

# Investigate High Availability performance problems

**adyen**

engineered  
for ambition



**HIGH AVAILABLE  
CONFIGURATIONS ARE VERY  
COMMON FOR POSTGRESQL.  
BUT HOW DO YOU  
INVESTIGATE PERFORMANCE  
PROBLEMS WHEN THE  
STANDBY CAN'T KEEP UP?**



engineered  
for ambition





# Boriss Mejias

Senior solution architect

EDB

adyen

engineered  
for ambition





# Derk van Veen

~~Pratiseuse specialist~~

## Adyen



**adyen**

engineered  
for ambition

# Agenda



A  
C

B



# Concepts

?



**Problem**



**Investigation**





**Solution**



**Lessons  
learned**

A

C

B

# Concepts

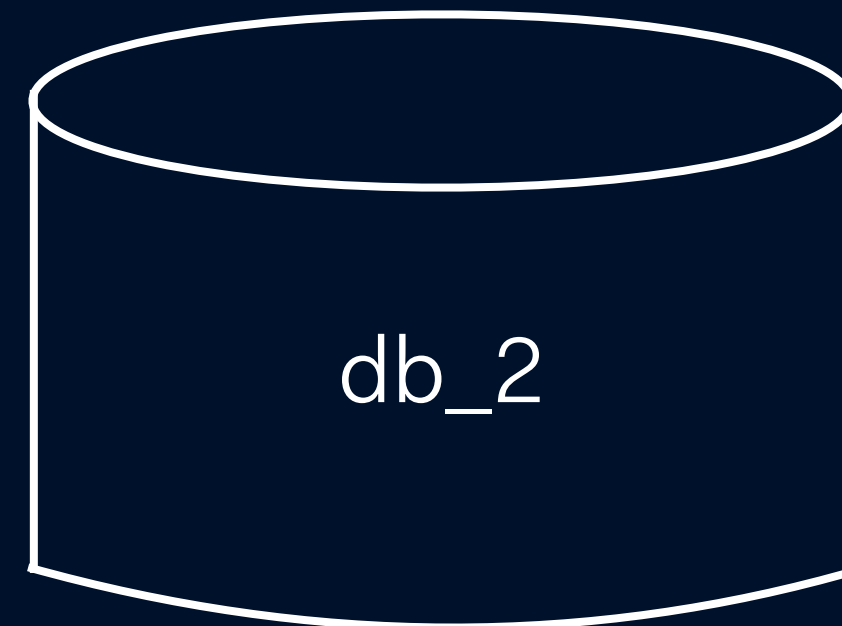
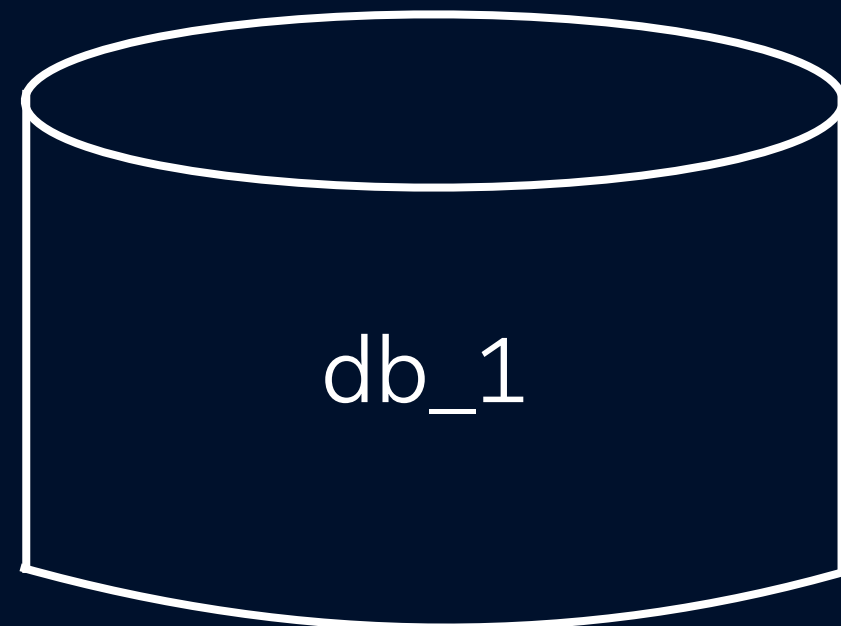
A

