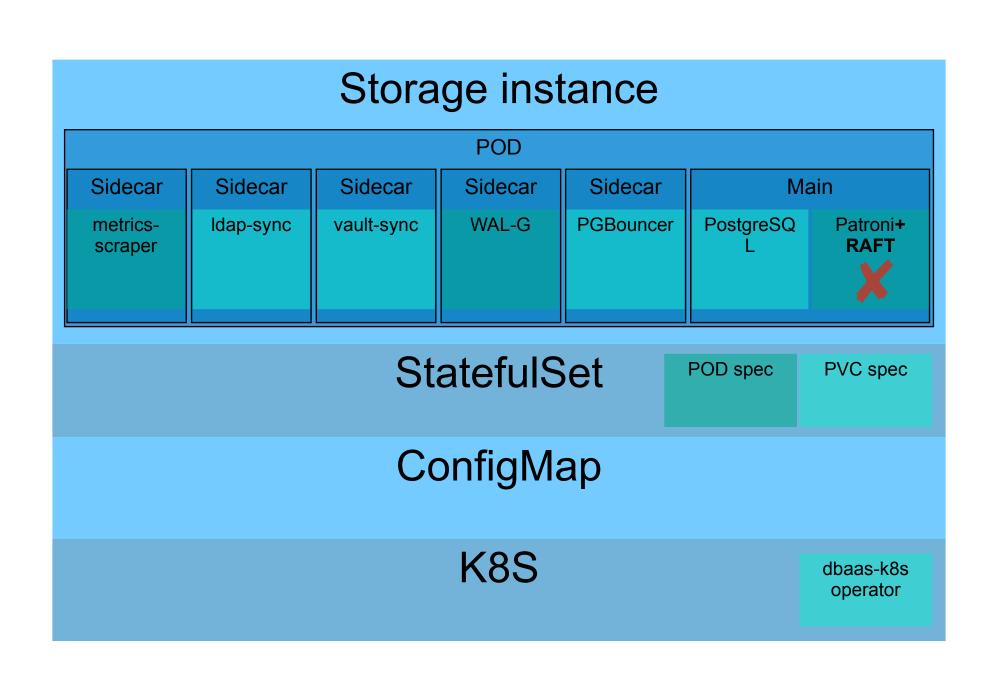


Проблема

Изменения конфиговне доезжают — некому сигналить процессам

Решение





Проблема

Изменения конфиговне доезжают — некому сигналить

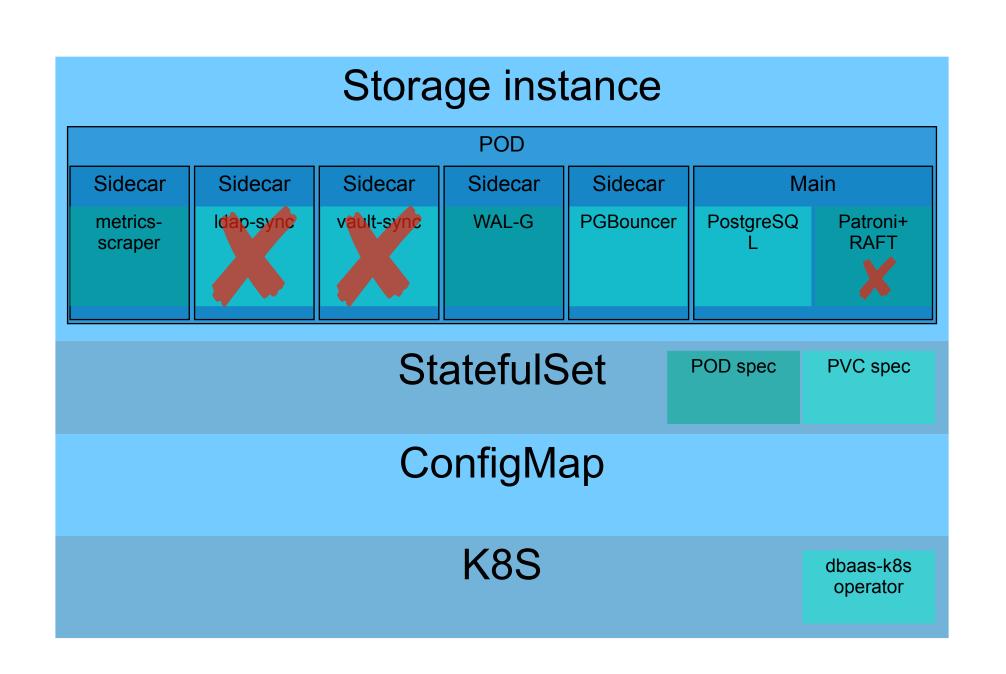


Ратгопі-гаїт оказался сыроват, а чуть позжеи вовсе deprecated

Решение







Проблема

Изменения конфиговне доезжают — некому сигналить



Ратгопі-гаїт оказался сыроват, а чуть позжеи вовсе deprecated



