



ОНЛАЙН-ОБРАЗОВАНИЕ

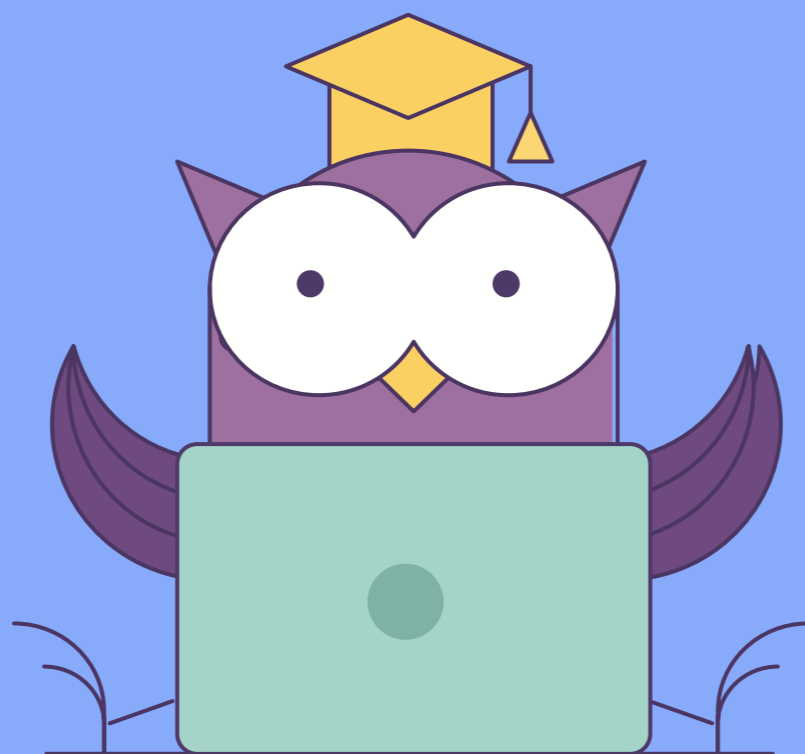
# PostgreSQL Cluster

Курс «Администратор Linux»

Занятие № 36



# Меня хорошо слышно && видно?



Напишите в чат, если есть проблемы!

Ставьте  + если все хорошо  
Ставьте  - если есть проблемы

Какие у нас есть варианты

Как делать/не делать Failover

Архитектура Patroni

Создание кластера

Как менять конфигурацию кластера

Немного про ETCD

Перенаправление клиентов на Master

Создание реплик и их реинициализация

- Встроенные решения
- Patroni
- Stolon:
  - Проксирует все запросы в мастер ноду. Нельзя давать нагрузку на реплики
  - Мастер выбирается самостоятельно при switchover-е
- permgr:
  - Нет фэнсинга из коробки (защита от двойного мастера)
  - Нет нужды в DCS - на мой взгляд это минус

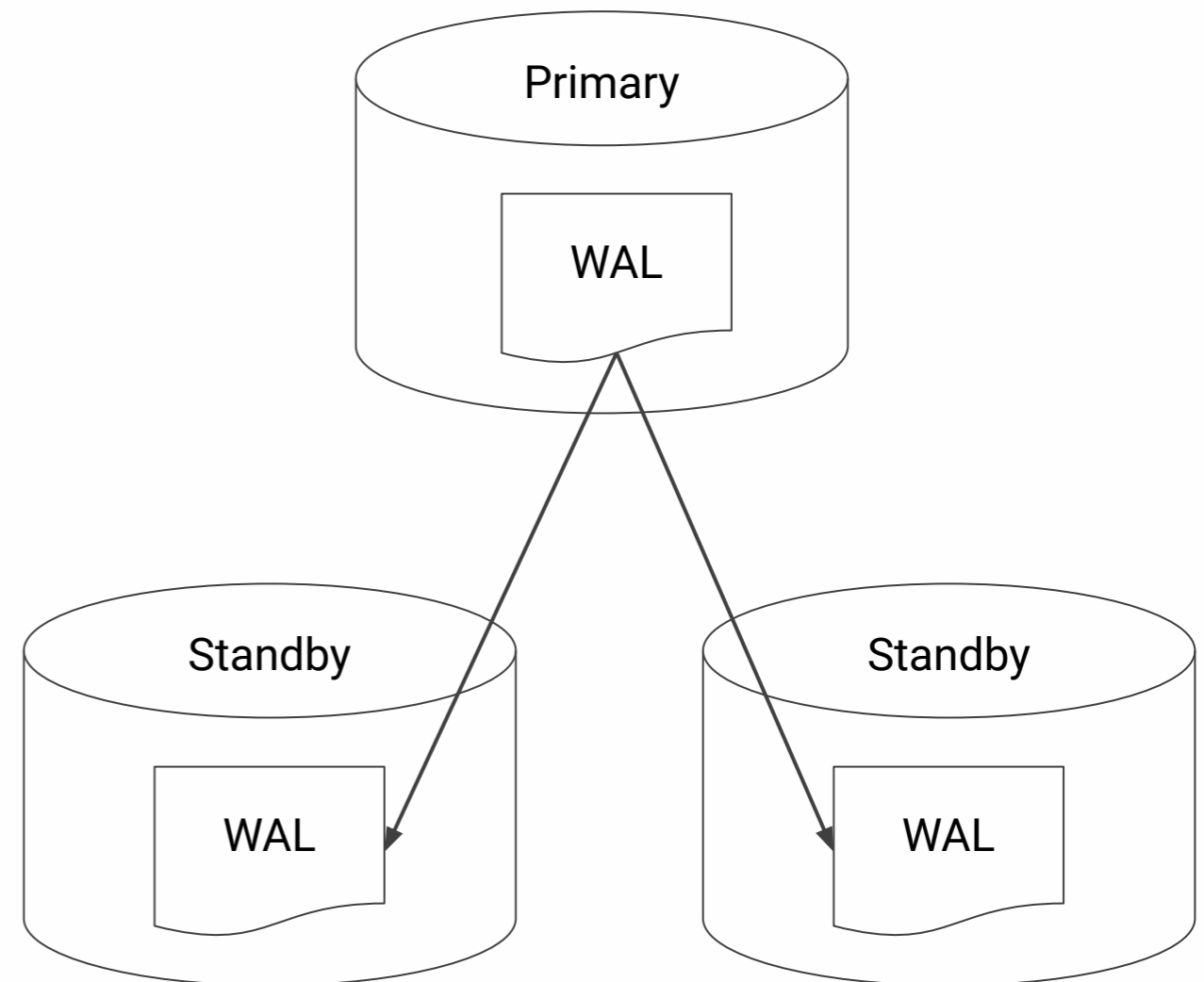
- Распределенное хранилище
  - NFS
  - DRBD
  - ISCSI (+ LVM)
- Мульти-мастер
  - BDR, Bucardo
- Логическая репликация
  - pglogical, slony, встроенная фича в postgresql 10
- Физическая репликация
  - В postgresql начиная с 9.0

