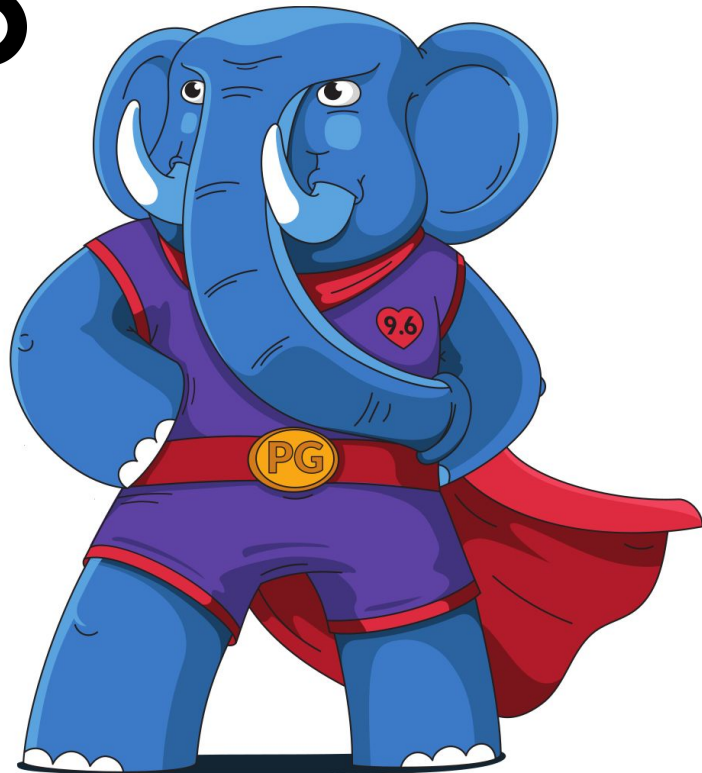


10 years of PostgreSQL in



Konstantin Evteev

Moscow 2020





2 500 000
housing ads



2 000 000
jobs



1 300 000
services



10 900 000
cars



31 000 000
goods for
sale

Avito helps people from all over the country to make deals

400 000 new
ads per day

59 702 904

120 deals
per minute

PostgreSQL in Avito 2020

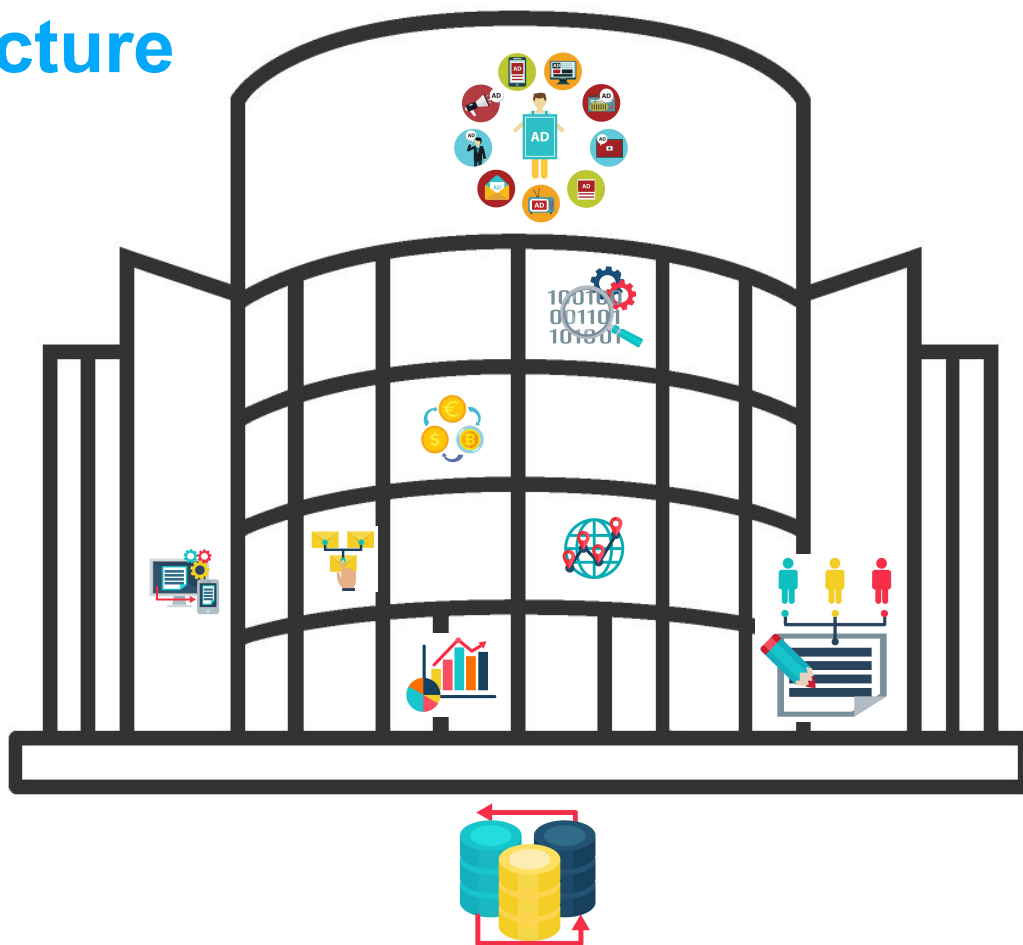
- > 100 physical servers
- > 250 databases
- 32Tb total size
- 150k RPS
- The size of backups is 93Tb



Plan

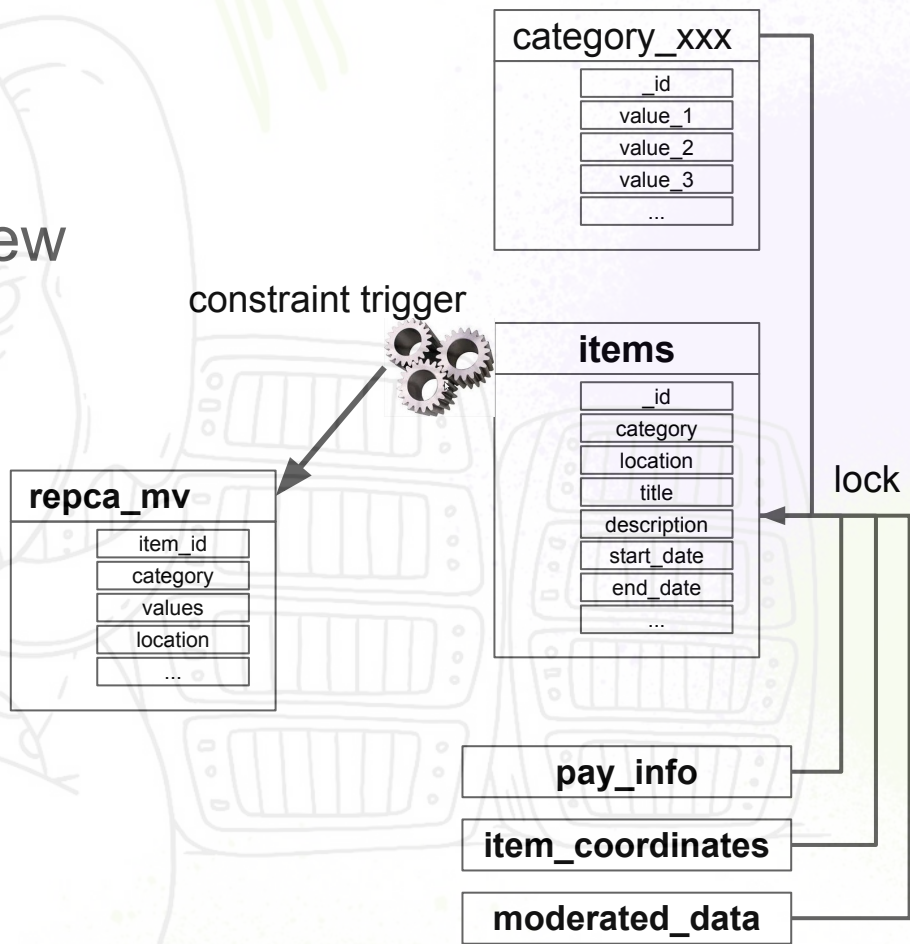
1. **Evolution of monolith architecture**
2. Migration to microservice architecture
3. Integration & communication
4. Dev tools and environment
5. Platform (DBaaS in 3 Datacenters)

Monolith architecture



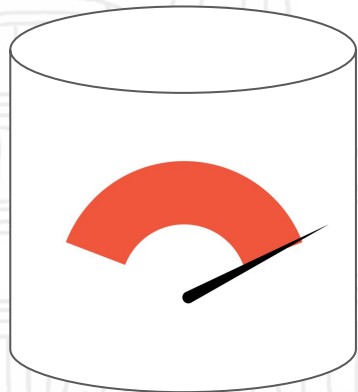
Turn on fast mode:

Denormalization: Materialized view



Turn on fast mode:

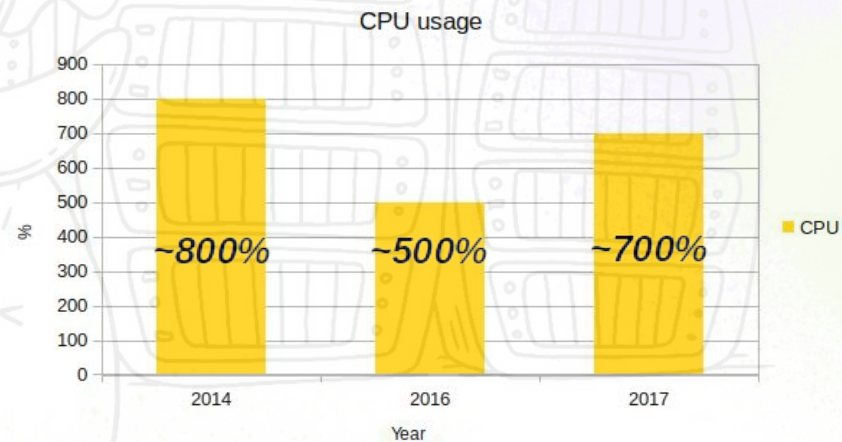
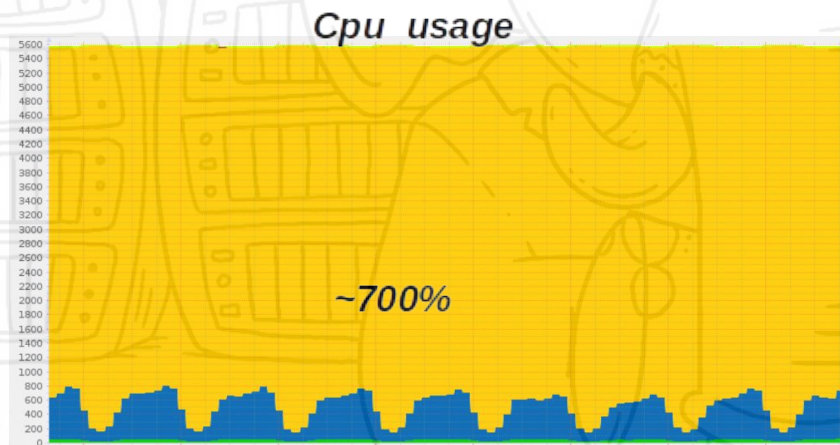
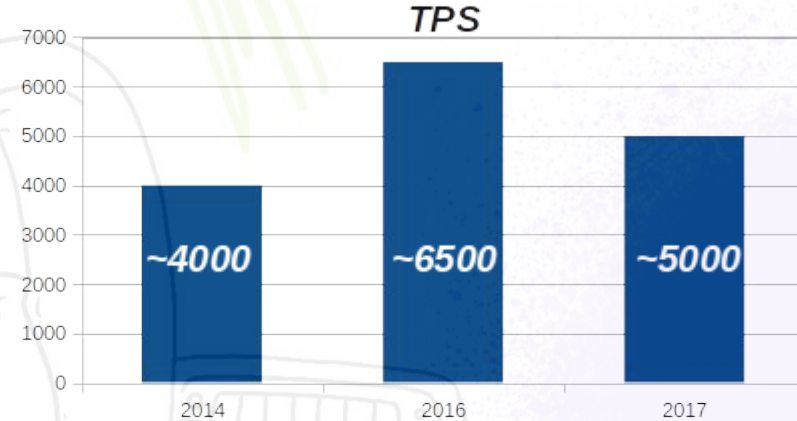
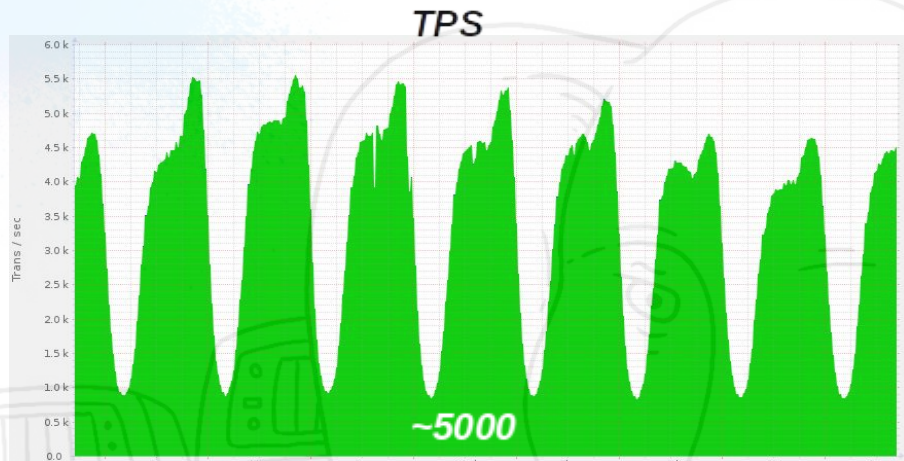
Materialized view in memory



```
mount -t tmpfs tmpfs /mnt/ramdisk
```

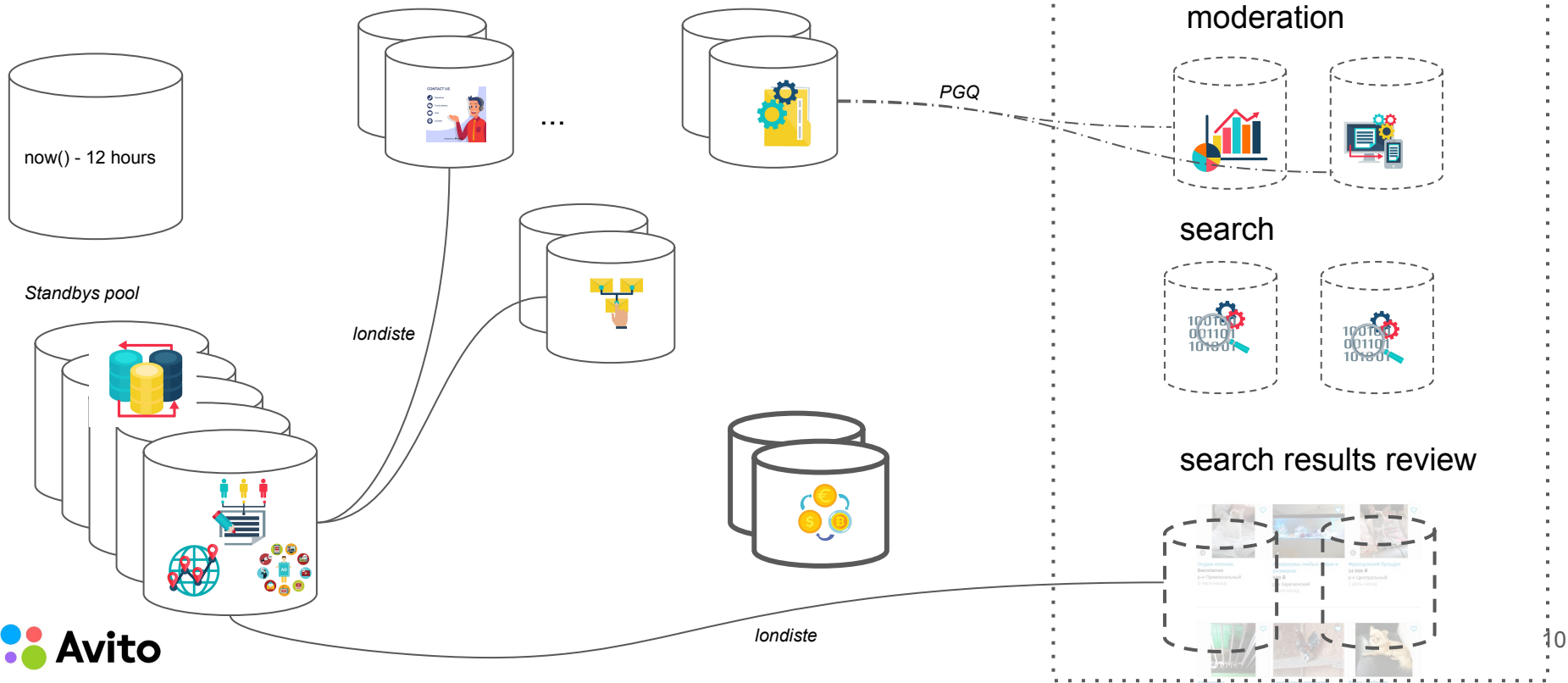
```
CREATE TABLESPACE fastspace LOCATION  
'/mnt/ramdisk/postgresql/data';
```


<https://www.slideshare.net/pavlushko/sphinx-10460333>

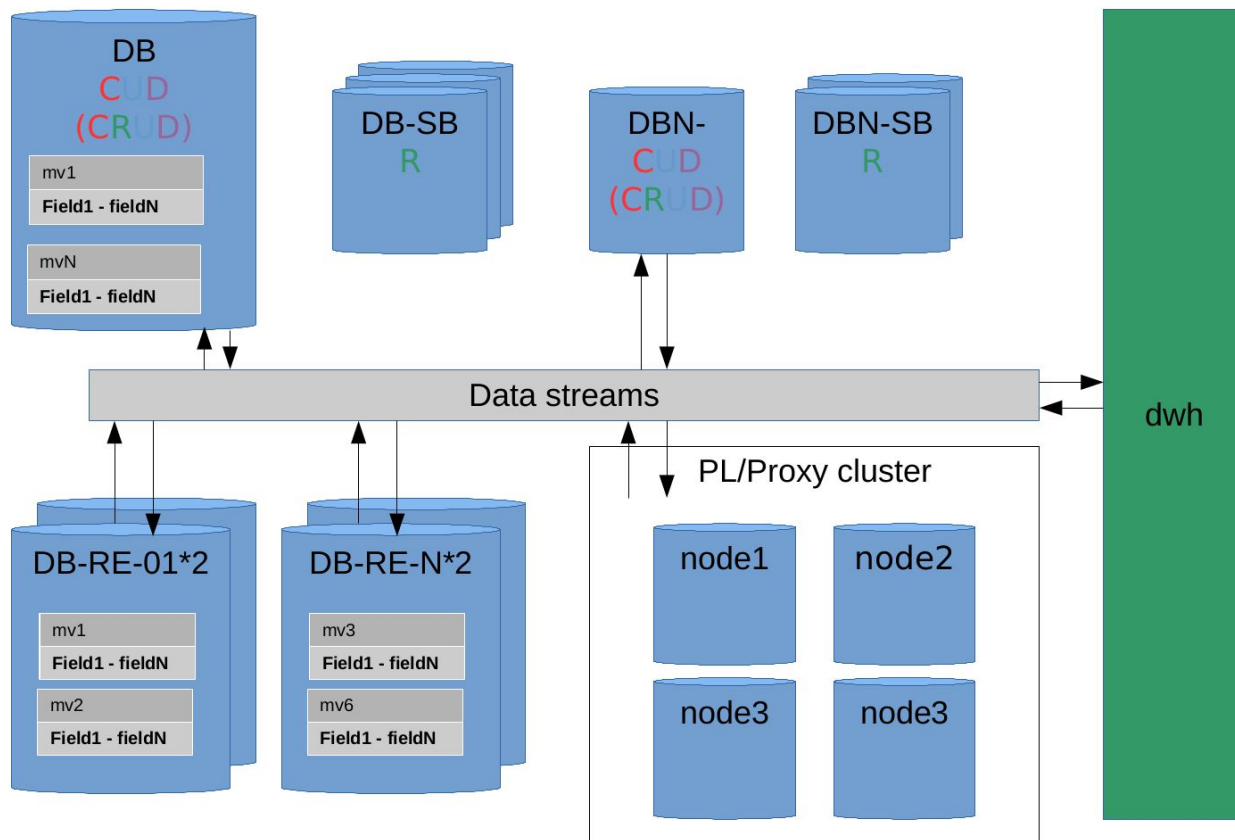
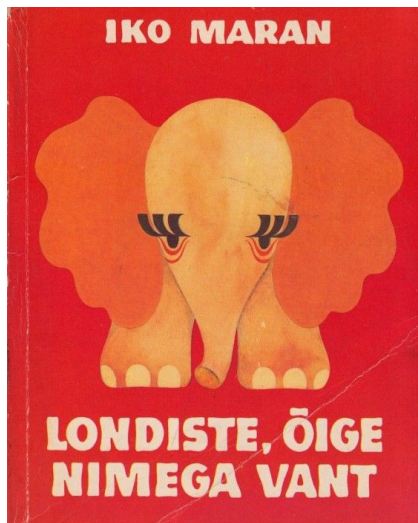


