



Standby in Production

Konstantin Evteev

Ottawa 2019

Avito has a nationwide audience

Every month more than 35 million people use Avito

Pageviews per month:

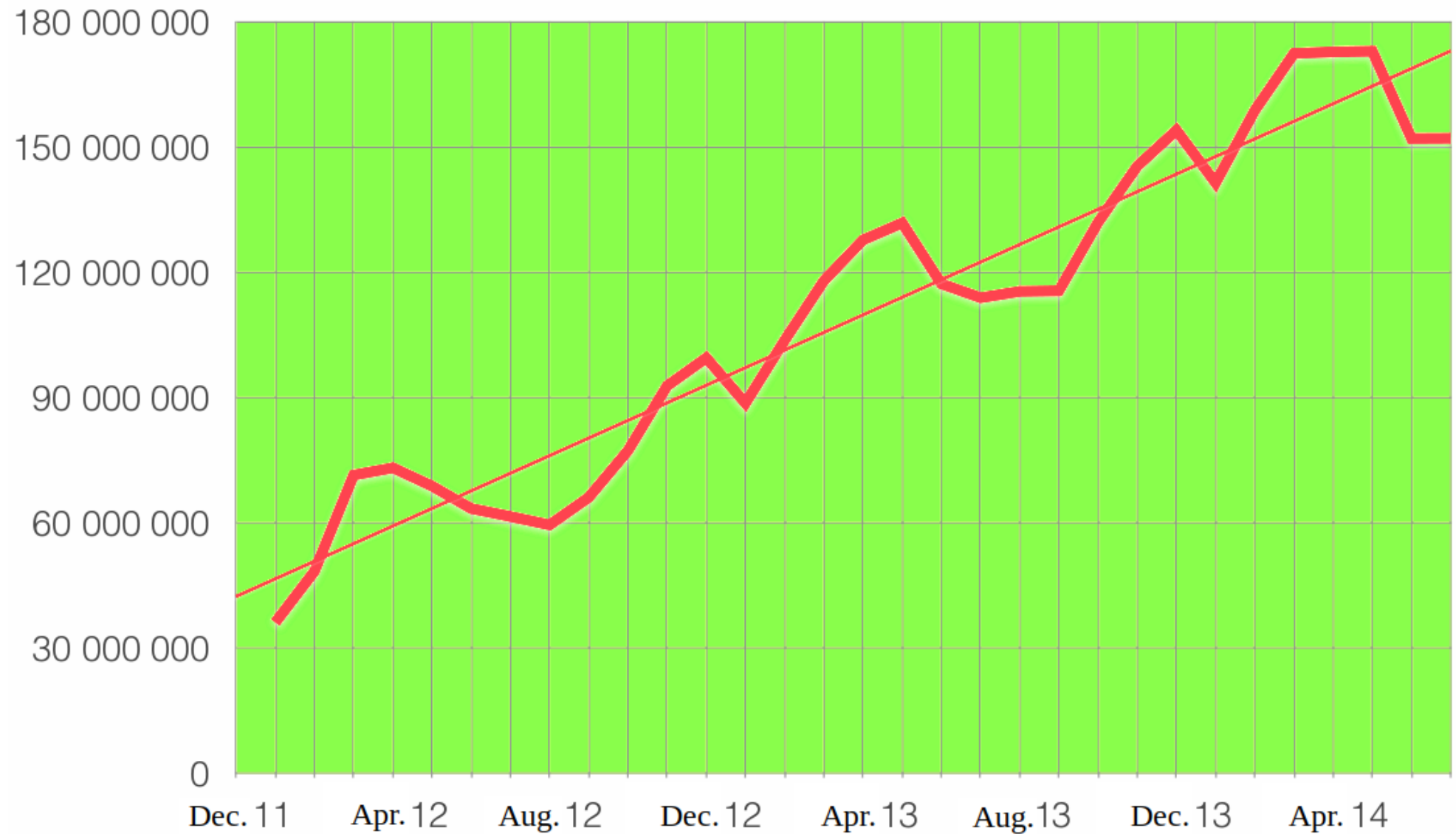


Statistics: February 2018

Some figures

- **600+ servers;**
- **4,5 Gbit/sec TX, 2 Gbit/sec RX without static;**
- **about a million queries per minute to the backend;**
- **270TB of images;**
- **>20 TB in Postgres on 100 nodes:**
 - **7 - 8K TPS on most nodes;**
 - **The largest - 20k TPS, 5 TB;**

pageviews

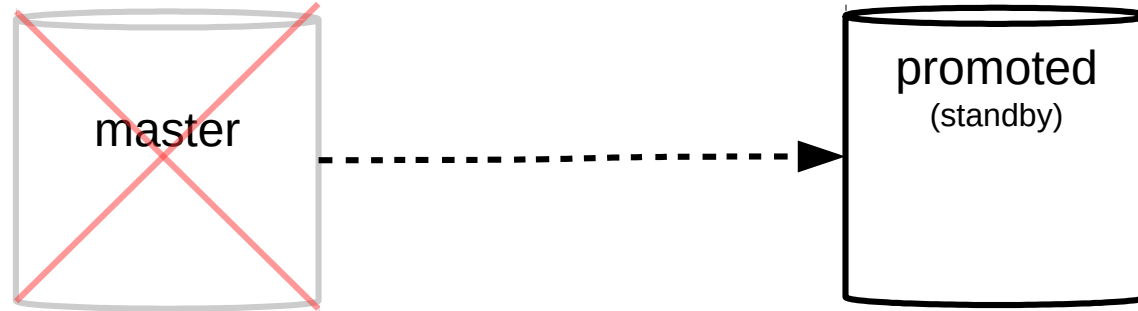


Contents

- a few words about standby and its history in general
- problems and solutions in replication based horizontal scale-out;
- Avito's solution to avoid stale reads from replica implementation;
- possible pitfalls while using standby with high request rate, applying DDL, receiving WAL files from archive;
- handling issues with technique of using several standbys in production and routing queries between them;
- logical replication based scaling example to compare
- conclusions and standby major upgrade features

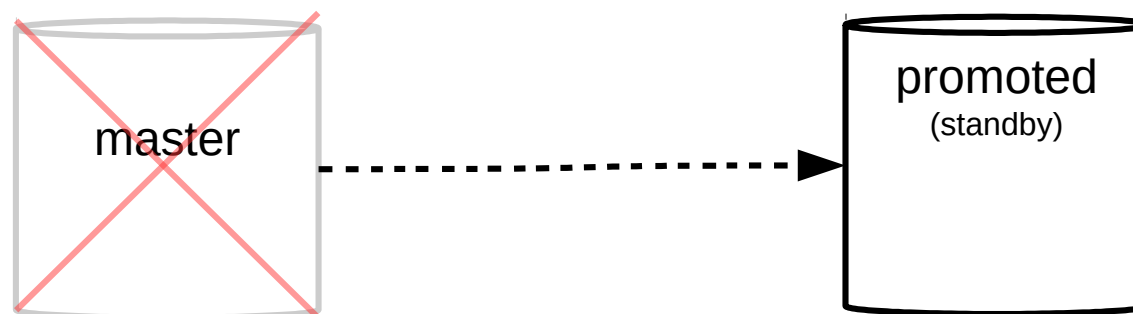
Standby

1 High Availability

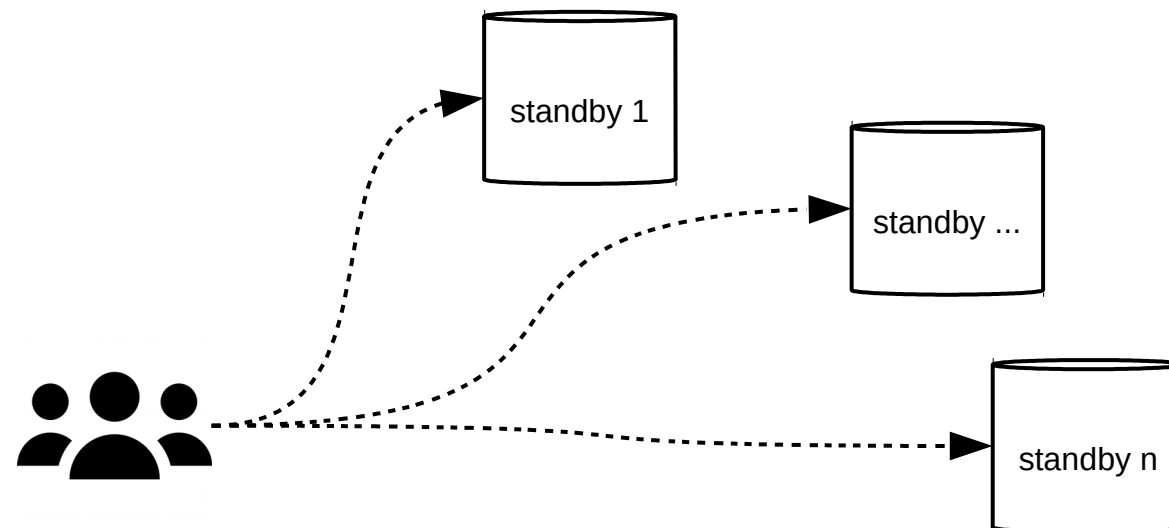


Standby

1 High Availability



2 Scaling



History

2000: RServ



History

2000: RServ

2001: PostgreSQL 7.1: write-ahead log



History

2000: RServ

2001: PostgreSQL 7.1: write-ahead log

2004: Slony



History

2000: RServ

2004: Slony

2001: PostgreSQL 7.1: write-ahead log

2005: PostgreSQL 8.0: point-in-time recovery



History

2000: RServ

2001: PostgreSQL 7.1: write-ahead log

2004: Slony

2005: PostgreSQL 8.0: point-in-time recovery

2007: SkyTools – PgQ & Londiste



History



2000: RServ

2001: PostgreSQL 7.1: write-ahead log

2004: Slony

2005: PostgreSQL 8.0: point-in-time recovery

2007: SkyTools – PgQ & Londiste

2008: 8.3 Easier administration of warm standby servers

PostgreSQL



[Download](#) · [Mailing Lists](#) · [Users Lounge](#) · [Hackers Lair](#) · [Documentation](#)

1. 2010: 9.0: hot standby, streaming replication

PostgreSQL



[Download](#) · [Mailing Lists](#) · [Users Lounge](#) · [Hackers Lair](#) · [Documentation](#)

1. 2010: 9.0: hot standby, streaming replication
2. 2011: 9.1: synchronous replication

PostgreSQL



[Download](#) · [Mailing Lists](#) · [Users Lounge](#) · [Hackers Lair](#) · [Documentation](#)

1. 2010: 9.0: hot standby, streaming replication
2. 2011: 9.1: synchronous replication
3. 2013: 9.3: sb can follow timeline switch

PostgreSQL



[Download](#) · [Mailing Lists](#) · [Users Lounge](#) · [Hackers Lair](#) · [Documentation](#)

1. 2010: 9.0: hot standby, streaming replication
2. 2011: 9.1: synchronous replication
3. 2013: 9.3: sb can follow timeline switch
4. 2014: 9.4: replication slots, logical decoding

PostgreSQL



[Download](#) · [Mailing Lists](#) · [Users Lounge](#) · [Hackers Lair](#) · [Documentation](#)

1. 2010: 9.0: hot standby, streaming replication
2. 2011: 9.1: synchronous replication
3. 2013: 9.3: sb can follow timeline switch
4. 2014: 9.4: replication slots, logical decoding
5. 2016: 9.6 multiple synchronous standbys, remote_apply

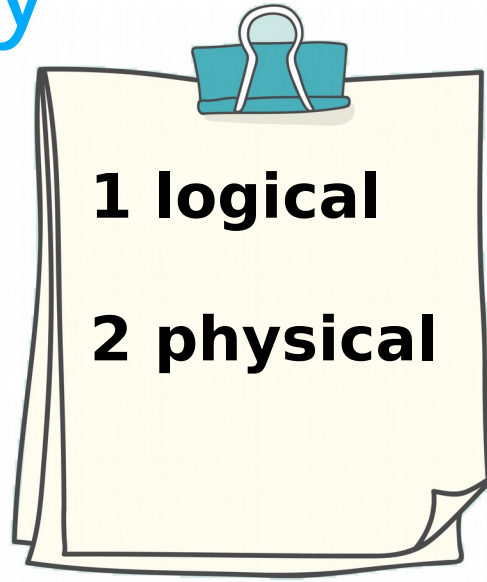
PostgreSQL



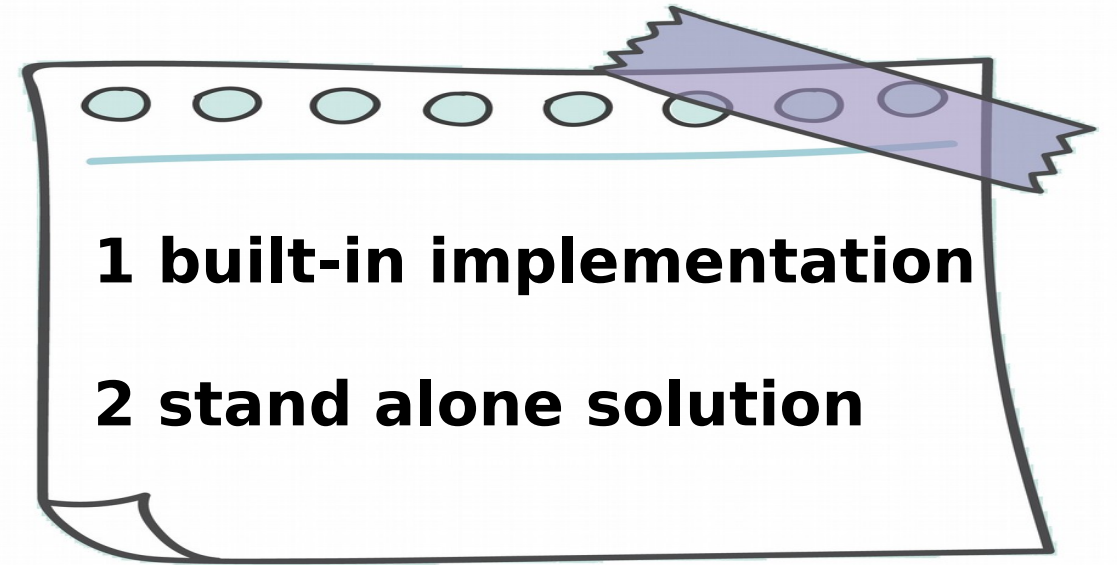
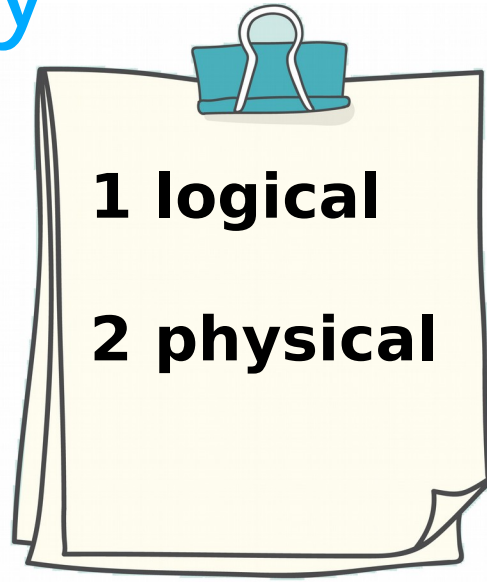
[Download](#) · [Mailing Lists](#) · [Users Lounge](#) · [Hackers Lair](#) · [Documentation](#)

1. 2010: 9.0: hot standby, streaming replication
2. 2011: 9.1: synchronous replication
3. 2013: 9.3: sb can follow timeline switch
4. 2014: 9.4: replication slots, logical decoding
5. 2016: 9.6 multiple synchronous standbys, remote_apply
6. 2017: 10: logical replication

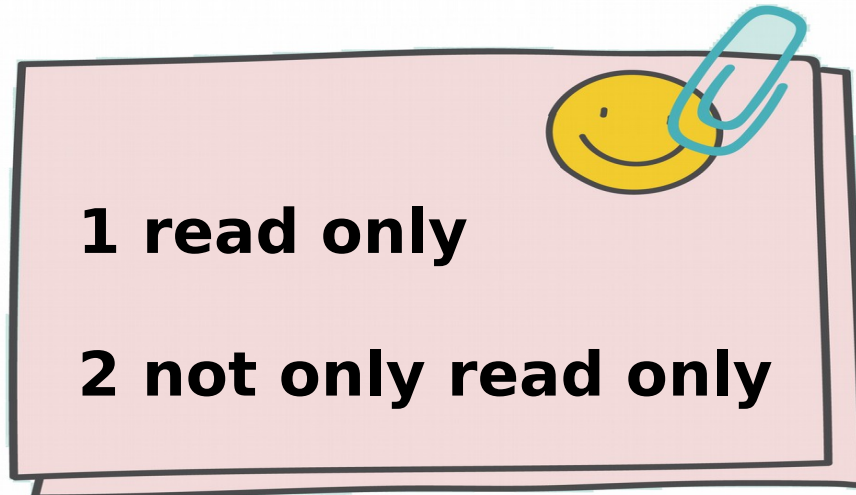
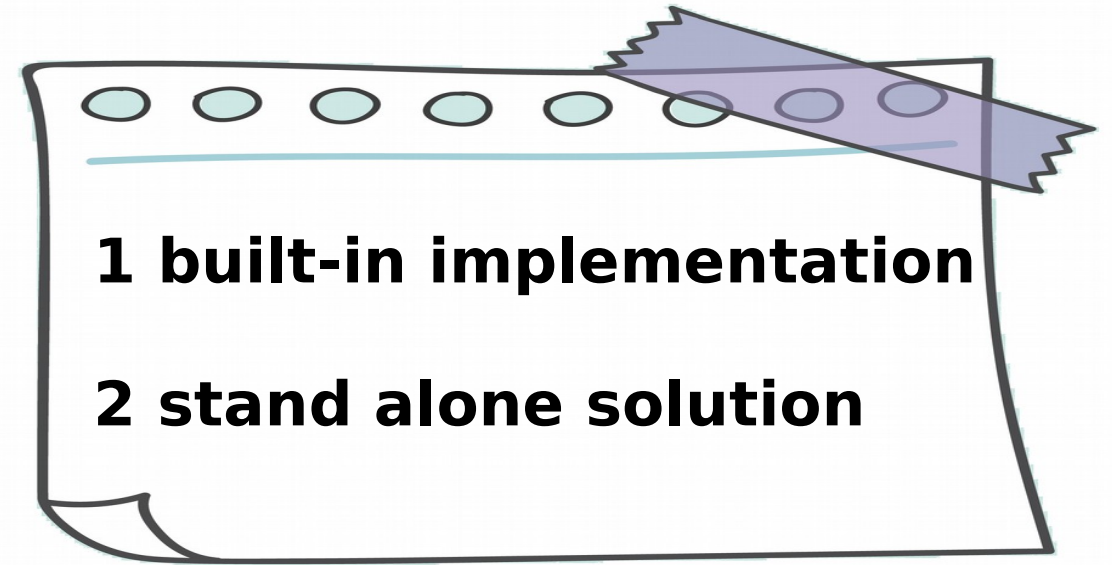
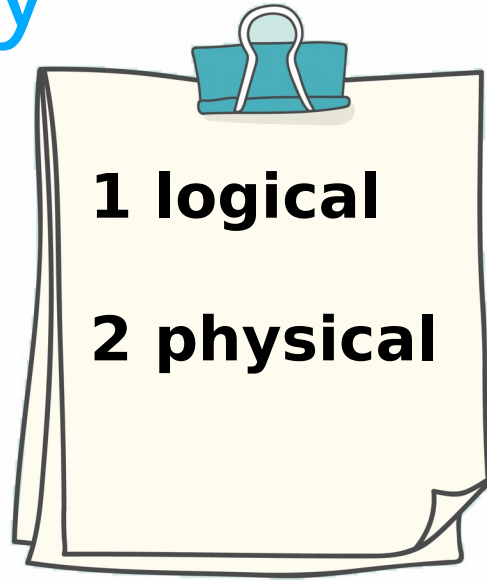
Standby