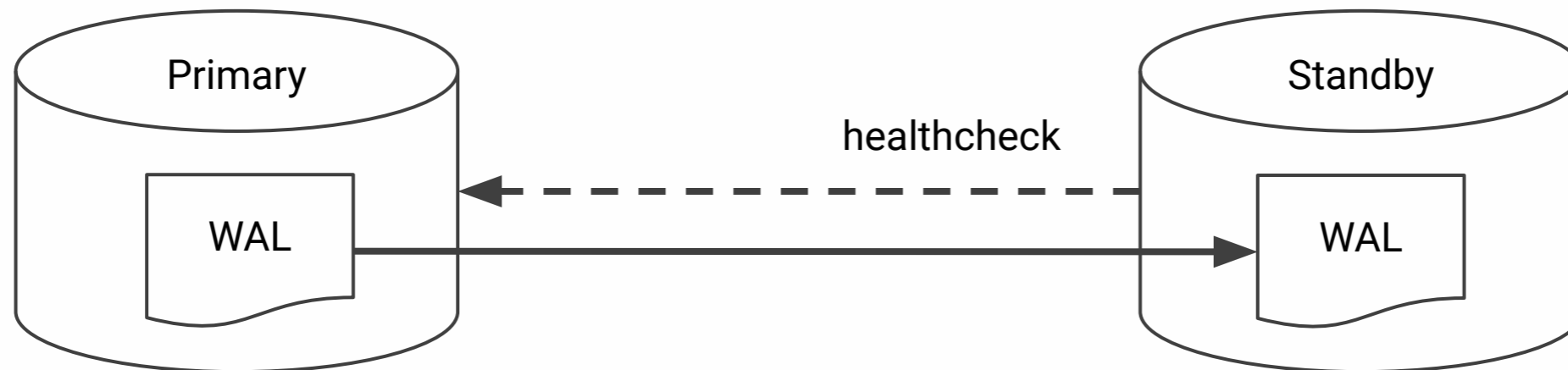
## Плюсы:

- Встроенная фича
- Минимальная задержка
- Идентичные копии

## Минусы:

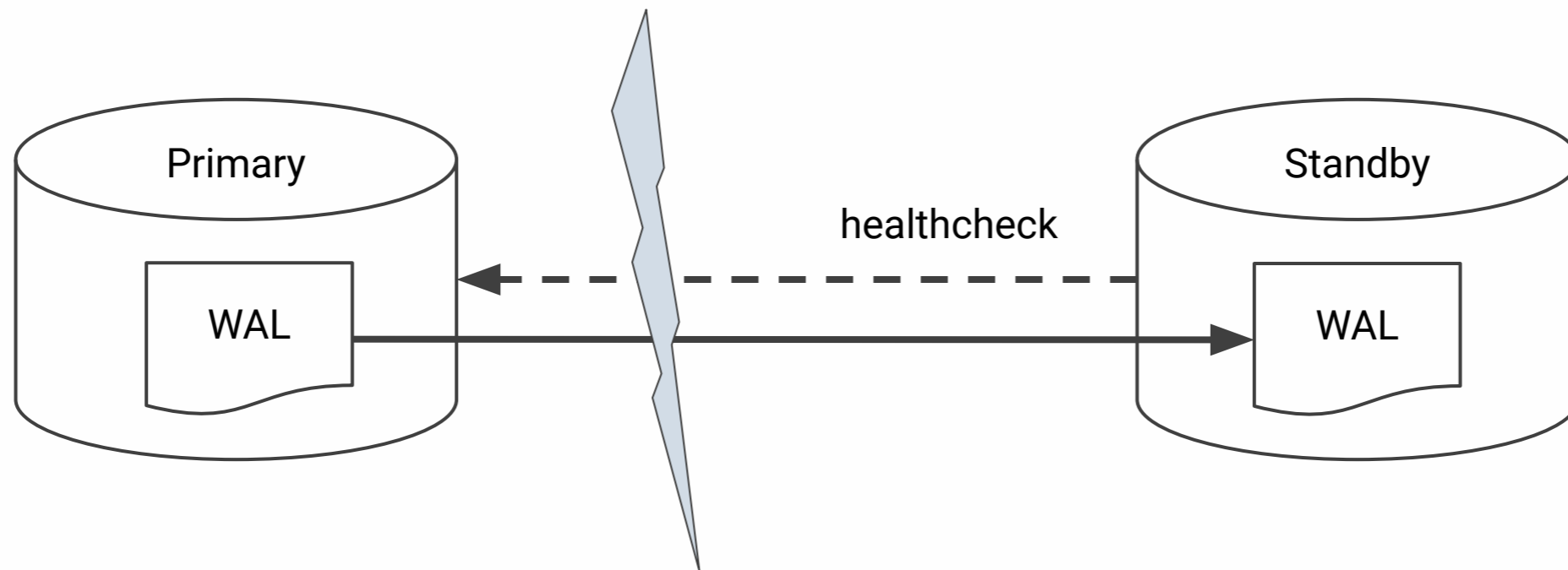
- Нужны одинаковый мажорные версии
- Нет автоматического failover