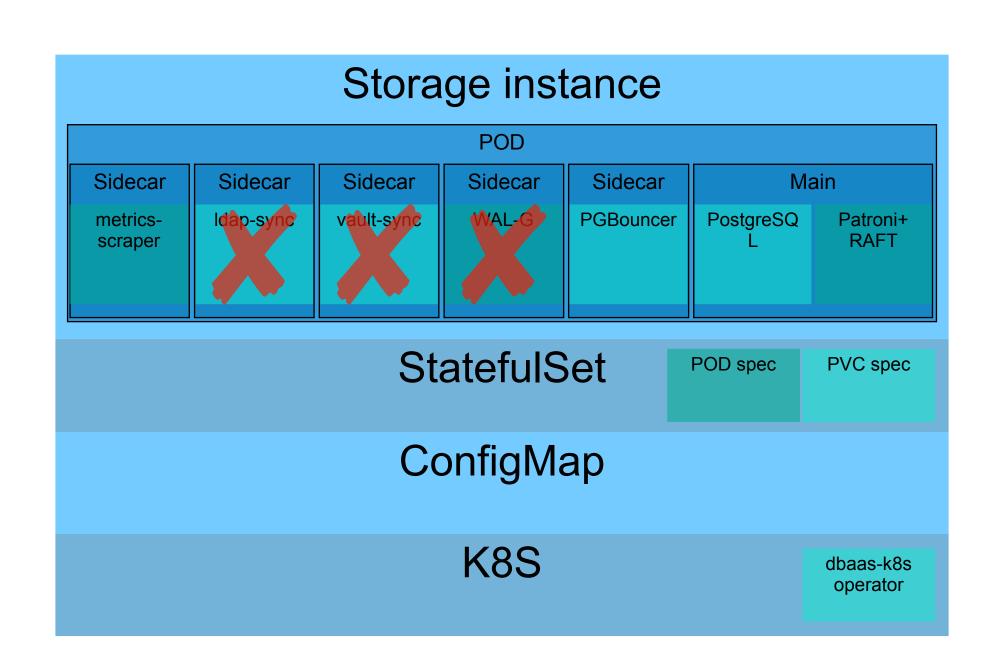
Много компонентов

Решени

PostgreSQL Agent

Встраиваем ETCDв PostgreSQL Agent

Taщим в PostgreSQL Agent всё, что можем



Проблема

Изменения конфиговне доезжают — некому сигналить

Ратгопі-гатт оказался сыроват, а чуть позжеи вовсе deprecated

Много компонентов

Нет бэкапов, т.к. в LXСмы полагались на общийсоnsullock

Решение

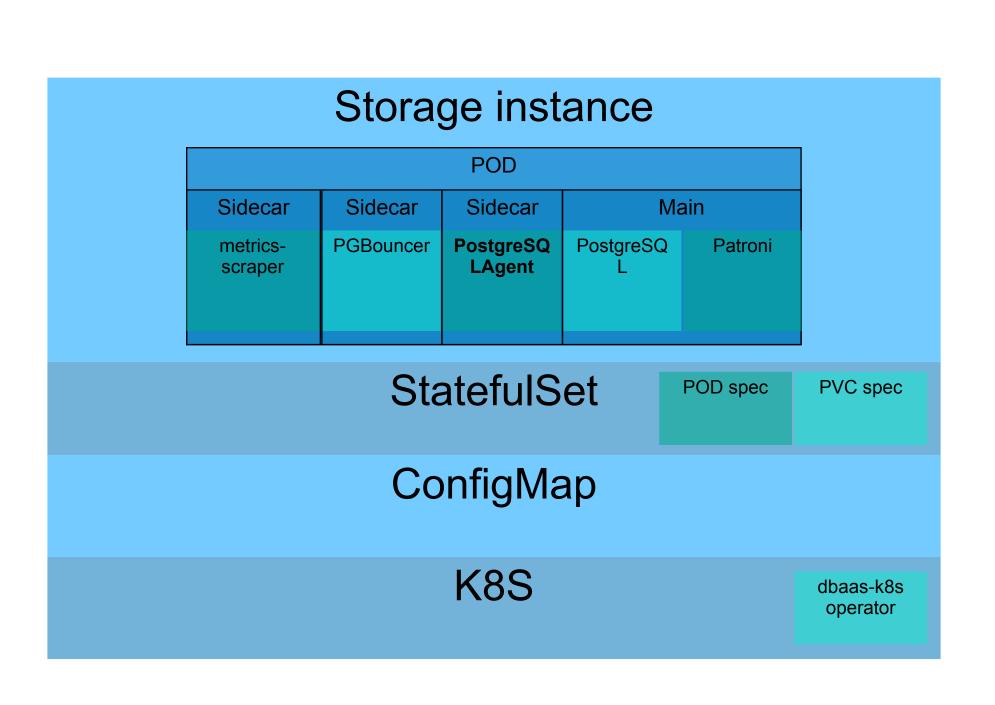