C

B

# Cluster

Shared buffers

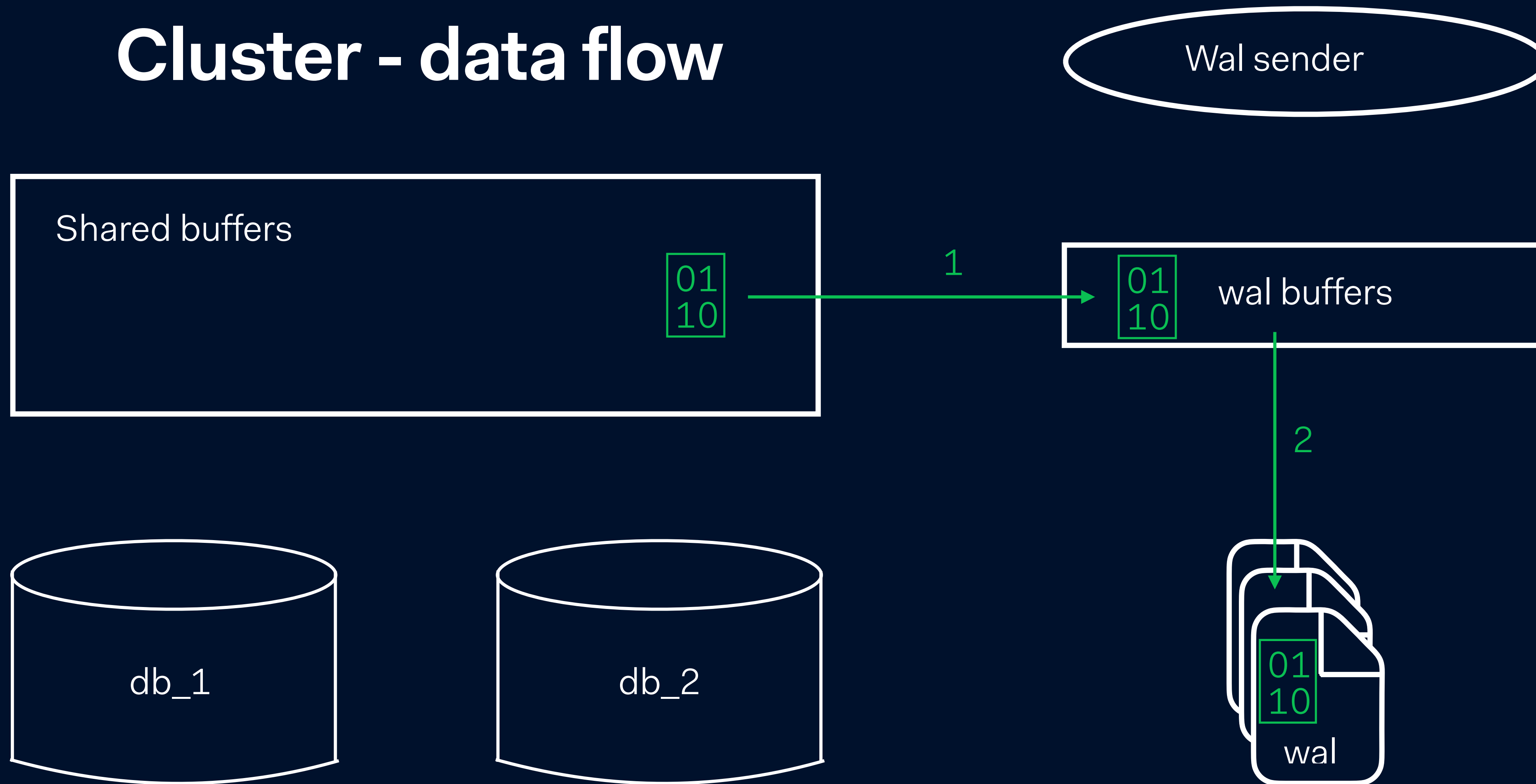


Wal sender

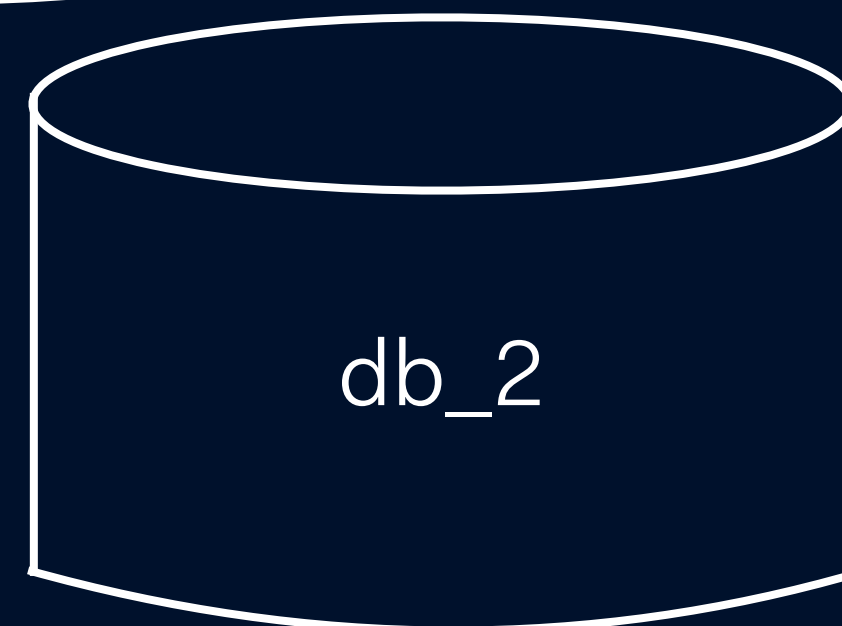