Monolith architecture 2016

functional, synced with the help of transactional queues

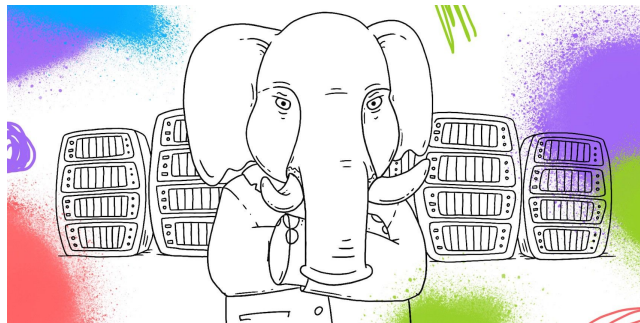


Avito 2016



Recovery

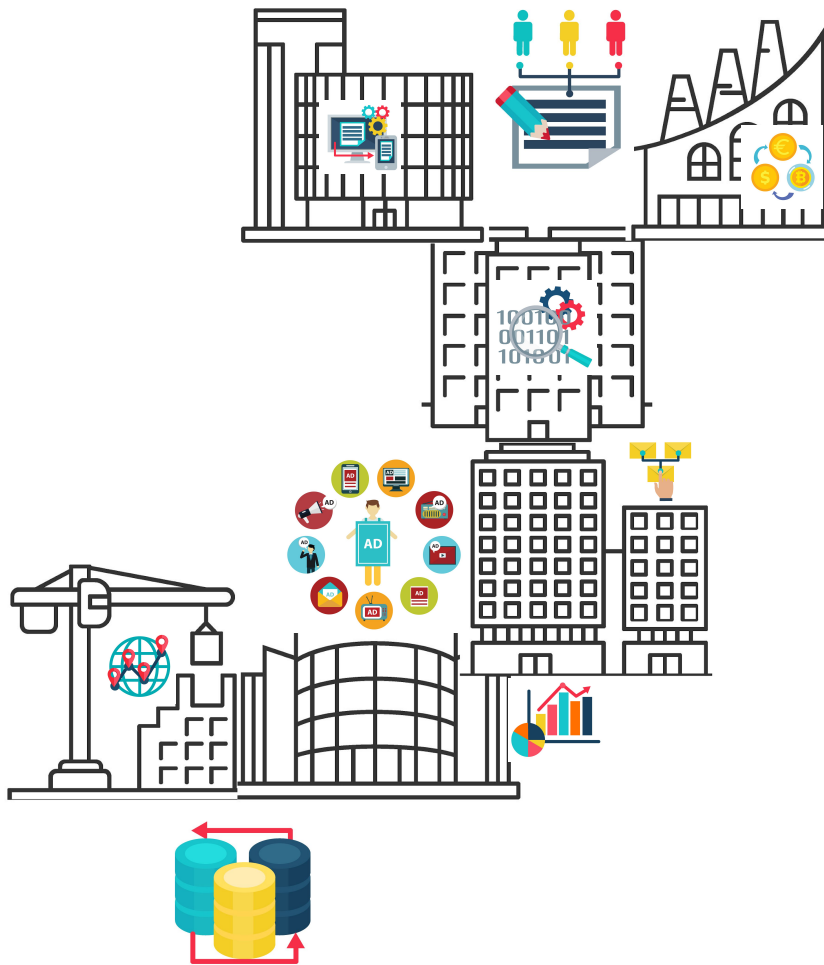
- (1) Reinitializing subscriber from another subscriber
- (2) UNDO recovery on the destination side
- (3) REDO - reposition source (subscriber's crash)
- (4) REDO 2 - on provider's side (provider's crash and switching to the provider's standby, subscriber is falling behind)



<https://medium.com/avitotech/recovery-use-cases-for-logical-replication-in-postgresql-10-a1e6bab03072>

Monolith architecture

1. Complex
2. Fragile
3. Low speed of improvements
4. Hard to scale

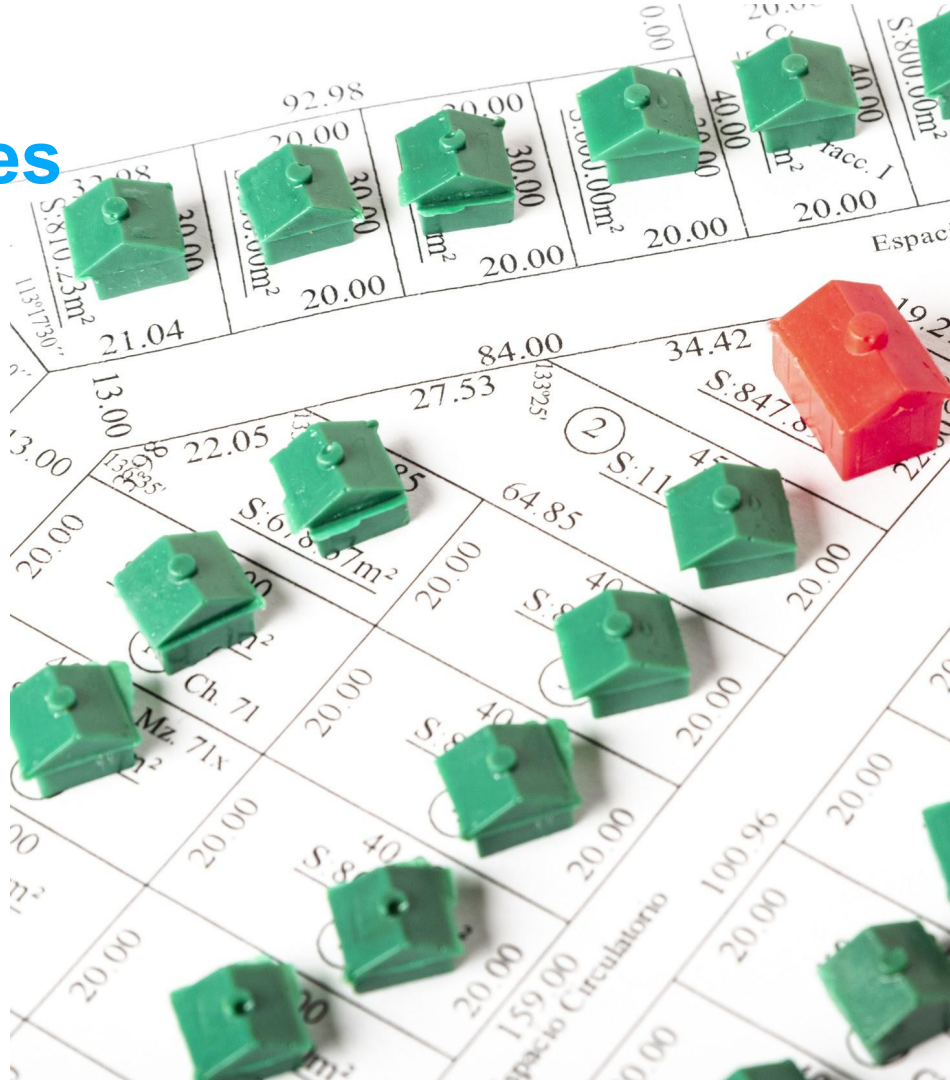


Plan

1. Evolution of monolith architecture
- 2. Migration to microservice architecture**
3. Integration & communication
4. Dev tools and environment
5. Platform (DBaaS in 3 Datacenters)

Migration to microservices

1. Architecture
2. Tooling
3. Integration tooling
4. Platform
5. Approaches & best practices



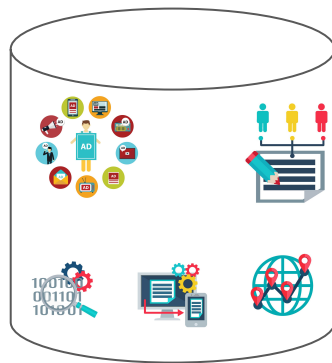
Shared Database is not an option

1. ...
2. Single point of failure
3. Can't find many db experts
4. High entry threshold
5. Stored procedures as code, CICD tools
- ...



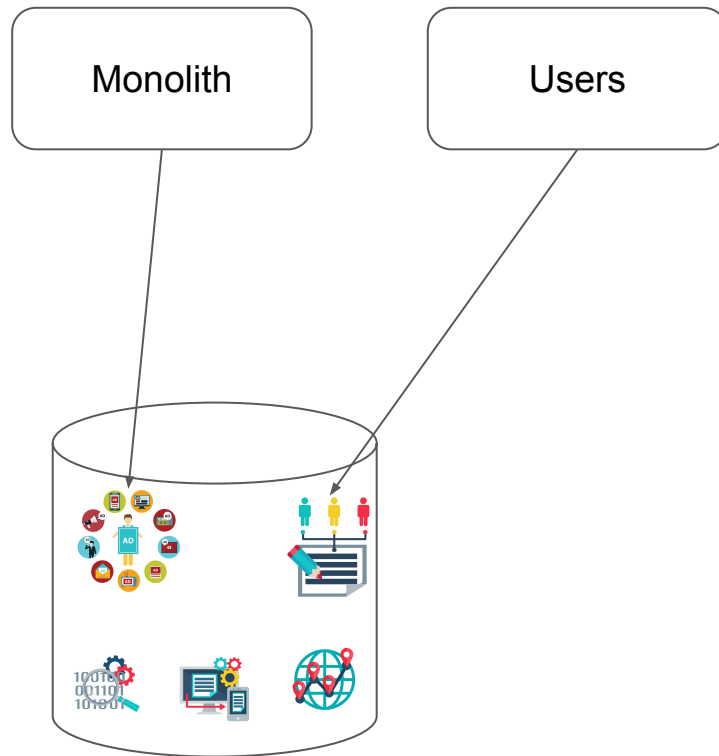
Breaking monolith DB into microservices DBs

1. Highlight the domain area



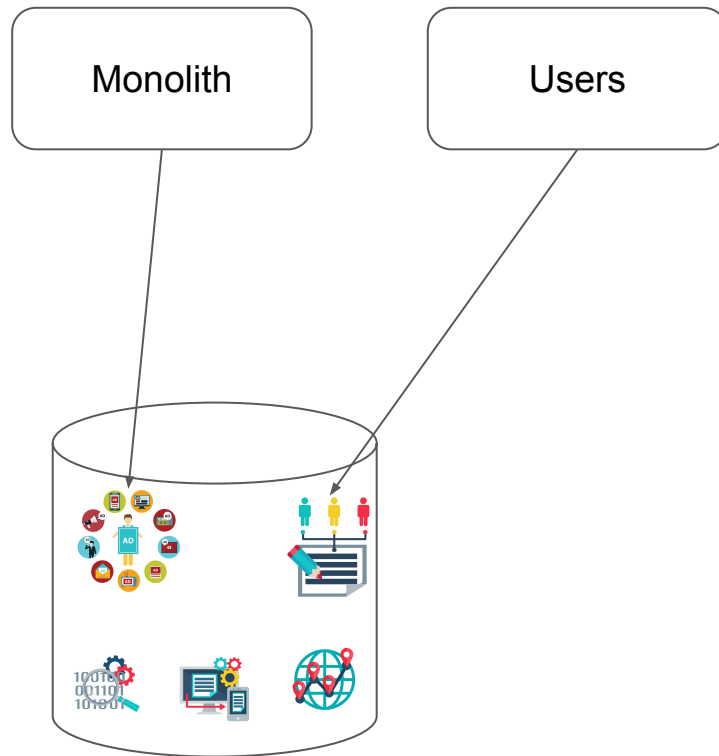
Breaking monolith DB into microservices DBs

1. Highlight the domain area
2. Split the code



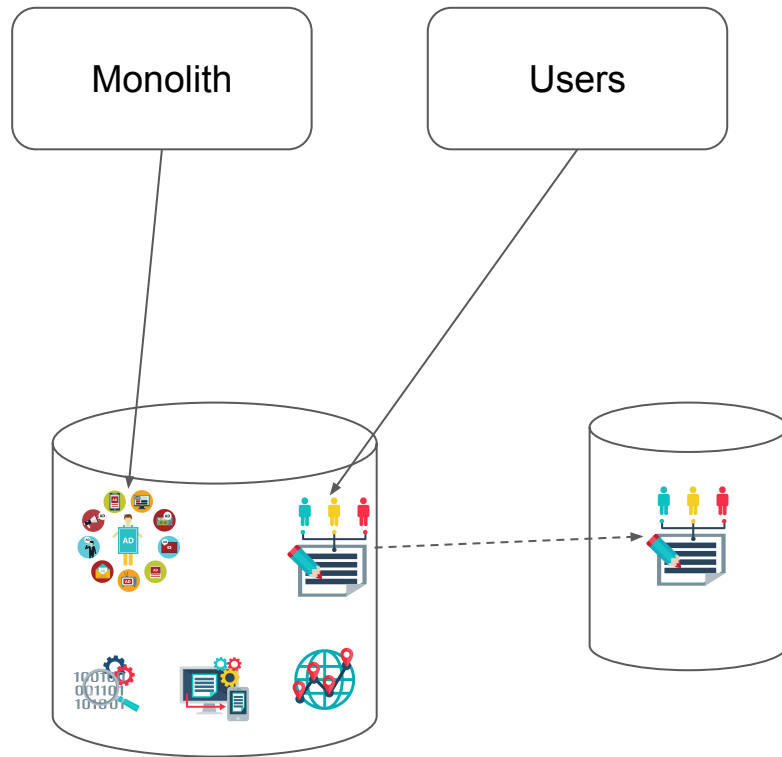
Breaking monolith DB into microservices DBs

1. Highlight the domain area
2. Split the code
3. Isolate the data



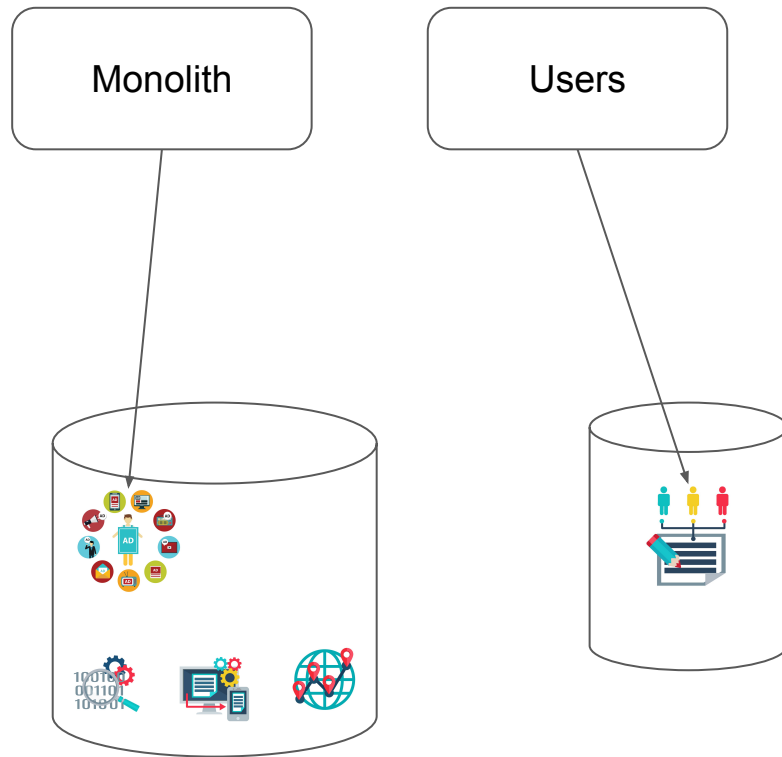
Breaking monolith DB into microservices DBs

1. Highlight the domain area
2. Split the code
3. Isolate the data
4. Switch to new DB
 - a. logical replication



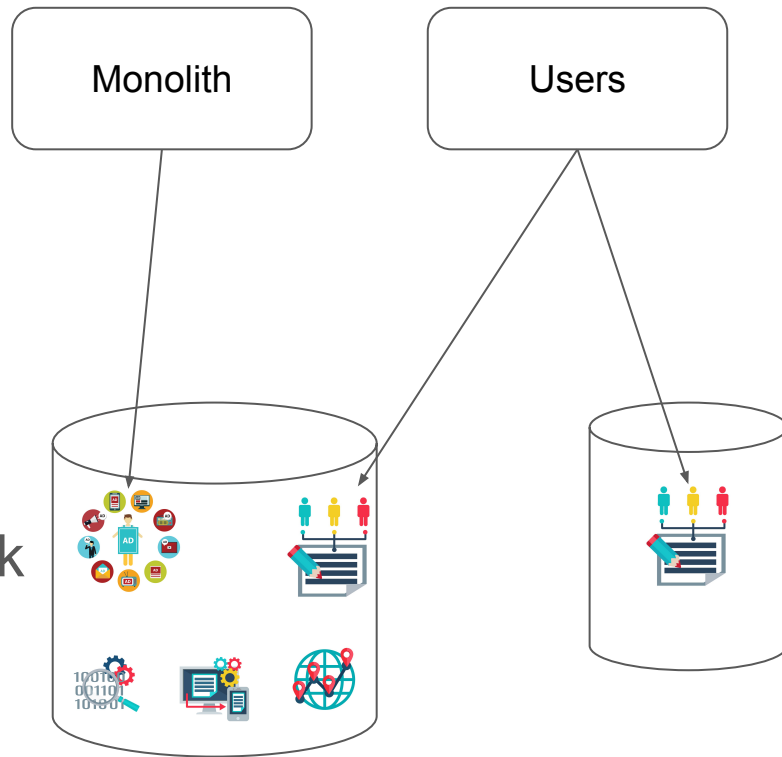
Breaking monolith DB into microservices DBs

1. Highlight the domain area
2. Split the code
3. Isolate the data
4. Switch to new DB
 - a. logical replication



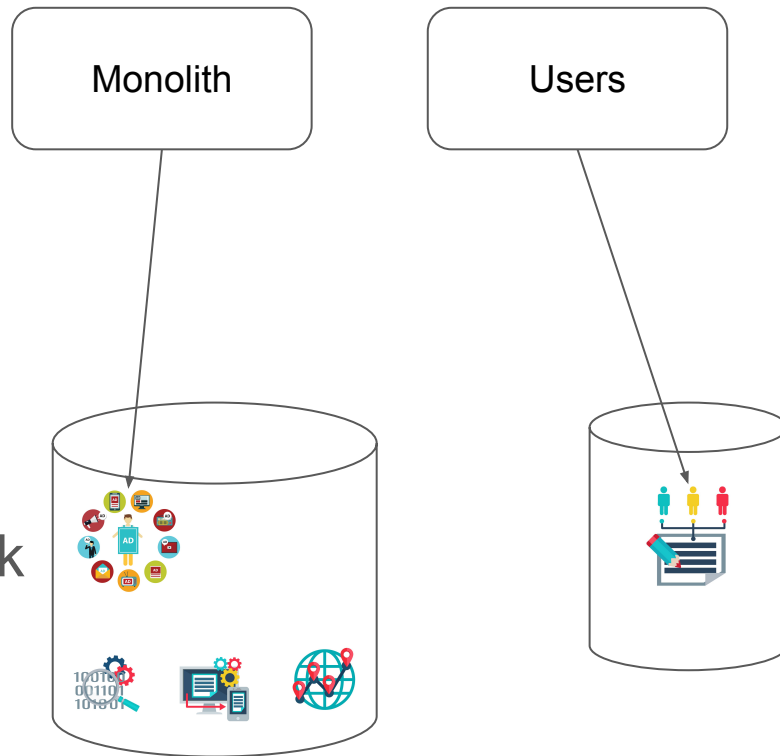
Breaking monolith DB into microservices DBs

1. Highlight the domain area
2. Split the code
3. Split the database
4. Switch to new DB
 - a. logical replication
 - b. duplicate changes and check

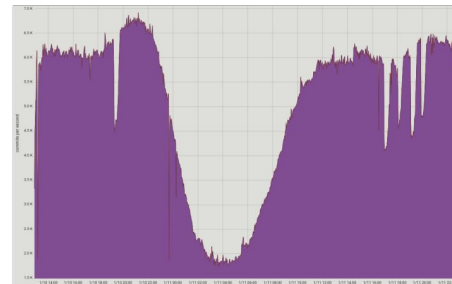


Breaking monolith DB into microservices DBs

1. Highlight the domain area
2. Split the code
3. Split the database
4. Switch to new DB
 - a. logical replication
 - b. duplicate changes and check

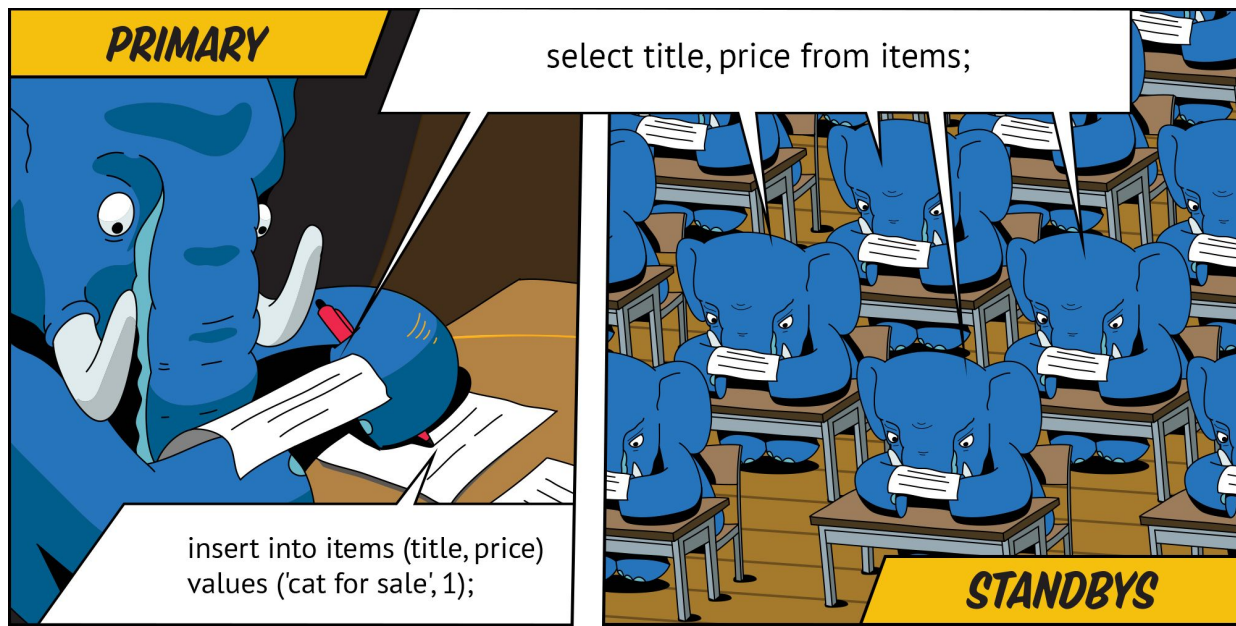


Scaling

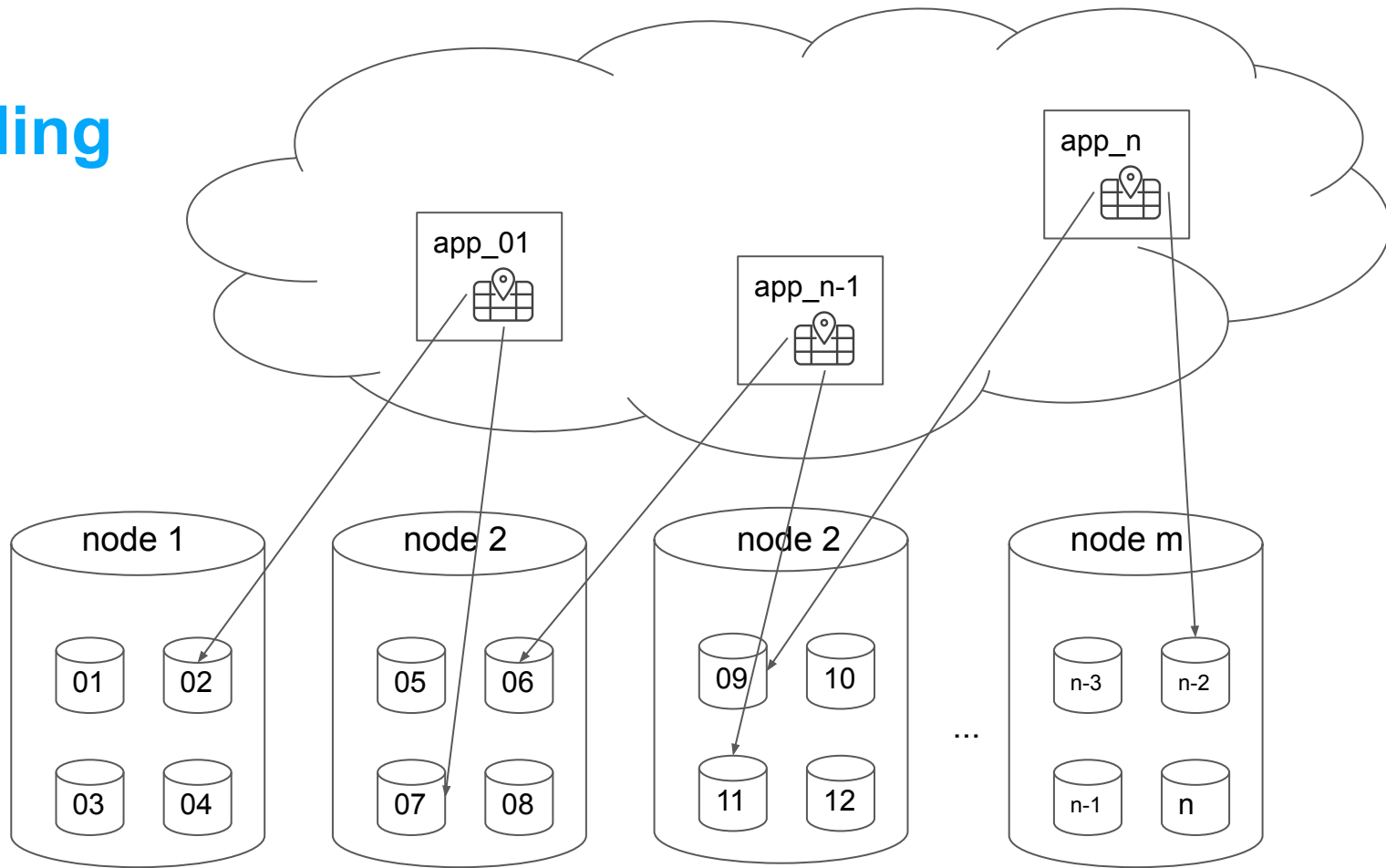


Standby in production: scaling application in the second largest classified site in the world