PostgreSQL Agent

Встраиваем ETCDв PostgreSQL Agent

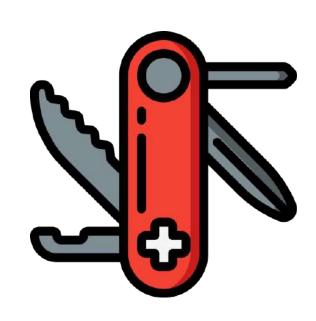
Taщим в PostgreSQL Agent всё, что можем

Тащим в PostgreSQL Agent Backup API

PostgreSQL B DBaaS/k8sPostgreSQL Agent: Swiss Army Knife

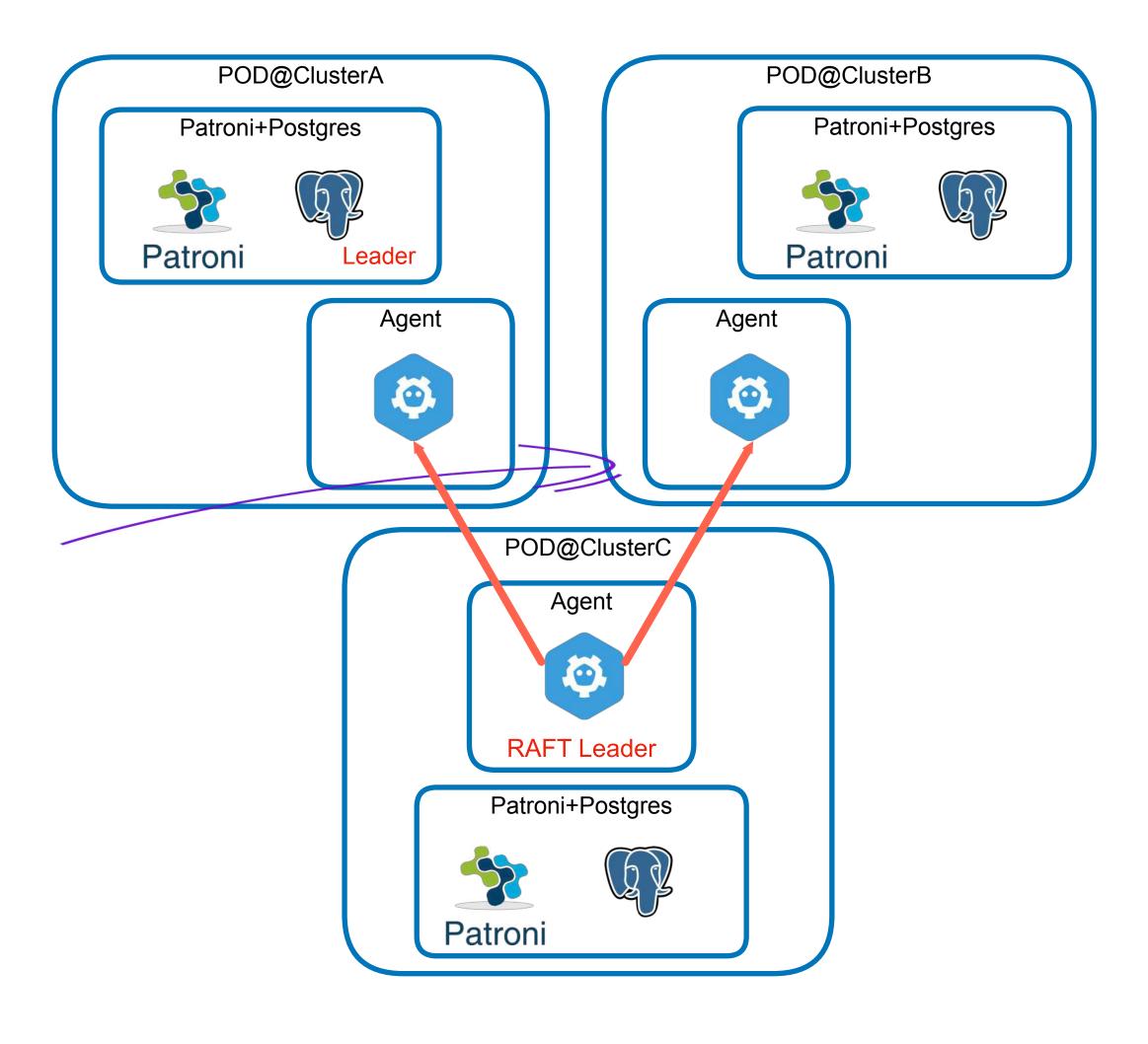


Решаемые задачи

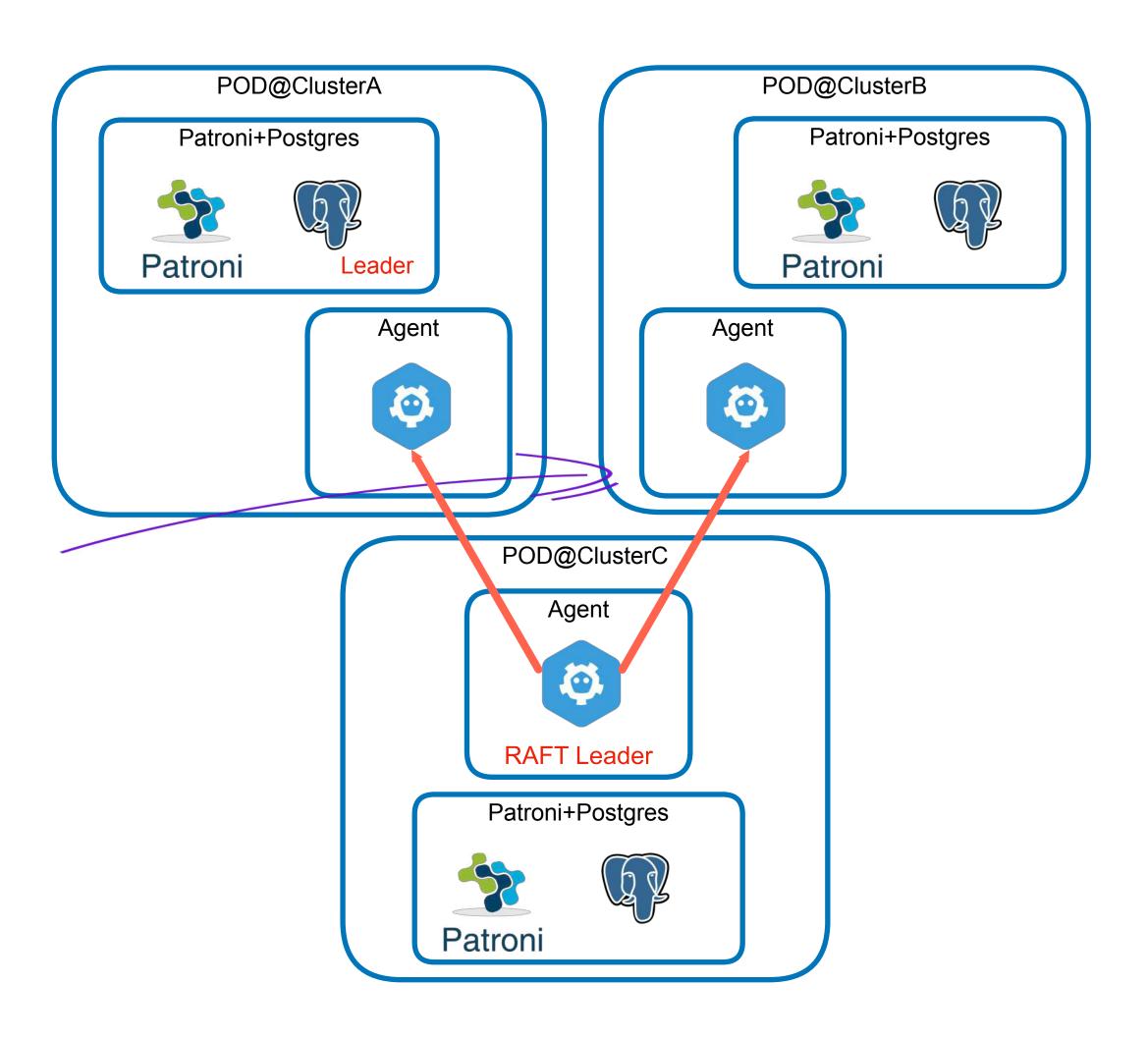


- Синхронизация
- Health-endpoints для k8s и
- контроллеров ETCD для Patroni
- Backup API