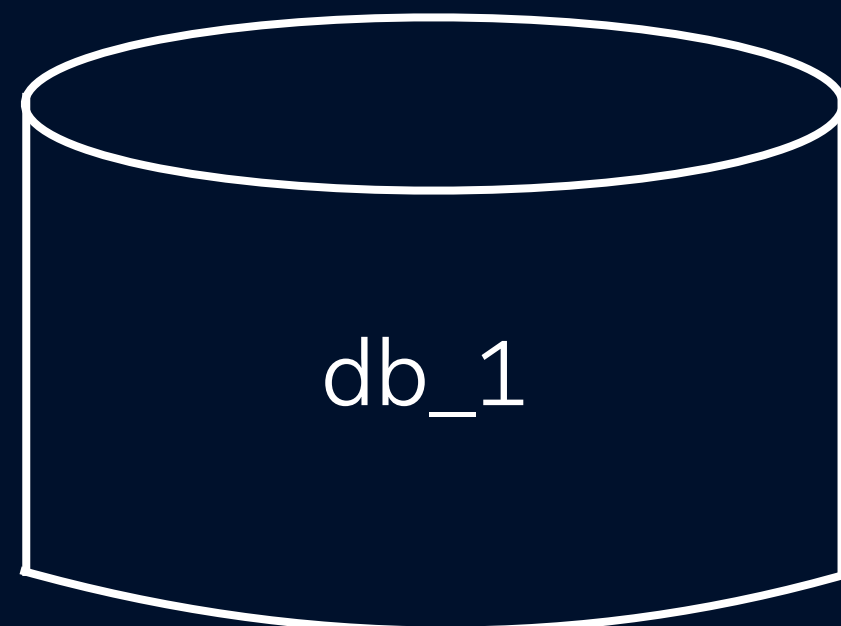
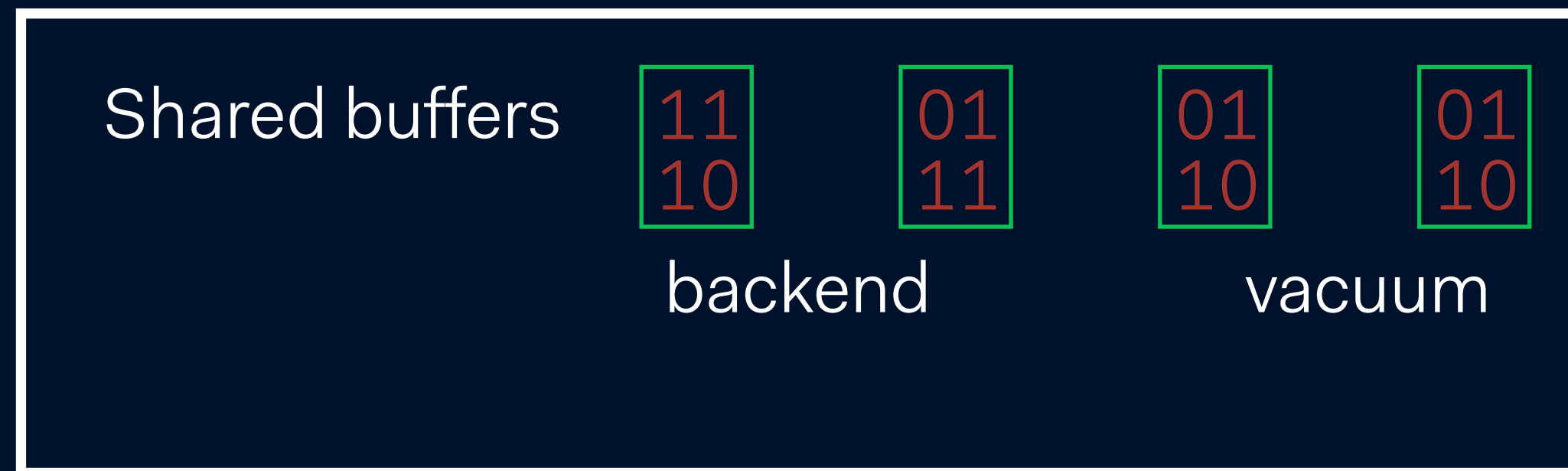
wal buffers