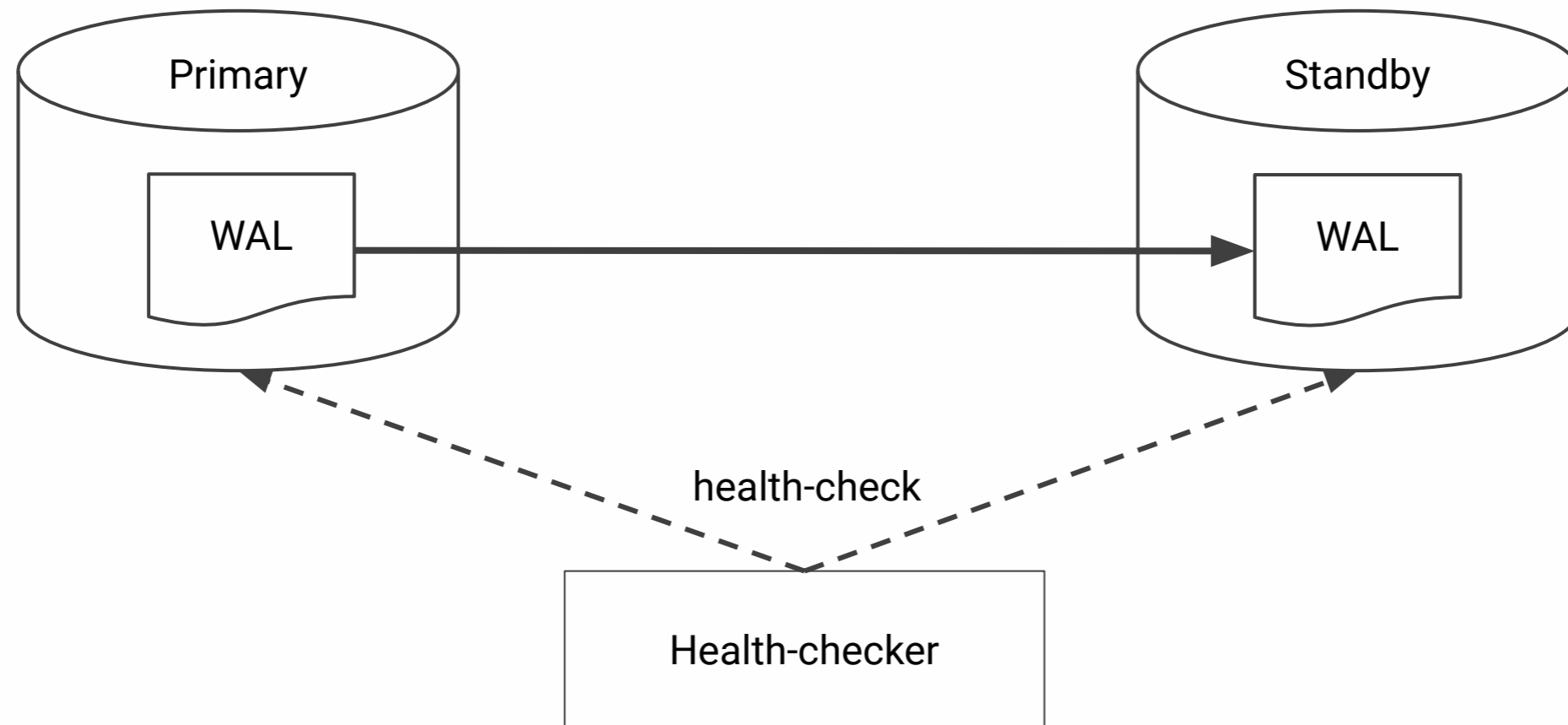


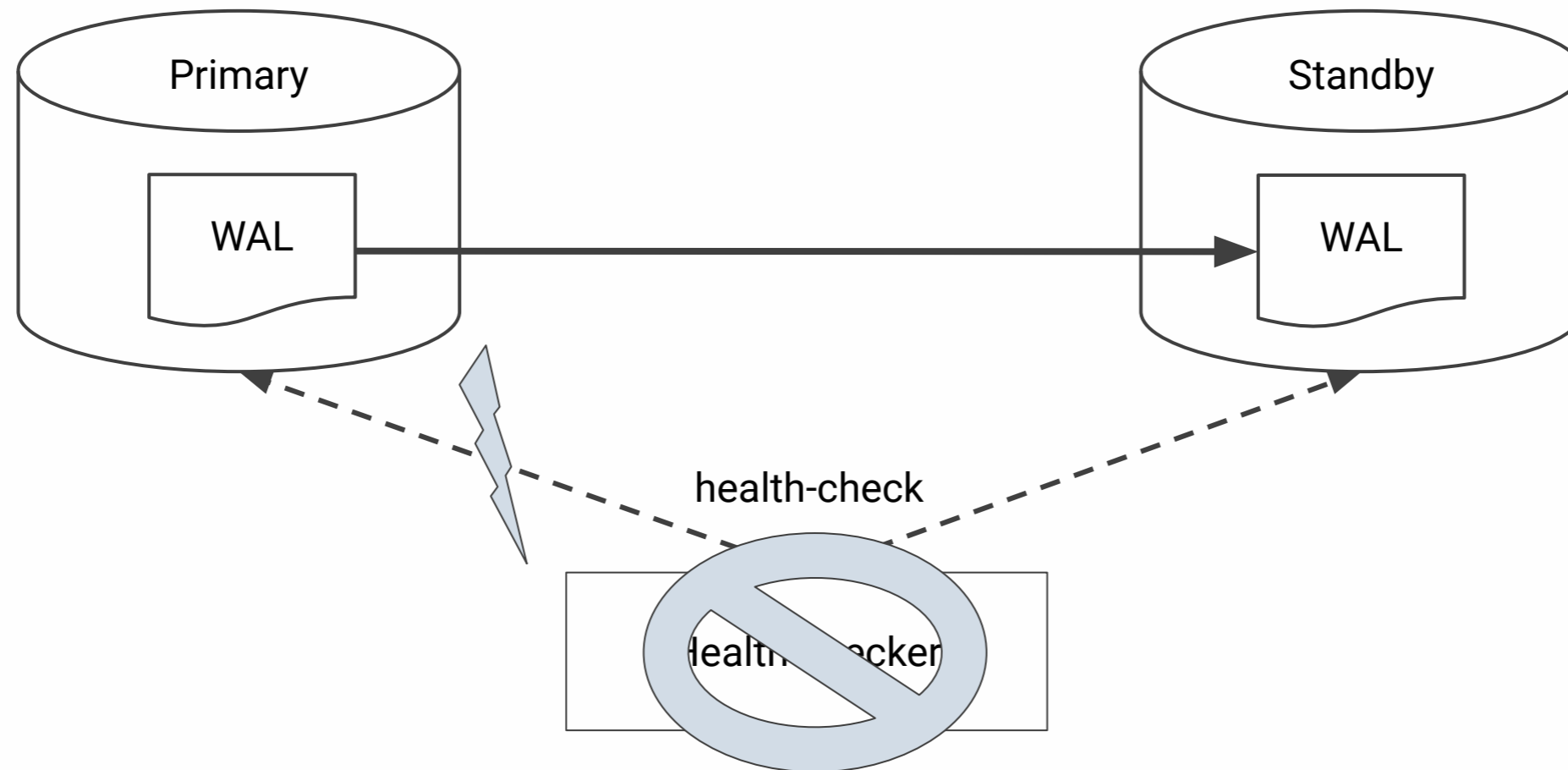


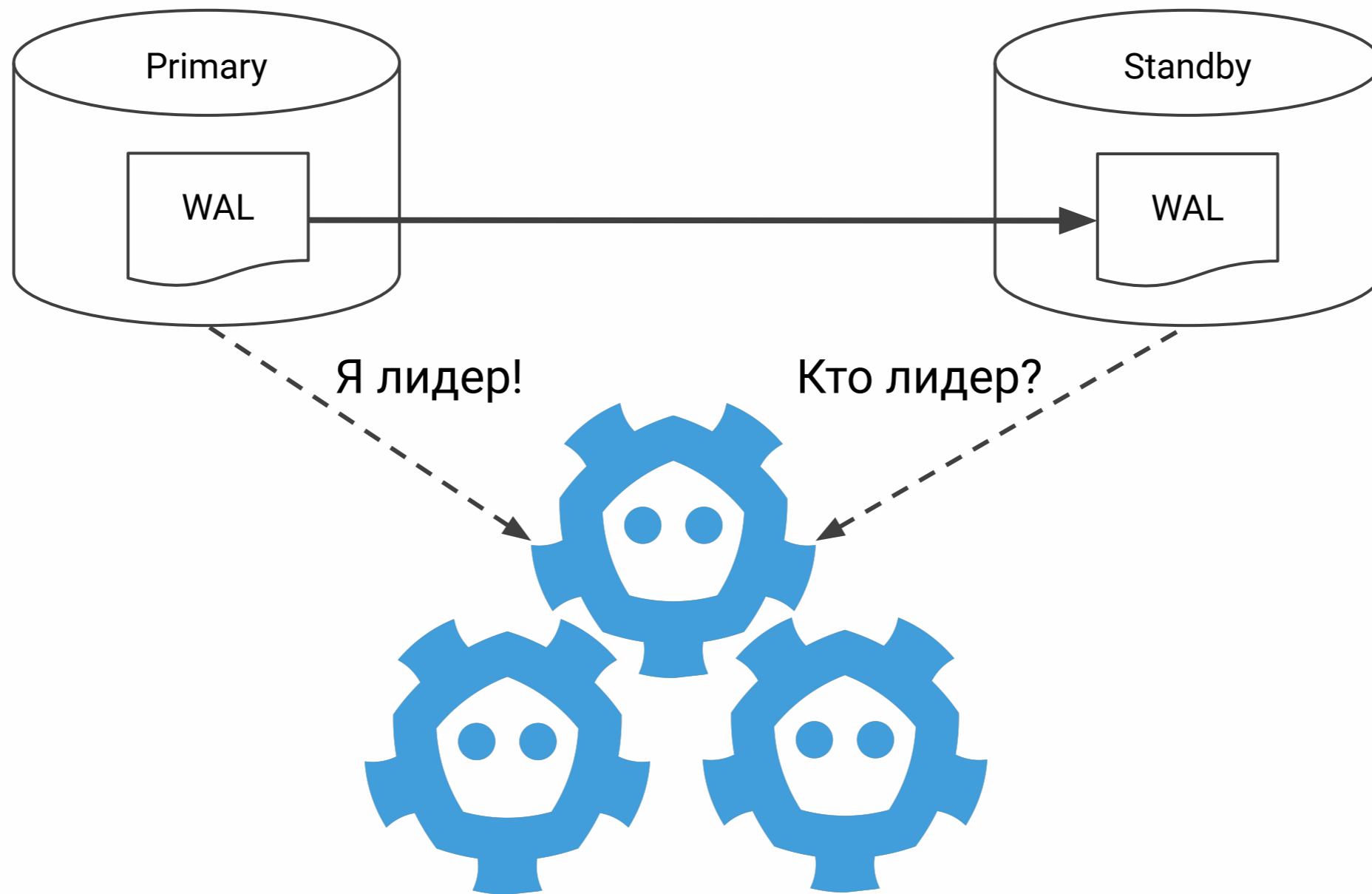
Запускаем healthcheck со стэндбая и при отрицательном ответе продвигаем (promote) его до Мастера

Split Brain!