PostgreSQL B DBaaS/k8sPostgreSQL Agent: Embedded ETCD



PostgreSQL B DBaaS/k8sPostgreSQL Agent: **Embedded ETCD**



Проблем a



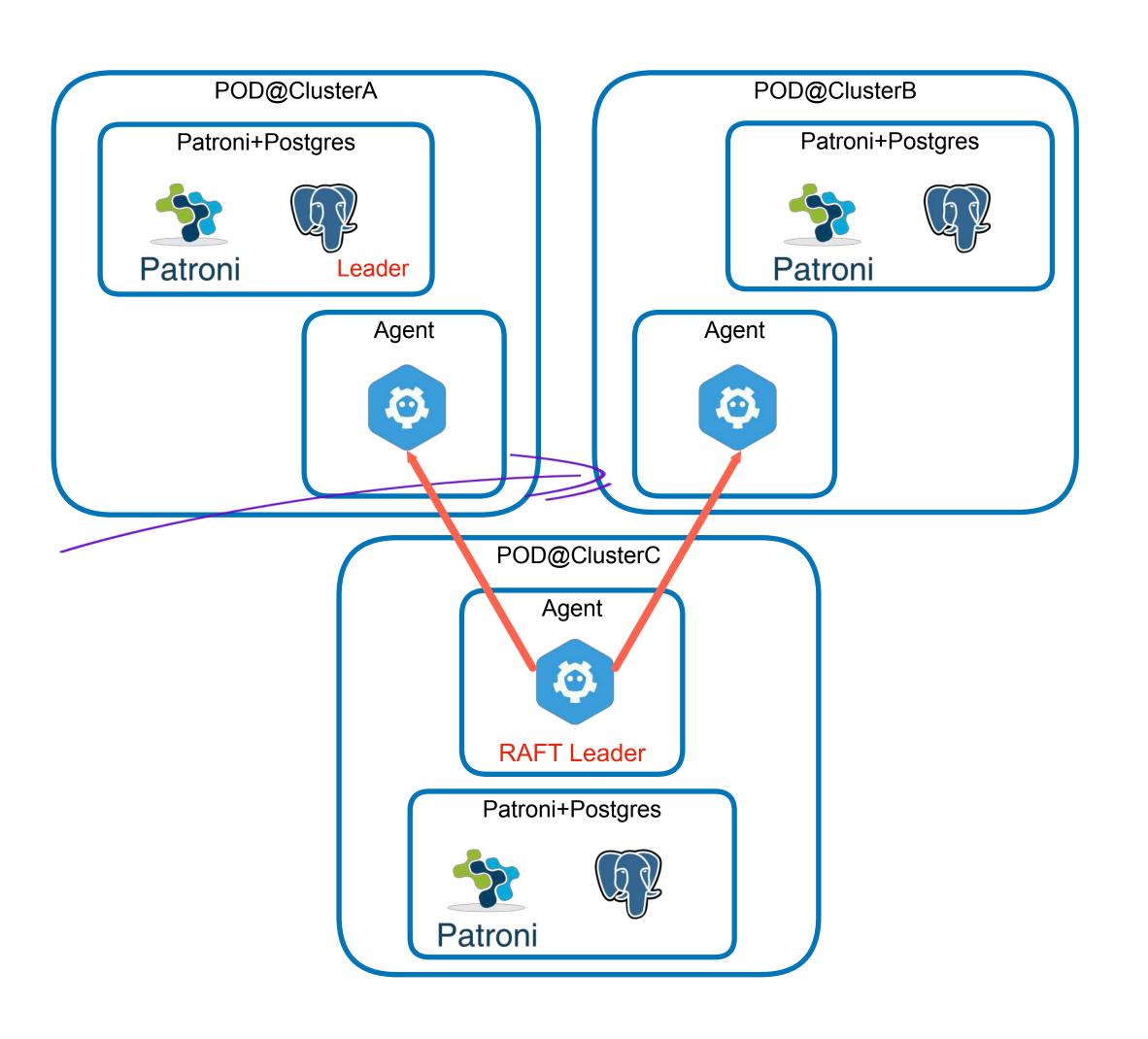
ETCD Bootstrap

Решени e



Реализуем на основеPostgreSQL Agent API