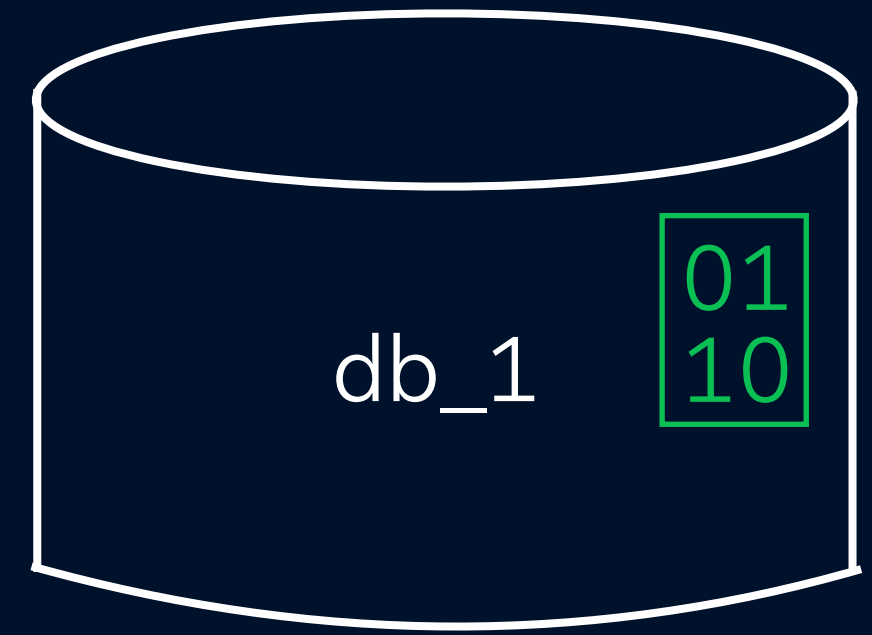
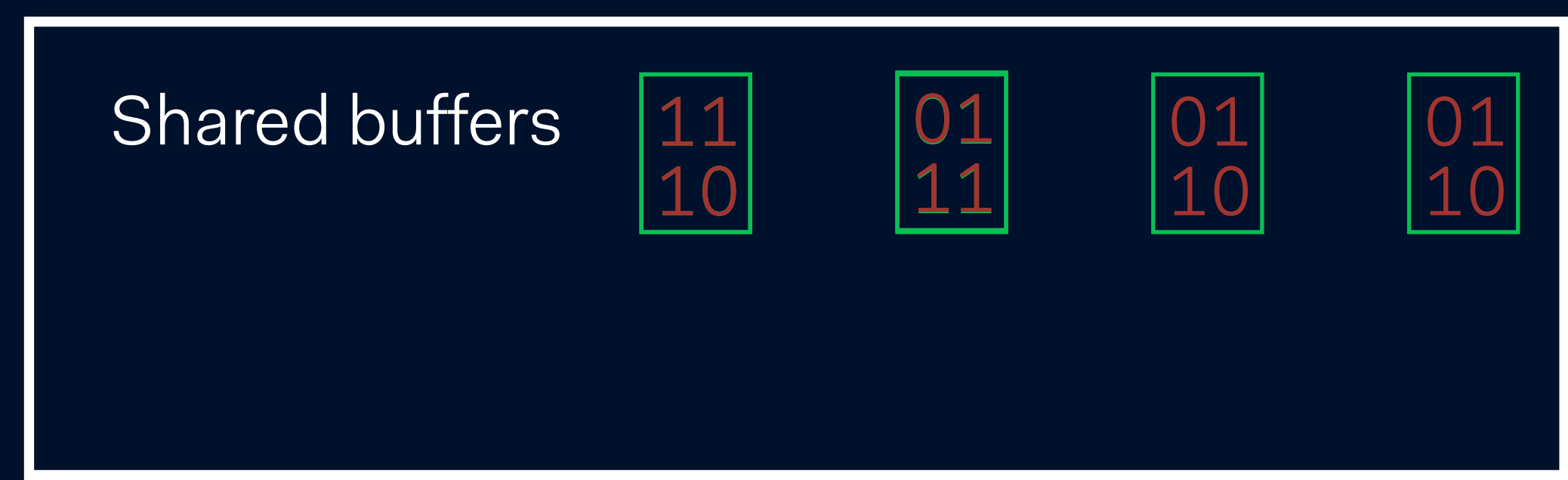
# Cluster - data flow