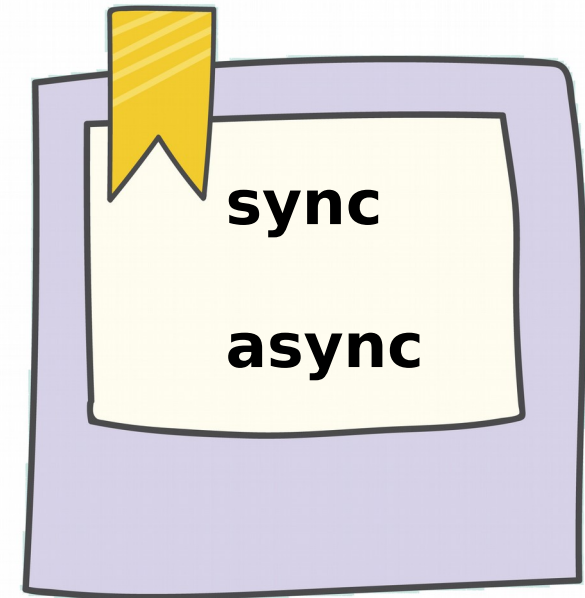
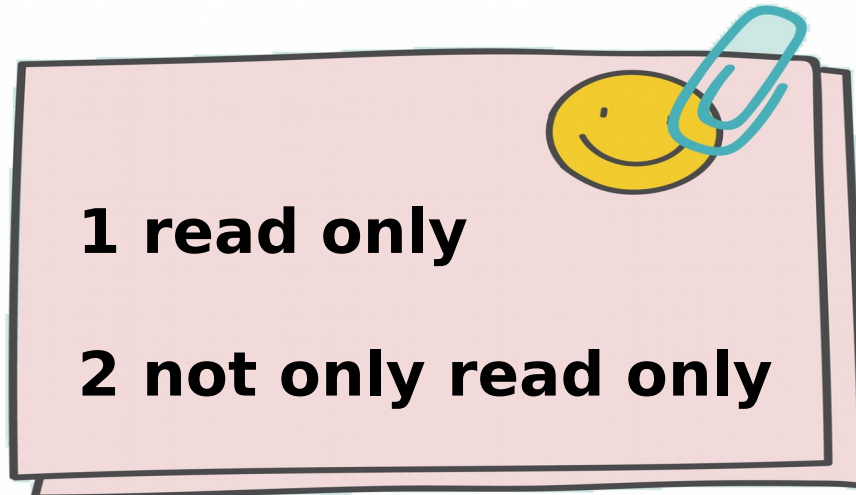
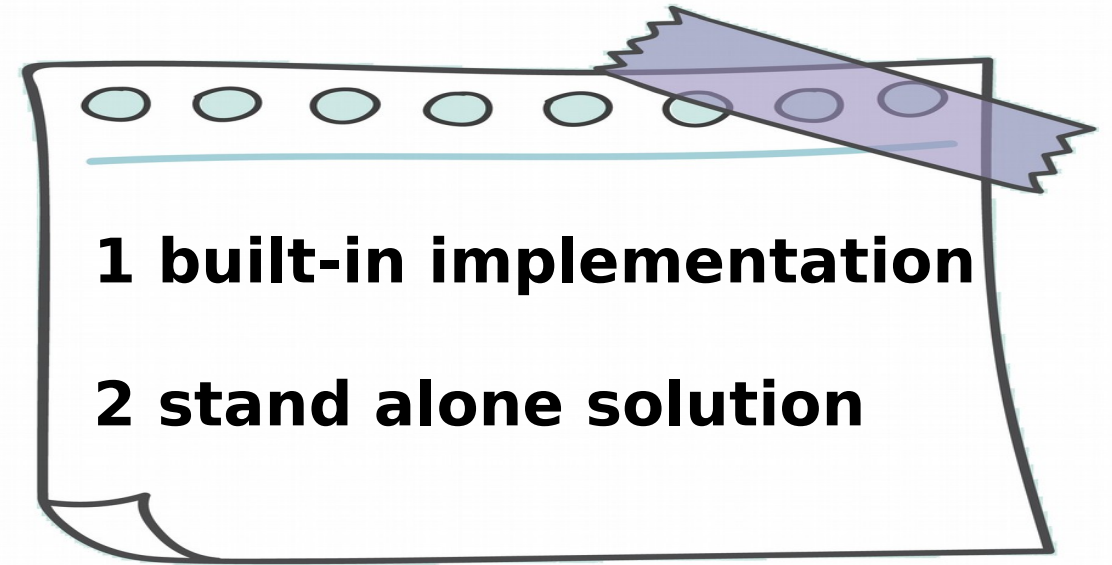
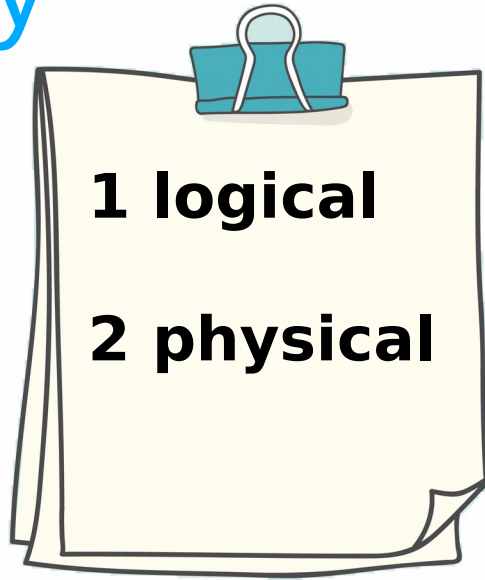
Standby



Standby

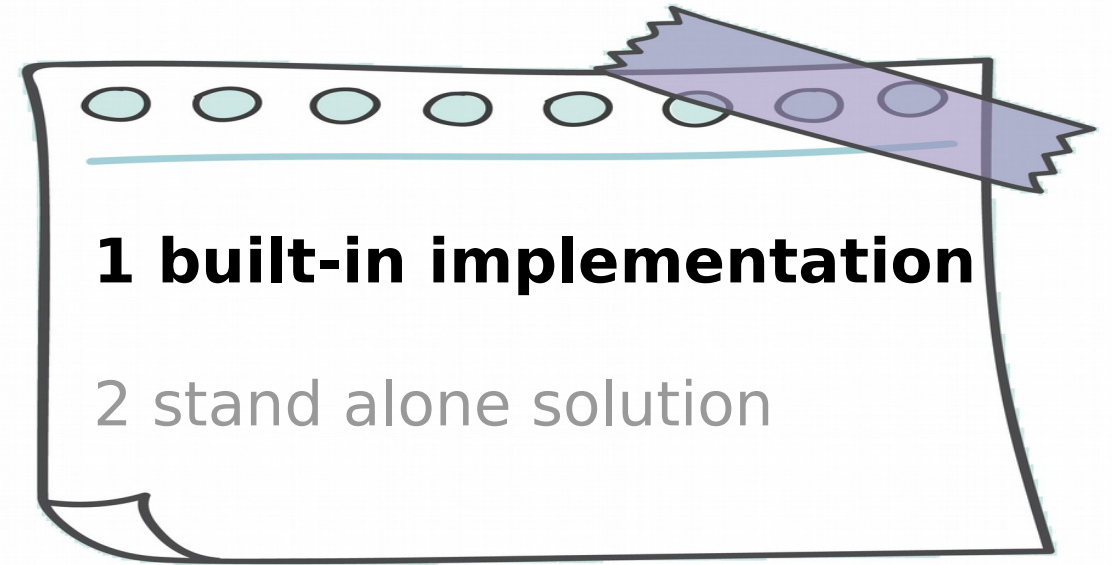
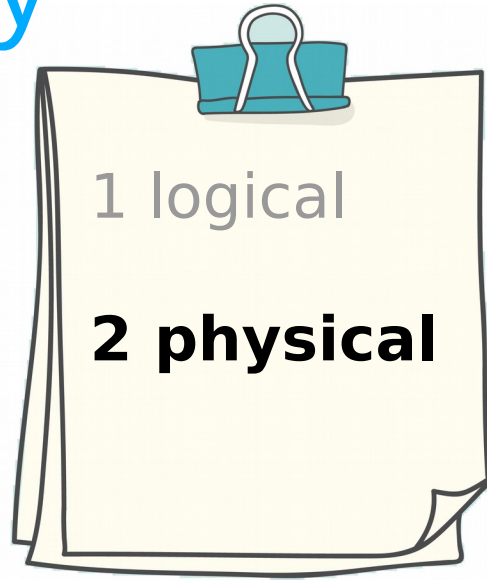


Standby



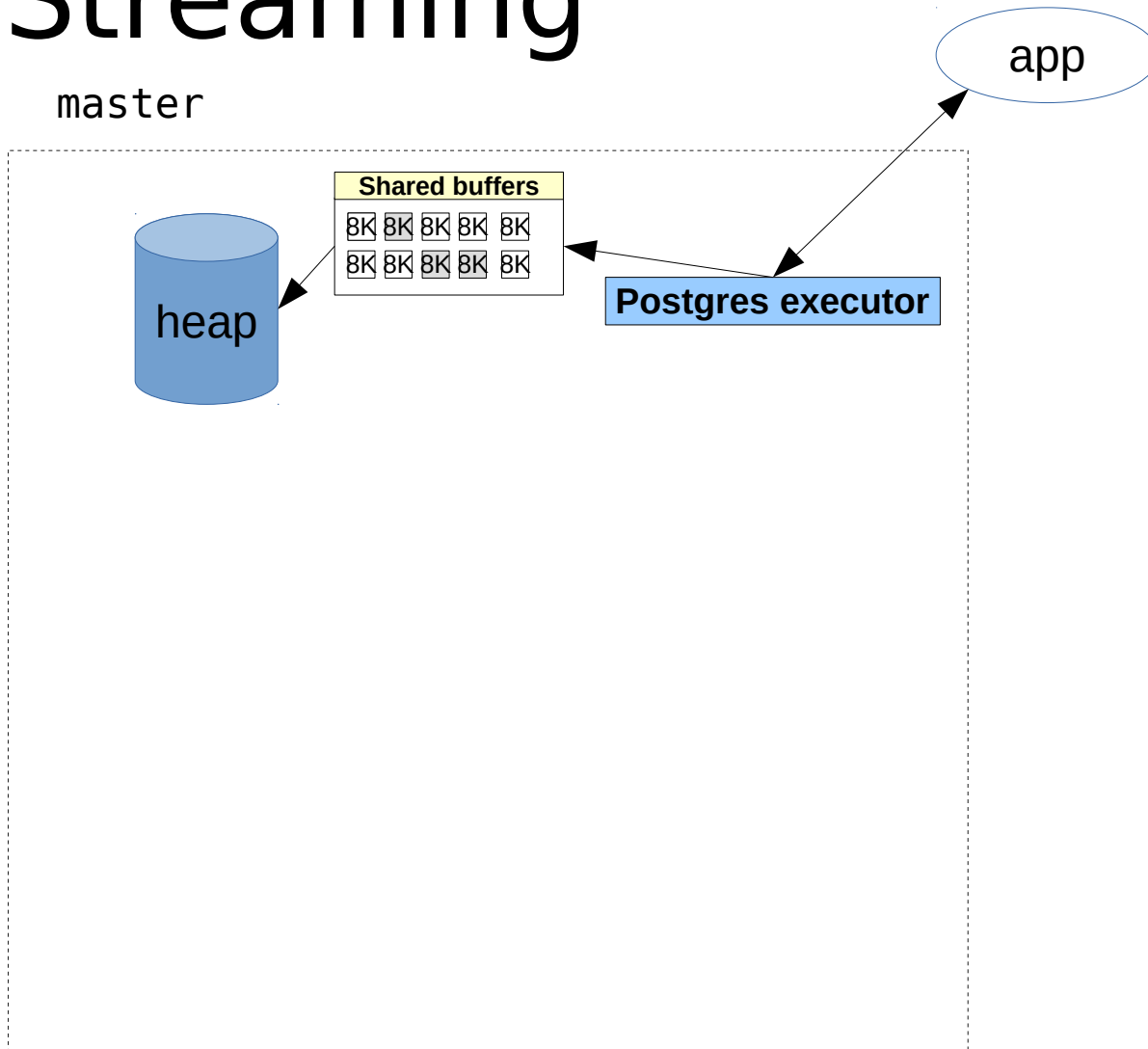
...