- (1) Deadlock on standby
- (2) DDL (statement_timeout and deadlock_timeout)
- (3) Vacuum replaying on standby and truncating data file
- (4) Restoring WAL from archive



Sharding

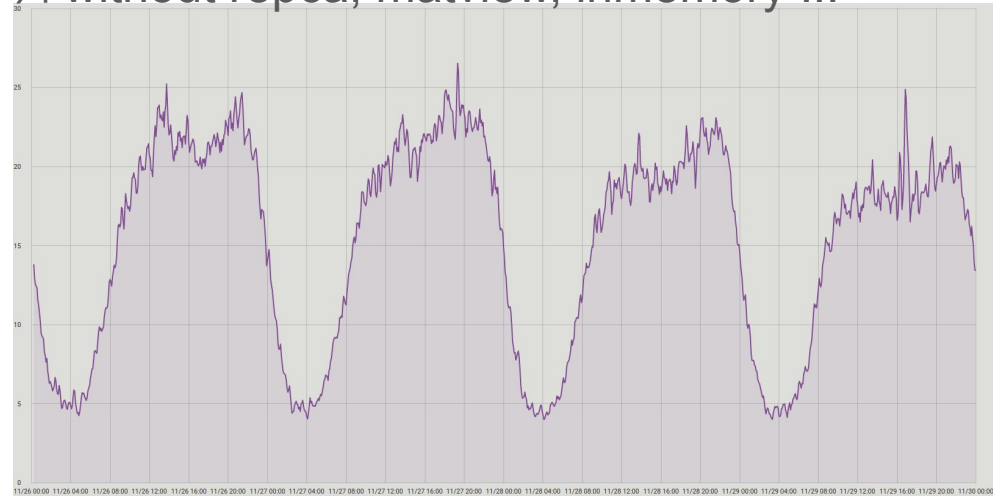


As a result

monolith



$\frac{1}{4}$ without repca, matview, inmemory ...



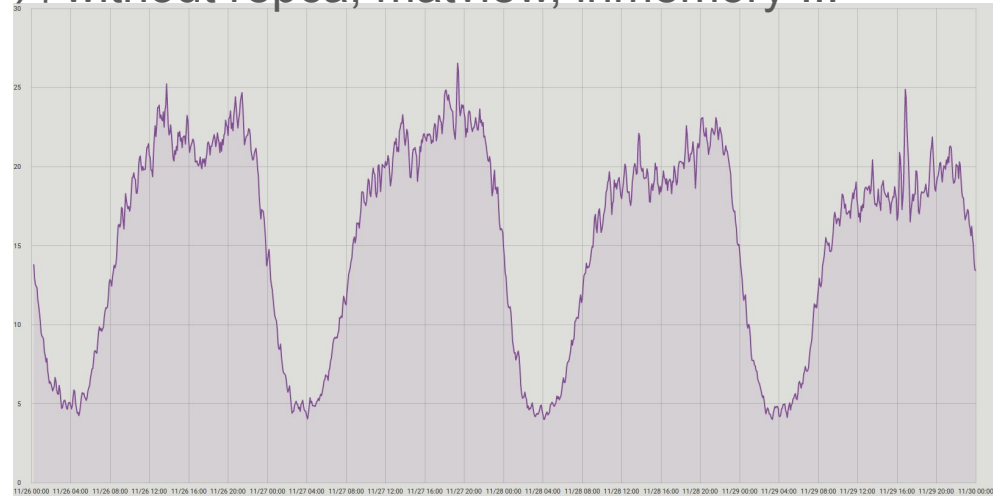
As a result

MATVIEW

monolith



$\frac{1}{4}$ without repca, matview, inmemory ...



As a result

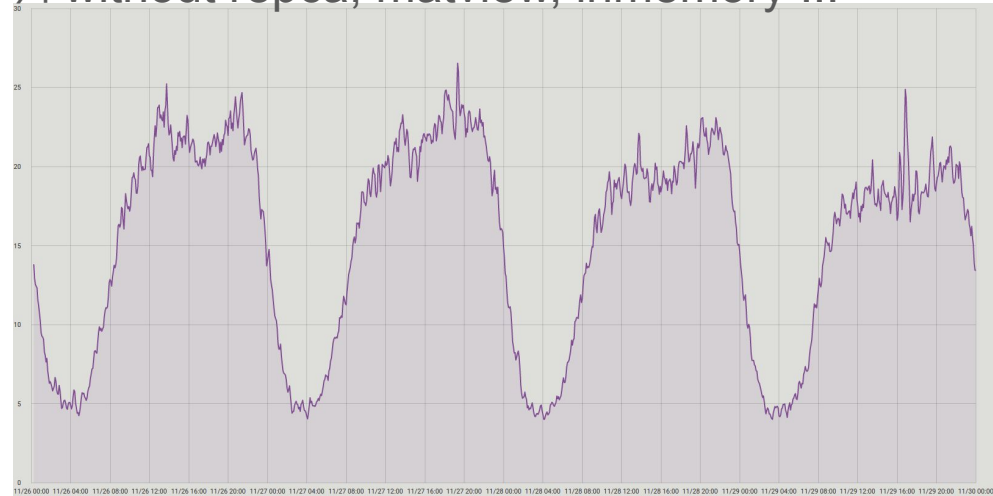
~~MATVIEW~~

~~LOGICAL REPLICAS~~

monolith



$\frac{1}{4}$ without repca, matview, inmemory ...



As a result

~~MATVIEW~~

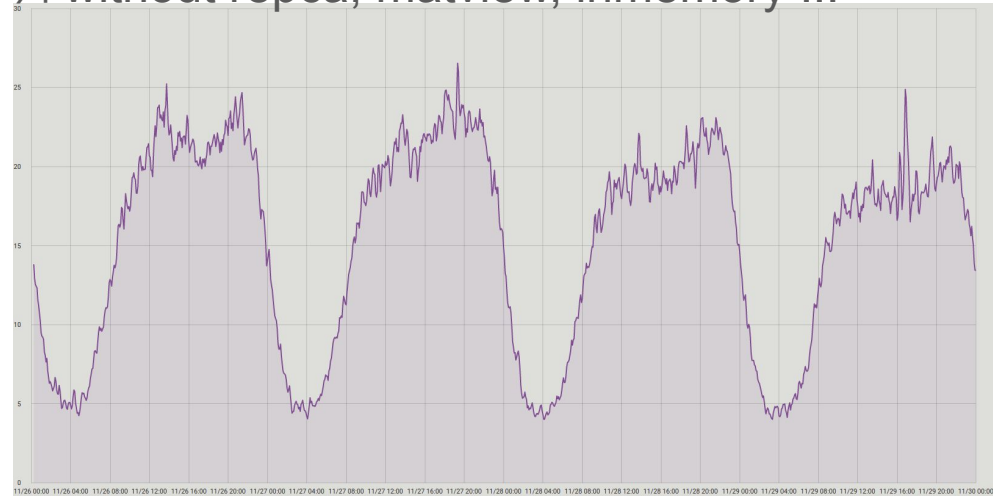
~~LOGICAL REPLICAS~~

~~IN-MEMORY TABLESPACES~~

monolith



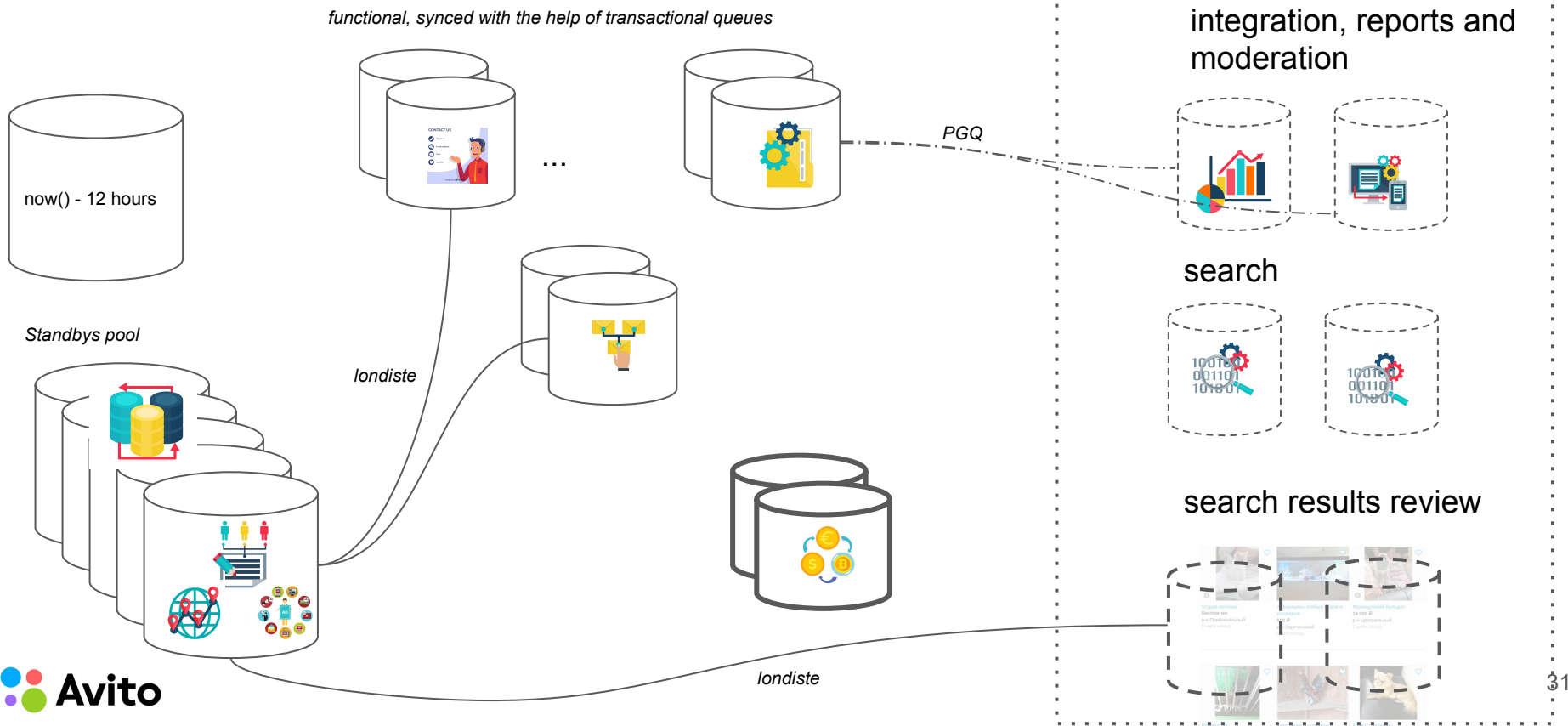
$\frac{1}{4}$ without repca, matview, inmemory ...



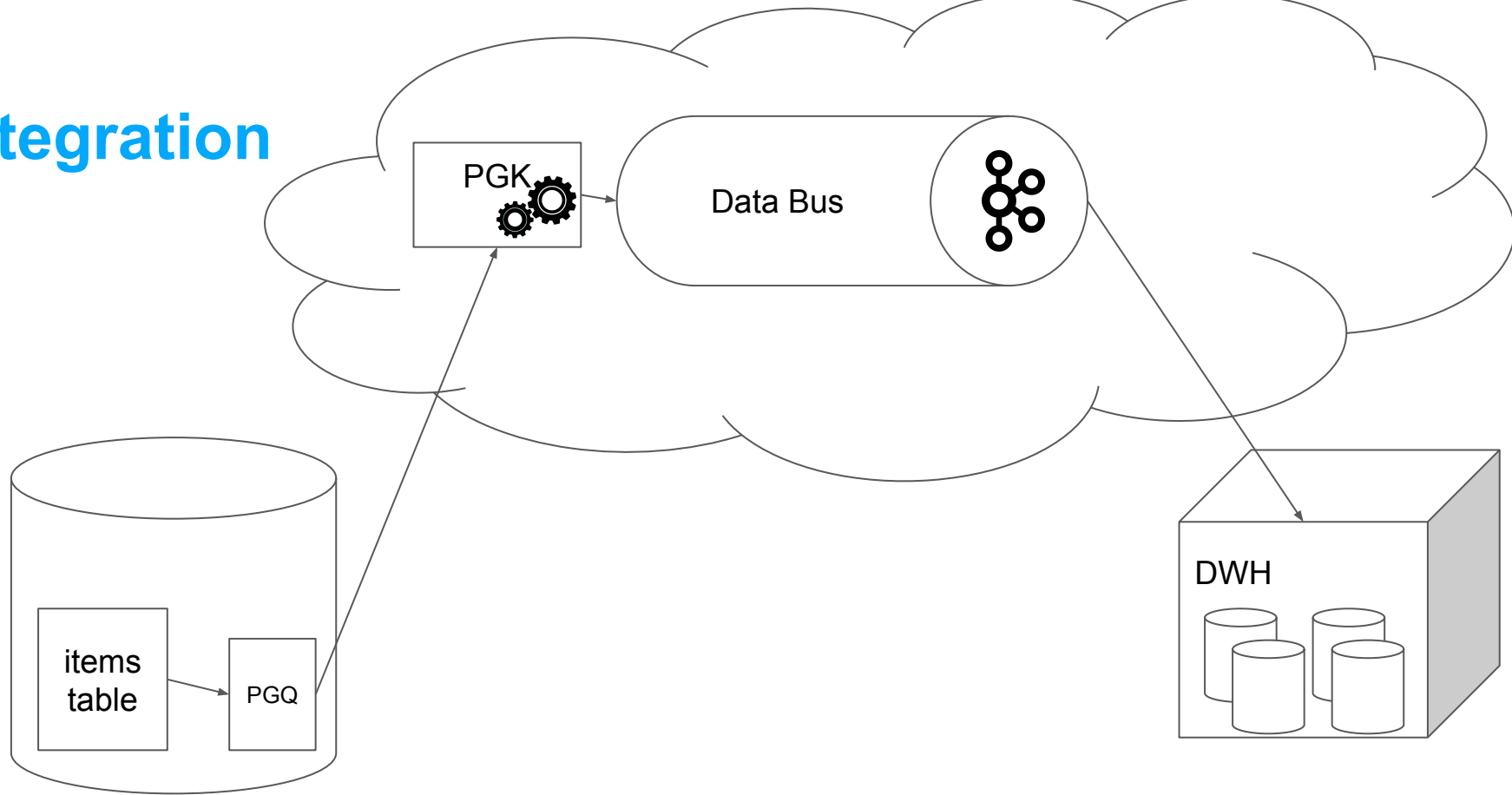
Plan

1. Evolution of monolith architecture
2. Migration to microservice architecture
- 3. Integration & communication**
4. Dev tools and environment
5. Platform (DBaaS in 3 Datacenters)