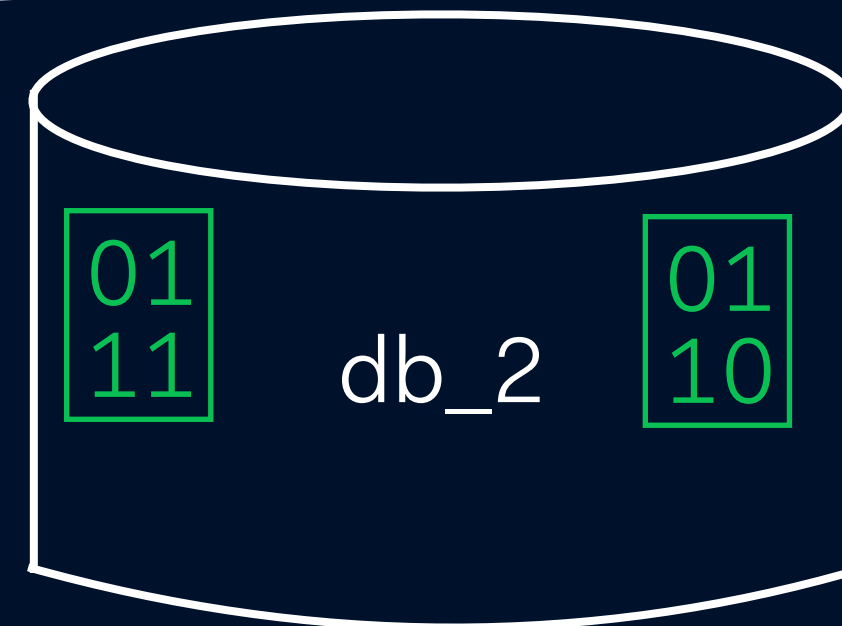
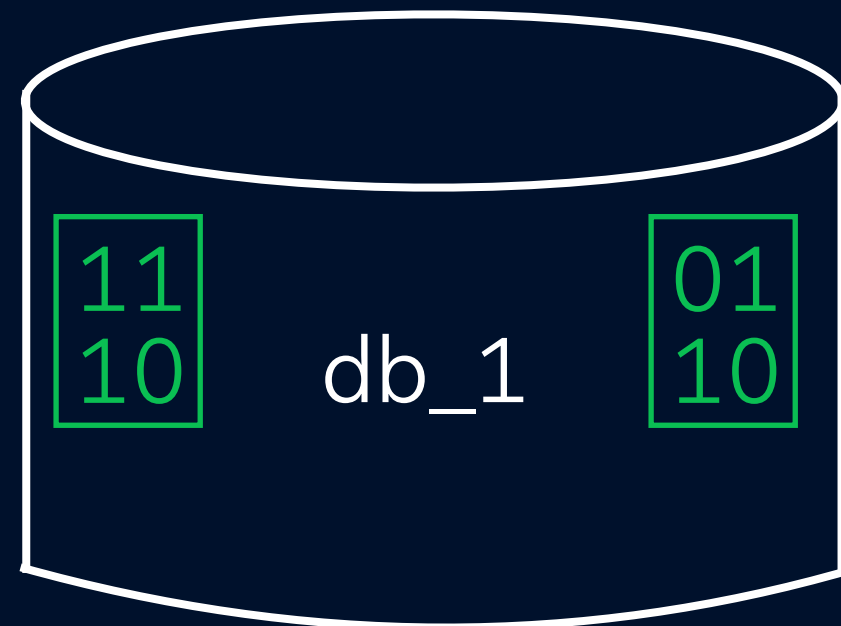
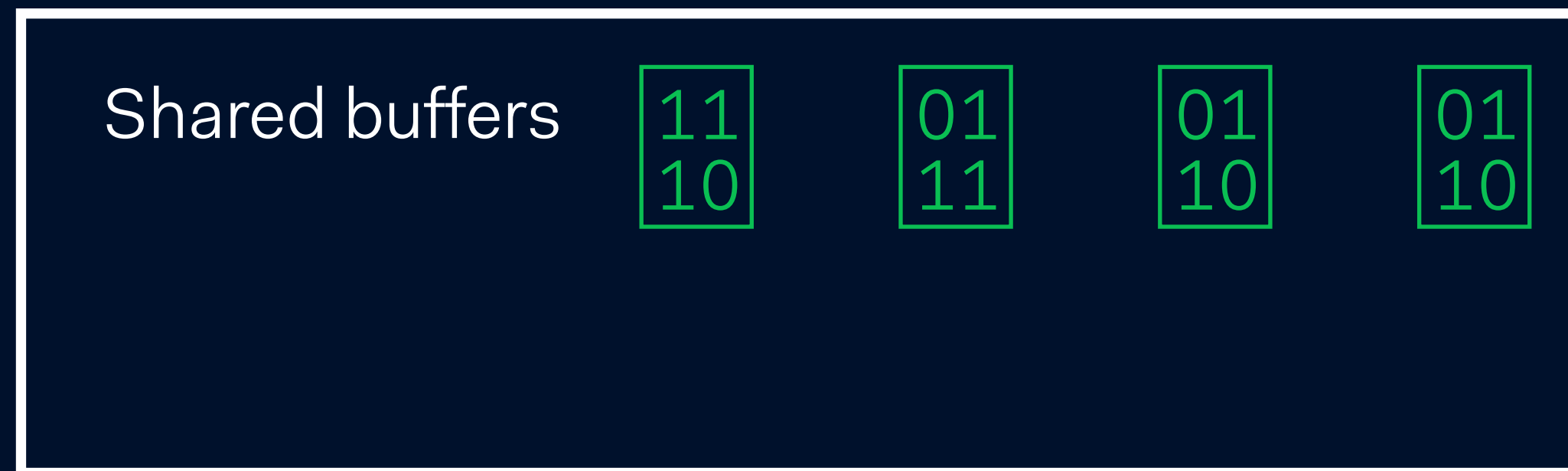
# Cluster - checkpoint