Standby



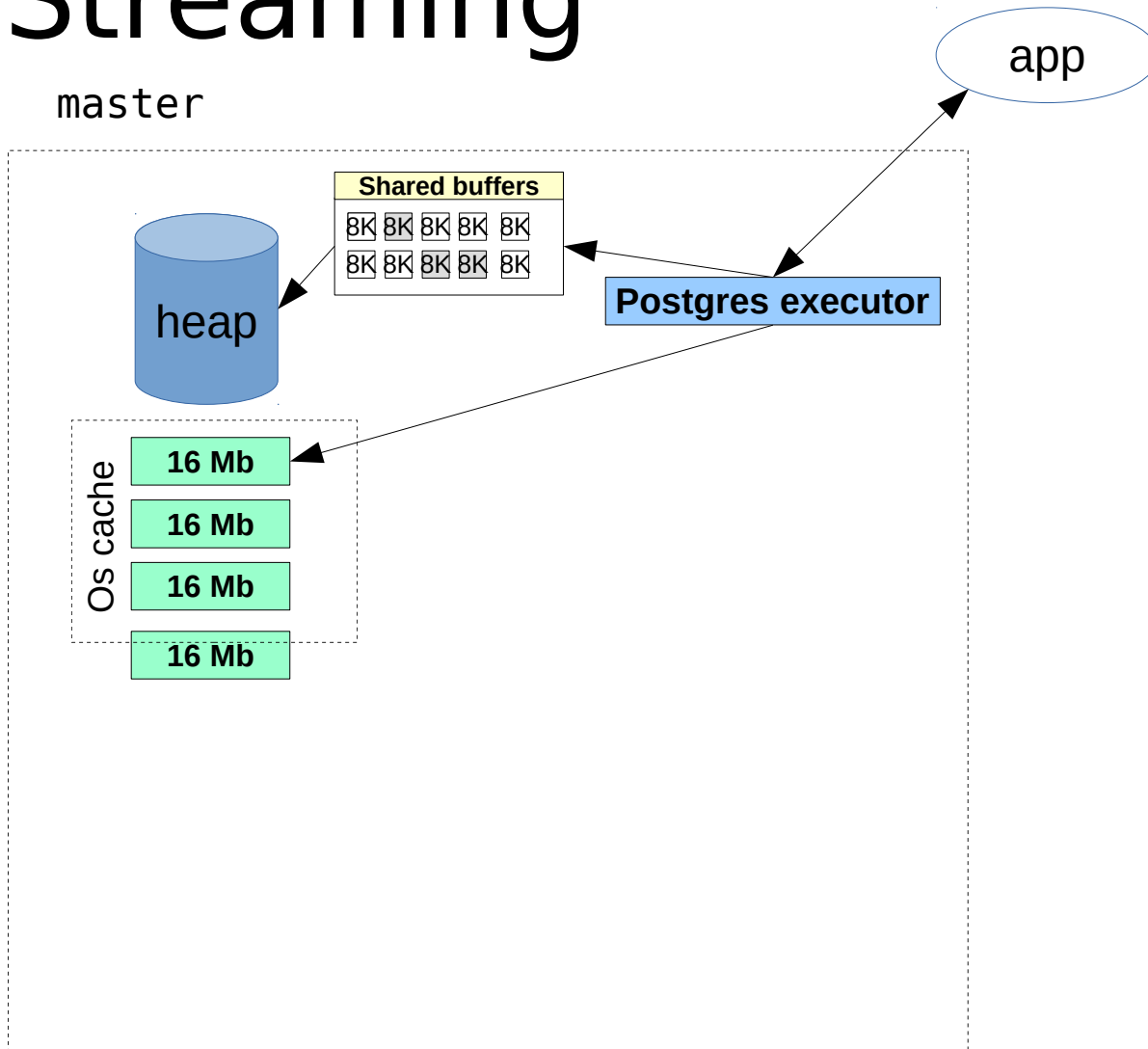
Streaming

master



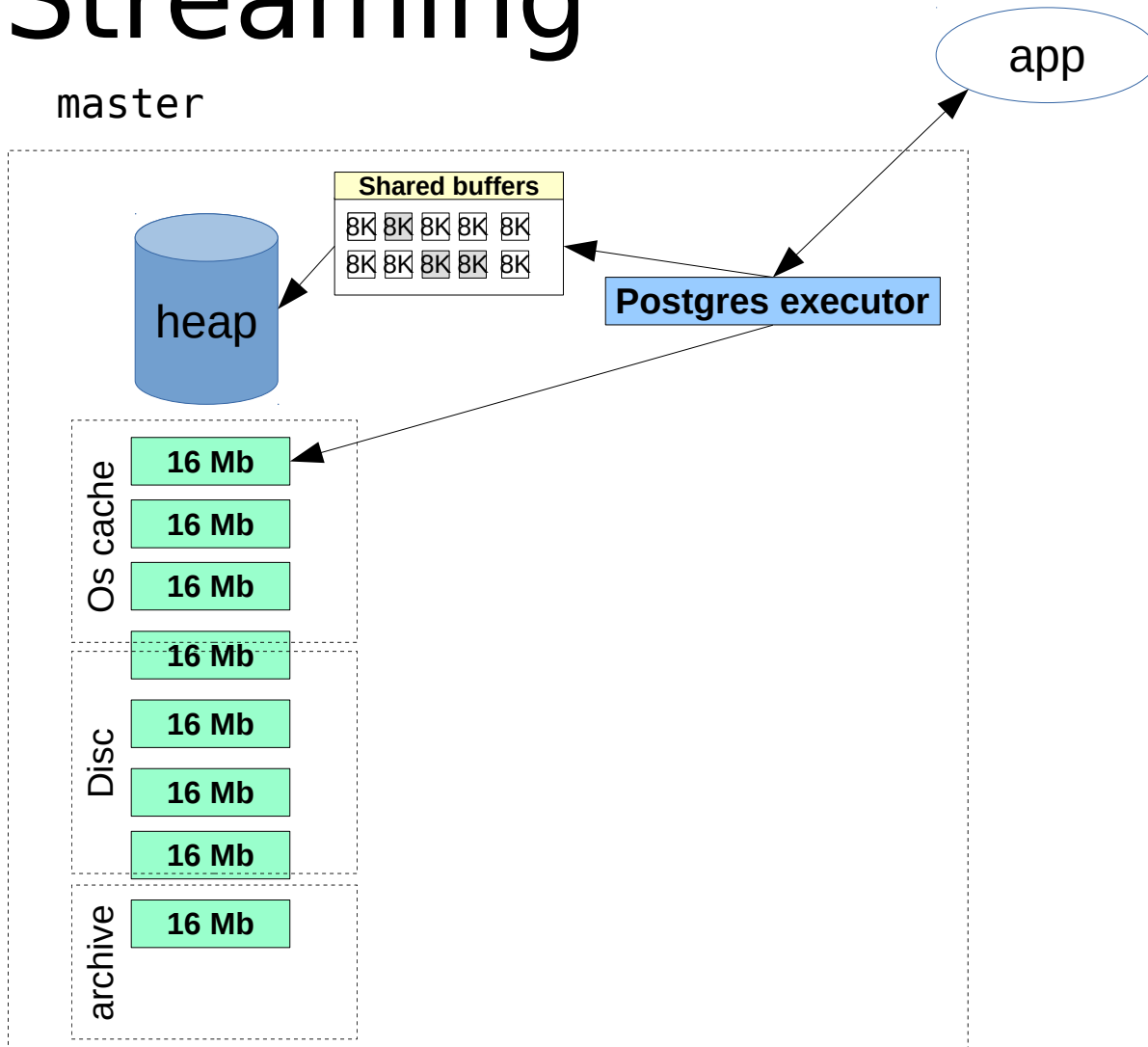
Streaming

master

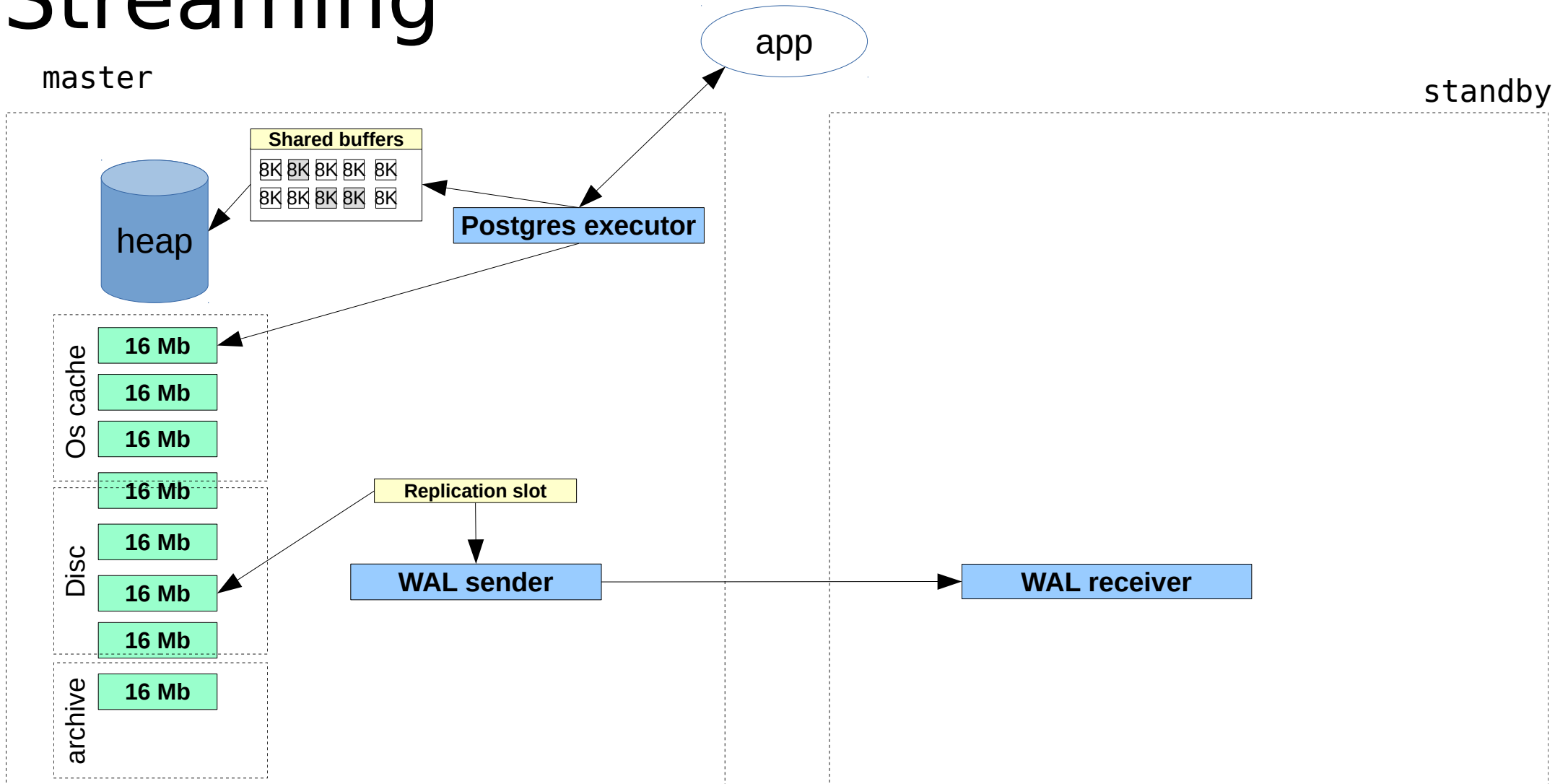


Streaming

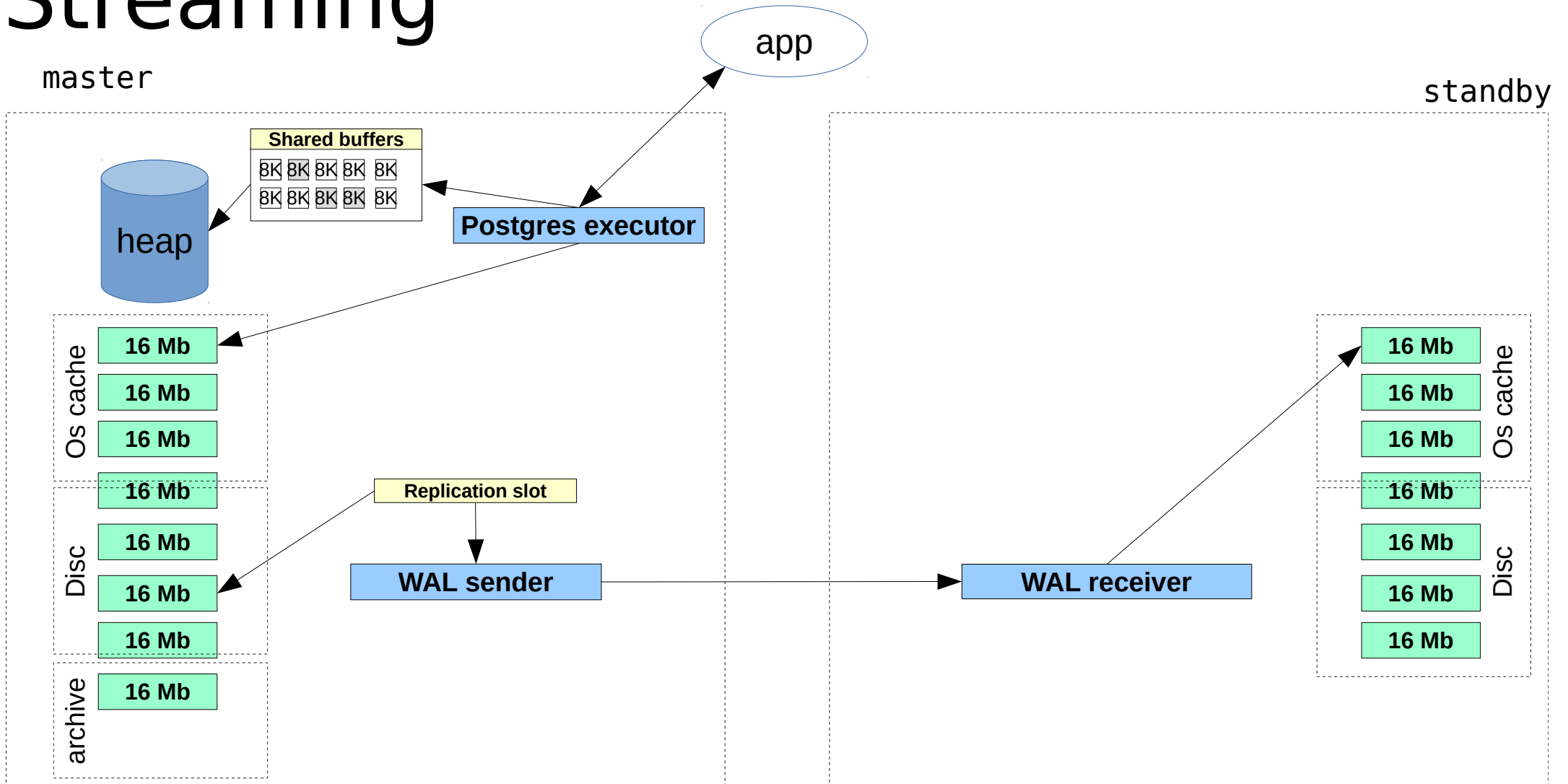
master



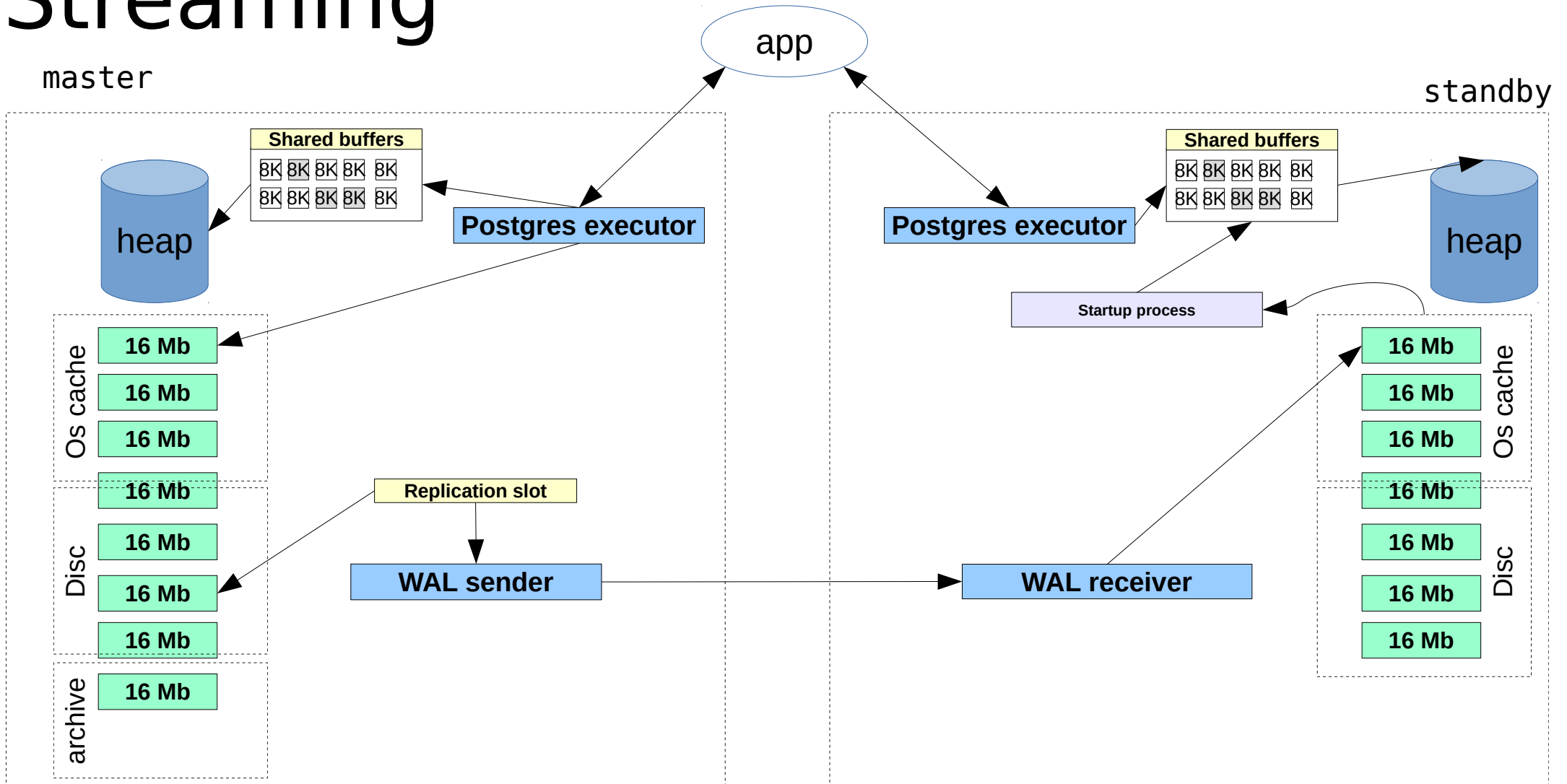
Streaming



Streaming

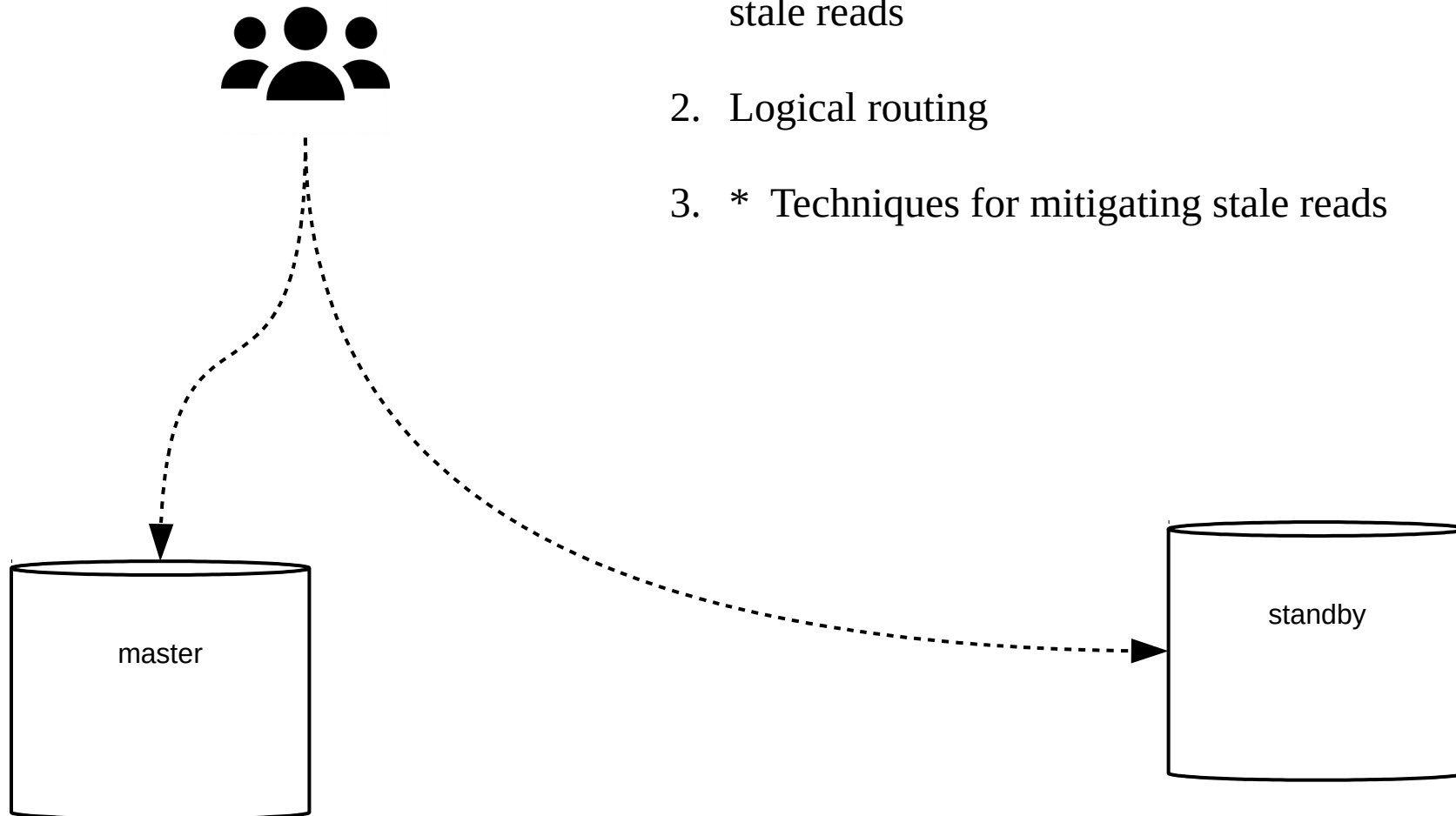


Streaming



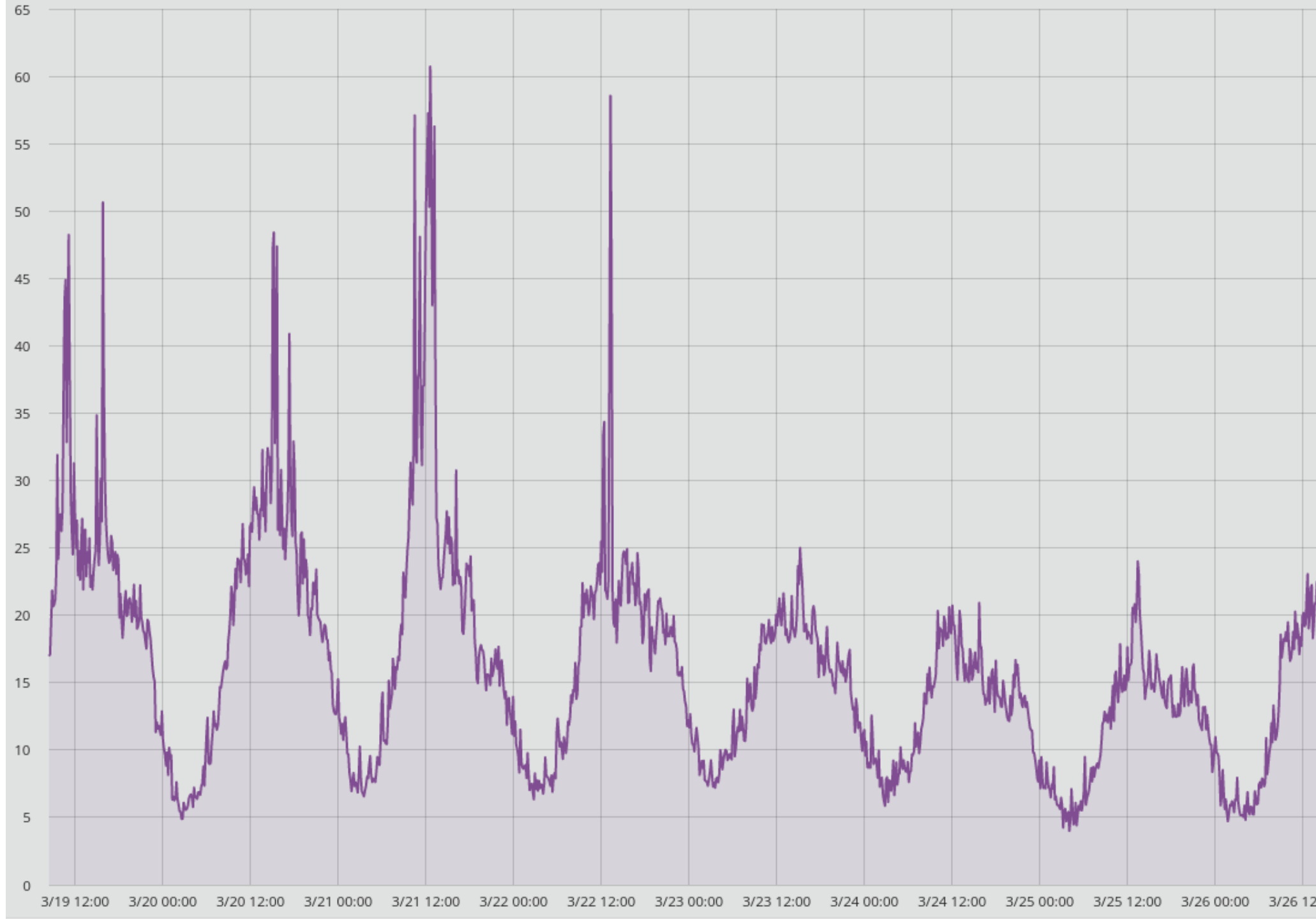
Scaling reads

1. Using read replicas without any techniques for mitigating stale reads
2. Logical routing
3. * Techniques for mitigating stale reads