PostgreSQL B DBaaS/k8sPostgreSQL Agent: **Embedded ETCD**



Проблем a



ETCD Bootstrap



He BCE Readiness Probes полезны

Решени e

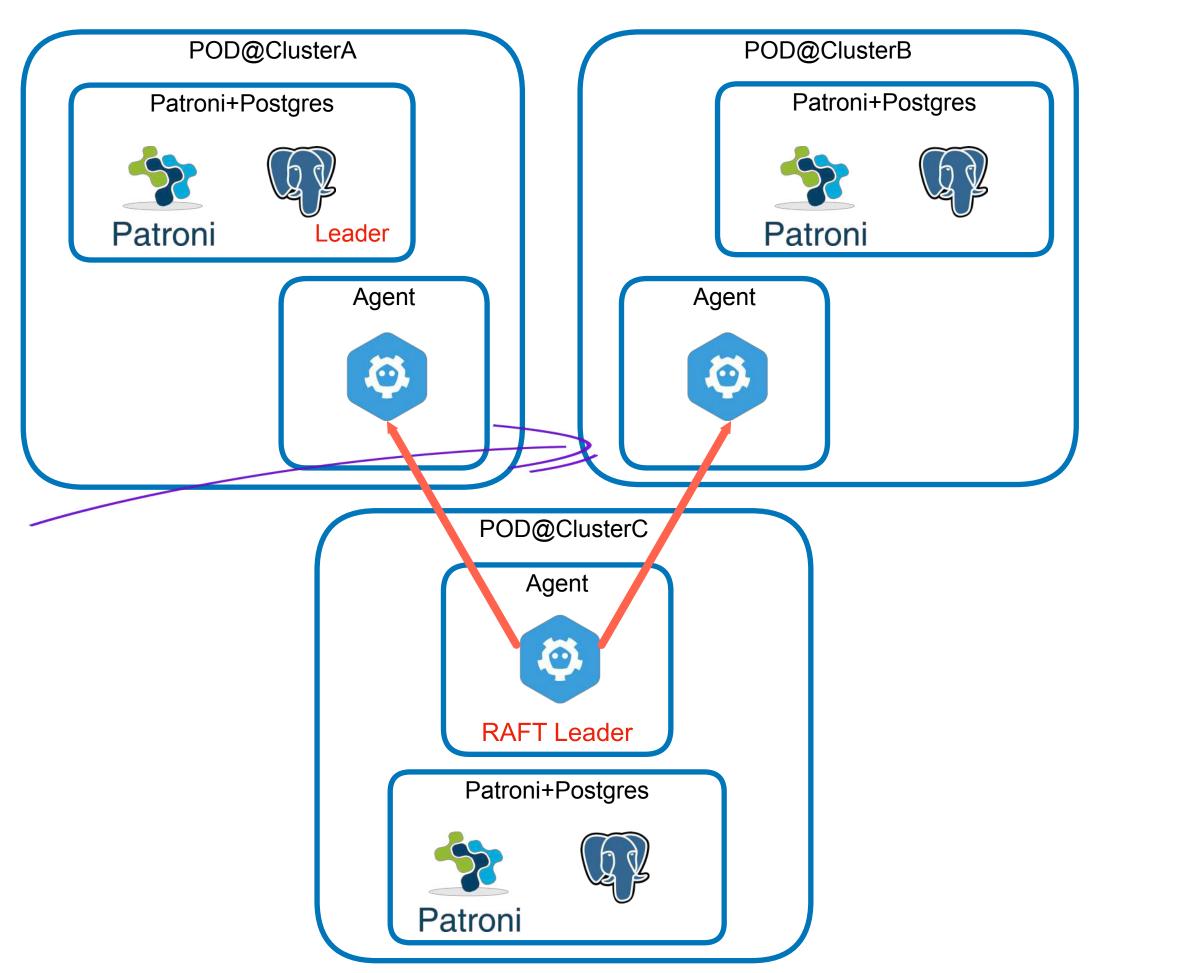


Реализуем на основеРоstgreSQL Agent API