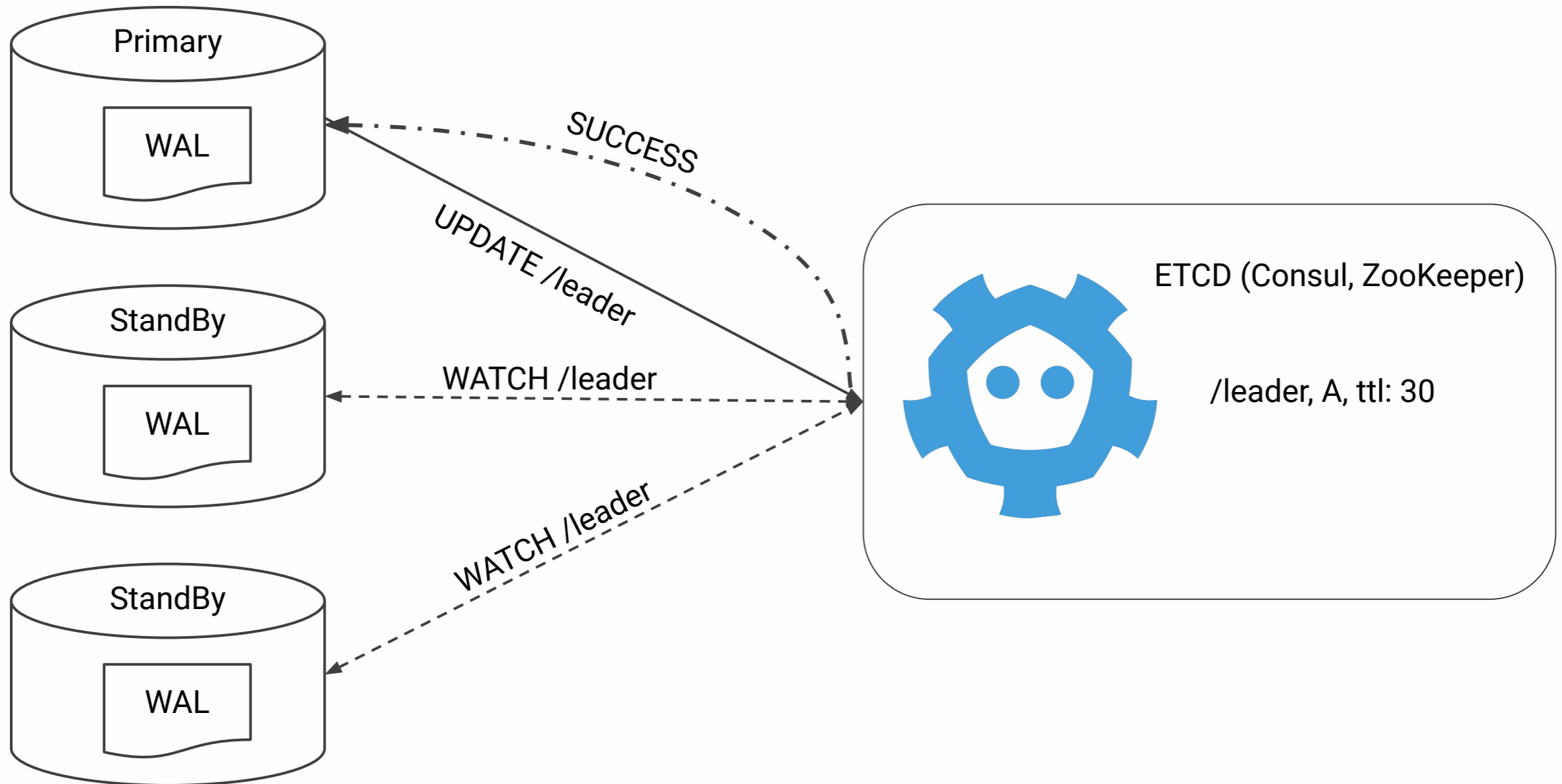


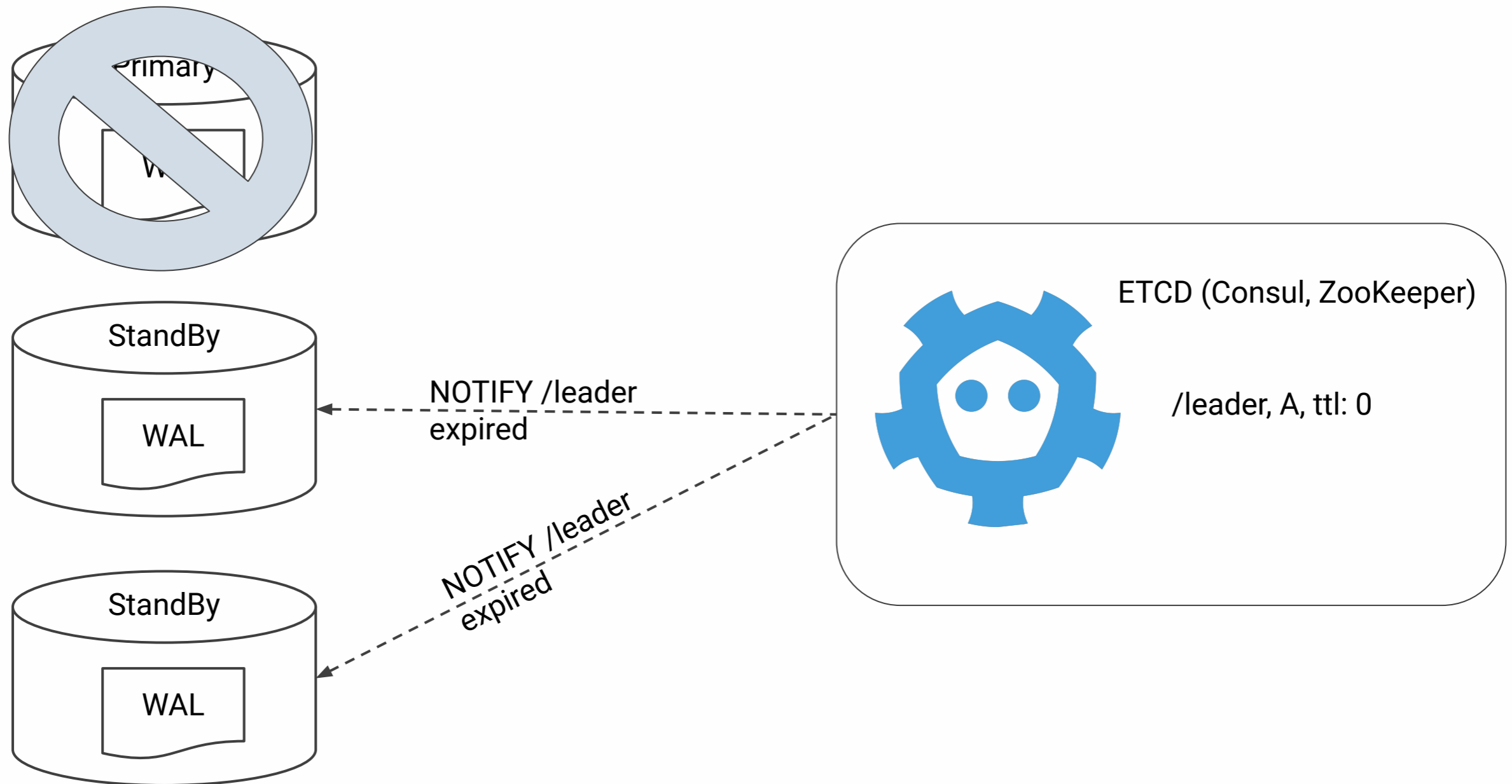
У постгреса нет  
какого либо решения  
по автоматическому  
фейловеру из коробки

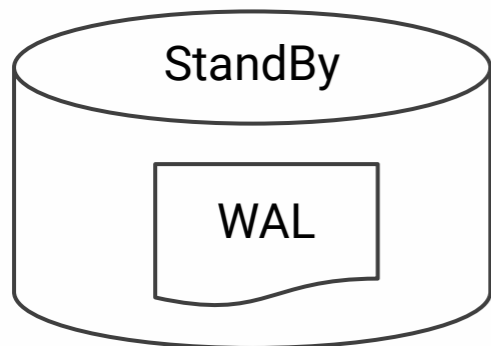
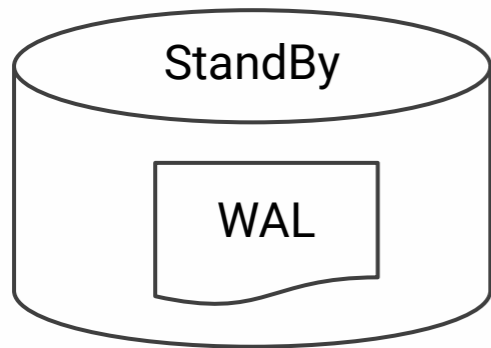
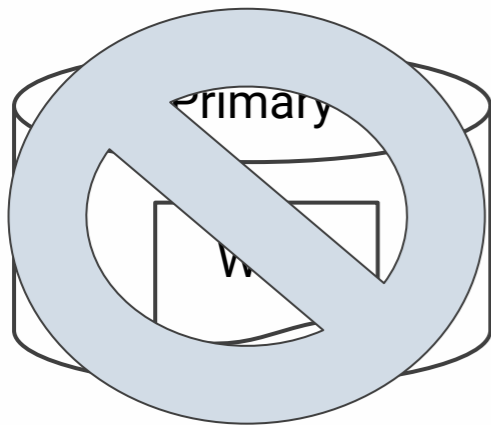