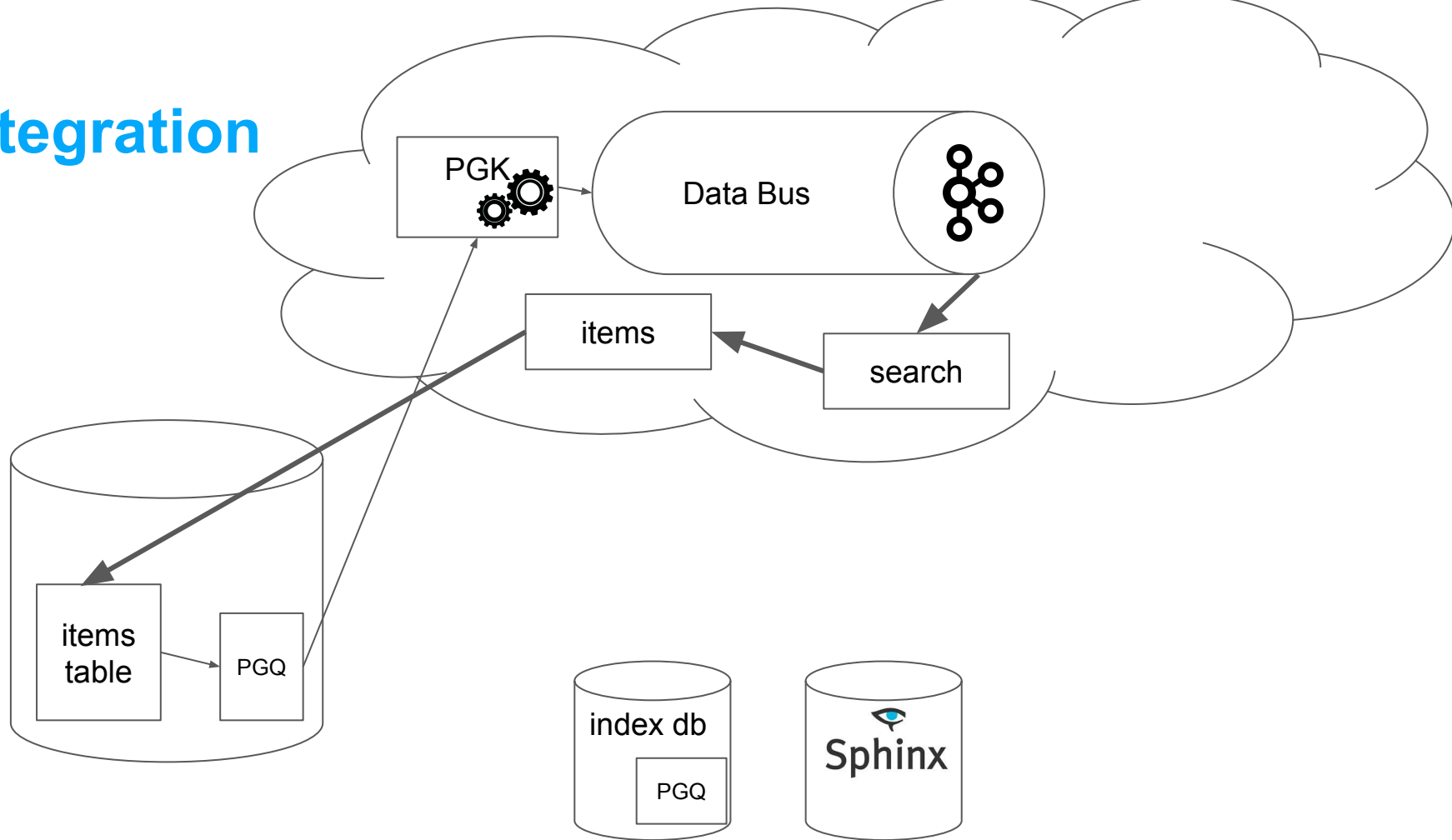
Integration



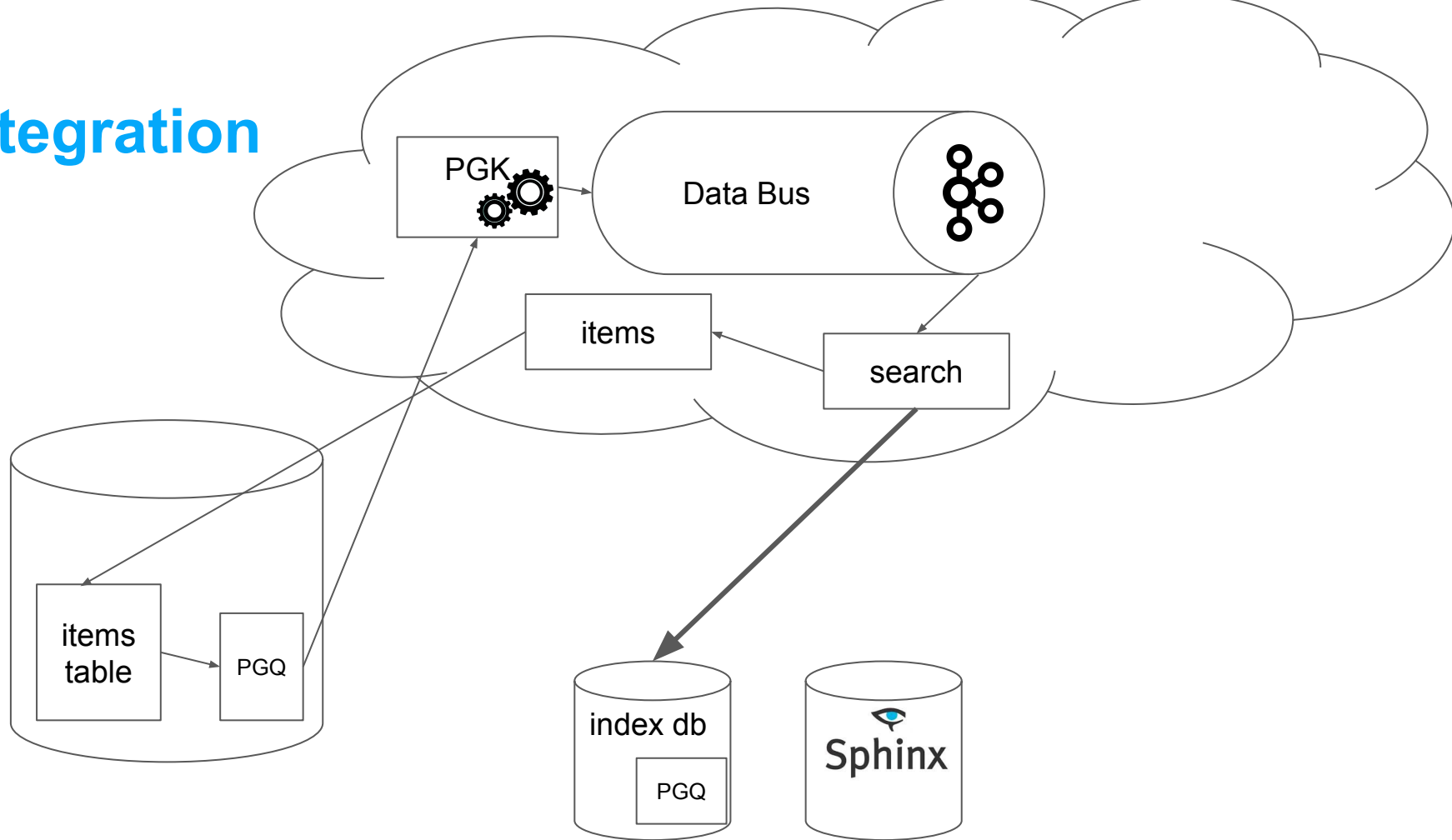
Integration



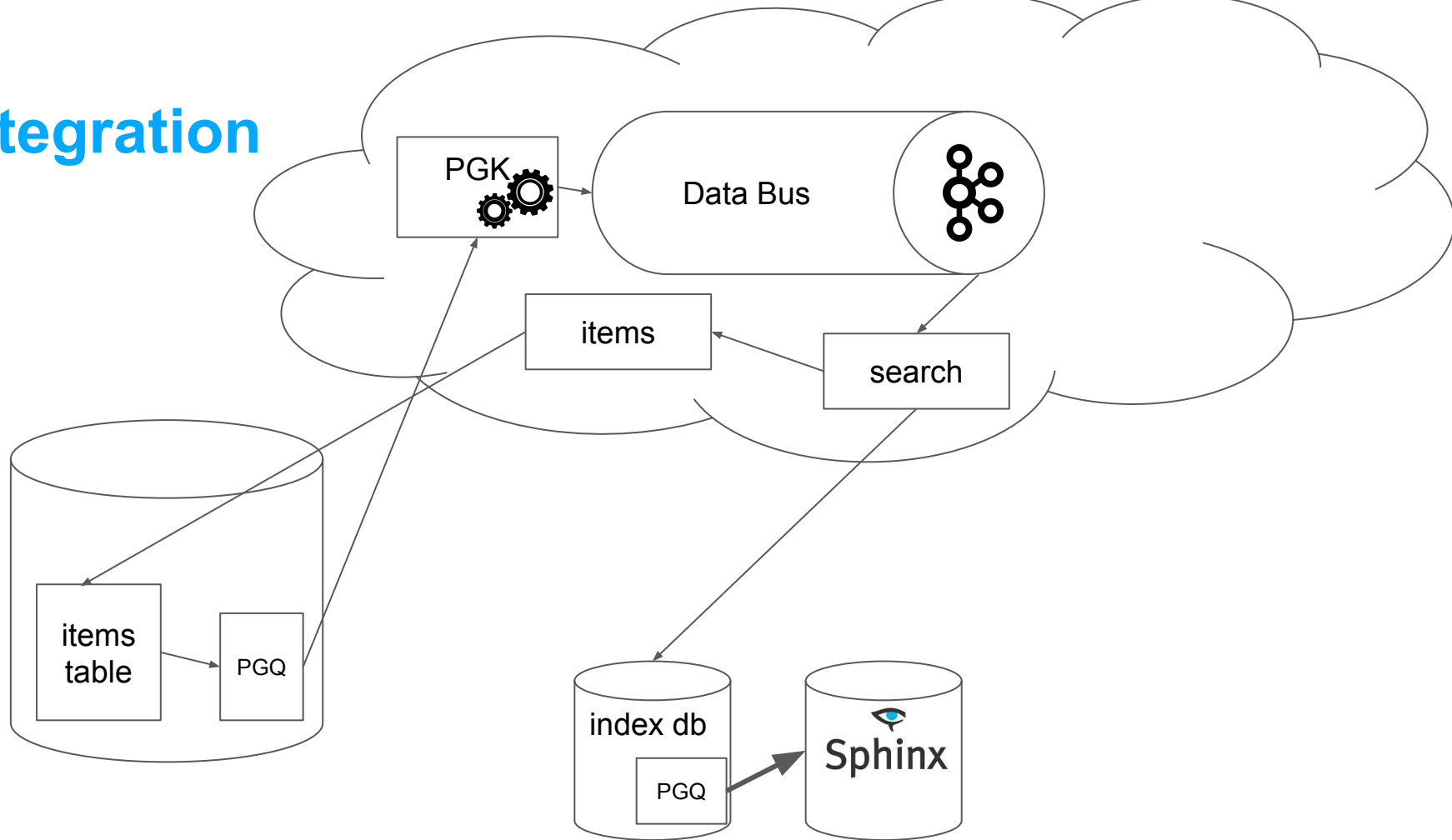
Integration



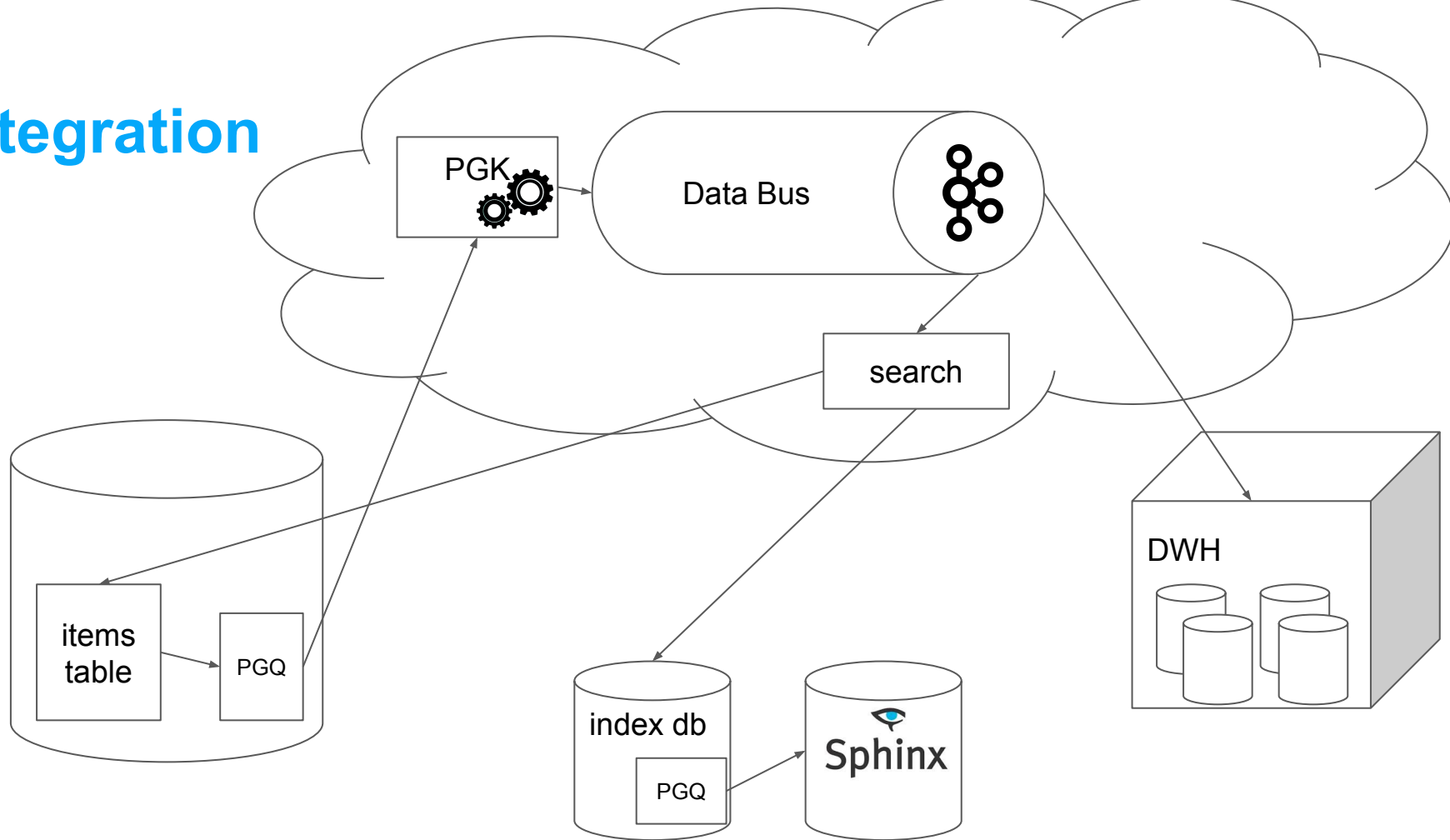
Integration



Integration

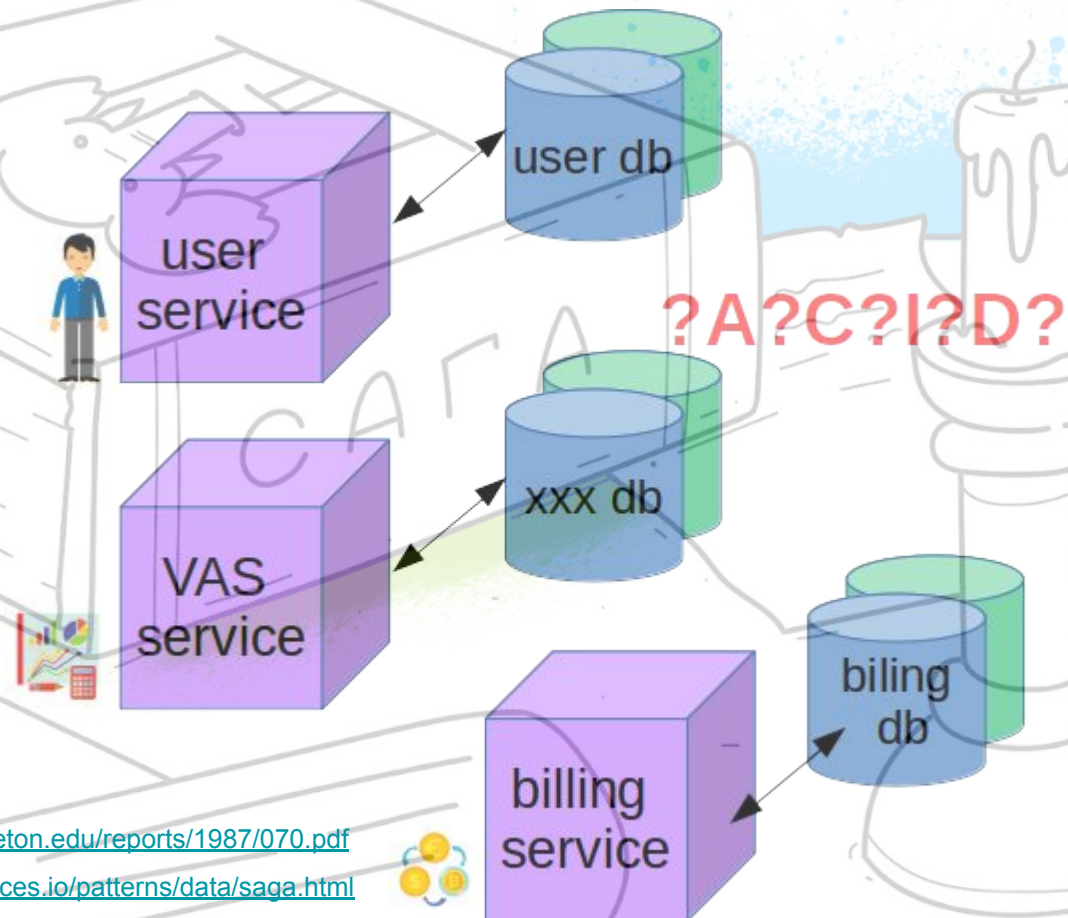
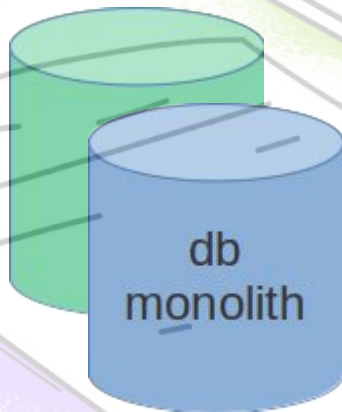


Integration



Saga

Monolith = shared database



<ftp://ftp.cs.princeton.edu/reports/1987/070.pdf>

<http://microservices.io/patterns/data/saga.html>

<https://habr.com/ru/company/avito/blog/426101/>

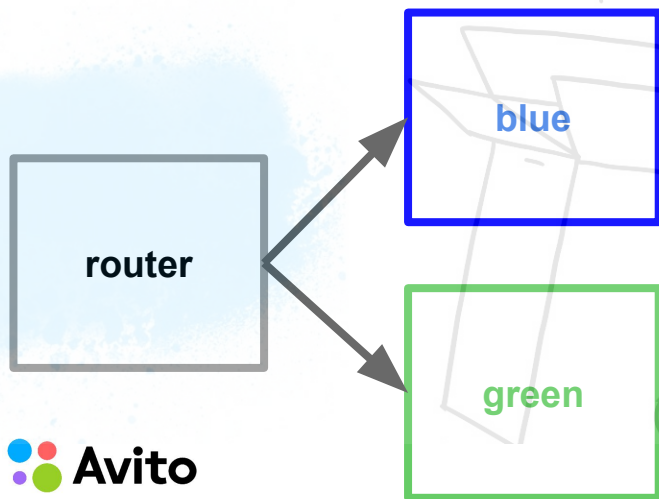
<https://habr.com/ru/company/oleg-bunin/blog/418235/>

Plan

1. Evolution of monolith architecture
2. Migration to microservice architecture
3. Integration & communication
- 4. Dev tools and environment**
5. Platform (DBaaS in 3 Datacenters)

Version control and code deploy

- migrators:
 - <https://github.com/yandex/pgmigrate/blob/master/doc/tutorial.md>
 - <https://flywaydb.org/>
 - <http://www.liquibase.org/>
 - <https://sqitch.org/>
- stored procedures versions stored in dictionary
- user-schema based deploy



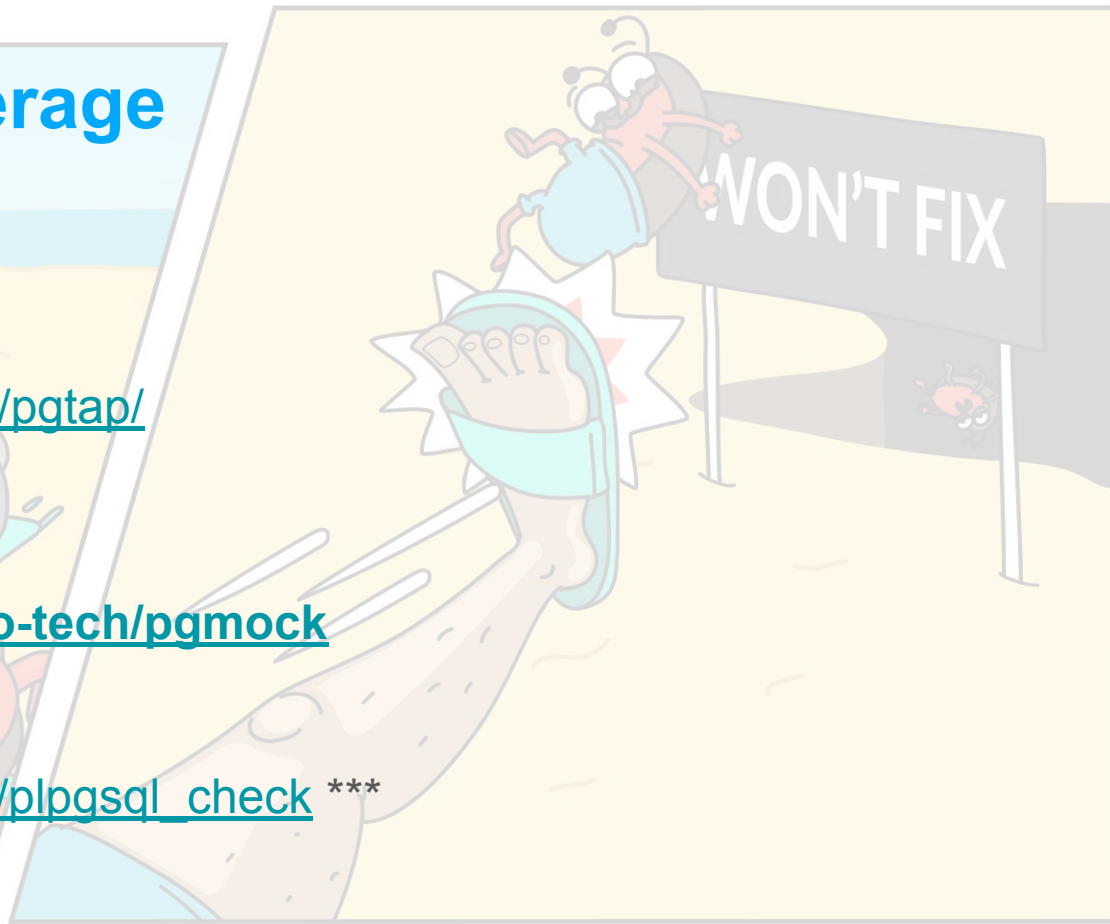
Tests and code coverage

- tests

- <https://github.com/theory/pgtap/>
- <https://pgtap.org/>
- <https://github.com/avito-tech/pgmock>

- coverage

https://github.com/okbob/plpgsql_check ***



Issues

1. Query to another services database
2. Production works from test environment
3. Run test and migrations in production environments
4. Drop table from IDE by chance
5. Security issues
6. ...

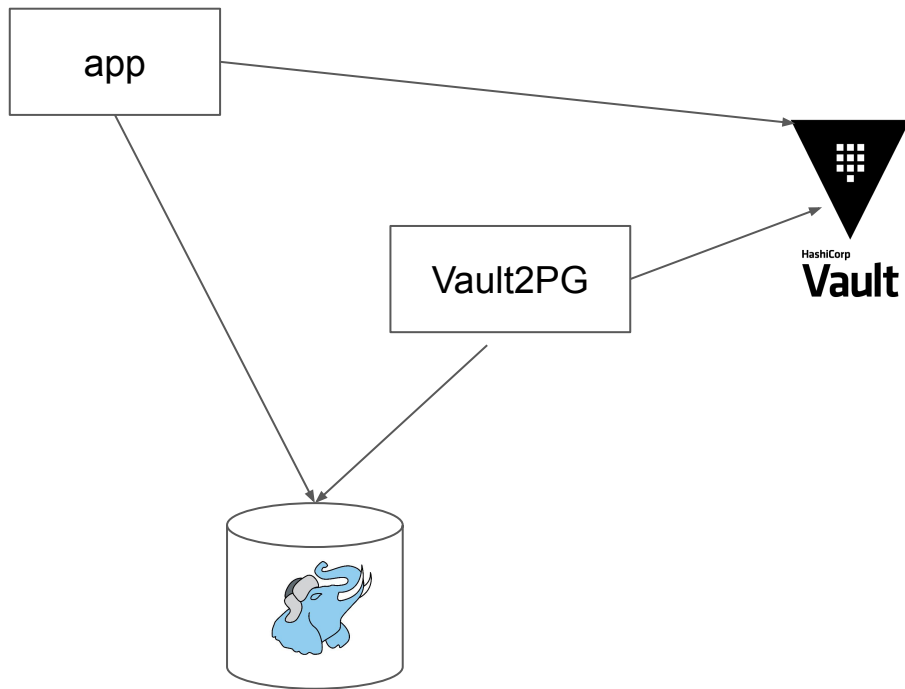
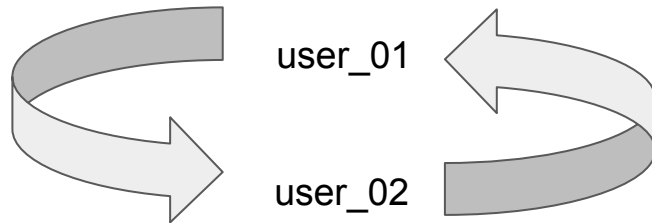