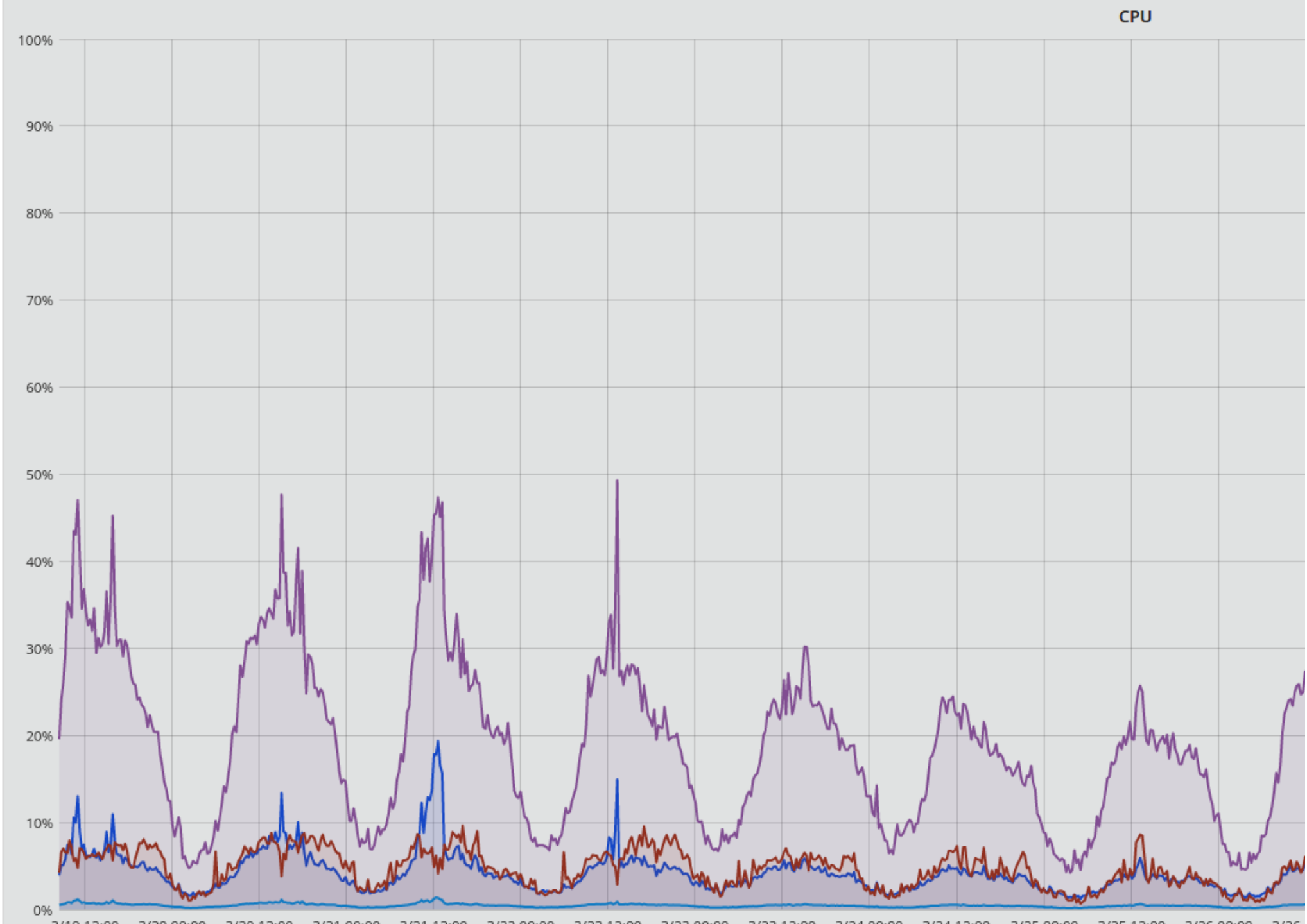
Load Average

* 28 physical cores

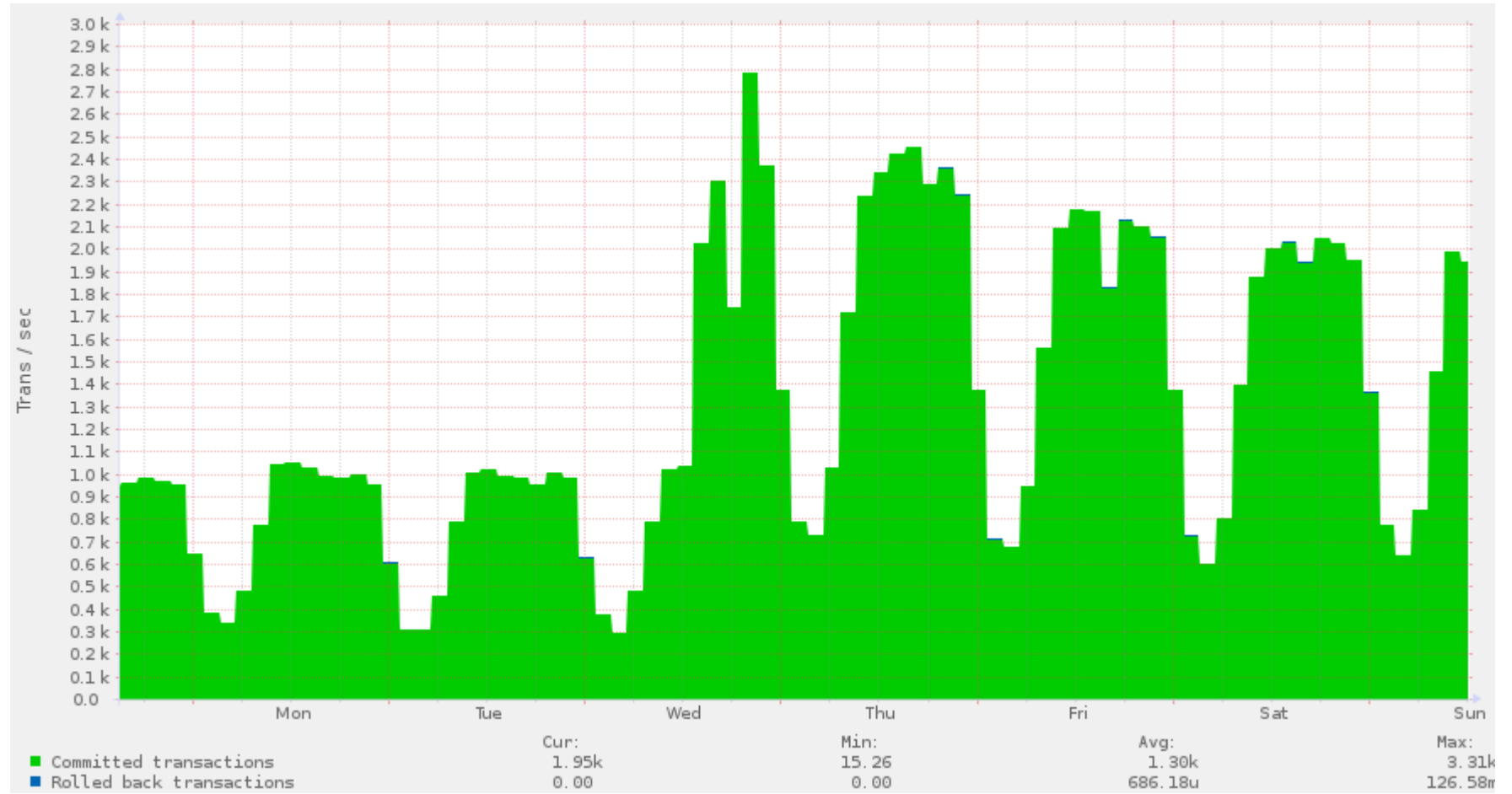


CPU

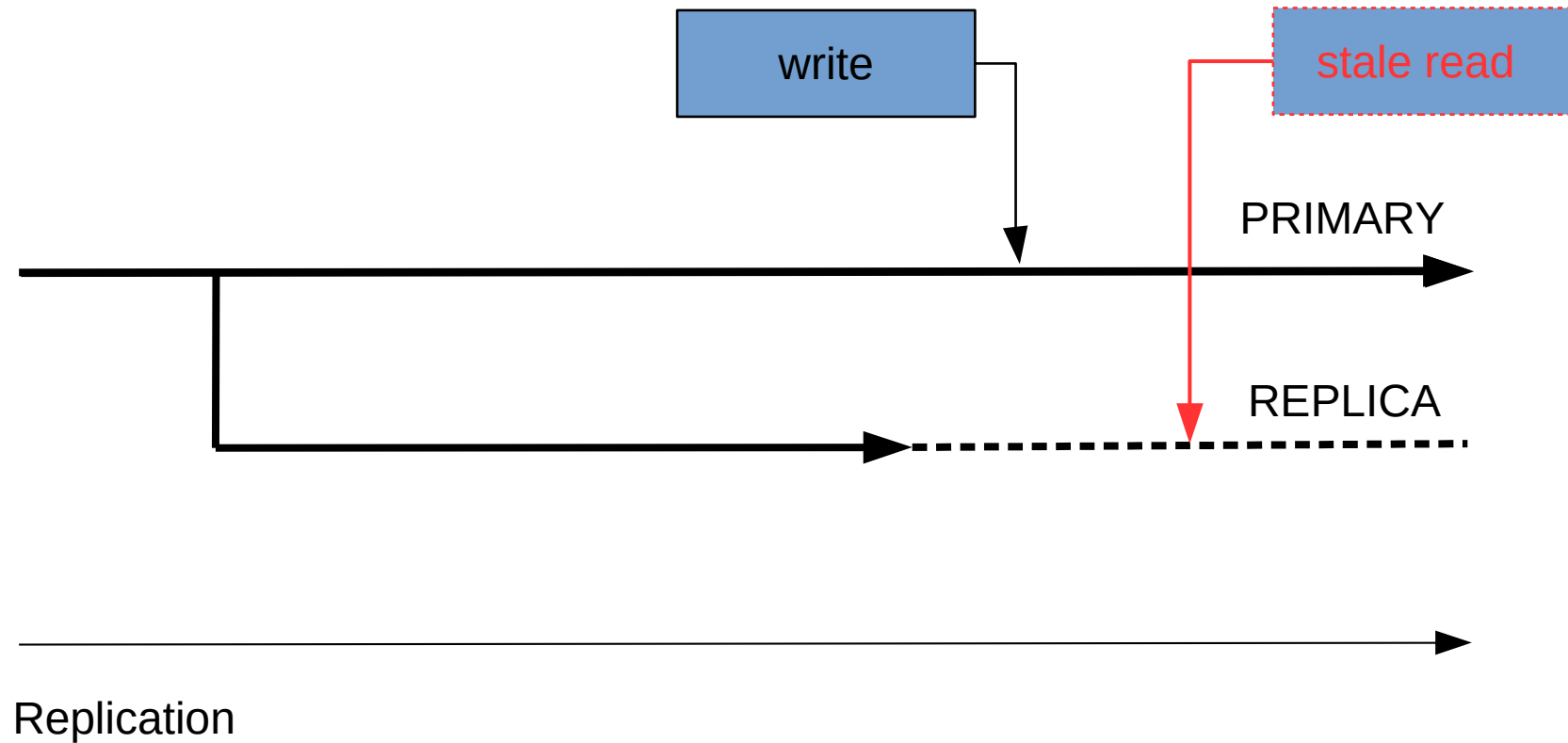
* 28 physical cores



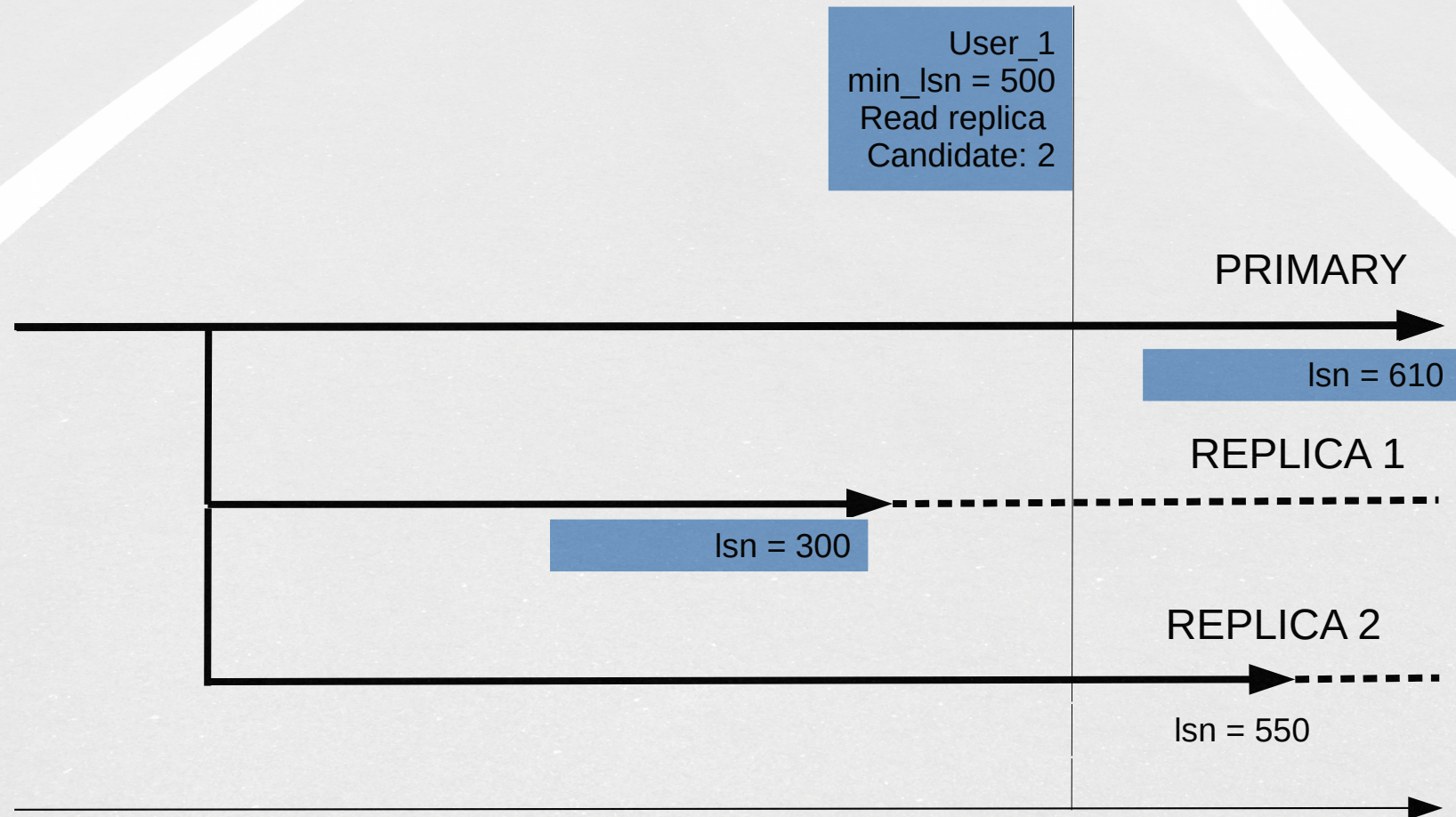
TPS



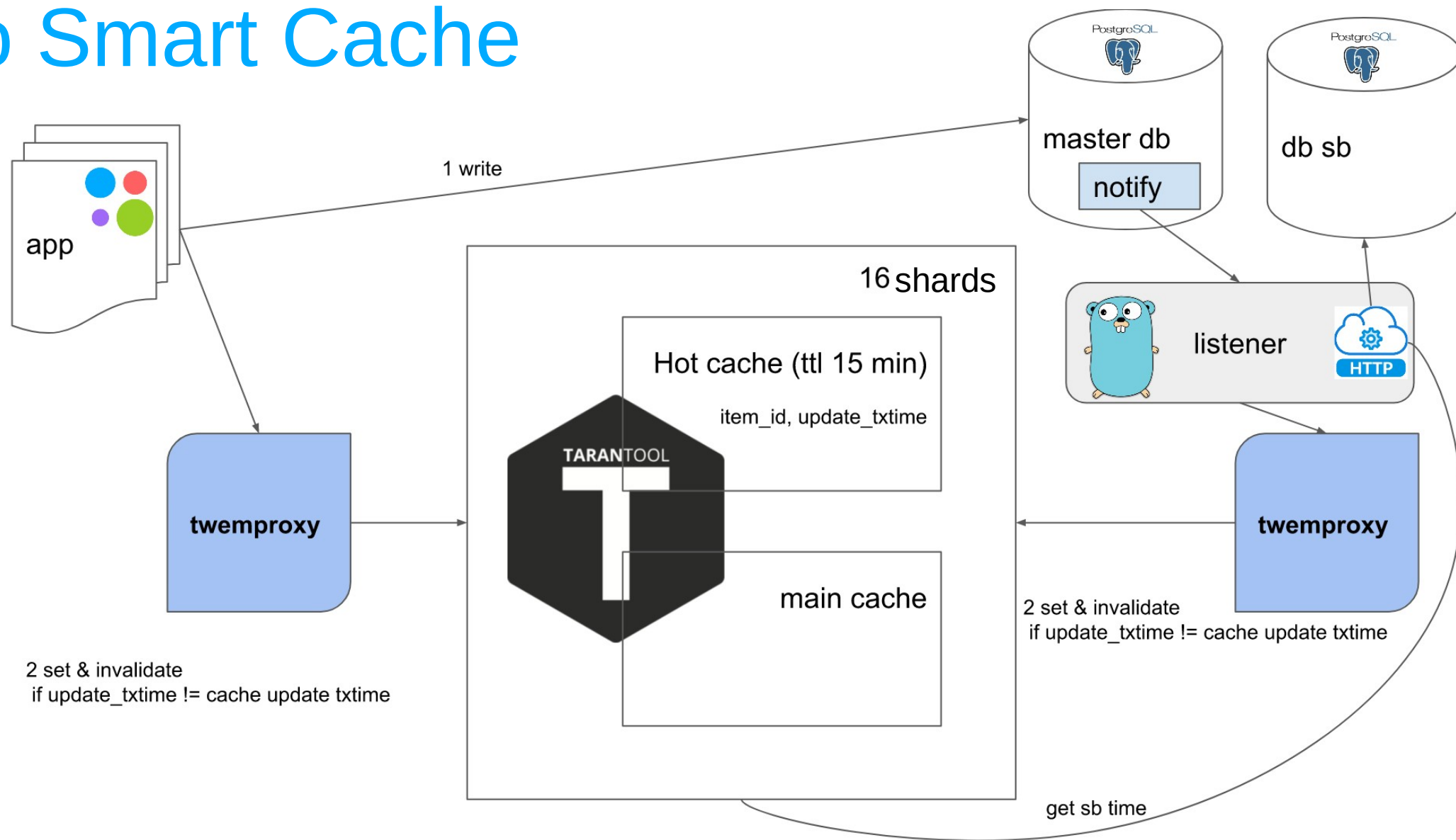
Stale Reads



Routing



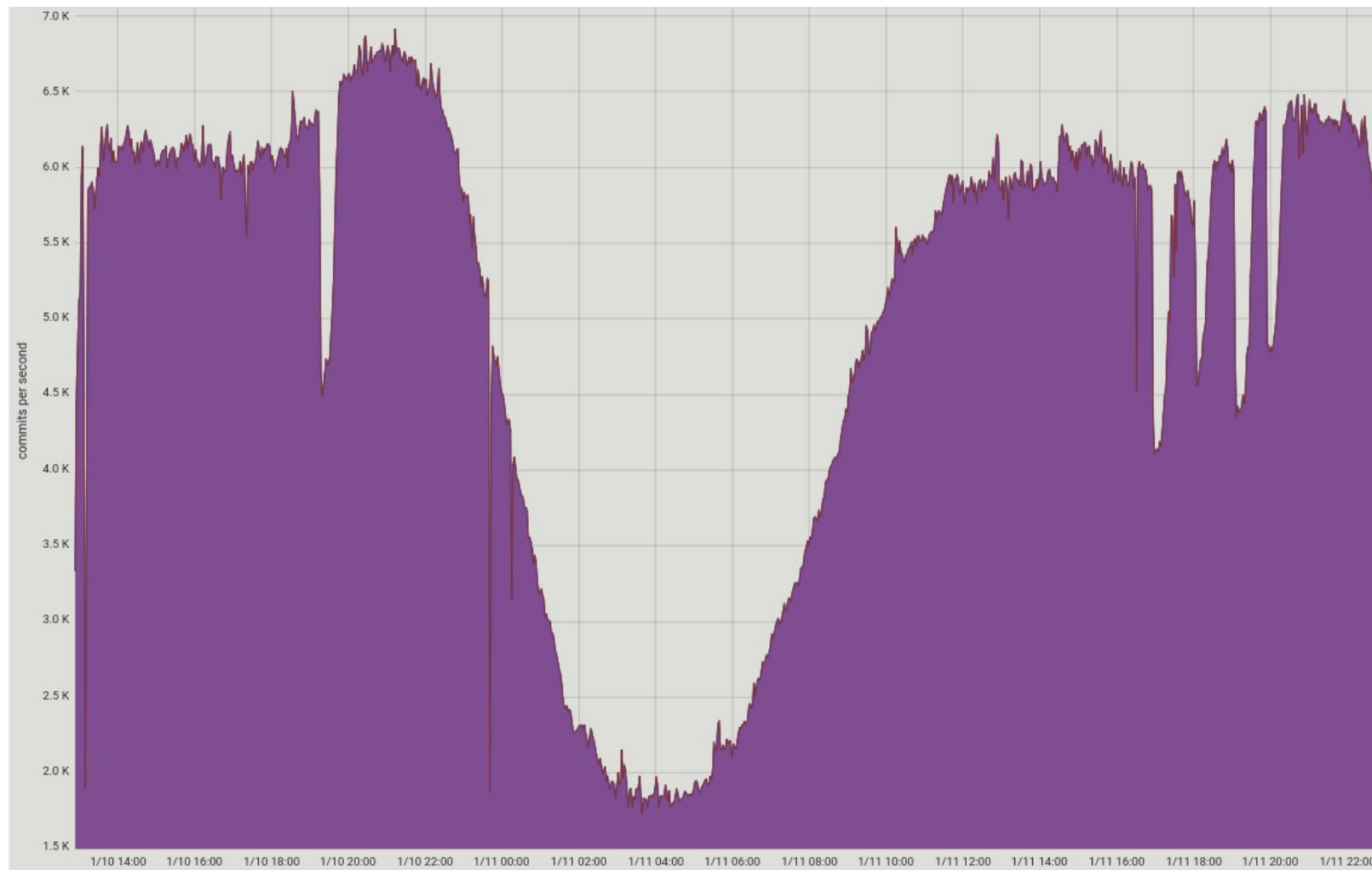
Avito Smart Cache



Two levels of cache

```
app.get_item(key):  
    data = main_cache.get(key) // try to get data from cache  
  
    if found then return  
  
    hot_cache.get(key) // get data from hot level of cache  
  
    if found then // if data was recently changed then route to master  
        data =db_master.get_item(key)  
        main_cache.set(key, data) // ttl 1 hour  
    else if sb too old then // if standby is falling behind – route to master  
        data =db_master.get_item(key)  
        main_cache.set(key, data) // ttl 24 hours  
    else // in other cases we can route to standby  
        data =db_slave.get_item(key)  
        main_cache.set(key, data) // ttl 24 hours  
    end if
```


Everything seems fine but ...



Pitfalls

(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



master

standby

```
create table items(item_id int);  
create table options(item_id int, v1 text);
```



master

standby

```
create table items(item_id int);  
create table options(item_id int, v1 text);
```

```
Begin;
```

```
alter table options add v2 int;
```



master

```
create table items(item_id int);  
create table options(item_id int, v1 text);
```

```
Begin;
```

```
alter table options add v2 int;
```



standby

```
Begin;
```

```
select * from items;
```



master



standby

```
create table items(item_id int);  
create table options(item_id int, v1 text);
```

```
Begin;
```

```
alter table options add v2 int;
```

```
alter table items add a text;