# Cluster - checkpoint



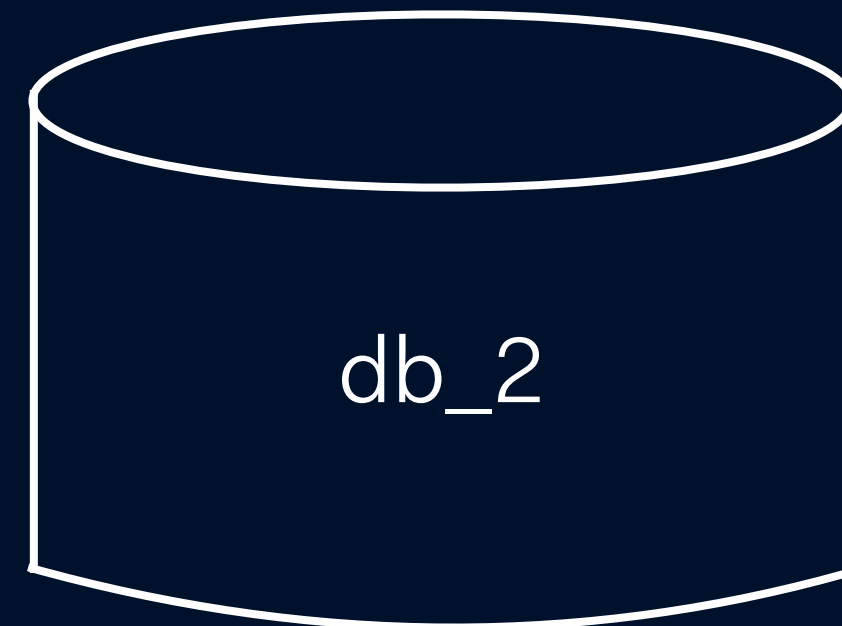
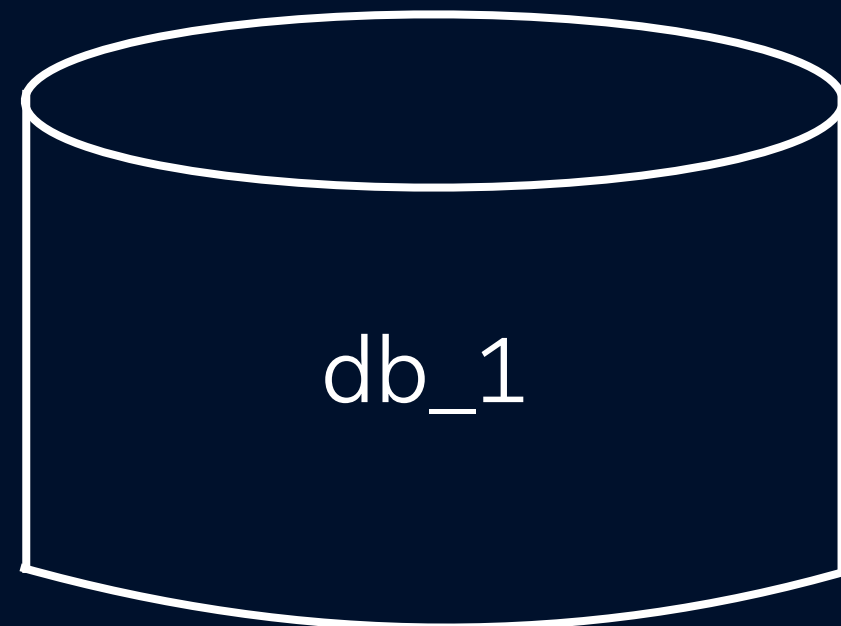
# Cluster - checkpoint