DB Access

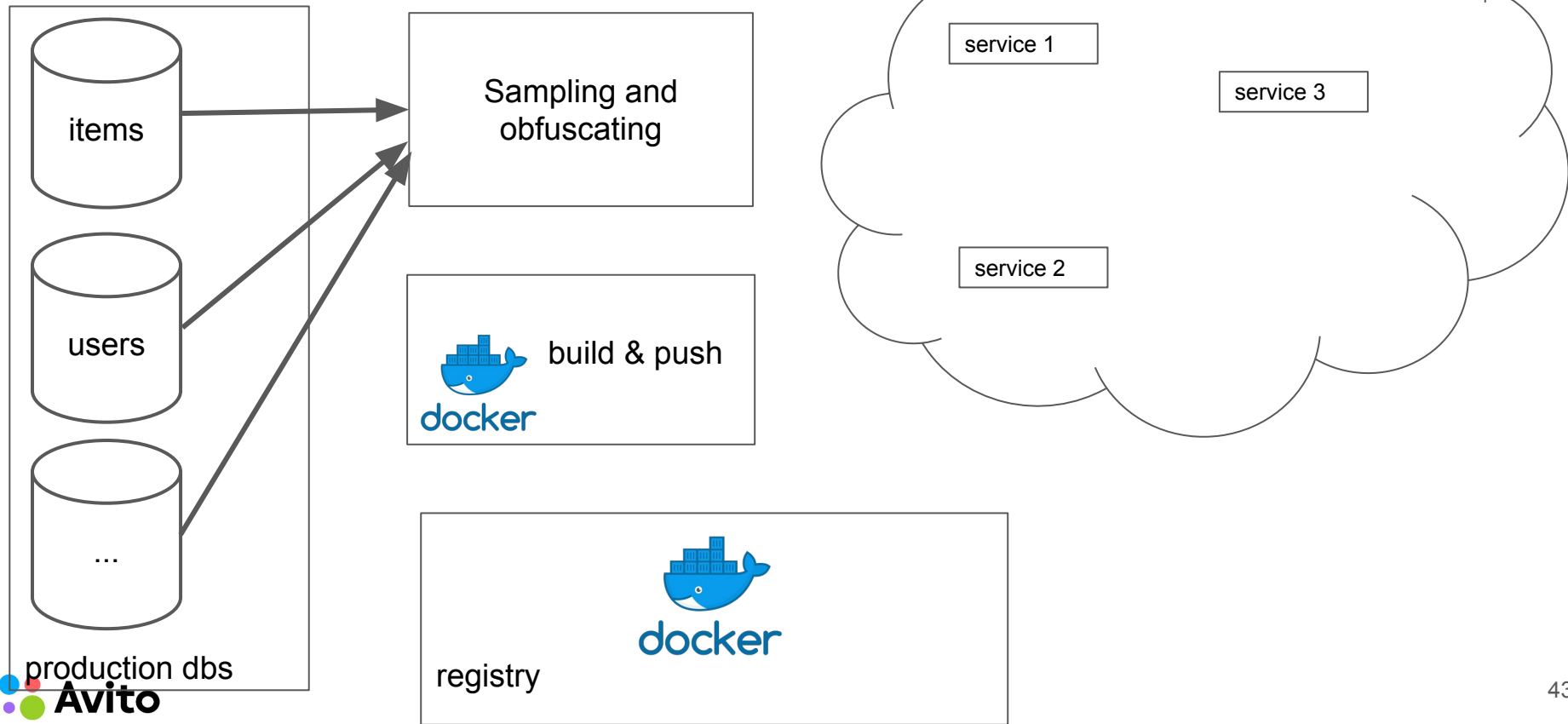
full access

read only

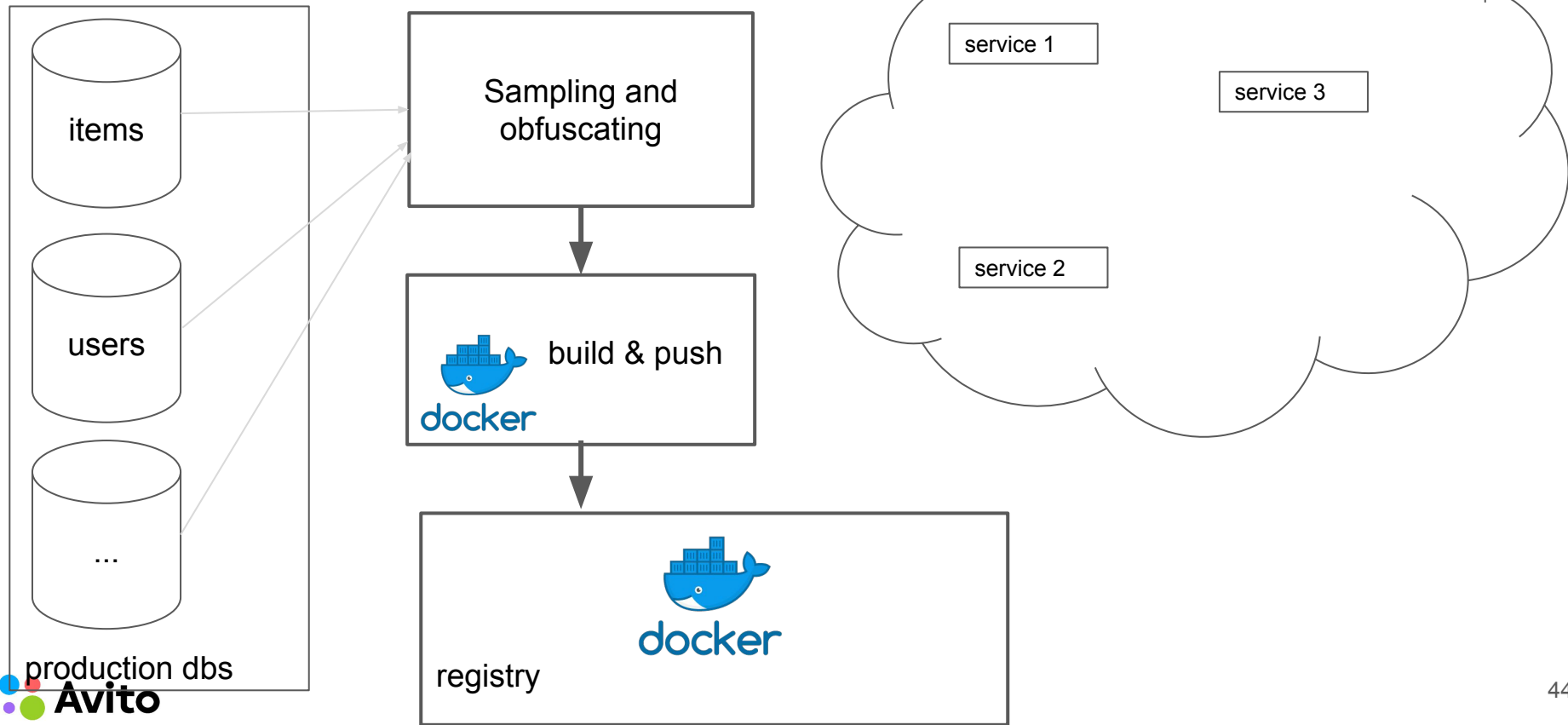
read write



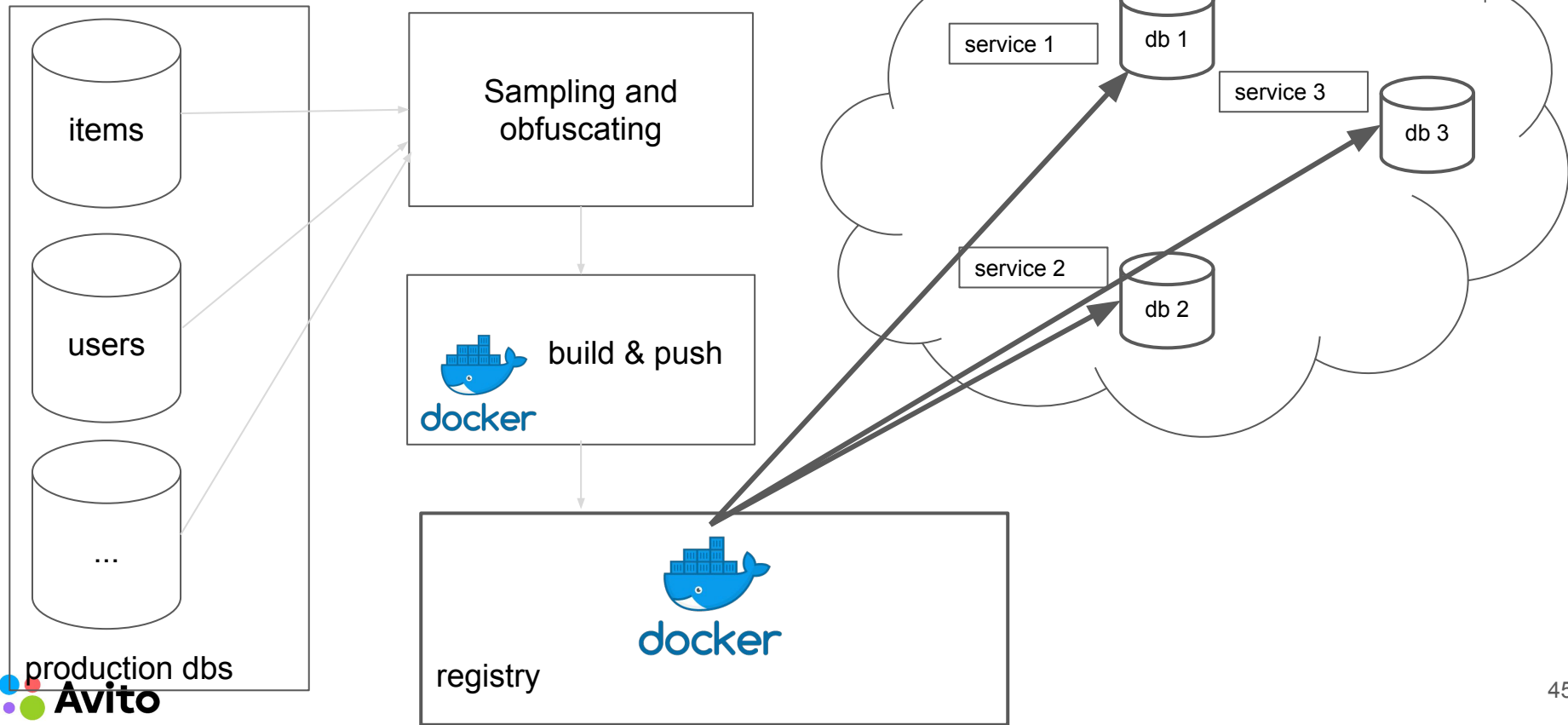
Dev and test environment



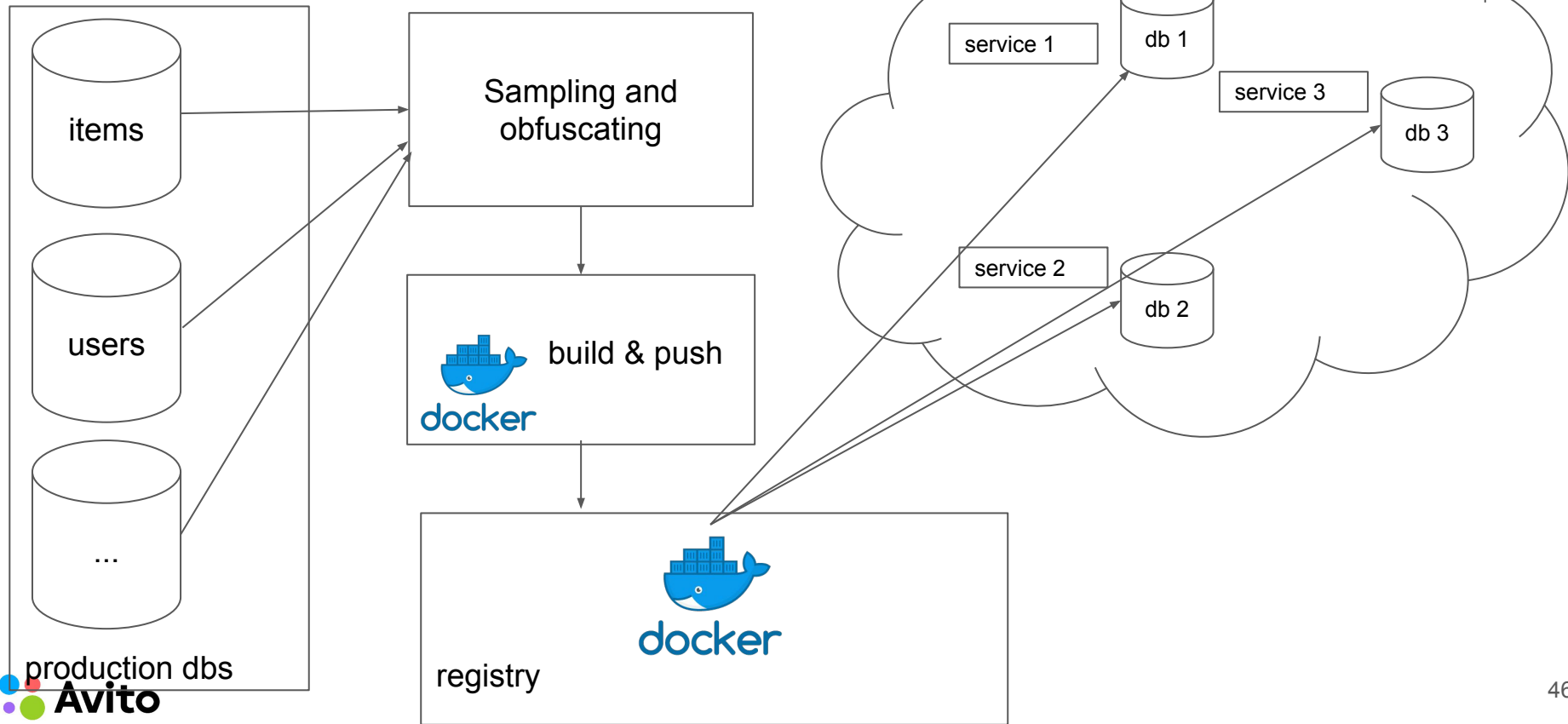
Dev and test environment



Dev and test environment



Dev and test environment

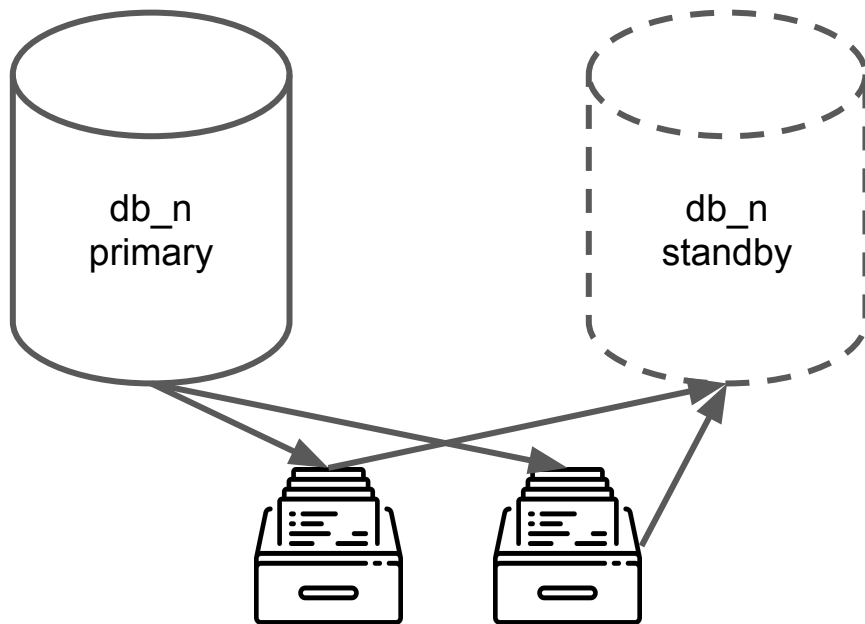


Plan

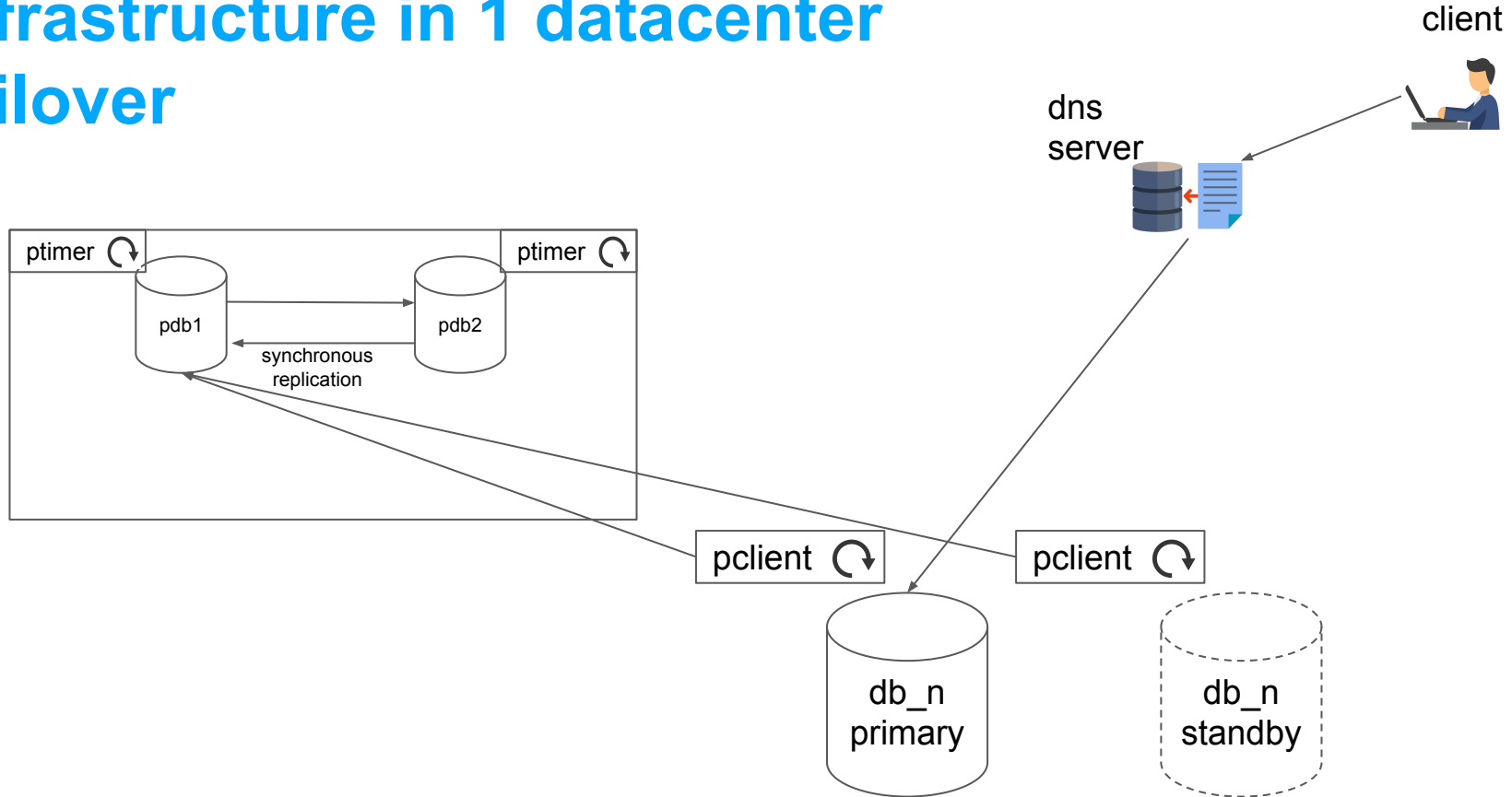
1. Evolution of monolith architecture
2. Migration to microservice architecture
3. Integration & communication
4. Dev tools and environment
5. **Platform (DBaaS in 3 Datacenters)**

Infrastructure in 1 datacenter archive & replication

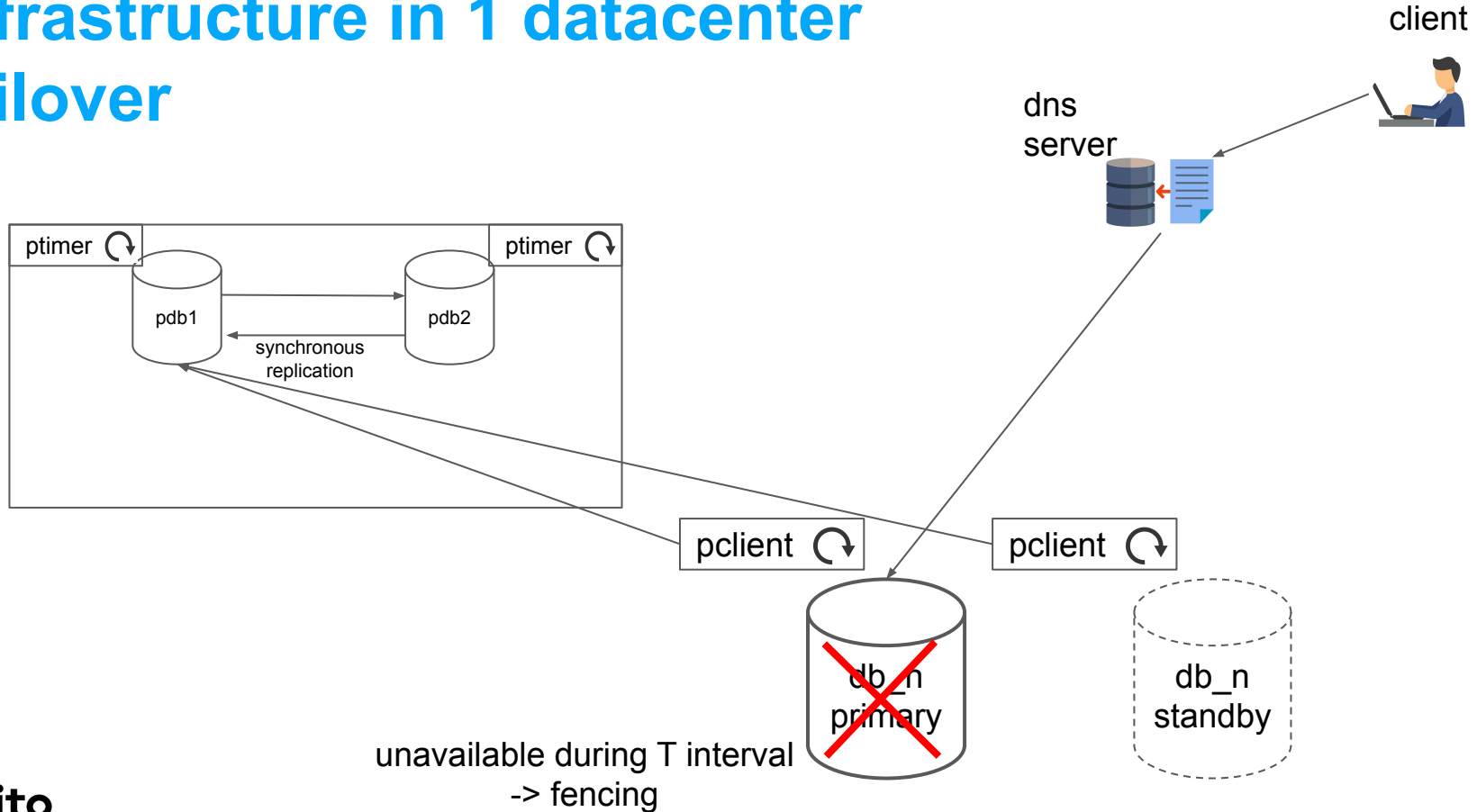
https://github.com/avito-tech/dba-utils/tree/master/pg_archive2