```

```
Begin;
```

```
select * from items;
```




master

```
create table items(item_id int);  
create table options(item_id int, v1 text);
```

```
Begin;
```

```
alter table options add v2 int;
```

```
alter table items add a text;
```



standby

```
Begin;
```

```
select * from items;
```

```
select * from options;  
--DEADLOCK is not detected
```

PostgreSQL 10

```
ERROR:  deadlock detected  
LINE 1: select * from options;  
                ^
```

```
DETAIL:  Process 25364 waits for  
AccessShareLock on relation 10000  
of database 9000; blocked by  
process 25322.
```

```
Process 25322 waits for  
AccessExclusiveLock on relation  
10000 of database 9000; blocked  
by process 25364.
```

```
HINT:  See server log for query  
details.
```

25322 is the PID of the apply process

Pitfalls

(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



master

standby

```
alter table options add v2 int;  
statement_timeout + deadlock_timeout
```




master



standby

```
alter table options add v2 int;  
statement_timeout + deadlock_timeout
```

```
select * from options;  
...  
select * from options;
```

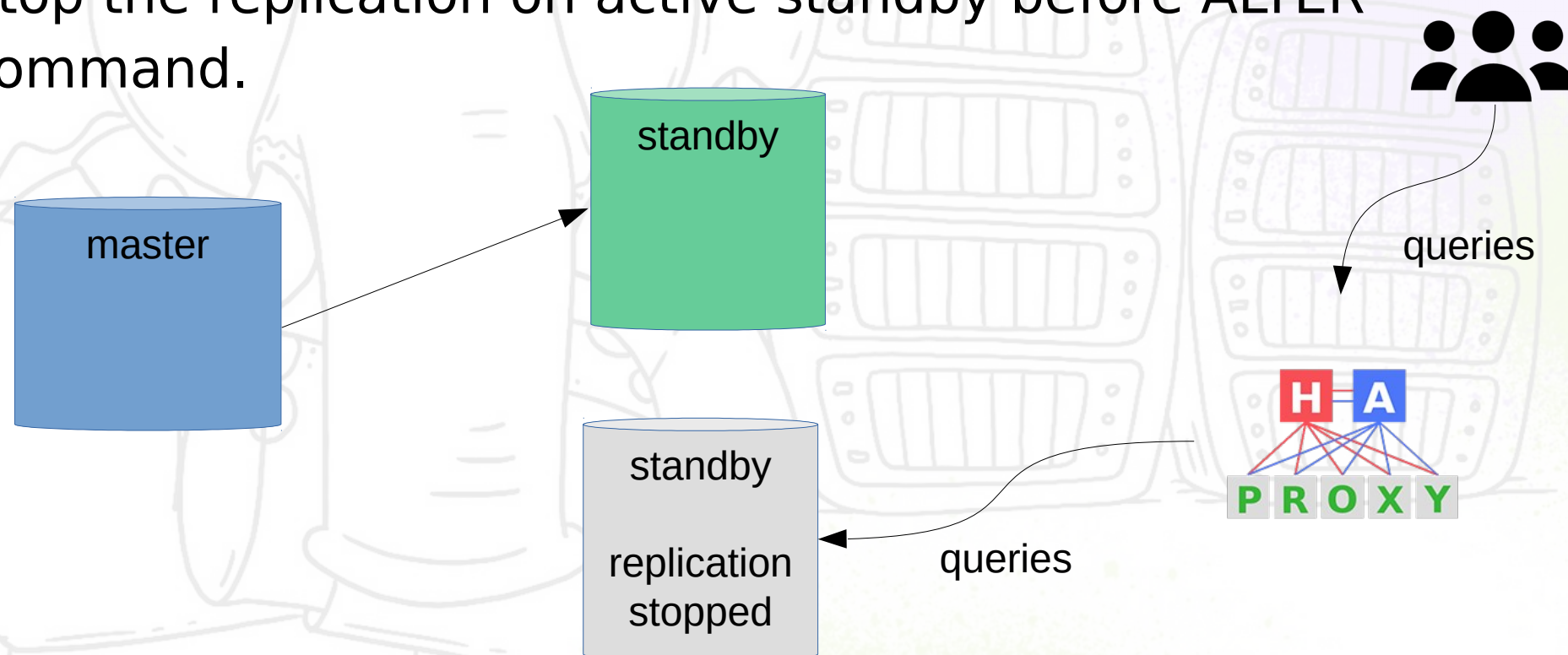
```
2018-01-12 16:54:40.208 MSK  
pid=20949,user=user_3,db=test,host=127.0.0.1:55763 LOG: process  
20949 still waiting for  
AccessShareLock on relation 10000  
of database 9000 after 5000.055  
ms
```

```
2018-01-12 16:54:40.208 MSK  
pid=20949,user=user_3,db=test,host=127.0.0.1:55763 DETAIL: Process  
holding the lock: 46091. Wait  
queue: 18639, 20949, 53445,  
20770, 10799, 47217, 37659, 6727,  
37662, 25742, 20771,
```

(2) DDL (statement_timeout and deadlock_timeout)

Script for HAProxy to implement external control
(switching your traffic from all nodes)

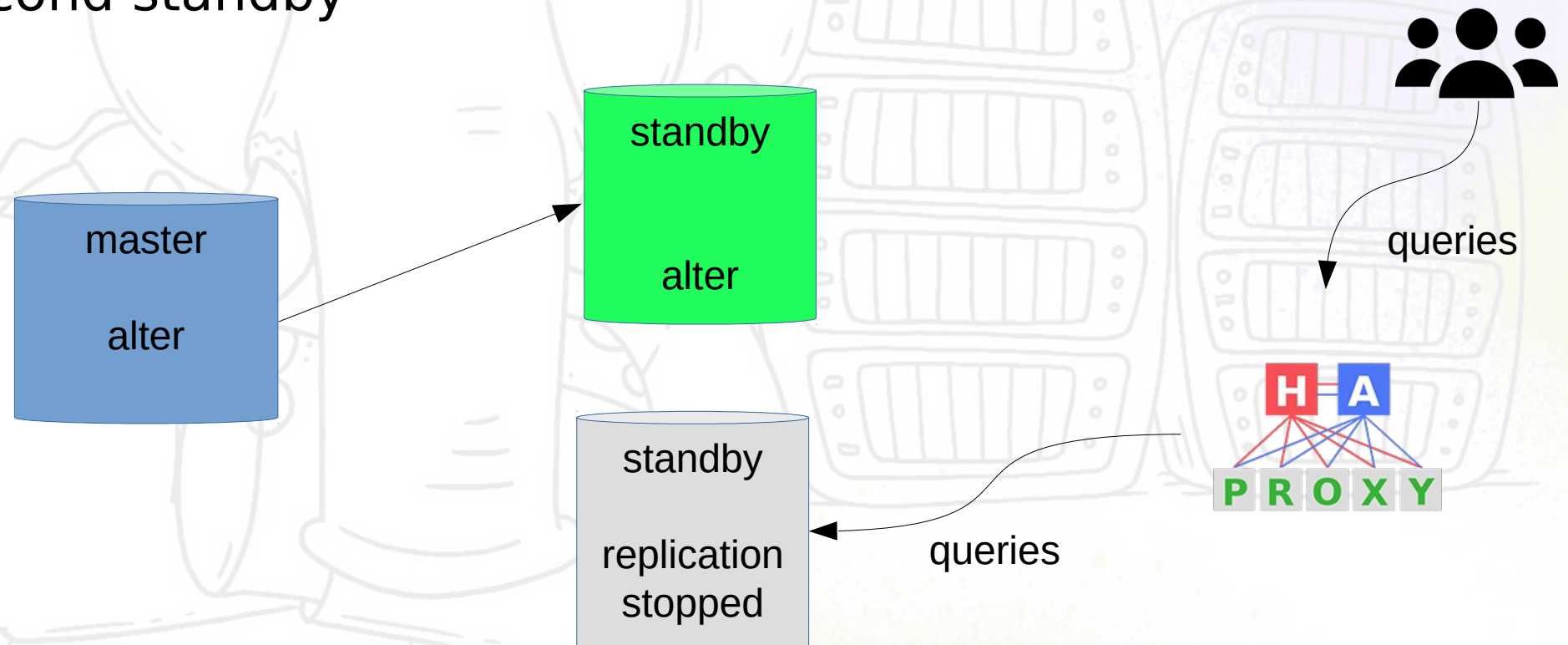
Stop the replication on active standby before ALTER
command.