# Cluster

Shared buffers



Wal sender

wal buffers



# Primary

Shared  
buffers

Wal sender

wal buffers

db\_1

wal

# HA setup

## Primary

Shared buffers

Wal sender

wal buffers

db\_1

wal

## Standby

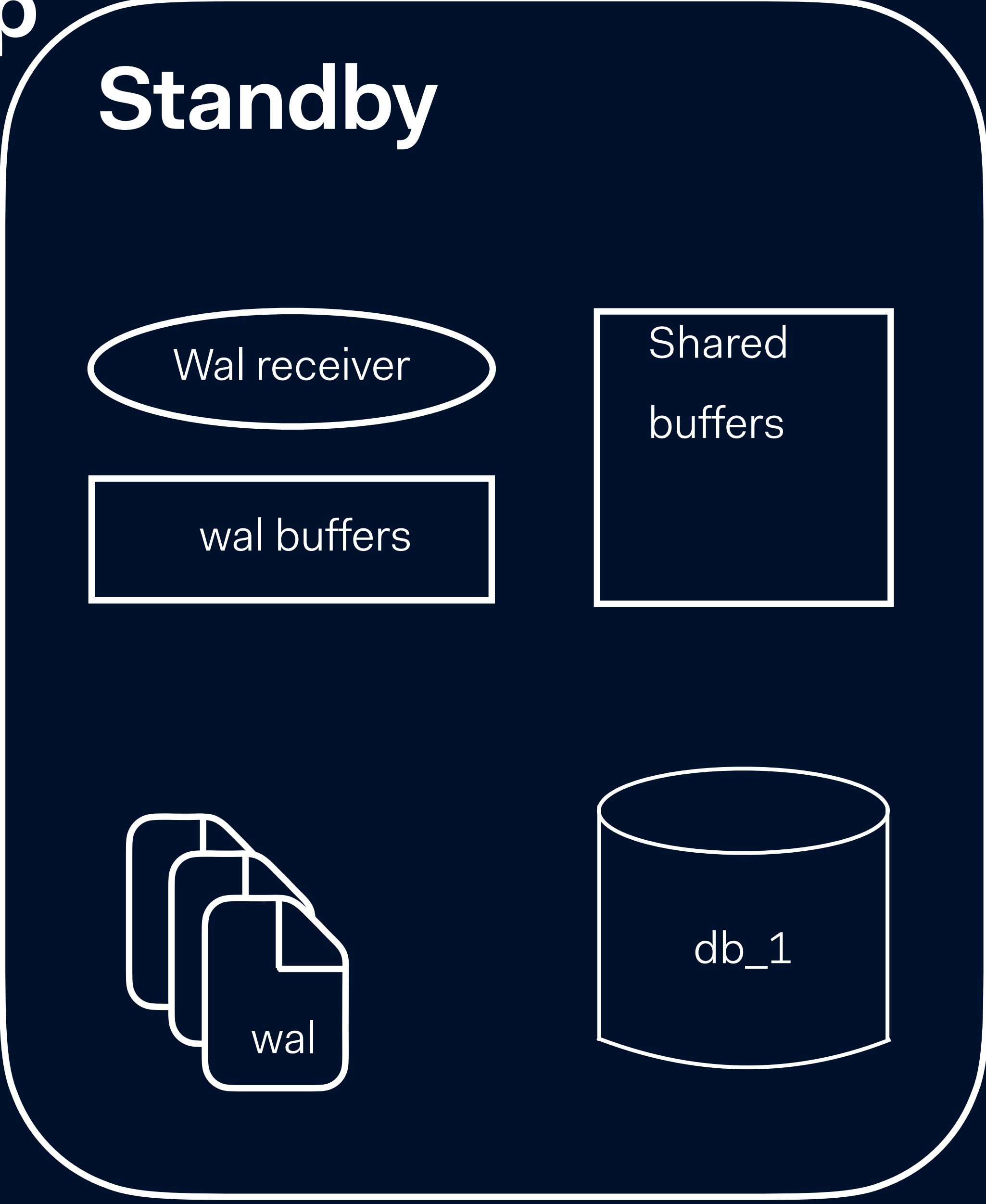
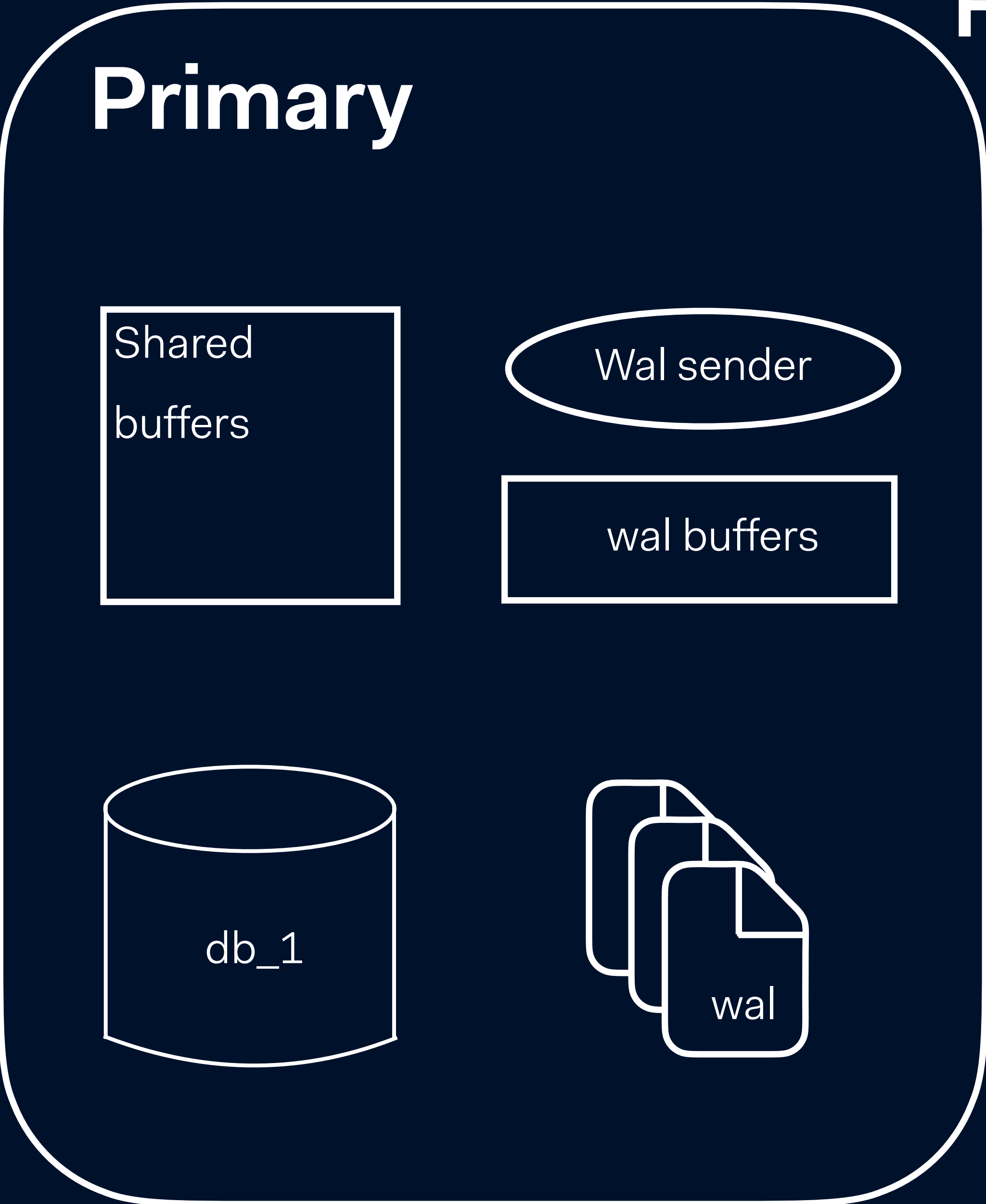
Wal receiver