- etcd (или Consul, Zookeeper) хранят инфу о том, кто сейчас лидер
- DCS хранит конфигурацию кластера
- помогает решить проблему с партиционированием сети
- убивает старые клиентские коннекты
- STONITH
- Неплохо бы иметь watchdog (Например, Nomad)

- PostgreSQL не умеет взаимодействовать с etcd
- Демон будет запущен рядом с PostgreSQL
- Демон умеет взаимодействовать с etcd
- Демон принимает решение promotion/demotion








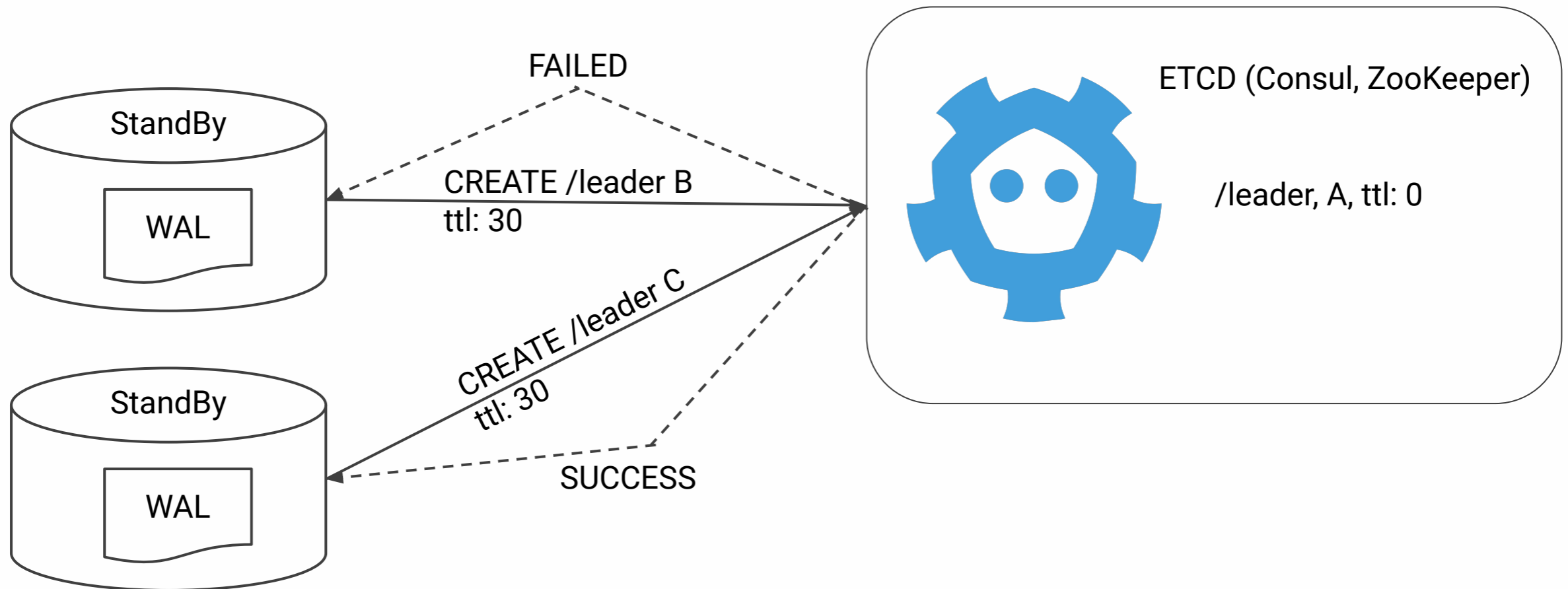
Node B:  
GET hostA:patroni -> Timeout  
GET hostB:patroni -> wal\_position: 200  
GET hostC:patroni -> wal\_position: 100

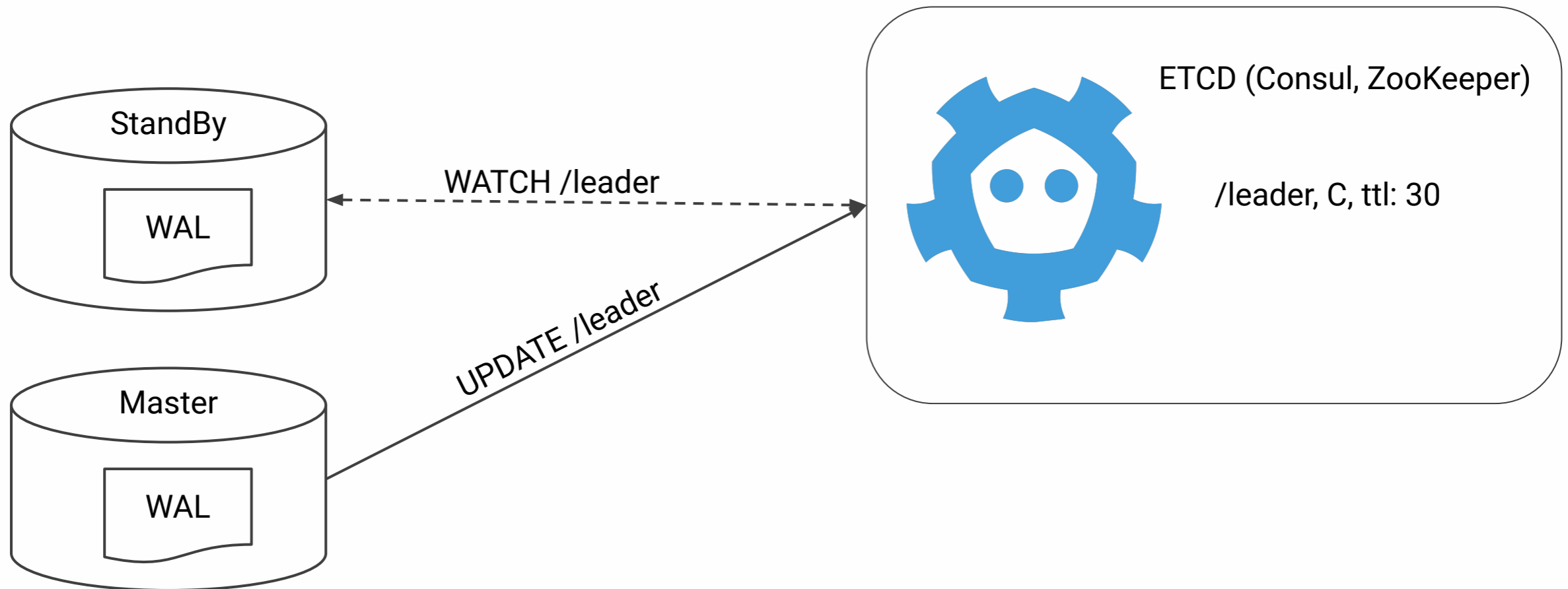
Node C:  
GET hostA:patroni -> Timeout  
GET hostB:patroni -> wal\_position: 200  
GET hostC:patroni -> wal\_position: 100