(2) DDL (statement_timeout and deadlock_timeout)

Run ALTER command

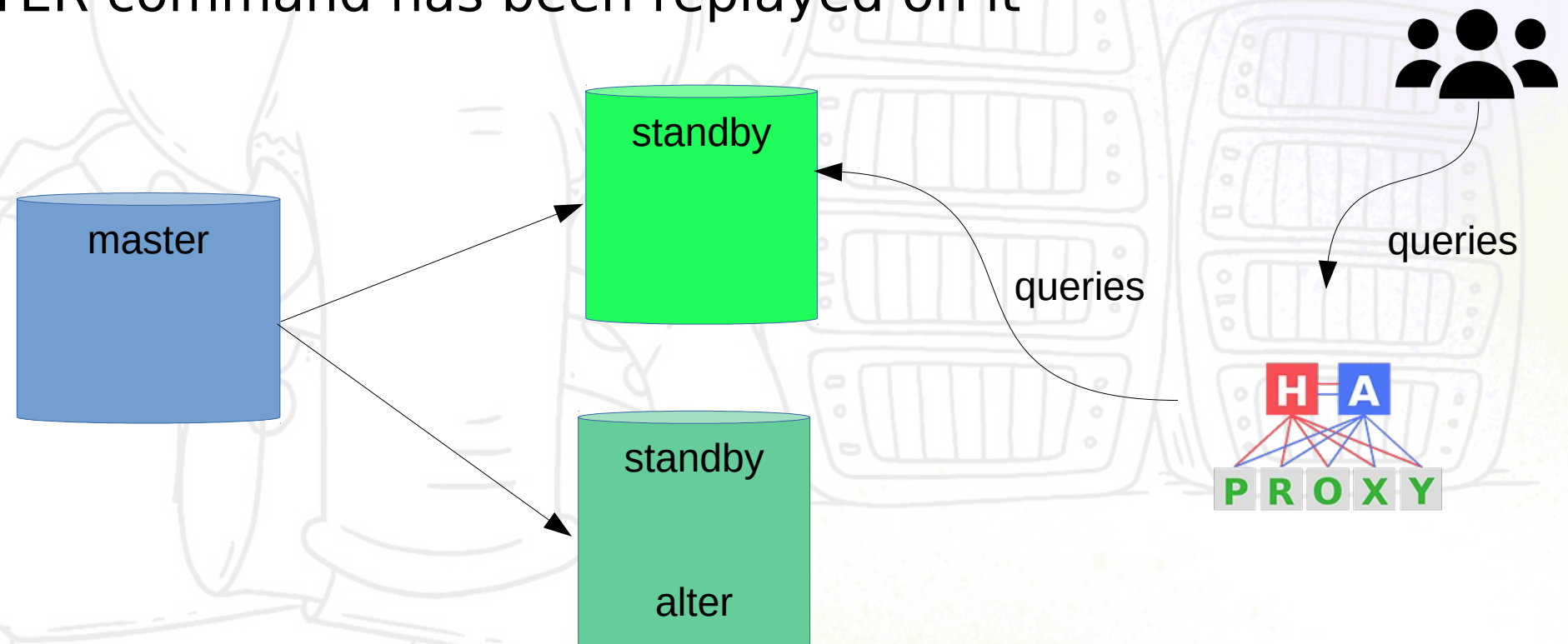
Wait till the ALTER command has been replayed on the second standby



(2) DDL (statement_timeout and deadlock_timeout)

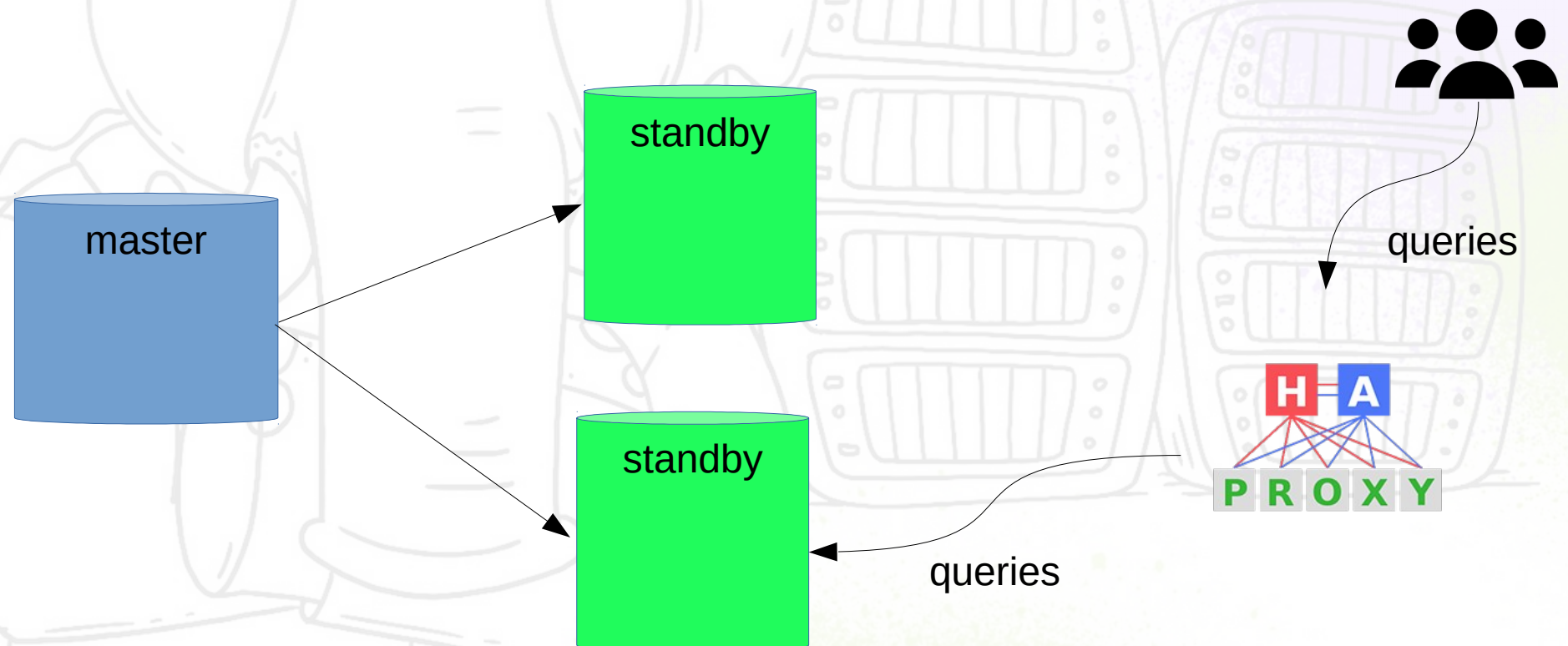
Switch traffic on the second standby

Start replication on the first standby and wait till the ALTER command has been replayed on it



(2) DDL (statement_timeout and deadlock_timeout)

Return the first standby to the pool of active standbys



Pitfalls

(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

(3) Vacuum replaying on standby and truncating data file

- Vacuum can truncate the end of data file — the exclusive lock is needed for this action. At this moment on the standby long locks between read only queries and recovery process occur
- It happens because some unlock actions are not written to WAL .
- On next slide you can see few AccessExclusive locks in one xid 920764691, and not a single unlock...
- Unlock happens much later. When standby replays commit

(3) Vacuum replaying on standby and truncating data file

```
tx: 920764691, lsn: 73CE0/10605980, desc: AccessExclusive locks: xid 920764691 db 16445 rel 3326466
tx: 920764691, lsn: 73CE0/10694568, desc: file truncate: base/16445/3326466 to 1965248 blocks
tx: 920764691, lsn: 73CE0/1105AB98, desc: AccessExclusive locks: xid 920764691 db 16445 rel 3326466
tx: 920764691, lsn: 73CE0/11116A88, desc: file truncate: base/16445/3326466 to 1965152 blocks
tx: 920764691, lsn: 73CE0/116C89C0, desc: AccessExclusive locks: xid 920764691 db 16445 rel 3326466
tx: 920764691, lsn: 73CE0/117211E0, desc: file truncate: base/16445/3326466 to 1965088 blocks
tx: 920764691, lsn: 73CE0/128DFF00, desc: AccessExclusive locks: xid 920764691 db 16445 rel 3326466
tx: 920764691, lsn: 73CE0/129A5DD0, desc: file truncate: base/16445/3326466 to 1964960 blocks
tx: 920764691, lsn: 73CE0/1315C4E8, desc: AccessExclusive locks: xid 920764691 db 16445 rel 3326466
tx: 920764691, lsn: 73CE0/134CF9E0, desc: file truncate: base/16445/3326466 to 1964832 blocks
```

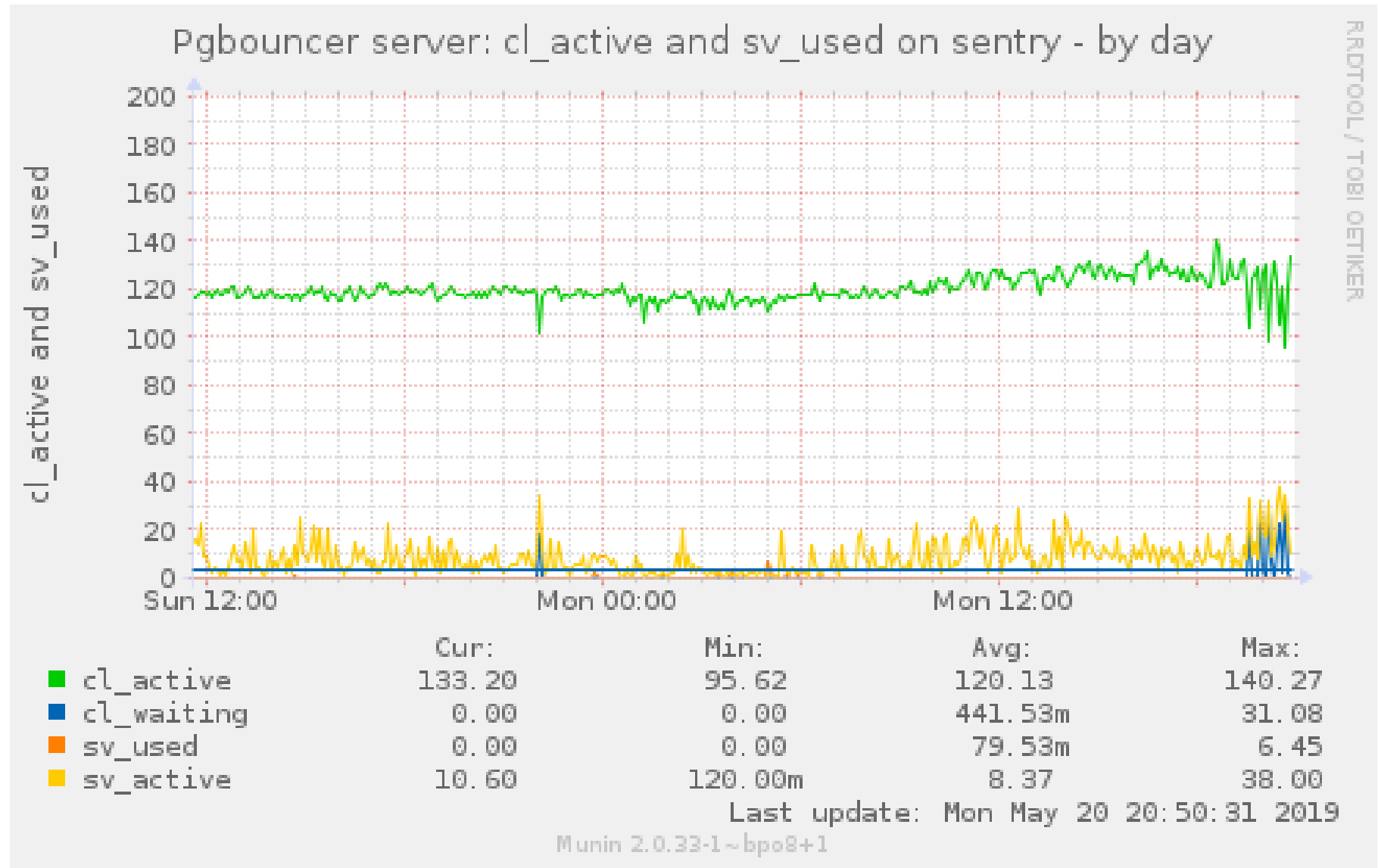
In our example there is 75-WAL-file interval between first lock and success truncate (unlock relation)

(3) Vacuum replaying on standby and truncating data file

The solution can be like:

- * <https://commitfest.postgresql.org/22/1981/>
Skip table truncation at VACUUM is coming in PostgreSQL 12
- * Decrease the number of locks on standby?
(Postgres Professional)

(3) Vacuum and truncating data file on primary

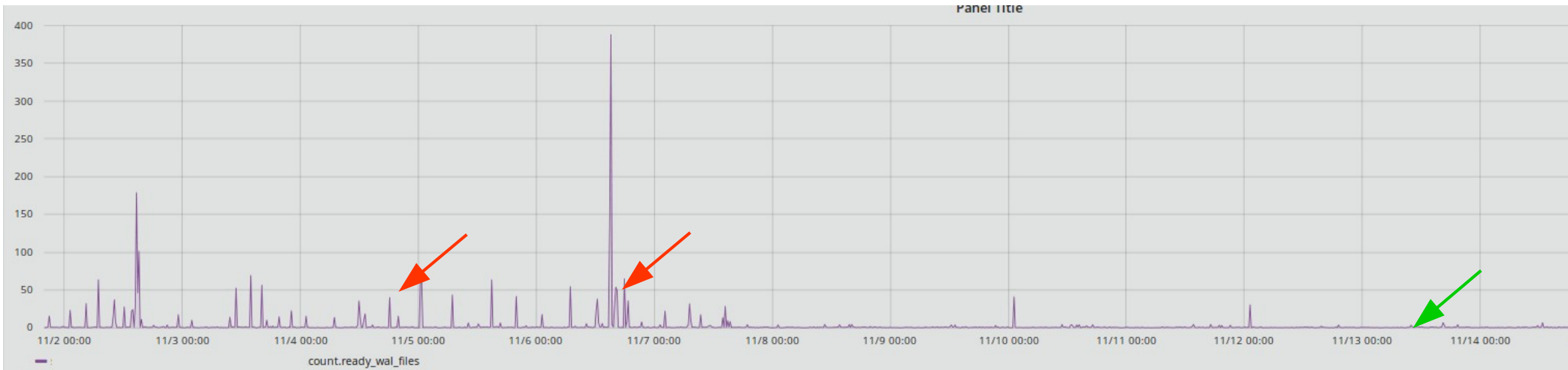
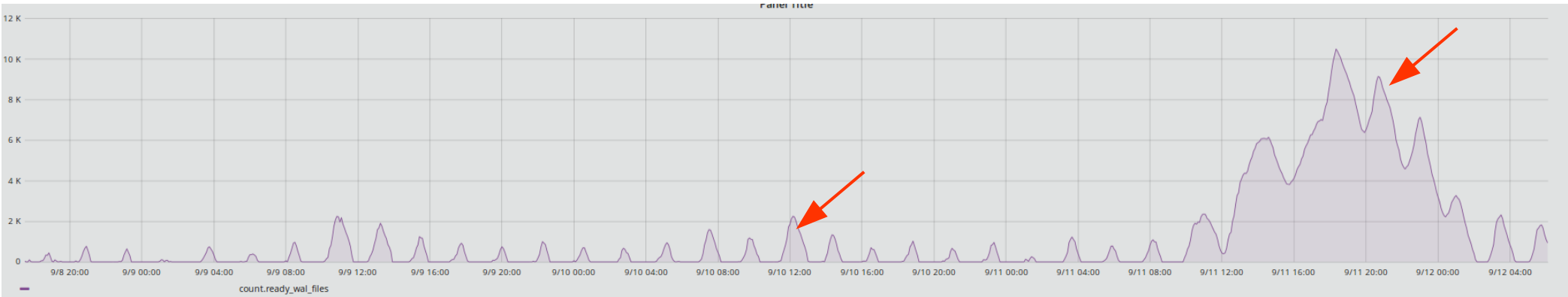


<https://sentry.io>

Pitfalls

- (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**

More and more WAL



Avito archive 2016

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

```
archive_command = '/usr/local/bin/archive_cmd HOSTNAME /postgresql/walldir/logs.complete %p %f'
```

```
usage: archive_cmd DST-HOSTNAME DST-DIR SRC-WAL-FILENAME-WITH-PATH SRC-WAL-FILENAME
```

```
DST-HOSTNAME          - for scp
```

```
DST-DIR               - archive directory for WALs
```

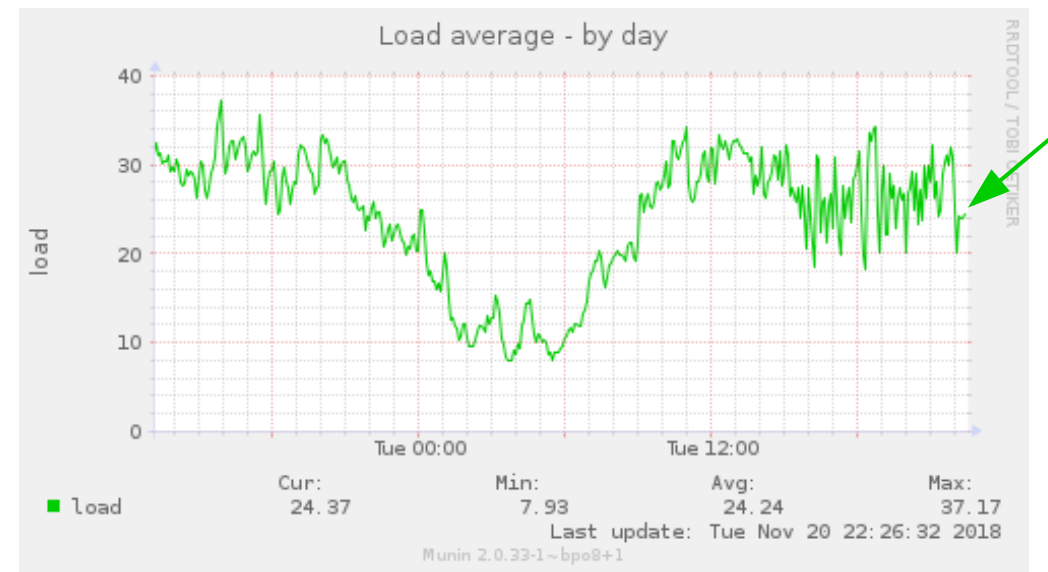
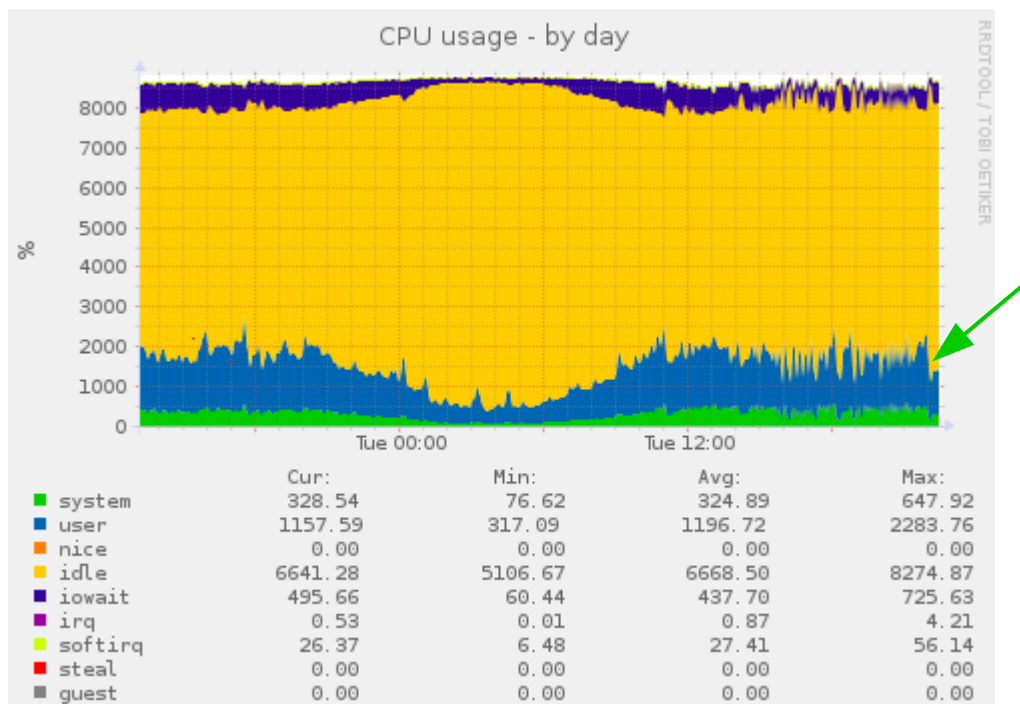
```
SRC-WAL-FILENAME-WITH-PATH - %p (file name with path)
```

```
SRC-WAL-FILENAME      - %f (file name)
```