Оставляем только дляPostgreSQL Agent,остальным — /bin/true

PostgreSQL B DBaaS/k8sPostgreSQL Agent: Embedded ETCD



Проблема



ETCD Bootstrap

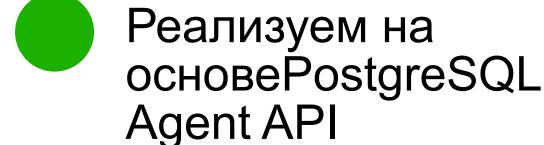


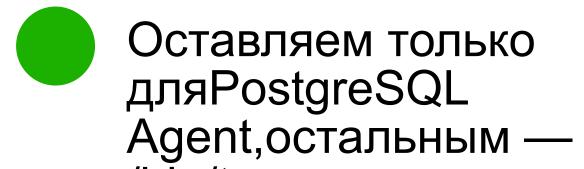
He все Readiness Probes полезны



Сложная конструкция

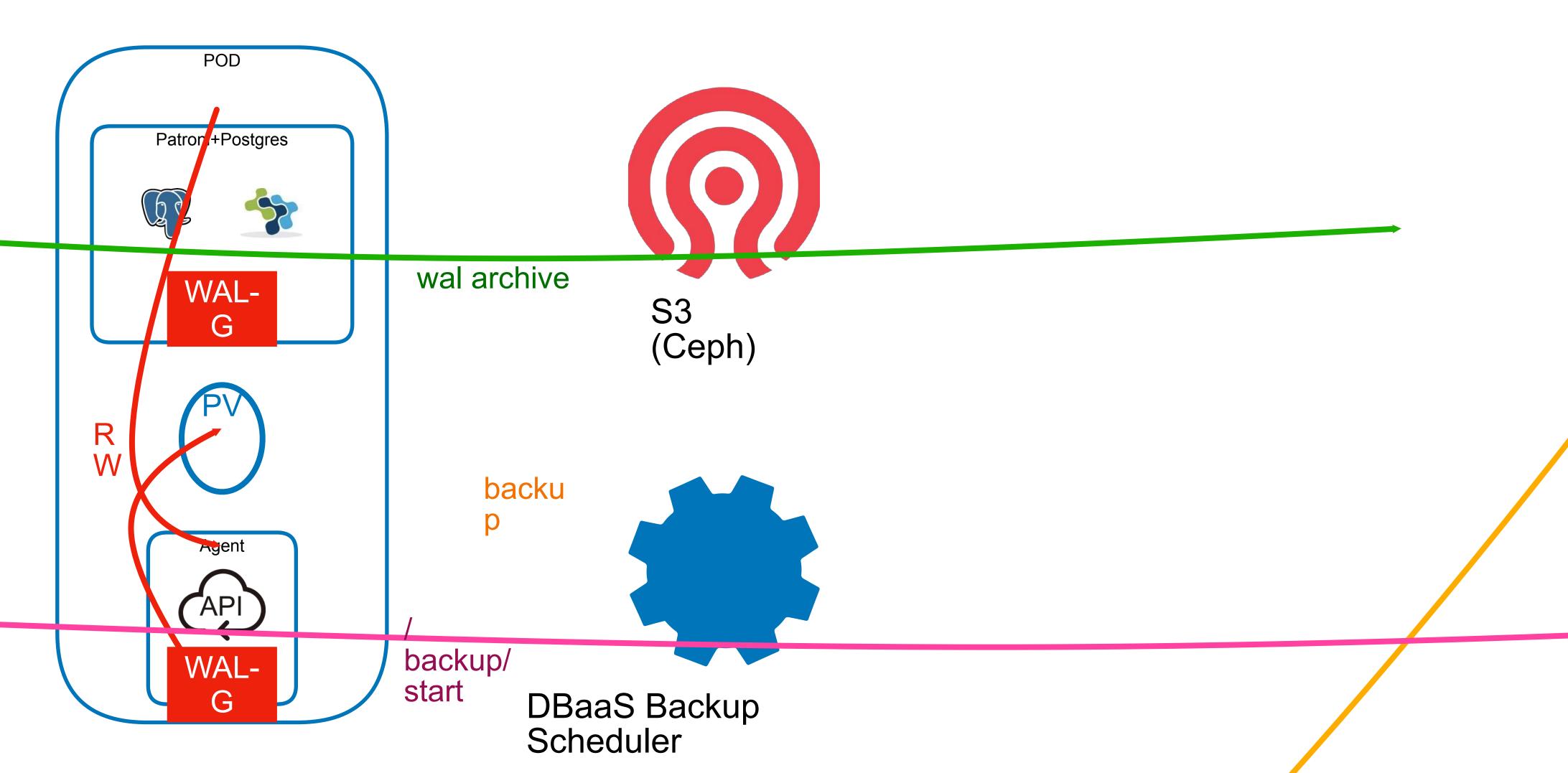
Решени



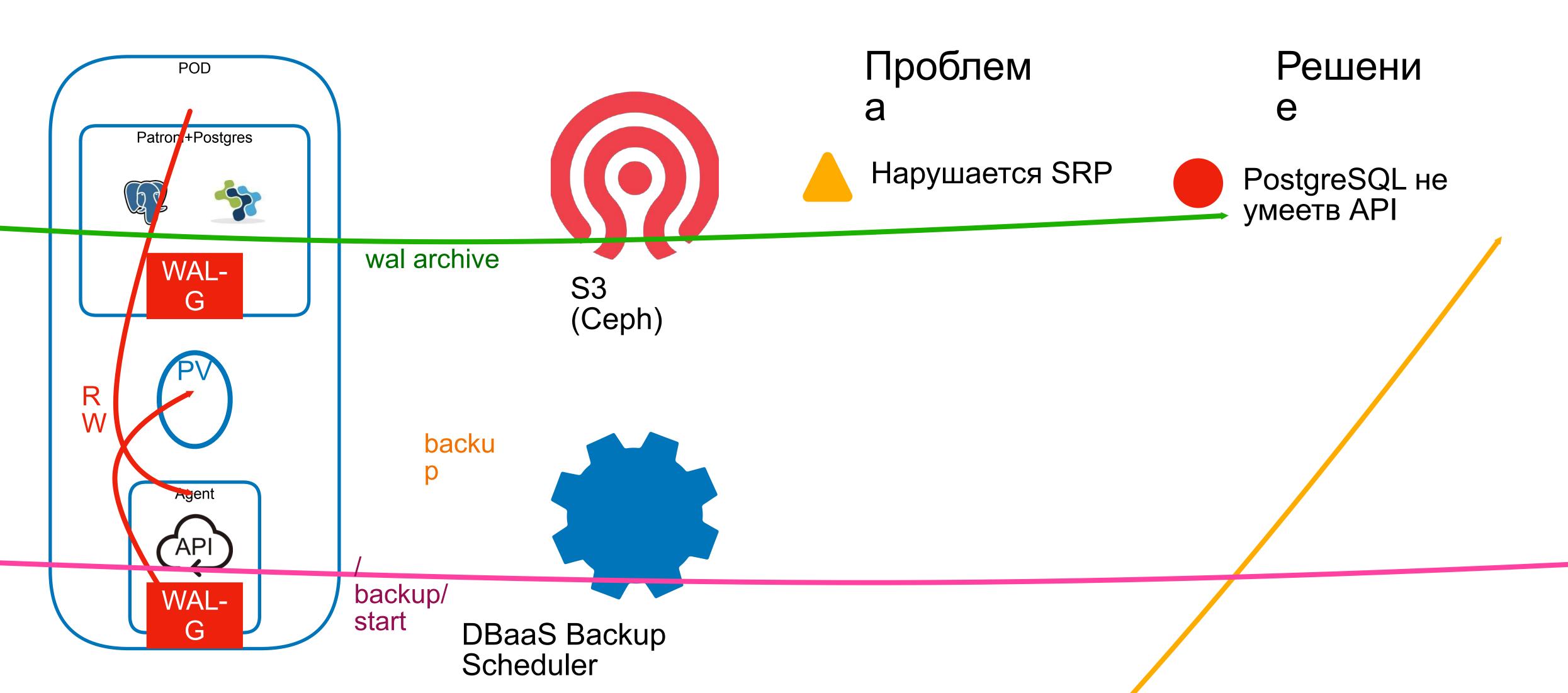


/bin/true Ничего прощене придумали

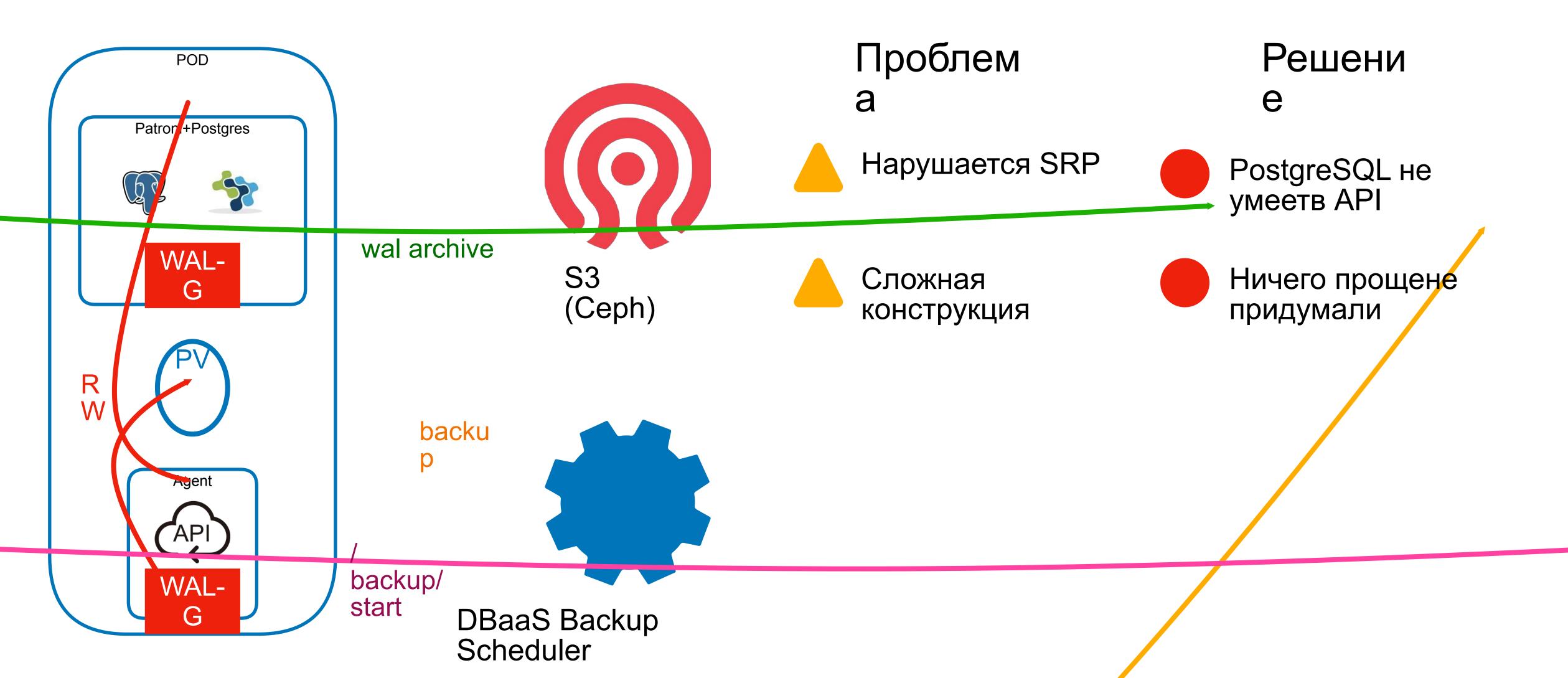
PostgreSQL B DBaaS/k8sPostgreSQL Agent: Backup API