ETCD (Consul, ZooKeeper)

/leader, A, ttl: 0





ETCD

HAProxy

pgbouncer (pgpool)

Нода Postgresql:

- Postgresql{9,10,11}-server
- pip install patroni и зависимости
- конфигурационный файл patroni.yml
- Дата директория - с правами для пользователя postgres

```
patroni /etc/patroni.yml OR systemctl start patroni.service
```

```
INFO: Selected new etcd server http://10.128.0.48:2379
```

```
INFO: Lock owner: None; I am pg01
```

```
trying to bootstrap a new cluster
```

```
LOG: listening on IPv4 address "10.128.0.49", port 5432
```

```
INFO: establishing a new patroni connection to the postgres cluster
```

```
INFO: Lock owner: pg01; I am pg01
```

```
INFO: no action. i am the leader with the lock
```

- patronictl - утилита для управления кластером
- patronictl -c /etc/patroni.yml list

```
[root@pg01 ~]# patronictl -c /etc/patroni.yml list
```

Cluster	Member	Host	Role	State	TL	Lag in MB
postgres	pg01	10.128.0.47		running	7	0.0
postgres	pg02	10.128.0.46	Leader	running	7	0.0
postgres	pg03	10.128.0.45		running	7	0.0

- **PATRONI\_CONFIGURATION** - путь до конфигурационного файла
- **PATRONI\_NAME** - имя текущей ноды. Должно быть уникально в контексте кластера
- **PATRONI\_SCOPE** - имя кластера
- **PATRONI\_LOG\_\*** - все что связано с логами

```
export PATRONI_CONSUL_HOST='192.168.11.100:8500'
```

```
export PATRONI_CONSUL_TOKEN=aabbccddeeff
```

```
# systemctl stop patroni
```

- 10 секунд по умолчанию на истечение ключа в ETCD
- После чего Patroni стучится на каждую ноду в кластере и спрашивает, не мастер ли ты, проверяет WAL логи, насколько близки они к мастеру. В итоге если WAL логи у всех одинаковые то, промоутится следующий по порядку
- Опрос нод идёт параллельно

etcd - это распределенная база данных для хранения важных/критичных данных в виде ключ-значение. В данном случае она нужна для того чтобы хранить данные о нодах в postgres кластере, а именно: мастер ключ кластера (ноды принадлежащие к одному кластеру, не смогут подключиться и работать в другом), информации о том, кто сейчас лидер в кластере, его статусах и статусах всех реплик.

```
yum install etcd -y
```

```
ETCD_DATA_DIR="/var/lib/etcd/default.etcd"
```

```
ETCD_LISTEN_PEER_URLS="http://10.128.0.48:2380"
```

```
ETCD_LISTEN_CLIENT_URLS="http://localhost:2379,http://10.128.0.48:2379"
```

```
ETCD_NAME="etcd0"
```

```
ETCD_INITIAL_ADVERTISE_PEER_URLS="http://10.128.0.48:2380"
```

```
ETCD_ADVERTISE_CLIENT_URLS="http://10.128.0.48:2379"
```

```
ETCD_INITIAL_CLUSTER="etcd0=http://10.128.0.48:2380"
```

```
ETCD_INITIAL_CLUSTER_TOKEN="cluster1" ETCD_INITIAL_CLUSTER_STATE="new"
```

```
[root@etcd ~]# etcdctl ls --recursive --sort -p /service/postgres
```

```
[root@etcd ~]# etcdctl get /service/postgres/leader
```

```
pg03
```

```
[root@etcd ~]# etcdctl get /service/postgres/members/pg01
```

```
{"conn_url":"postgres://10.128.0.47:5432/postgres","api_url":"http://10.128.0.47:8008/  
patroni","timeline":8,"state":"running","role":"replica","xlog_location":150997240}
```

```
[root@pg02]# patronictl -c /etc/patroni.yml edit-config
```

```
---
```

```
+++
```

```
@@ -2,5 +2,6 @@
```

```
maximum_lag_on_failover: 1048576
```

```
postgresql:
```

```
  use_pg_rewind: true
```

```
+ parameters:
```

```
+ maintenance_work_mem: 256MB
```

```
retry_timeout: 10
```

```
ttl: 30
```

```
Apply these changes? [y/N]:
```

```
Mar 21 09:59:50 pg03 patroni: 2019-03-21 09:59:50,666 INFO: Changed maintenance_work_mem from 65536 to 256MB
```

```
Mar 21 09:59:50 pg03 patroni: 2019-03-21 09:59:50,667 INFO: PostgreSQL configuration items changed, reloading configuration.
```

Попробуем поменять параметр требующий перезагрузки: **max\_connections**

```
Mar 21 10:04:10 pg03 patroni: 2019-03-21 10:04:10,665 INFO: Changed  
max_connections from 100 to 200 (restart required)
```

```
[root@pg02] http http://10.128.0.45:8008
```

Ручной Switchover:

```
patronictl -c /etc/patroni.yml switchover --master pg03 --candidate pg01
```

Что делать если нужно поменять конфигурацию PostgreSQL только локально.

- etcd
- patroni.yml
- postgresql.base.conf
- ALTER SYSTEM SET - имеет наивысший приоритет

Некоторые параметры, такие как: `max_connections`, `max_locks_per_transaction`, `wal_level`, `max_wal_senders`, `max_prepared_transactions`, `max_replication_slots`, `max_worker_processes` не могут быть переопределены локально - Patroni их перезаписывает.