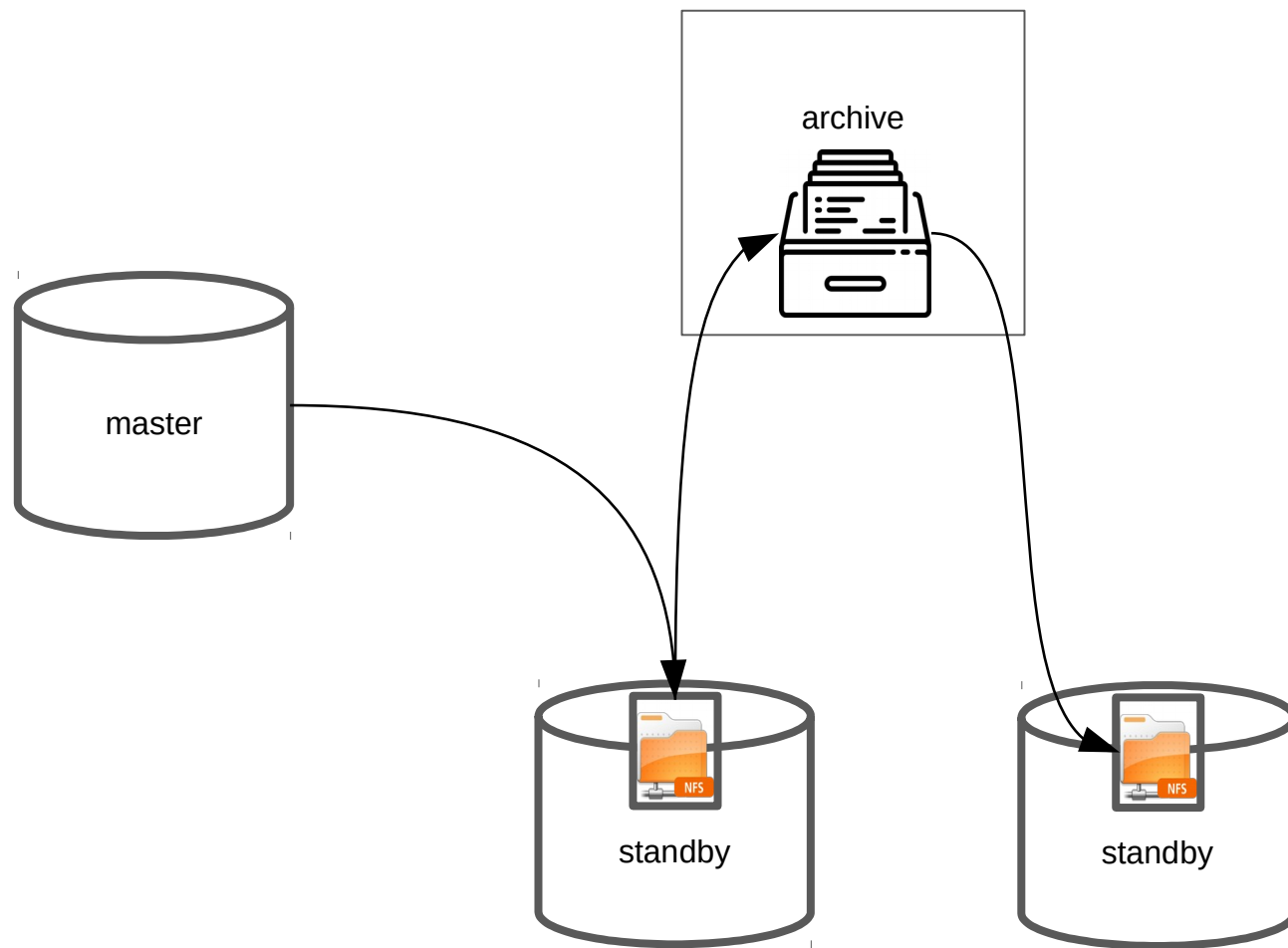
```
# archive in one thread if
```

```
# - ready WAL files lower than threshold ready_wals_for_parallel
```

Archiving 1 WAL ~ 60ms

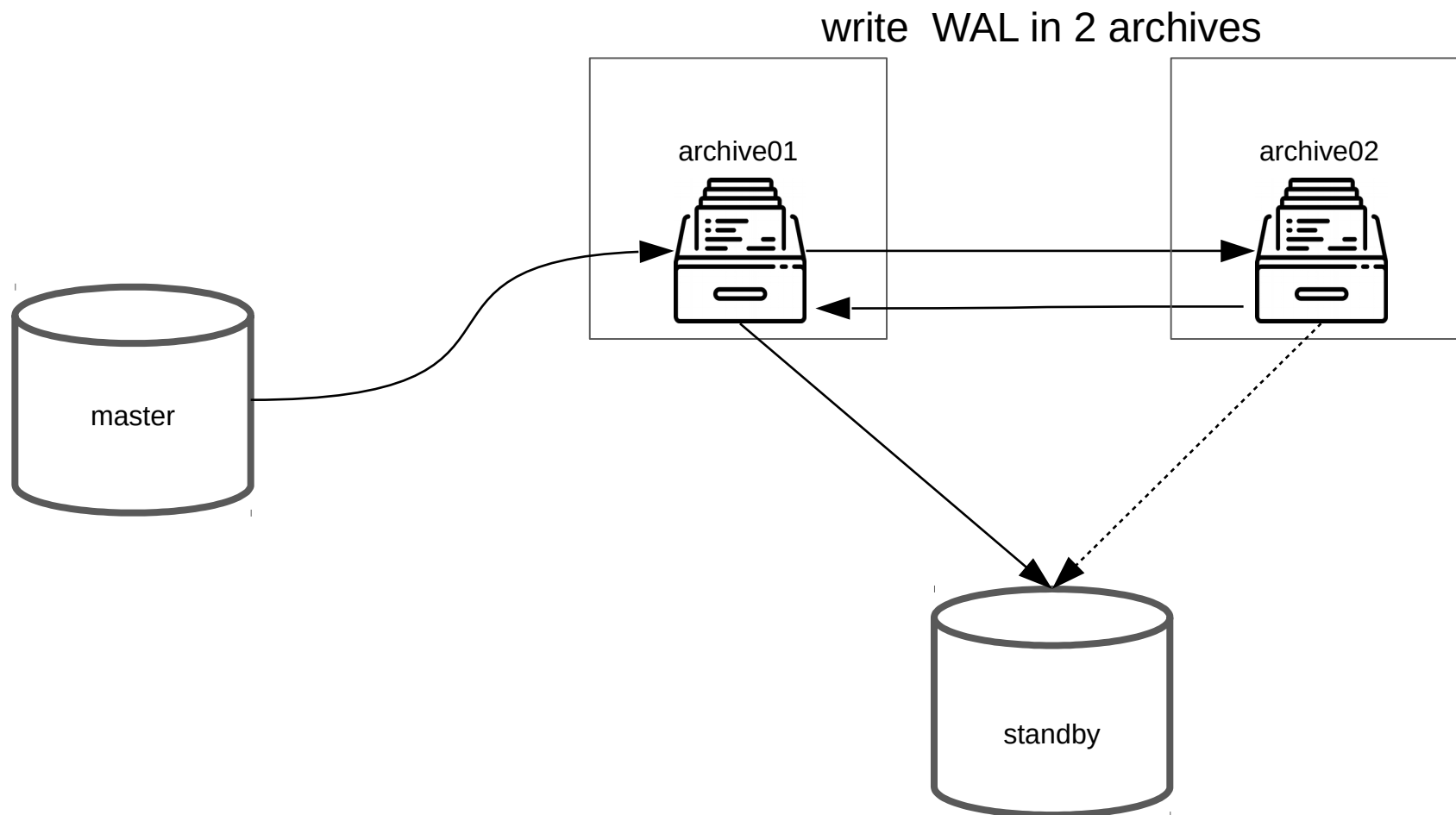


Old archive schema



New archive schema

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

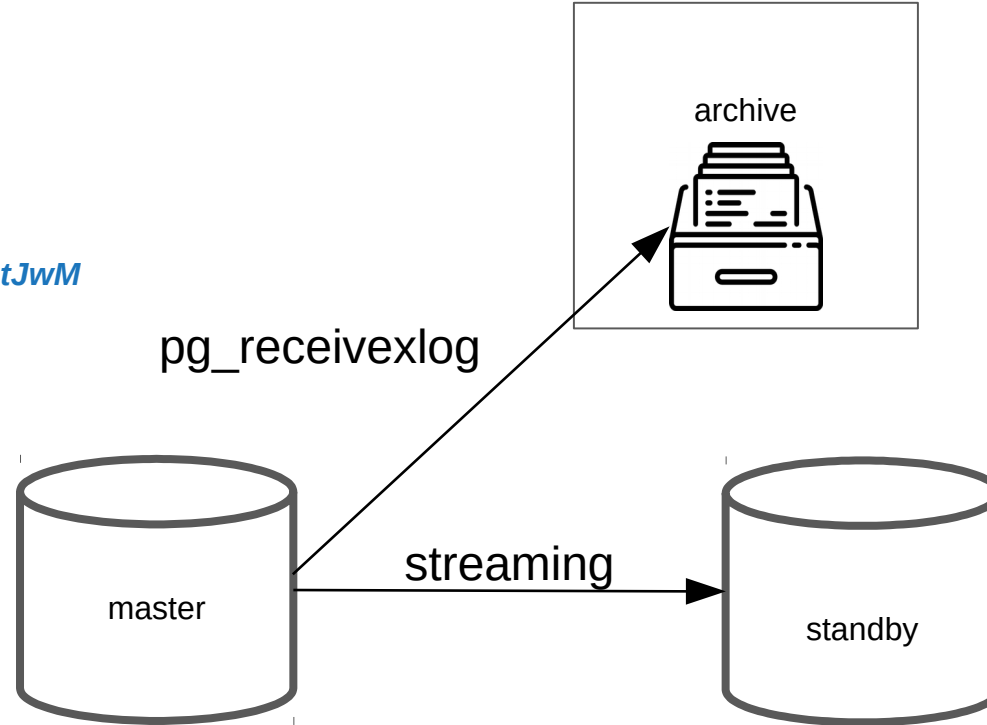


Streaming

WARM standby done right
Heikki Linnakangas

<https://pgday.ru/ru/2015/papers/8>

<https://www.youtube.com/watch?v=mIQ90MntJwM>

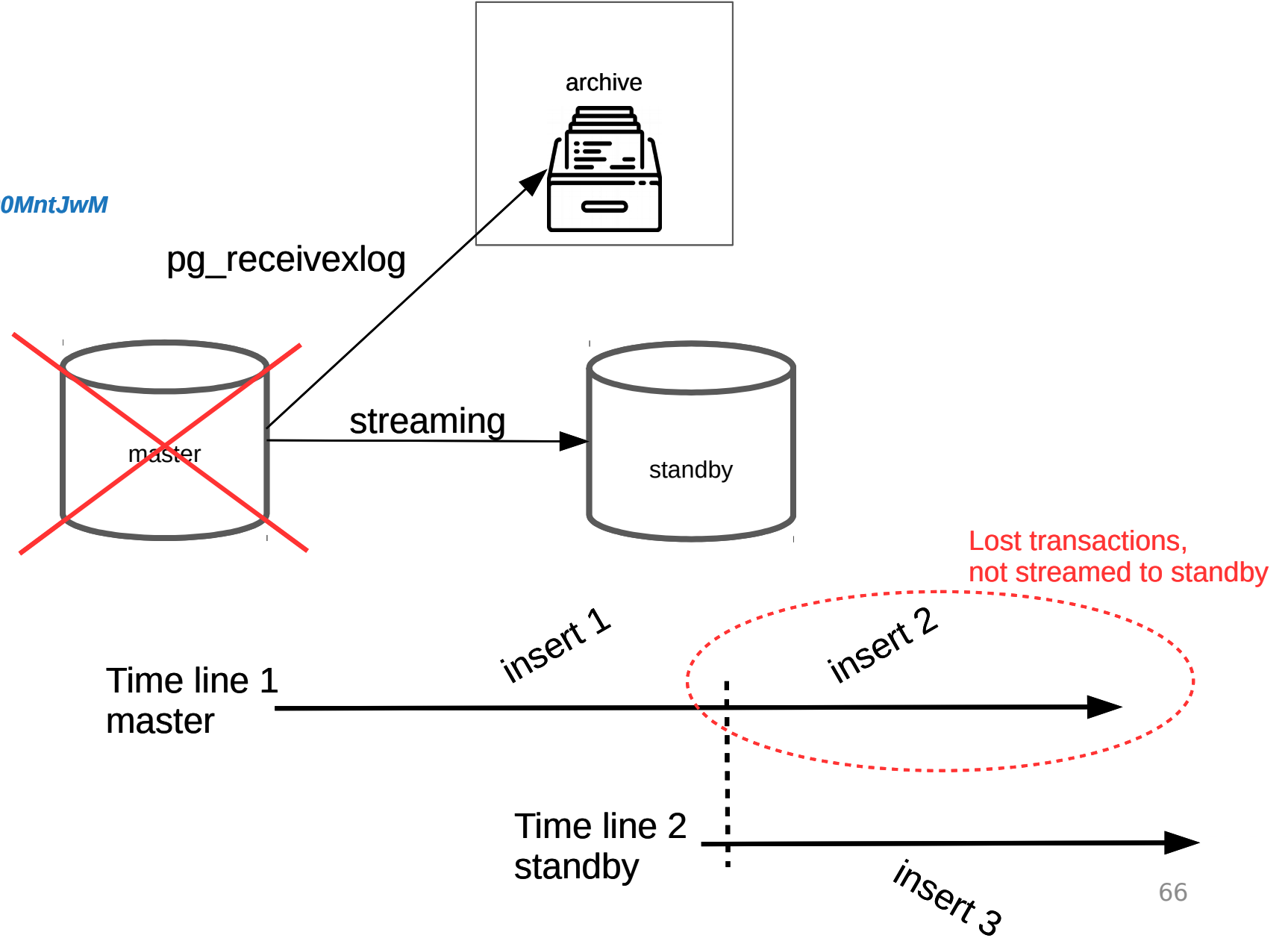


Streaming

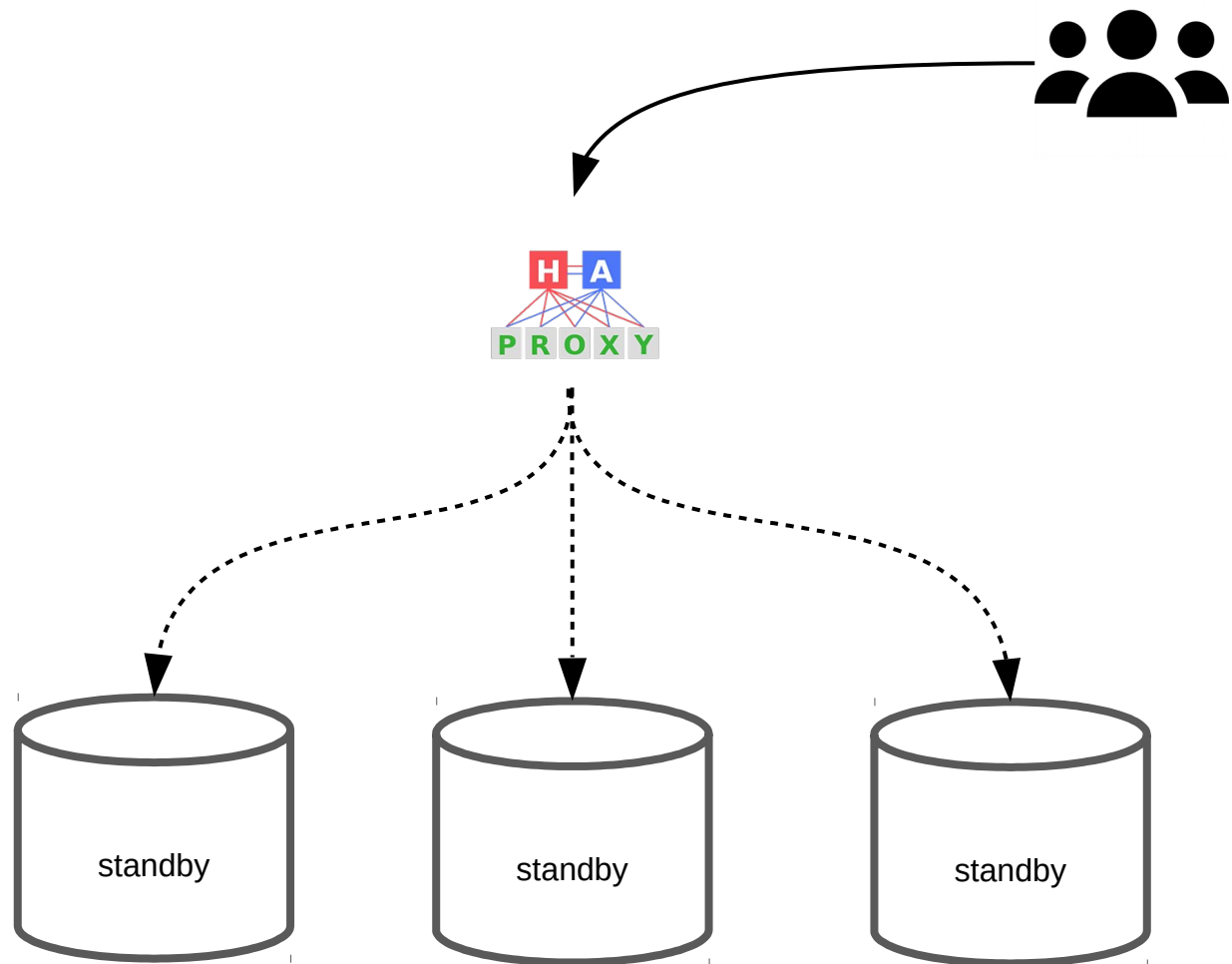
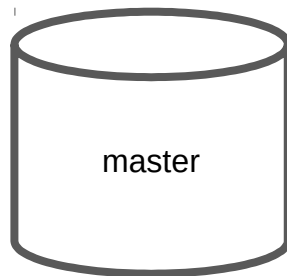
WARM standby done right
Heikki Linnakangas