Infrastructure in 1 datacenter failover

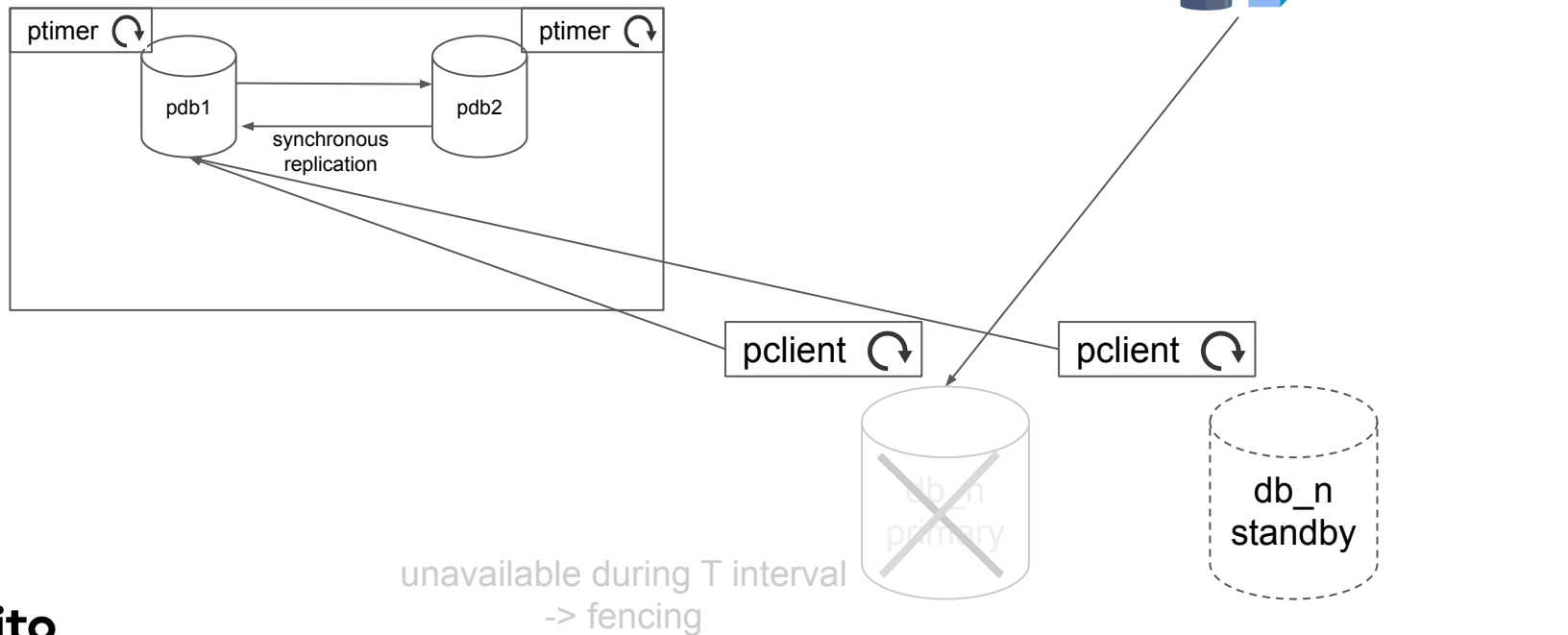


Infrastructure in 1 datacenter failover

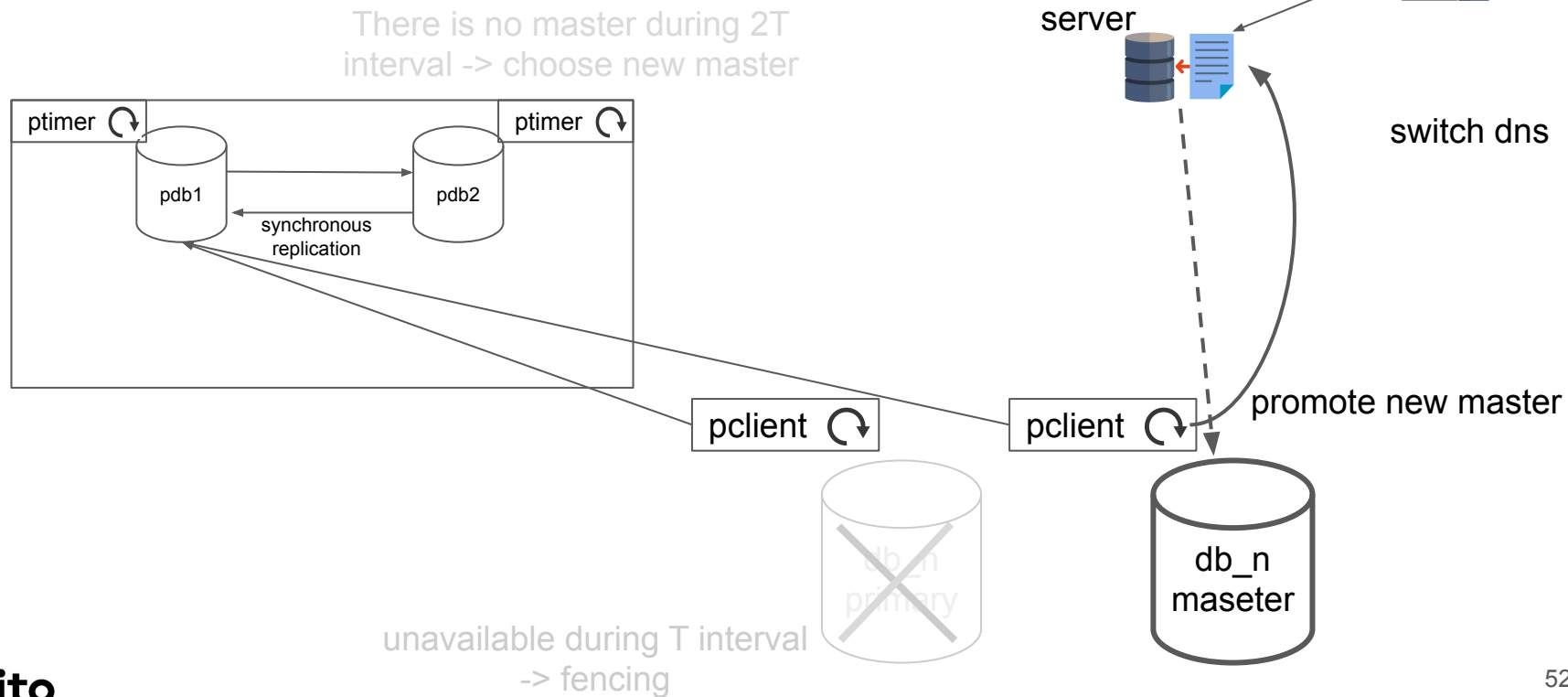


Infrastructure in 1 datacenter failover

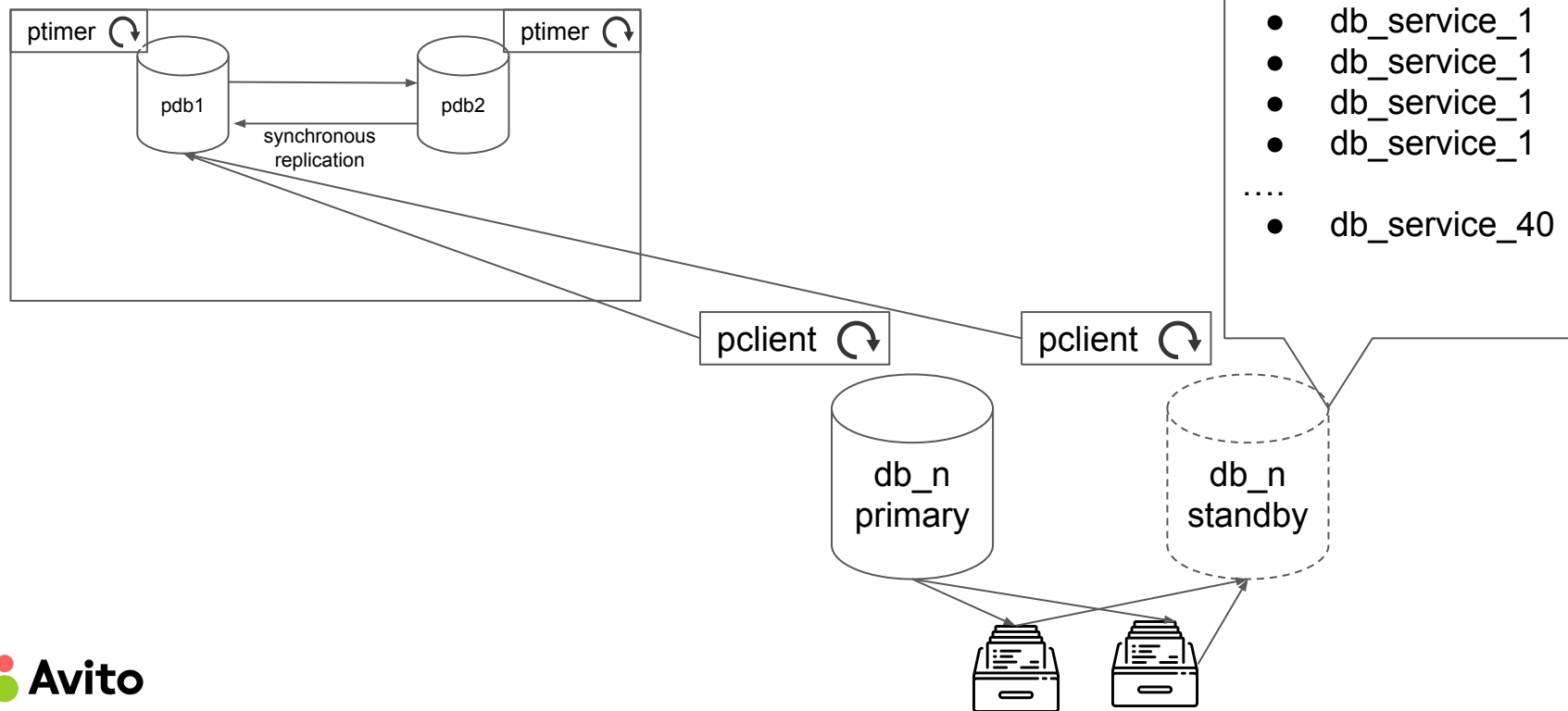
There is no master during 2T
interval -> choose new master



Infrastructure in 1 datacenter failover



Infrastructure in 1 datacenter



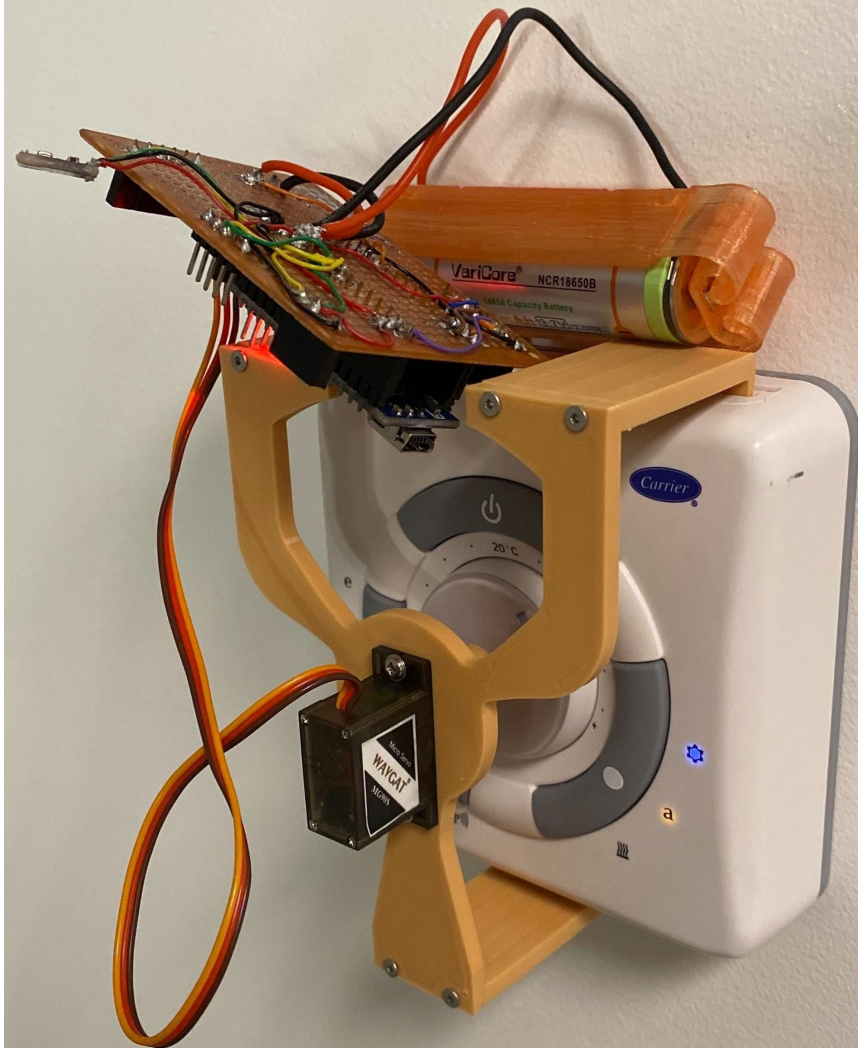
Problems with noisy neighbours

1. OOM
2. Hot data set out of cache
3. CPU/network/IO bound queries
4. Long transactions
5. ...

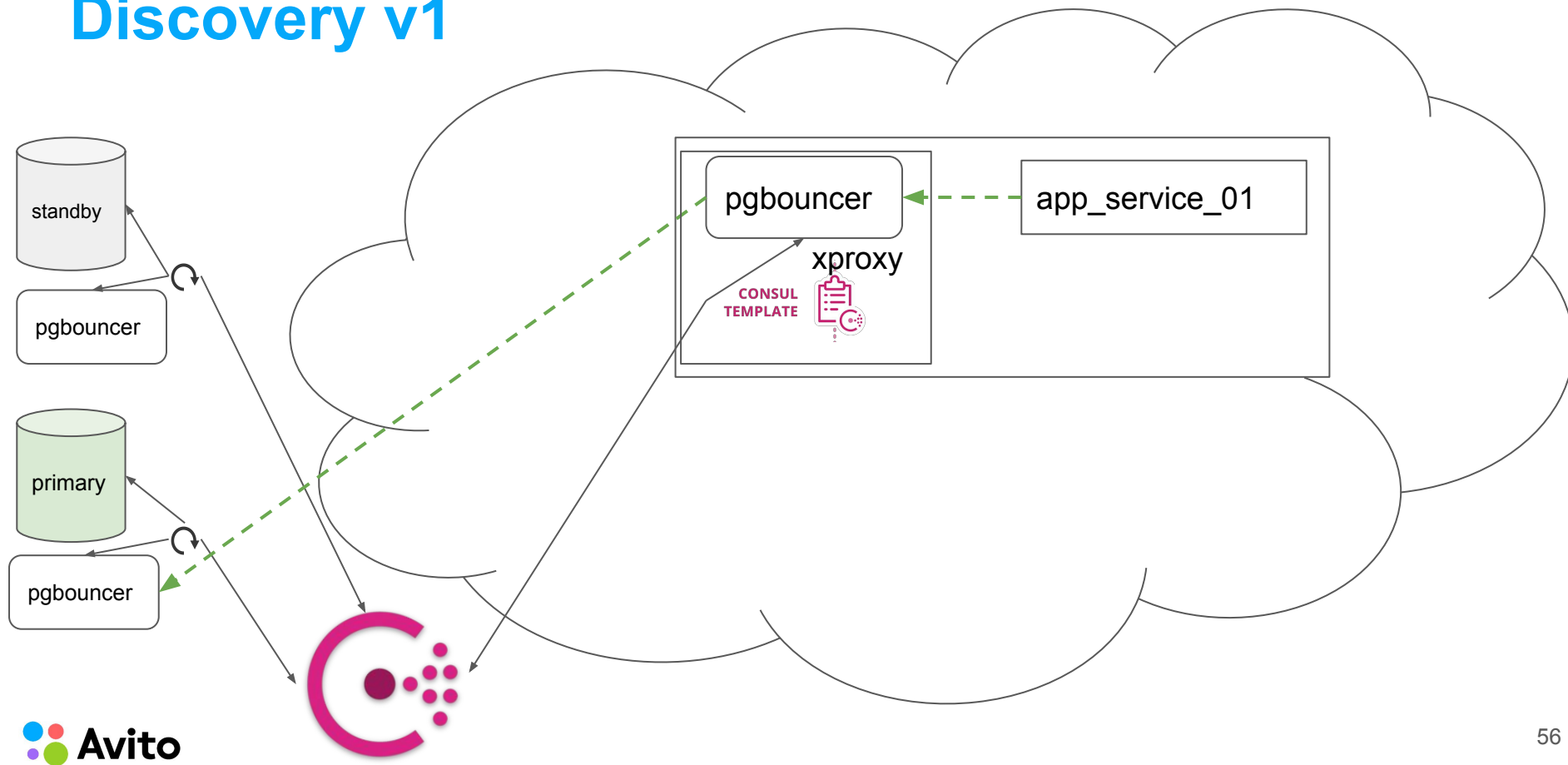


DBaaS

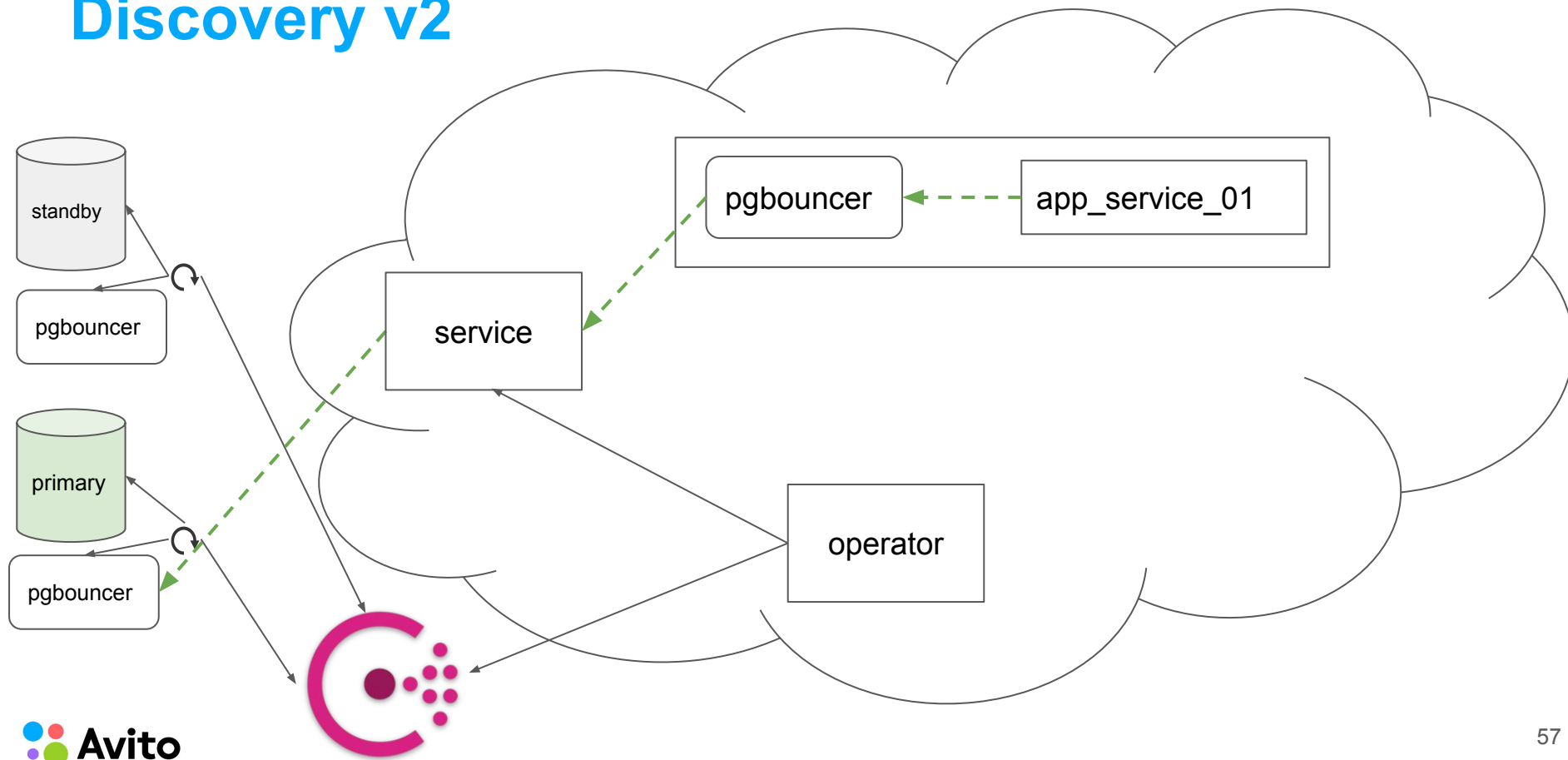
1. Database discovery
2. Archive
3. Autofailover
4. Limits and fully-guaranteed resources.