Проверка запущен ли PostgreSQL мастер:

- GET /master - должно возвращать 200 ТОЛЬКО для одной ноды

Проверка работают ли реплики

- GET /patroni с мастера должно возвращать replication:{{state: streaming}} для всех реплик

Запущен ли сам PostgreSQL:

- GET /patroni должен возвращать state:running для каждой ноды

Отставание реплики:

- GET /patroni - xlog: location с реплик не должен быть далеко от этого же параметра на мастере

- HAProxy. haproxy.toml
- Pgbouncer - решит проблему с дисконнектом у клиентов

postgresql:

callbacks:

on\_start: /opt/pgsql/pg\_start.sh

on\_stop: /opt/pgsql/pg\_stop.sh

on\_role\_change: /opt/pgsql/pg\_role\_change.sh

- `nofailover (true/false)` - в положении `true` нода никогда не станет мастером
- `noloadbalance (true/false)` - `/replica` всегда возвращает код 503
- `clonefrom (true/false)` - `patronictl` выберет предпочтительную ноду для `pgbasebackup`
- `nosync (true/false)` - нода никогда не станет синхронной репликой
- `replicatefrom (node name)` - указать реплику с которой снимать реплику

- Switchover
  - Переключение роли Мастера на новую ноду. Делается вручную, по сути плановые работы
- Failover
  - Экстренное переключение Мастера на новую ноду
  - Происходит автоматически
  - Ручной вариант - manual failover - только когда не система не может решить на кого переключать, или не настроен автомат

- `patronictl switchover cluster_name`
- отложенный switchover

- `patronictl -c /etc/patroni.yml restart postgres pg02`
  - Применение новых параметров требующих обязательной перезагрузки
  - Обновление Postgres

- `patronictl -c /etc/patroni.yml reinit postgres pg03`
  - Реинициализирует ноду в кластере. Т.е. по сути удаляет дата директорию и делает `pgbasebackup`

- Отключается автоматический failover
- Ставится глобальная пауза на все ноды
- Проведение плановых работ, например с etcd или обновление PostgreSQL

Тем не менее:

- Можно создавать реплики
- Ручной switchover возможен
- `patronictl -c /etc/patroni.yml pause|resume`

- **synchronous\_mode:** true/false - не делает failover ни на какую реплику кроме синхронной
- **synchronous\_mode\_strict:** true/false - если синхронная реплика пропала, то мастер не принимает новые записи пока она не вернется

Полные и инкрементные бэкапы создаются кастомными скриптами по плану (cron)/barman/wal-g/wal-e/etc

- Роль узла в кластере можно узнать запросом к DCS
- Архивные транзакционные логи (WAL):
  - сегментами в 16 Мб с мастер узла (**archive\_command=on**)
  - потоком по протоколу физической репликации (**pg\_receive\_wal**)

- Возможность восстановиться из бэкапа на любую точку по:
  - времени
  - id транзакции (xid)
  - Lsn транзакционной записи в журнале
  - именной записи в журнале

bootstrap:

method: `probackup`

probackup:

command: `"pg_probackup restore -B /path/to/backup --instance <scope> -D <datadir> --time='2019-09-08 00:00:00 UTC' \ --recovery-target-action=promote"`

recovery\_conf:

recovery\_target\_timeline: latest

restore\_command: `pg_probackup-11 archive-get -B /var/backup --instance <scope> --remote-user=dbbackup --wal-file-path %p --wal-file-name %f --remote-host=AA.BB.CC.DD`

```
: pg_probackup archive-get from /var/backup/wal/db-mt/00000013000000000000000076 to
/var/data/base/pg_wal/RECOVERYXLOG
ERROR: Source WAL file "/var/backup/wal/db-mt/00000013000000000000000076" doesn't exist
2019-08-28 10:06:57.782 UTC [23] LOG: redo done at 0/75000198
INFO: pg_probackup archive-get from /var/backup/wal/db-mt/00000013000000000000000075 to
/var/data/base/pg_wal/RECOVERYXLOG
INFO: pg_probackup archive-get completed successfully
2019-08-28 10:07:01.015 UTC [23] LOG: restored log file "00000013000000000000000075" from
archive
INFO: pg_probackup archive-get from /var/backup/wal/db-mt/00000014.history to
/var/data/base/pg_wal/RECOVERYHISTORY
ERROR: Source WAL file "/var/backup/wal/db-mt/00000014.history" doesn't exist
2019-08-28 10:07:01.639 UTC [23] LOG: selected new timeline ID: 20
2019-08-28 10:07:01.677 UTC [23] LOG: archive recovery complete
INFO: pg_probackup archive-get from /var/backup/wal/db-mt/00000013.history to
/var/data/base/pg_wal/RECOVERYHISTORY
INFO: pg_probackup archive-get completed successfully
2019-08-28 10:07:02.280 UTC [23] LOG: restored log file "00000013.history" from archive
2019-08-28 10:07:02.389 UTC [21] LOG: database system is ready to accept connections
```

- По умолчанию реплика создается с помощью утилиты `pg_basebackup`
- Это поведение можно переопределить параметром `create_replica_methods`
- Важно, обязательно нужно указать `basebackup`, иначе если из бекапа не получится, то реплика не заведется.

postgresql:

create\_replica\_methods:

- `probackup`
- `basebackup`

probackup:

command: `"ssh dbbackup@10.23.2.163 'bash /var/backup/pg_restore.sh'"`

no\_params: `True`

basebackup:

max-rate: `'100M'`

2019-08-20 14:17:51,986 INFO: Removing data directory: /var/data/base

INFO: Validating backup PWJ0PZ

INFO: Backup PWJ0PZ data files are valid

INFO: Backup PWJ0PZ WAL segments are valid

INFO: Backup PWJ0PZ is valid.

INFO: Restore of backup PWJ0PZ completed.

2019-08-20 14:17:56,150 INFO: replica has been created using probackup

2019-08-20 14:17:56,153 INFO: bootstrapped from leader 'AA.BB.CC.DD'

- Docker образ для минимального запуска
- Скрипт с восстановлением из бекапа
  - Минимальный конфиг для старта
  - pg\_hba.conf с trust доступами (для упрощения)
- `docker exec pgvalid pg_dump -h localhost -U postgres > /dev/null`
- Amcheck
  - `CREATE EXTENSION amcheck;`
  - `pg_probackup-11 checkdb --amcheck --heapallindexed ...`

Custom recovery.conf:

- `recovery_min_apply_delay = '12h'`

Tags:

- `nosync`
- `nofailover`
- `nobalance`

# Домашнее задание

- Развернуть кластер PostgreSQL из трех нод. Создать тестовую базу - проверить статус репликации
  - Сделать switchover/failover
  - Поменять конфигурацию PostgreSQL + с параметром требующим перезагрузки
- \* Настроить клиентские подключения через HAProxy

**Ваши вопросы?**



**Заполните, пожалуйста,  
опрос в ЛК о занятии**

**Спасибо  
за внимание!**

**До встречи в Slack и на вебинаре**