<https://pgday.ru/ru/2015/papers/8>

<https://www.youtube.com/watch?v=mlQ90MntJwM>



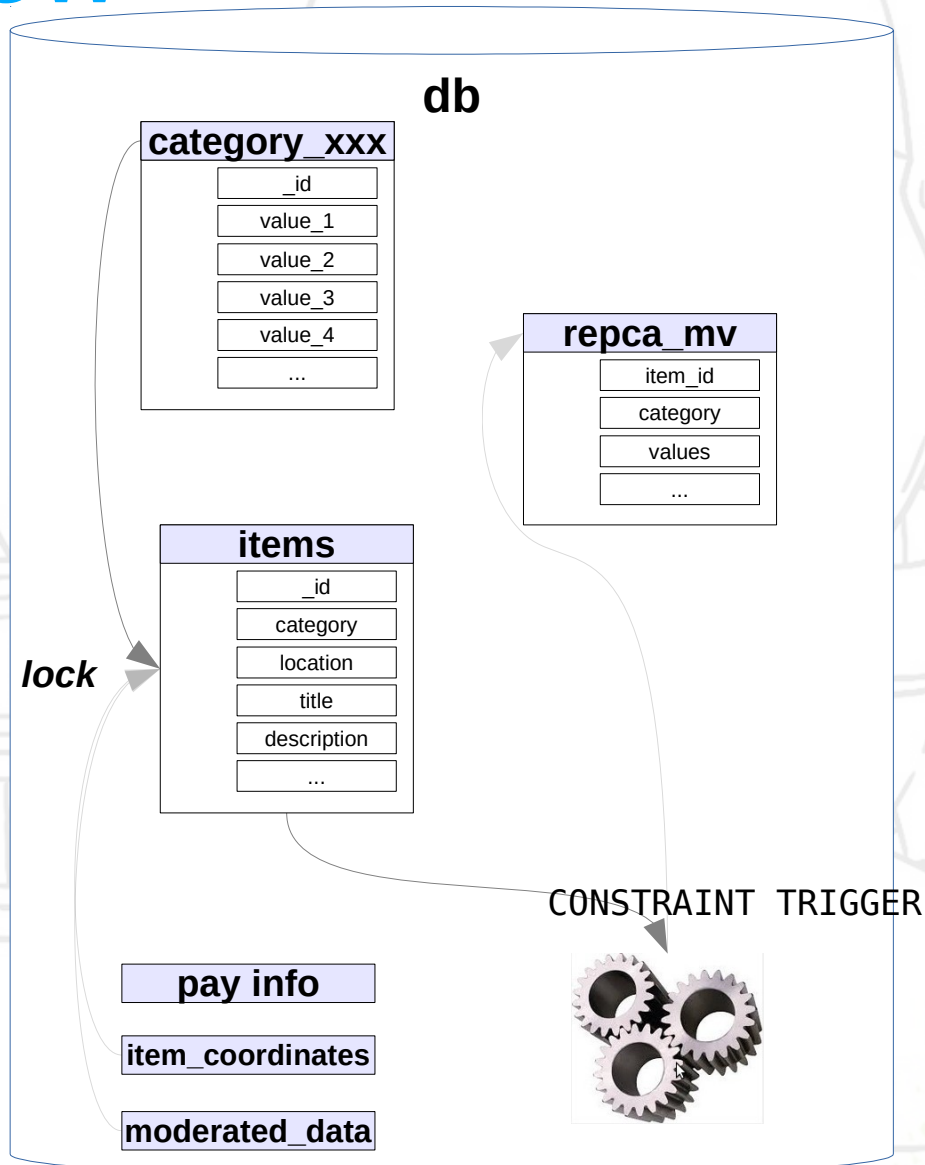
Standbys pool



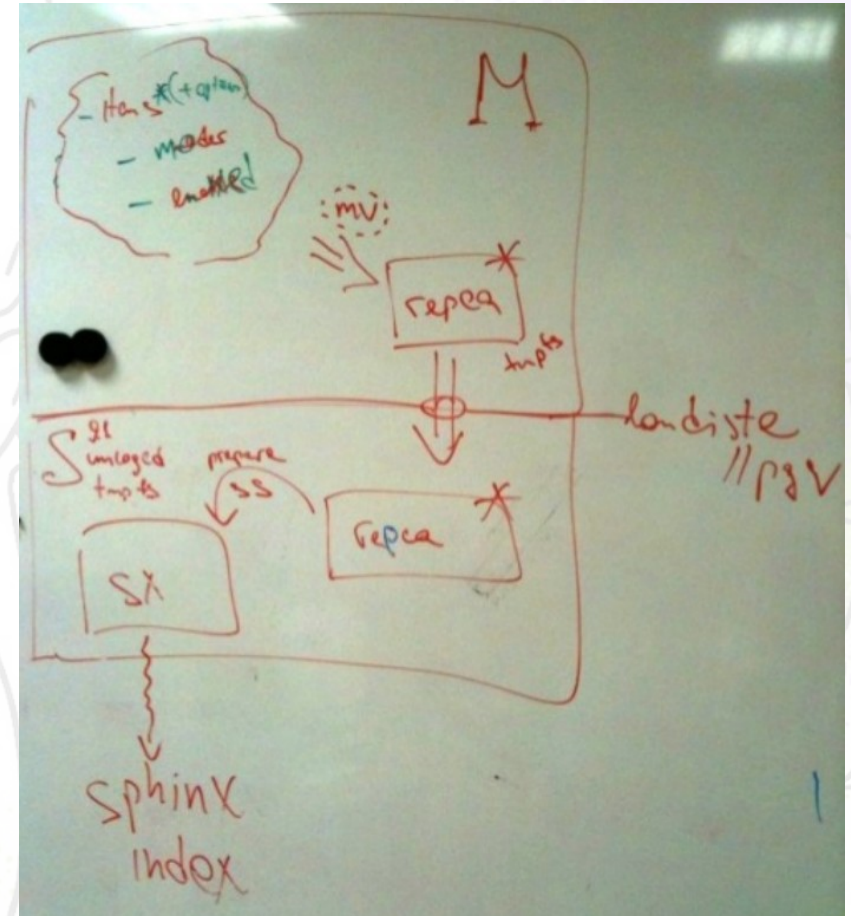
HAProxy check function

```
if master
    then false
if lag > max
    then create file and return false
if lag > min and file exists
    then return false
if lag < min and file exists
    then remove file and return true
else
    true
```

MatView

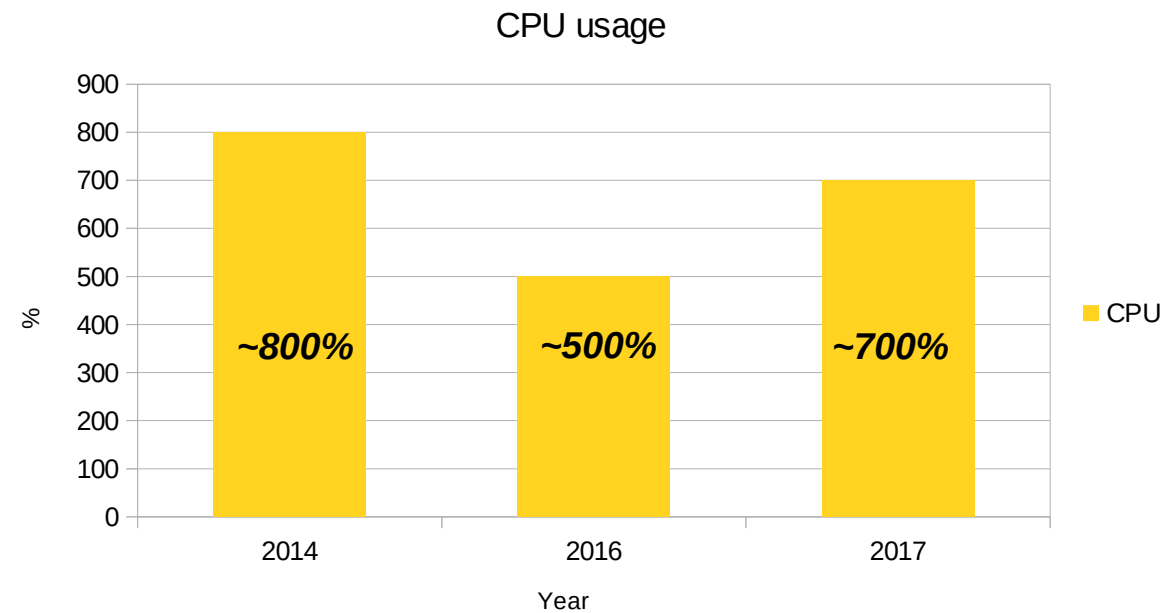
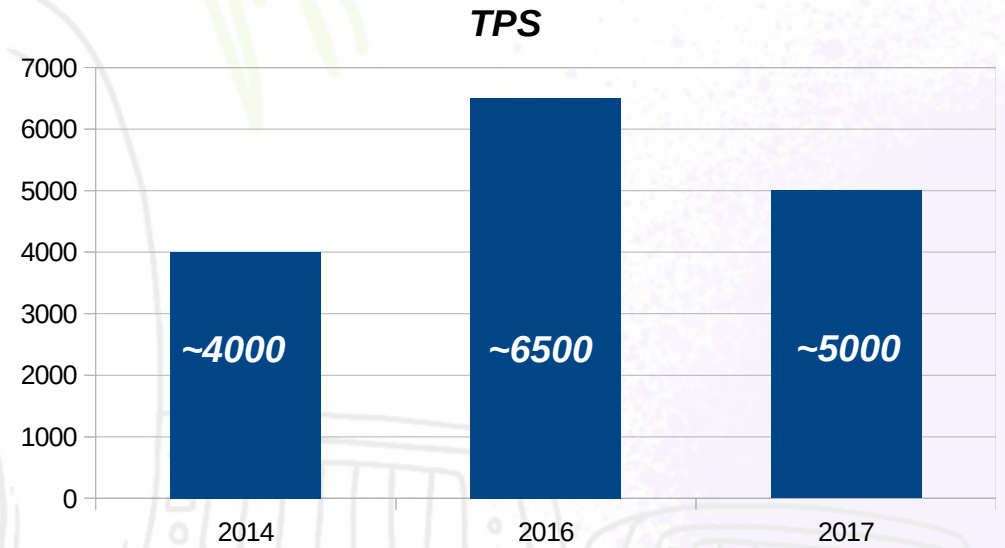
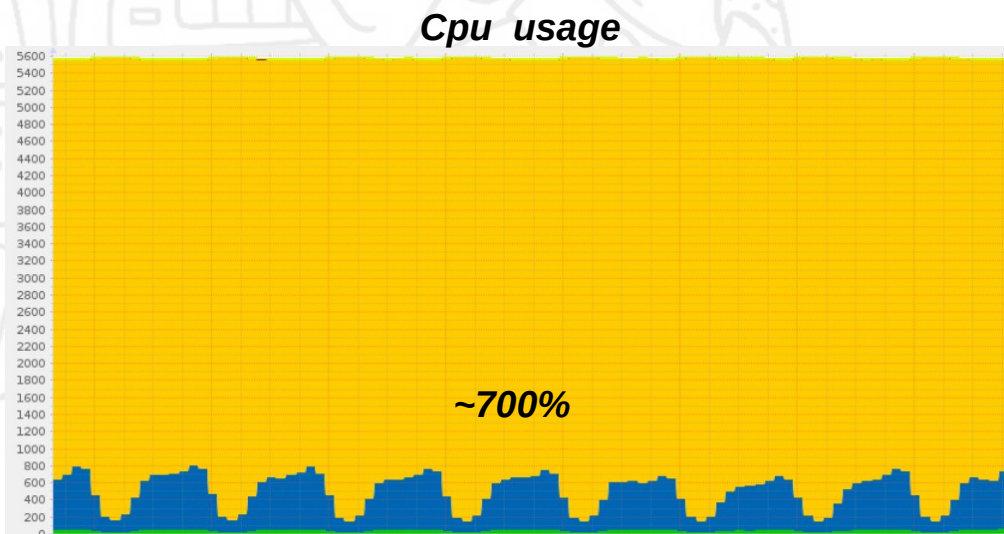
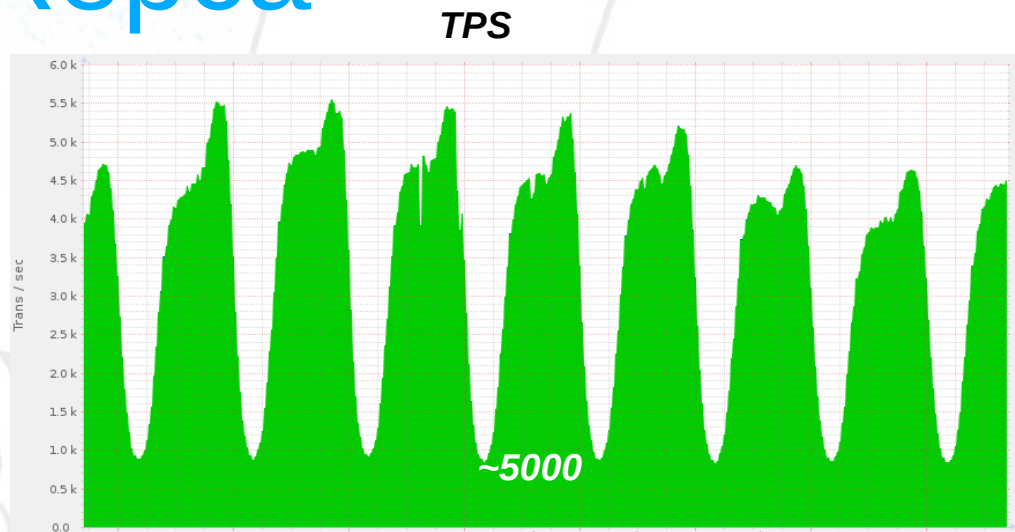


!2009 Michael Tyurin



<https://goo.gl/sXABqK>

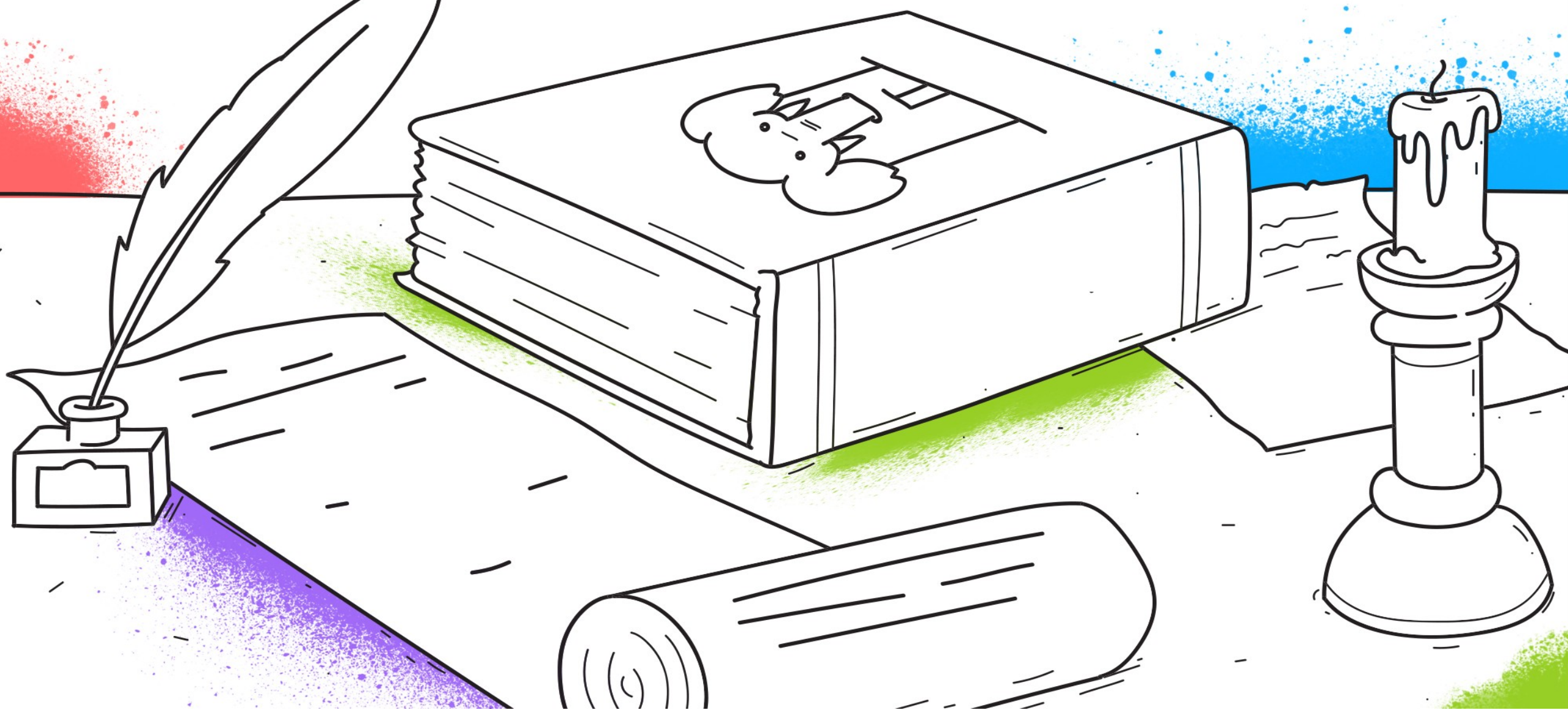
Repca



Conclusion



- There are a few kinds of standbys
- We can scale reads with the help of standby
 - ignore stale reads
 - logical routing
 - hot cache
- There are some pitfalls with standby in production
- Archive and backup depends on your DRP
- Major upgrade with standby also needs advanced manipulations



<https://www.avito.ru/company/job/dp-eng>