Discovery v1

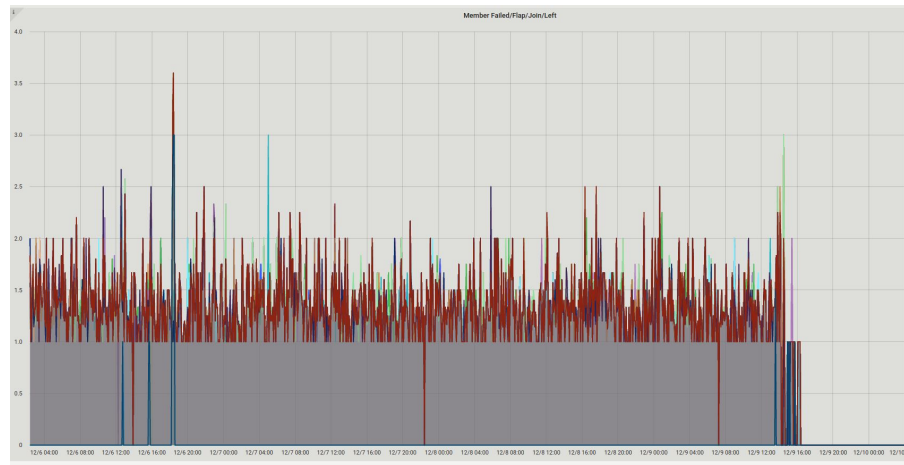
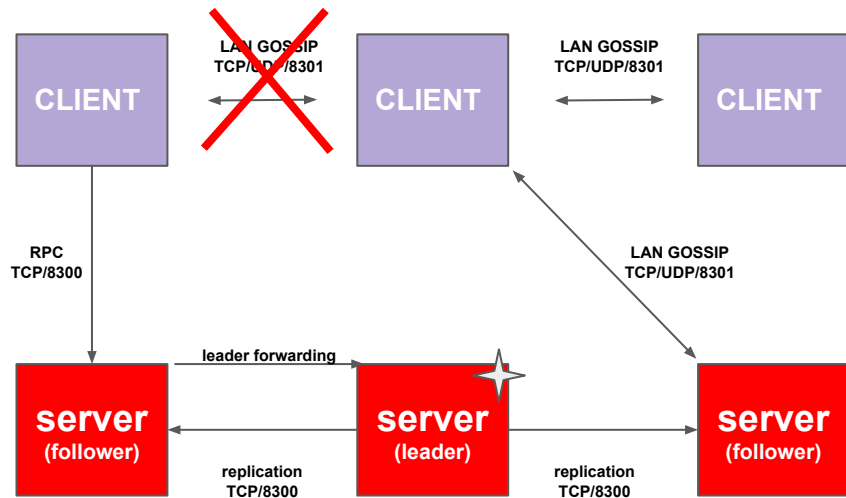


Discovery v2



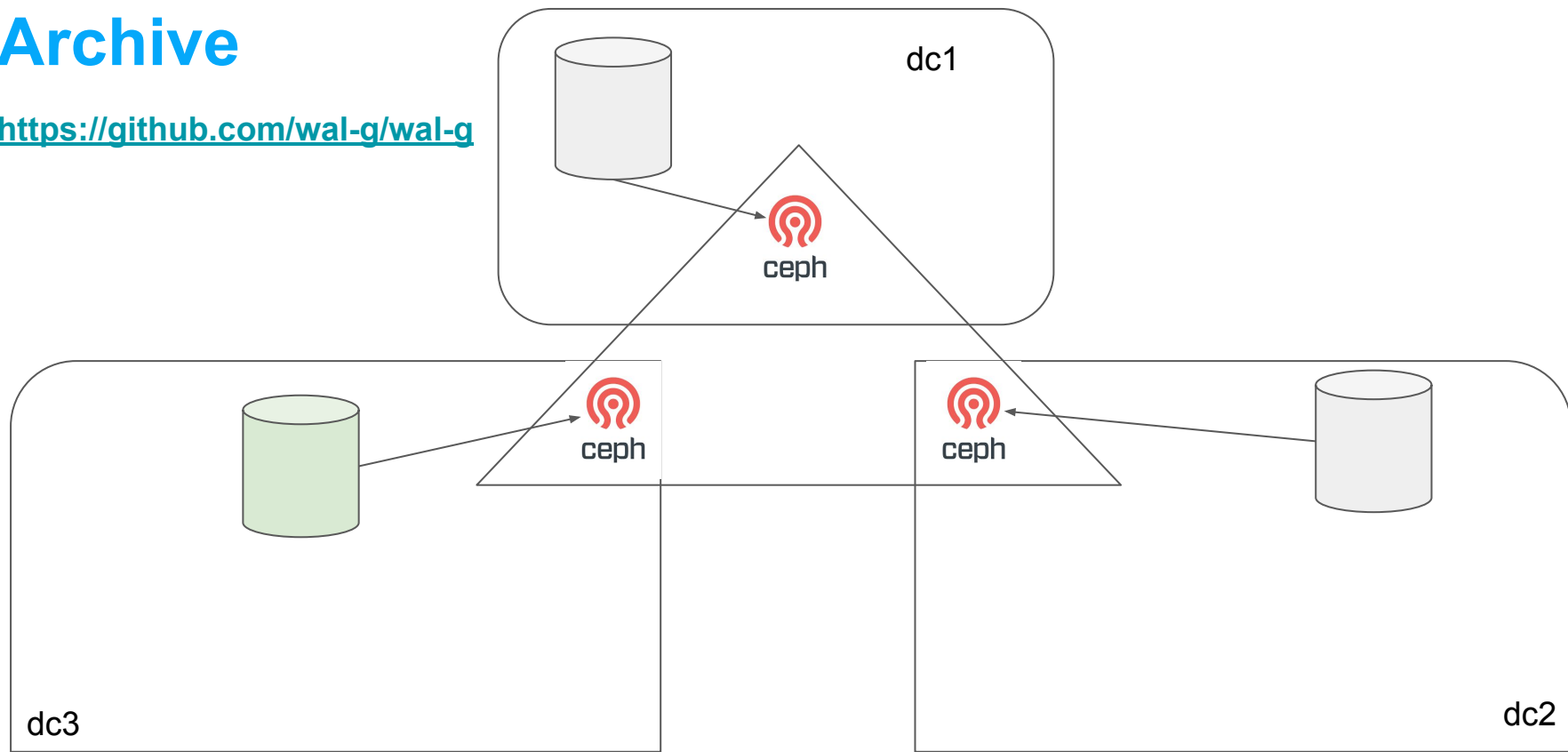
Caveats

even 2 nodes can successfully play the mafia game



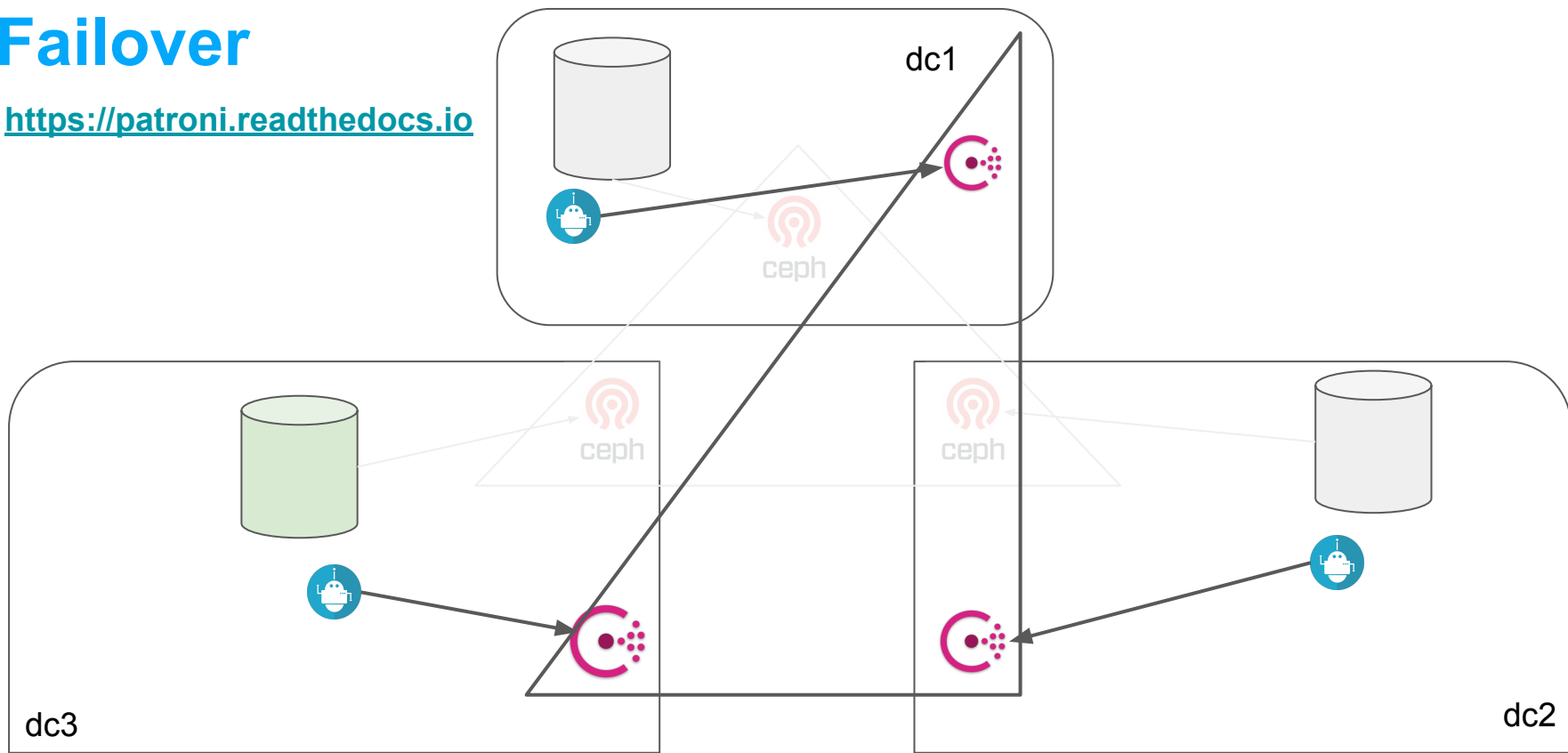
Archive

<https://github.com/wal-g/wal-g>



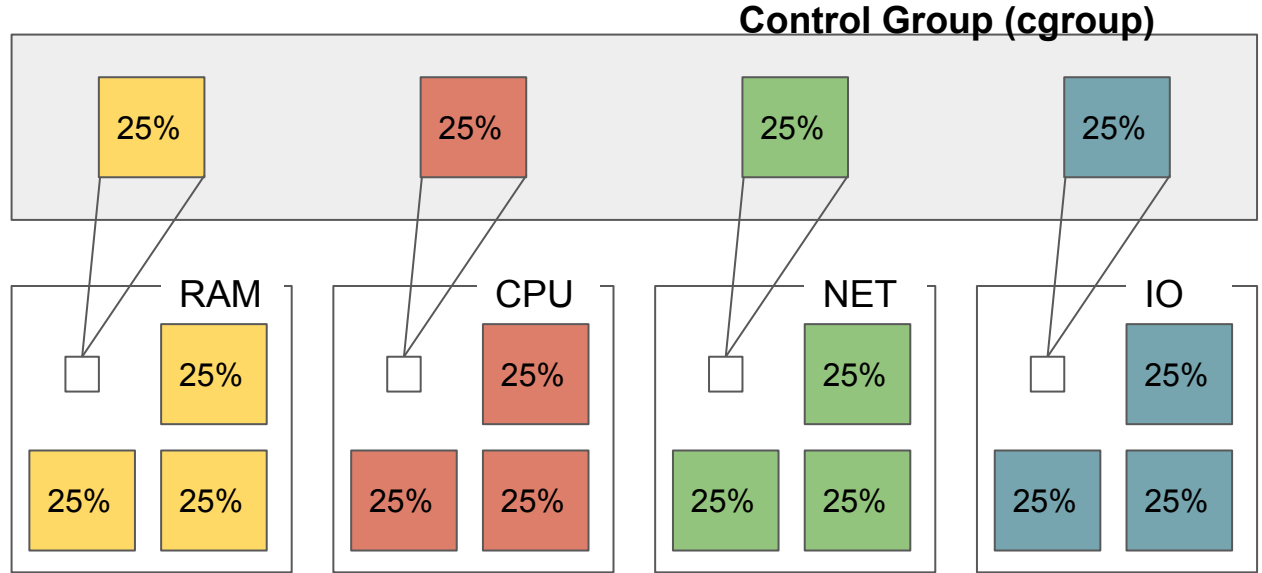
Failover

<https://patroni.readthedocs.io>

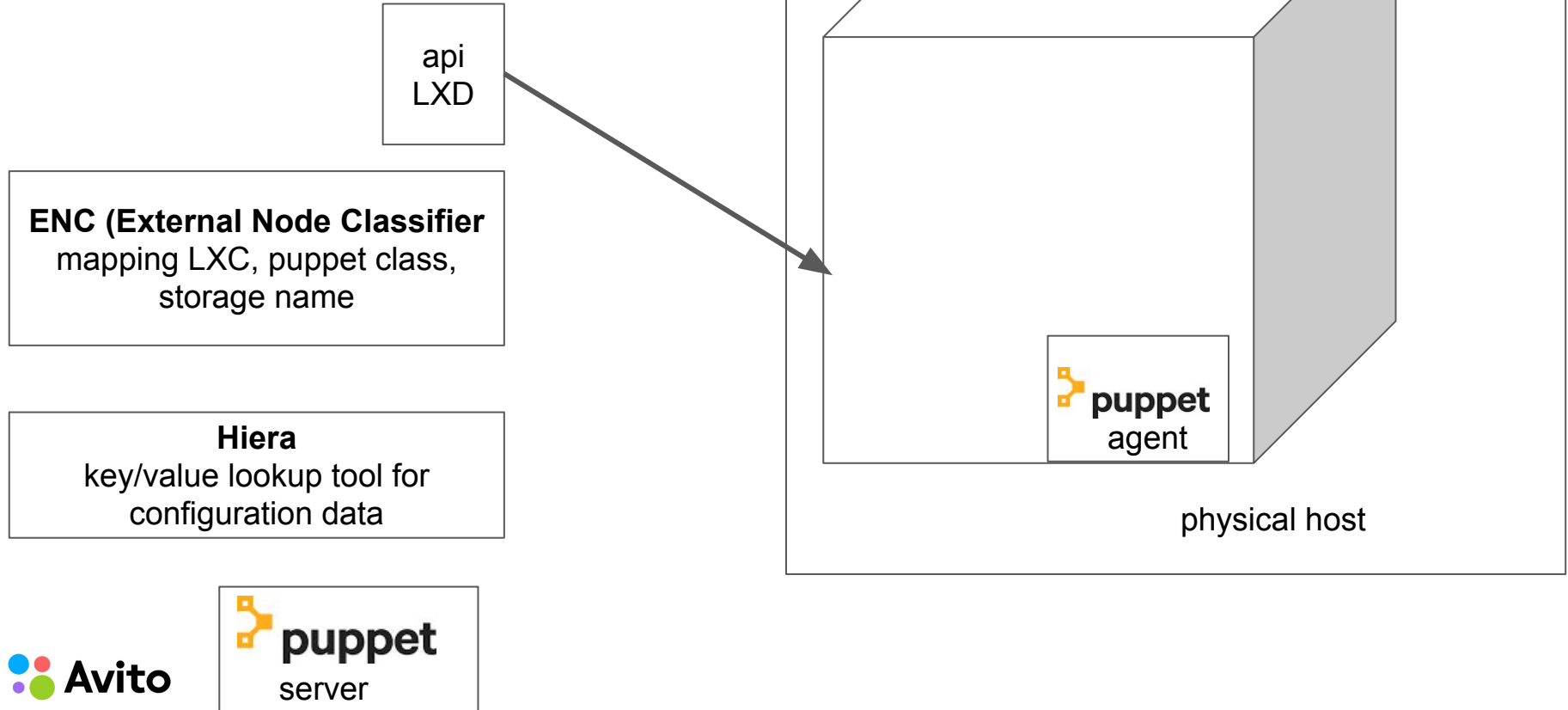


Limits and fully-guaranteed resources cgroups

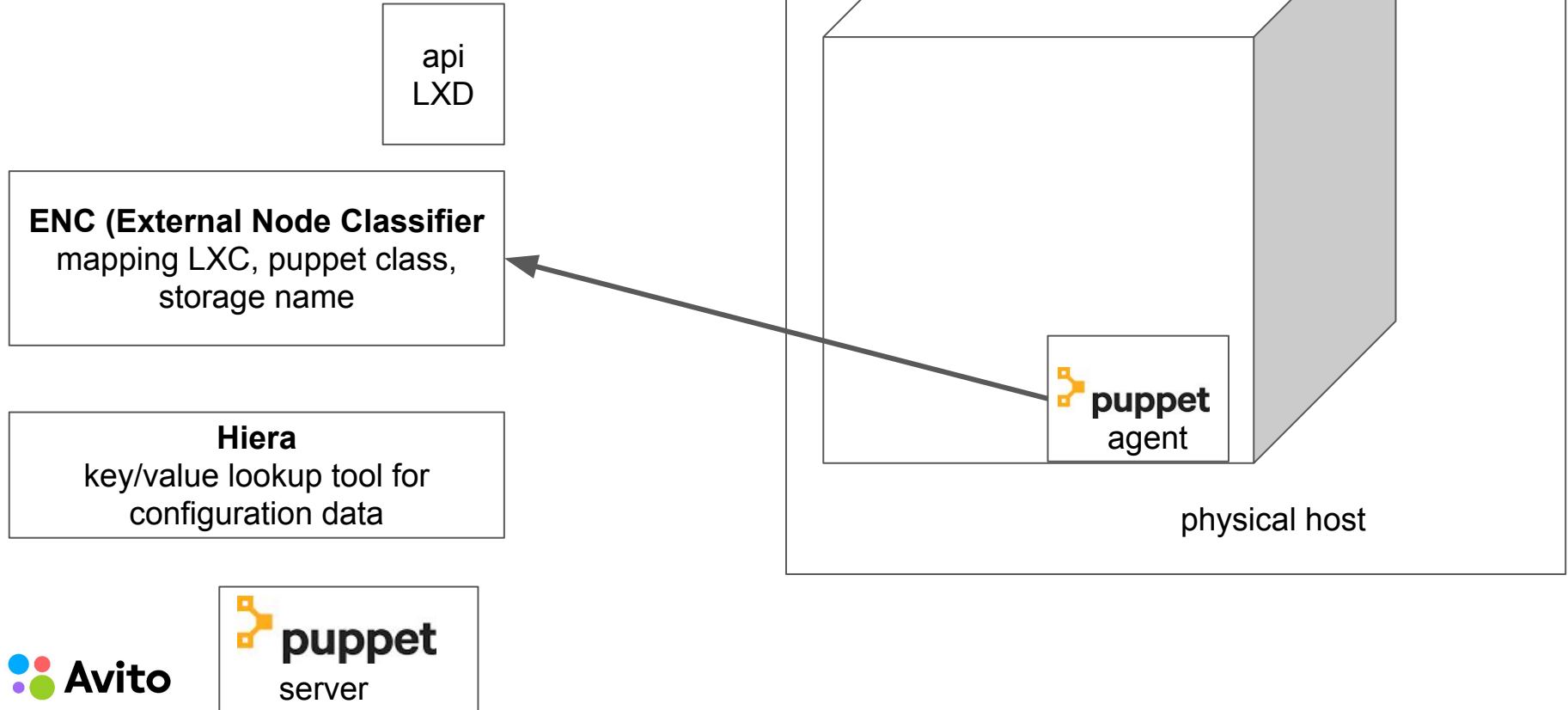
- CPU
- storage size
- network throughput
- IO limiting



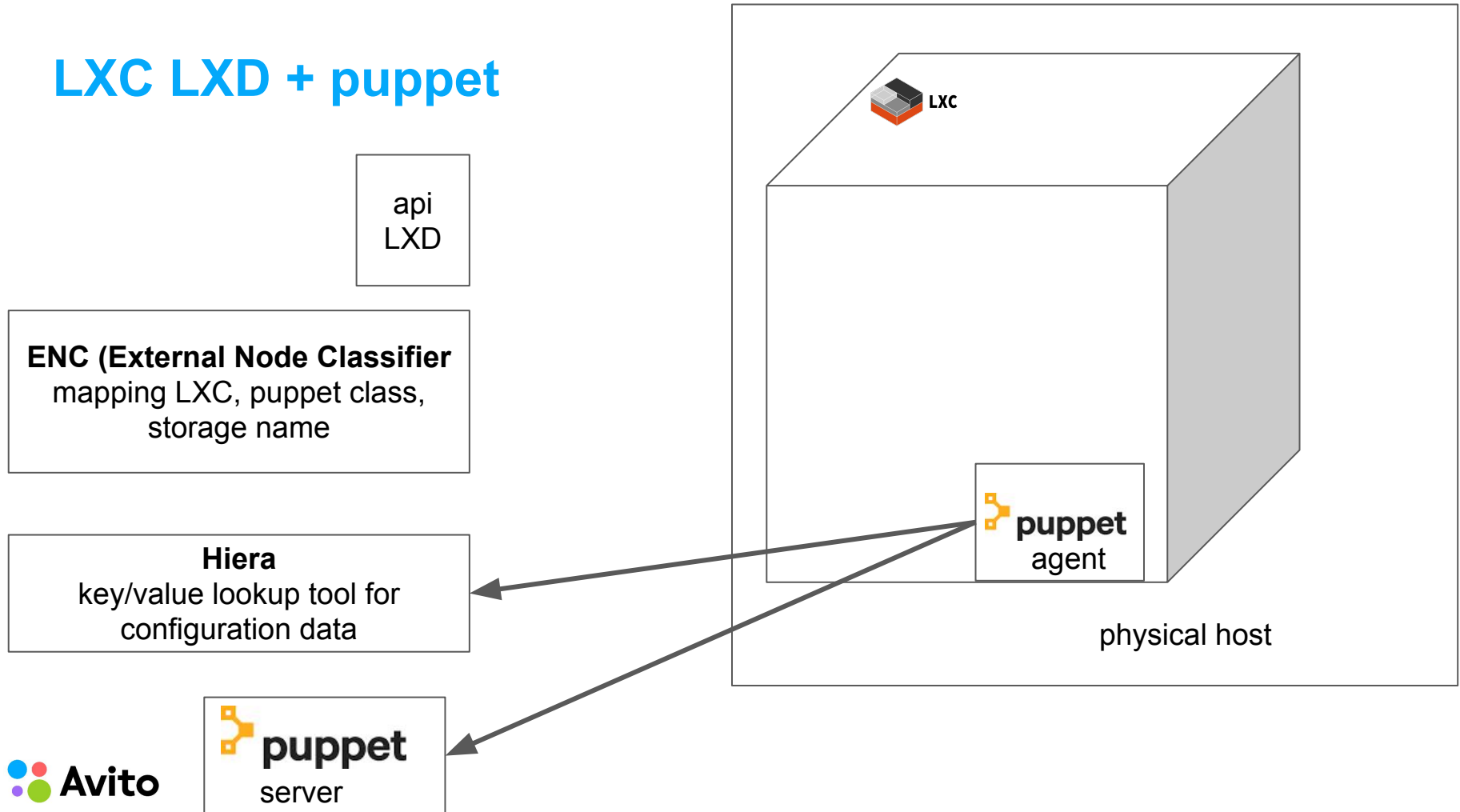
LXC LXD + puppet



LXC LXD + puppet



LXC LXD + puppet

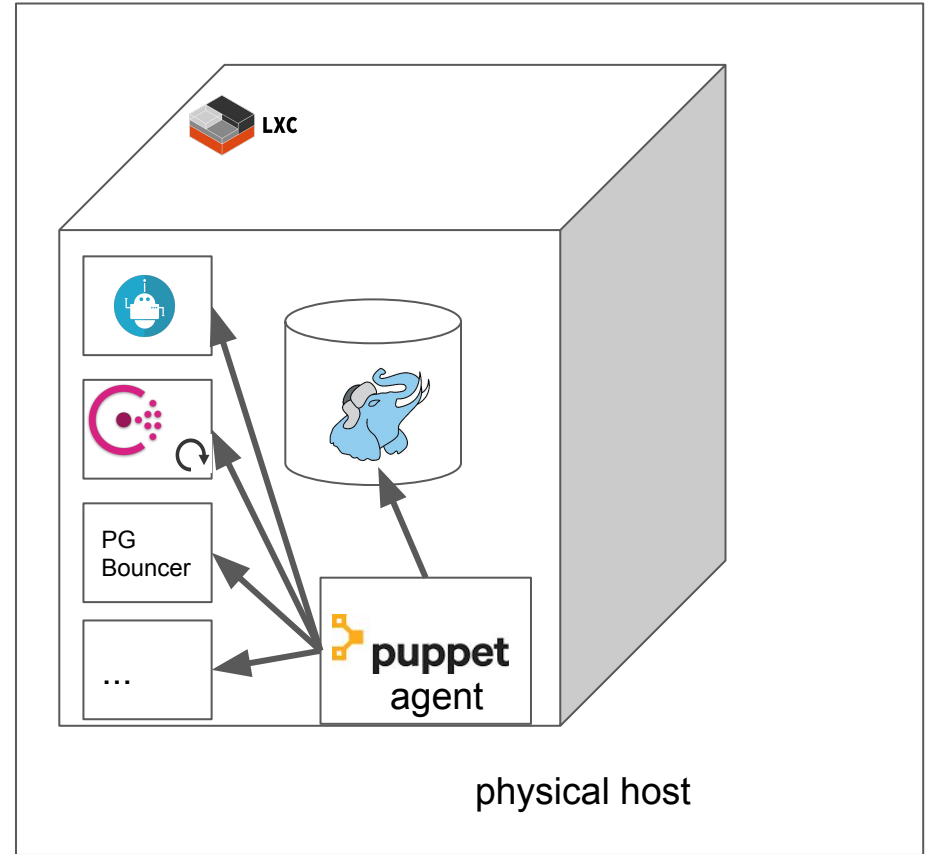


LXC LXD + puppet

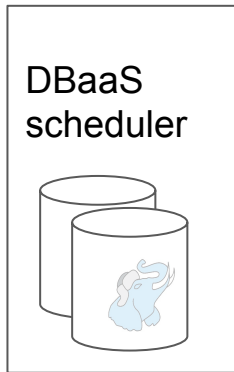
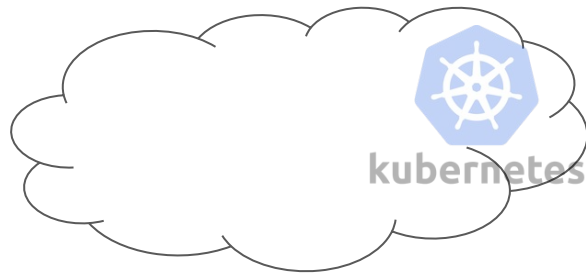
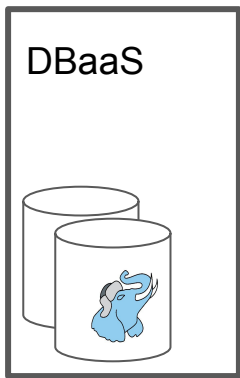
api
LXD