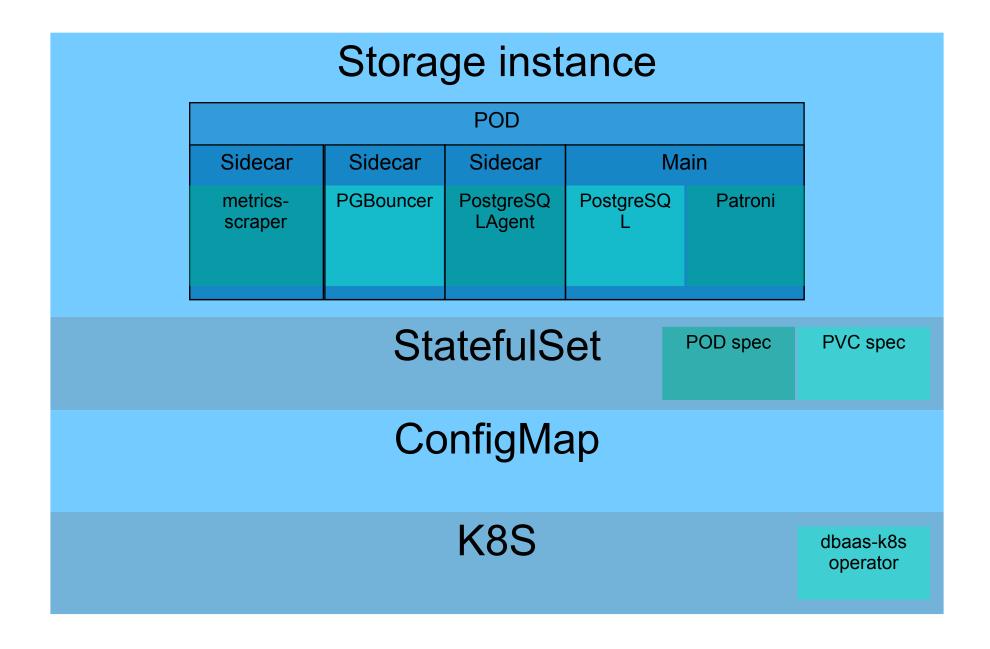
PostgreSQL в DBaaS/k8sPostgreSQL Agent: Backup API



PostgreSQL в DBaaS/k8sPostgreSQL Agent: Backup API



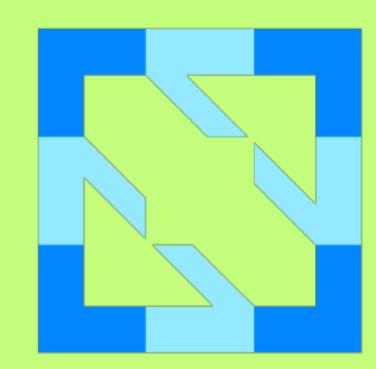
PostgreSQL в DBaaS/k8sЧто в итоге?



PostgreSQL на платформе DBaaS/k8s Автоматизированн

- Автоматизированн ый жизненный цикл
- Время создания базы 2-4 минуты
- Много разработки, но меньше
- операционки Достаточно сложная
- реализация PostgreSQL всё ещё плохо масштабируется

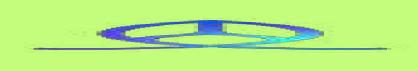
Что там с Cloud Native?



1

3

4

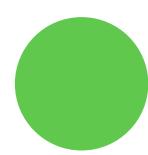


Что там с Cloud Native?CockroachDB «из коробки»



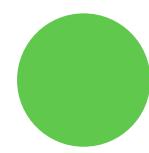
High Availability & Fault Tolerance

Native: RAFT, Replication



ObservabilityNative: Prometheus-

exporter, Json-лог



DurabilityNative: Backup to S3



ScalabilityNative: вертикальное (с оговорками)Native: Горизонтальное

Что там у Cloud Native?Адаптация CockroachDB на платформе DBaaS/k8s



Андрей Лярский

Инженер в SQL/DBA





@oxumorron

He все базы одинаковок8s-friendly.