ENC (External Node Classifier)
mapping LXC, puppet class,
storage name

Hiera
key/value lookup tool for
configuration data



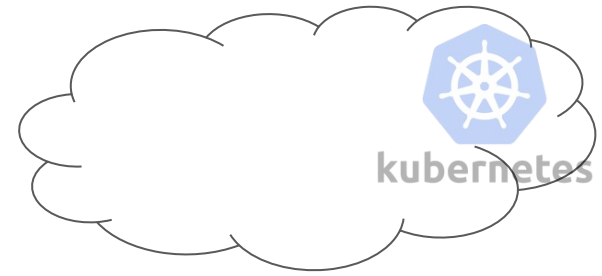
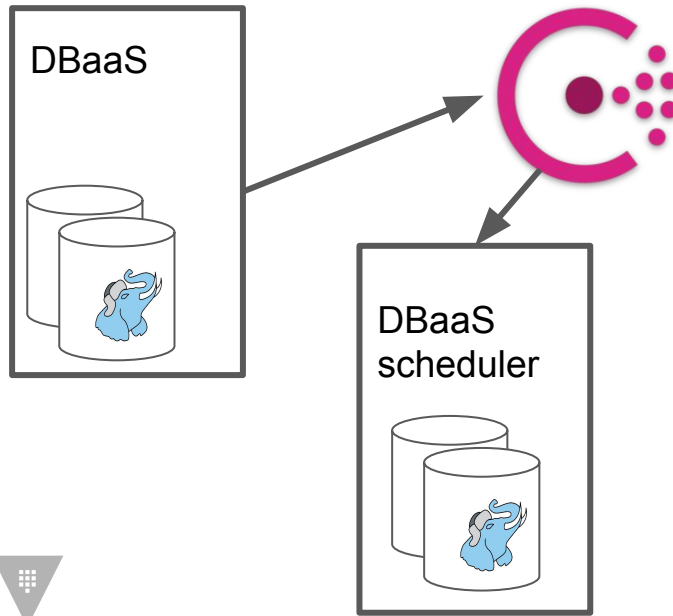
DBaaS



Hiera
key/value
lookup tool
for
configuratio
n data

**ENC (External
Node
Classifier**
mapping LXC,
puppet class,
storage name

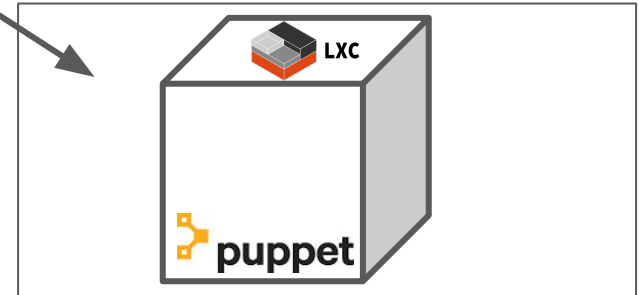
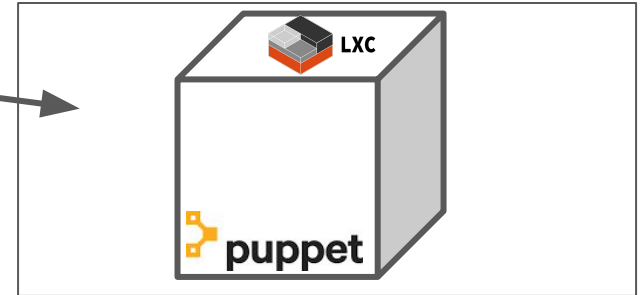
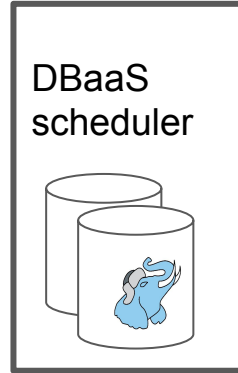
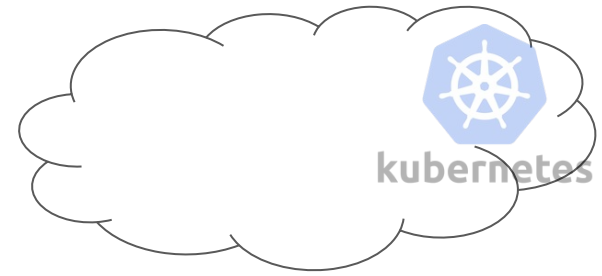
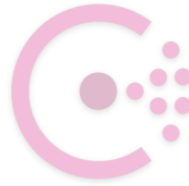
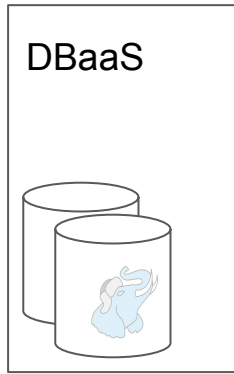
DBaaS



Hiera
key/value
lookup tool
for
configuratio
n data

**ENC (External
Node
Classifier**
mapping LXC,
puppet class,
storage name

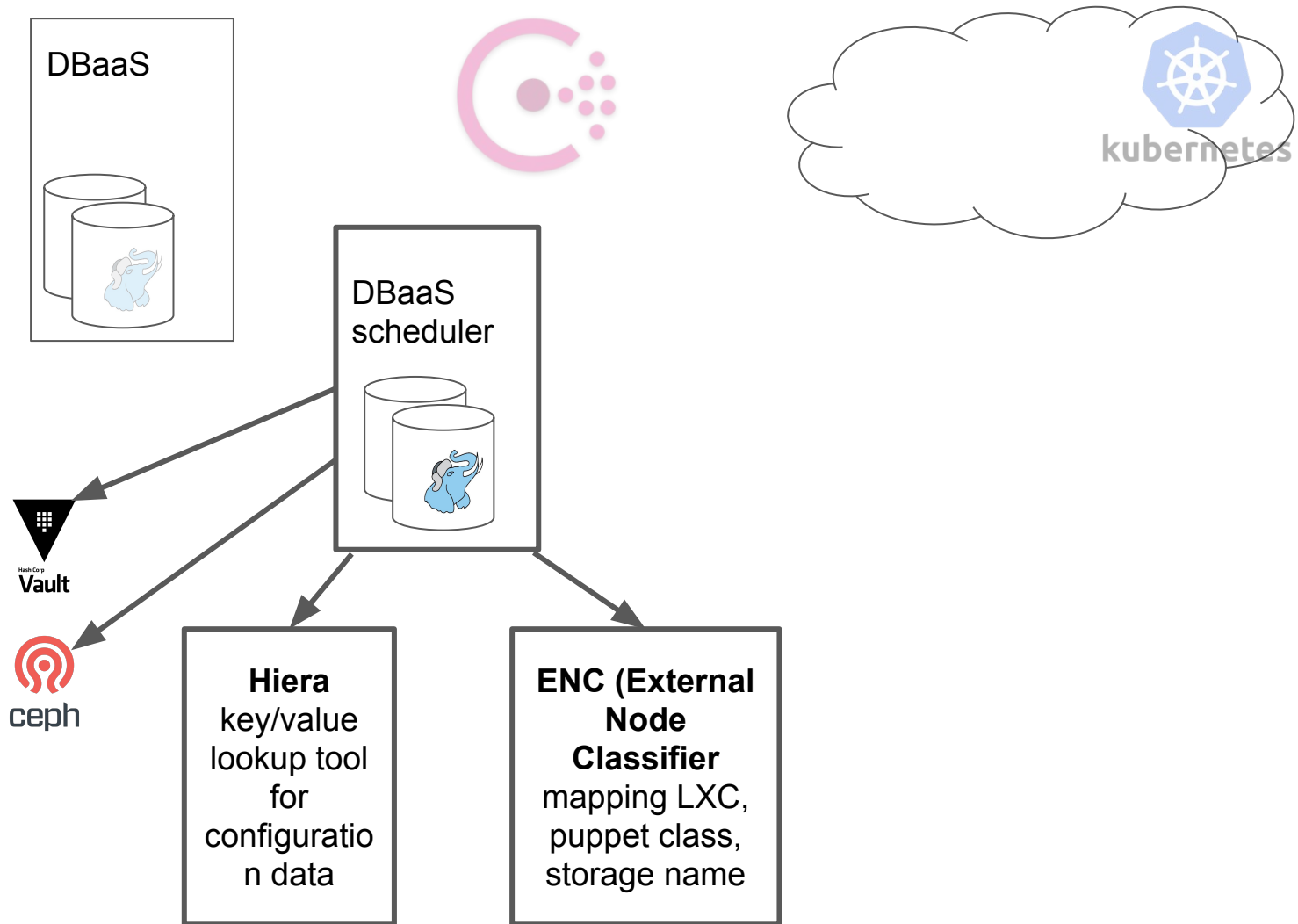
DBaaS



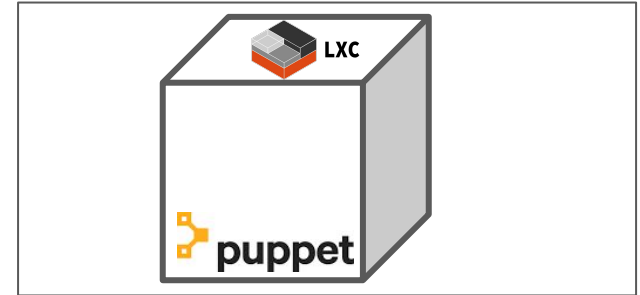
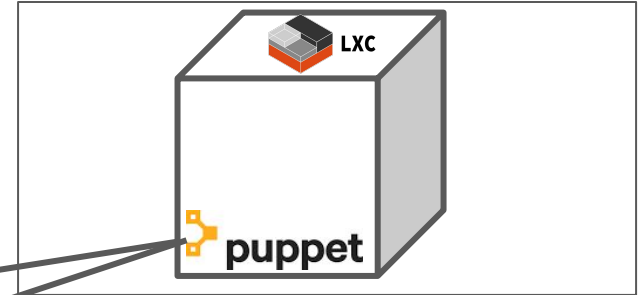
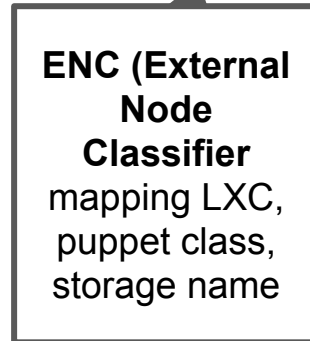
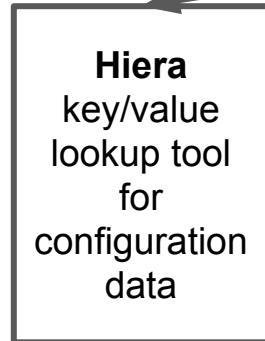
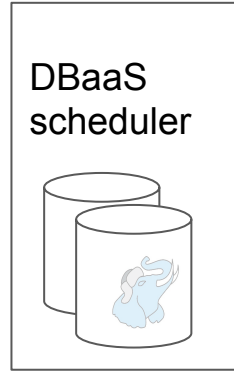
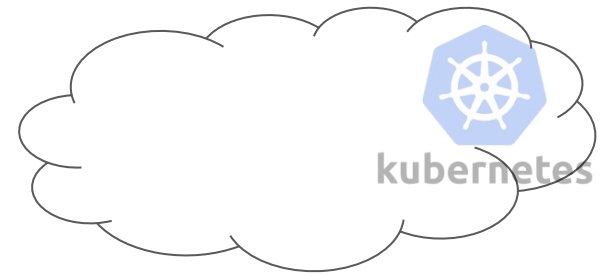
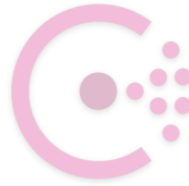
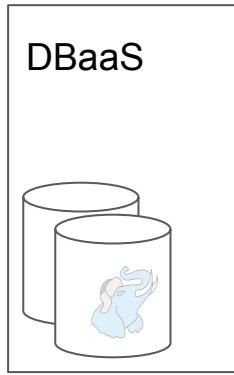
Hiera
key/value
lookup tool
for
configuration
data

**ENC (External
Node
Classifier**
mapping LXC,
puppet class,
storage name

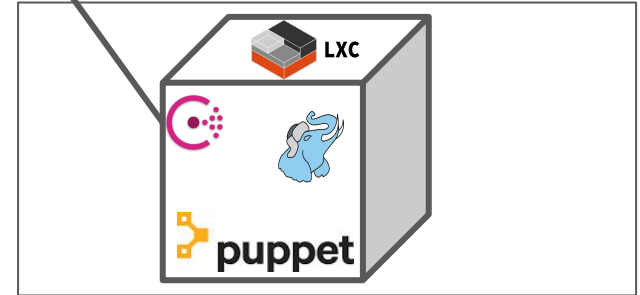
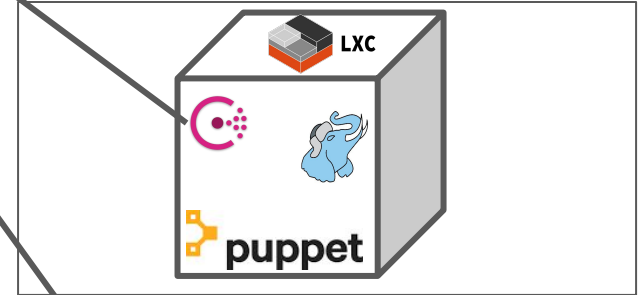
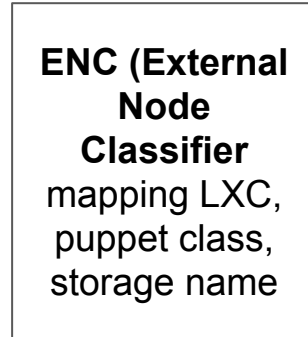
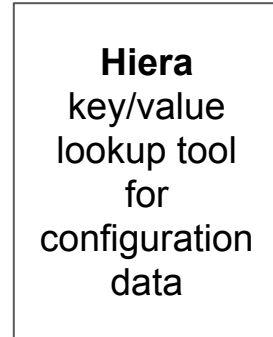
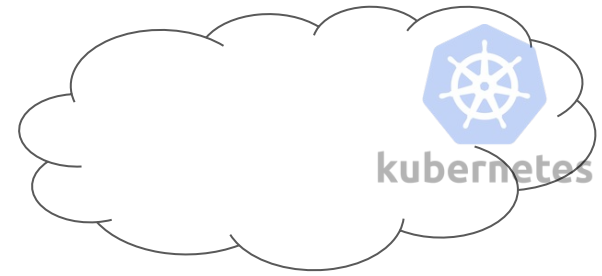
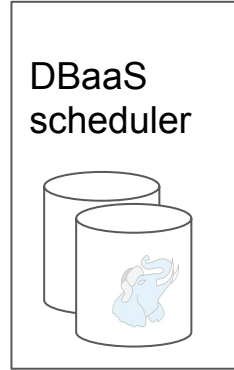
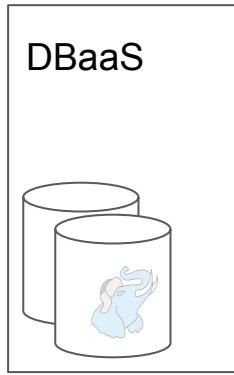
DBaaS



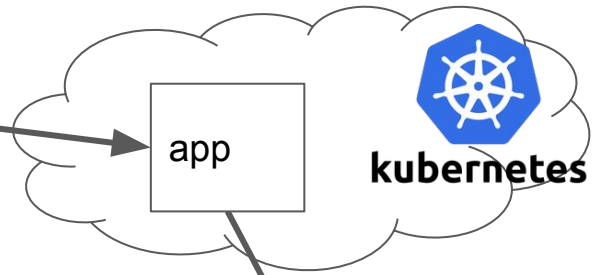
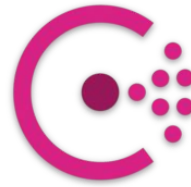
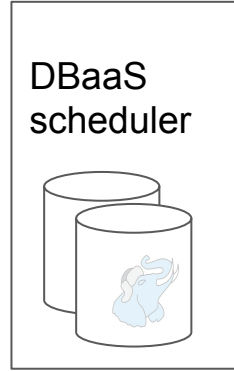
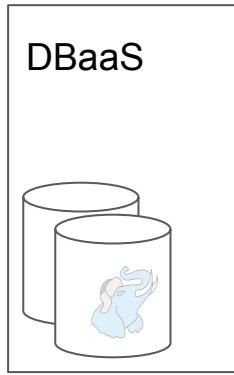
DBaaS



DBaaS

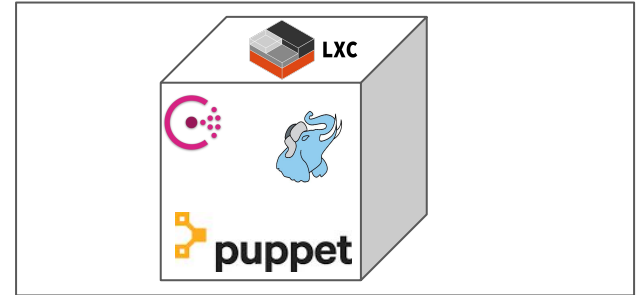
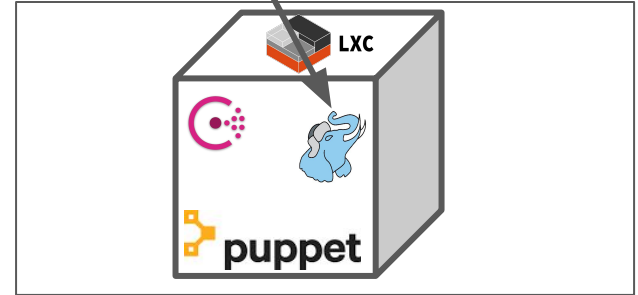


DBaaS



Hiera
key/value
lookup tool
for
configuration
data

**ENC (External
Node
Classifier**
mapping LXC,
puppet class,
storage name



Team evolution

OLTP Team:

- business logic
- integration
- code review
- code rules
- platform tasks
- database administration routines



Team evolution

OLTP Team:

- business logic
- integration
- code review
- code rules
- platform tasks
- database administration routines

PG SWAT team:

- monolith business logic
- integration
- code review
- code rules

DBA team:

- platform tasks
- database administration routines



Team evolution

OLTP Team:

- business logic
- integration
- code review
- code rules
- platform tasks
- database administration routines

PG SWAT team:

- monolith business logic
- integration
- code review
- code rules

DBA team:

- platform tasks
- database administration routines

PostgreSQL experts
in cross-functional
teams

PG SWAT team:

- core services
- integration
- code review
- code rules

DBA team:

- platform tasks
- database administration routines
- **DBaaS**



Avito

monolith

splitting monolith

microservices

Team evolution

OLTP Team:

- business logic
- integration
- code review
- code rules
- platform tasks
- database administration routines

PG SWAT team:

- monolith business logic
- integration
- code review
- code rules

DBA team:

- platform tasks
- database administration routines

PostgreSQL experts
in cross-functional
teams

PG SWAT team:

- core services
- integration
- code review
- code rules

DBA team:

- platform tasks
- database administration routines
- **DBaaS**



Wishlist

Multi DC k8s / multi cluster operator

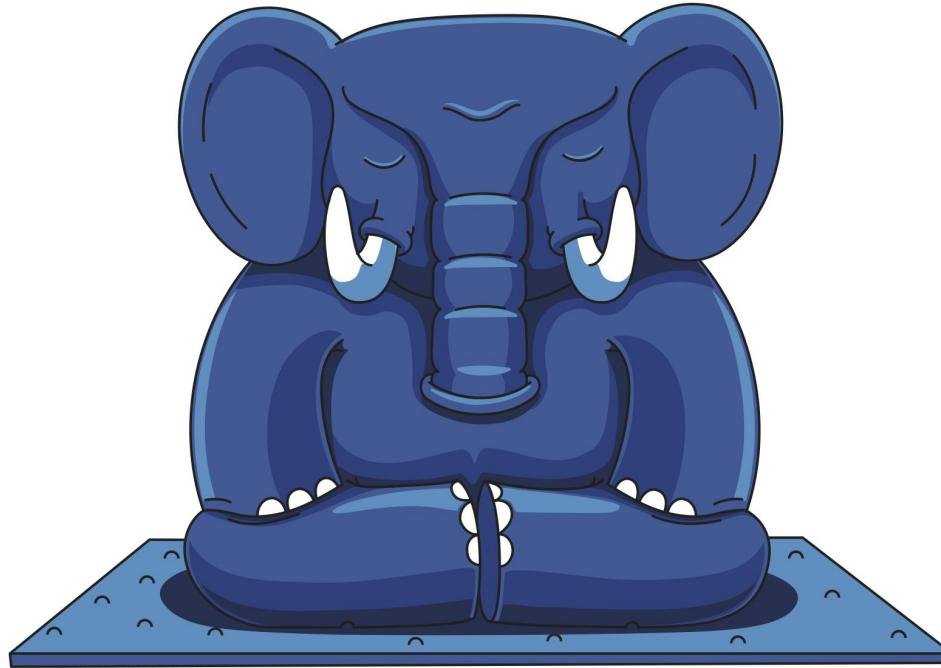
Code coverage

CICD tools and improvements

Standby improvements

Logical replication improvements (recovery, long transactions, parallel apply, ddl

Thank you!



<https://tech.avito.ru>

<https://github.com/avito-tech>

<https://www.facebook.com/evteev.k.s>