Shared buffers

wal buffers

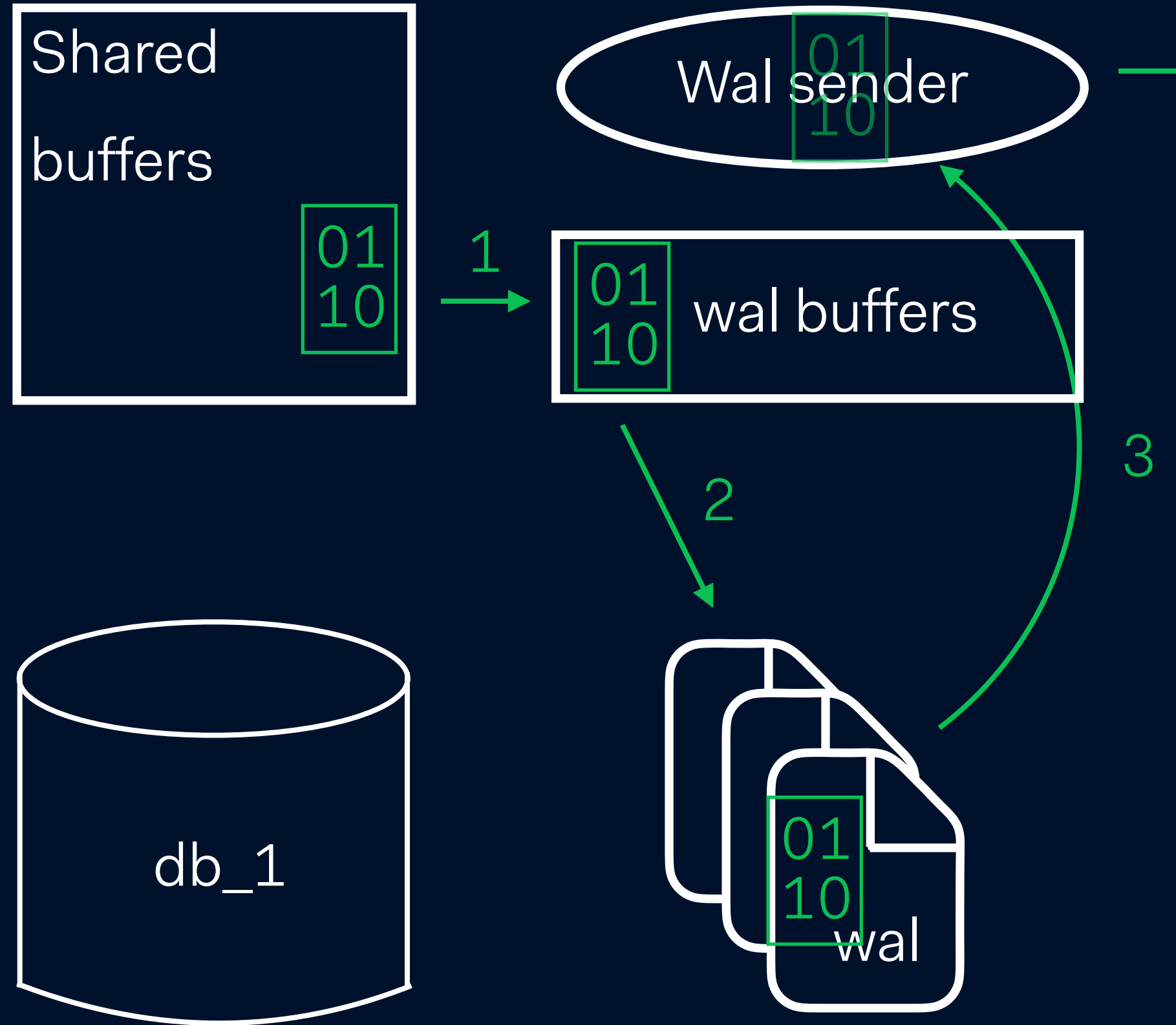
wal

db\_1

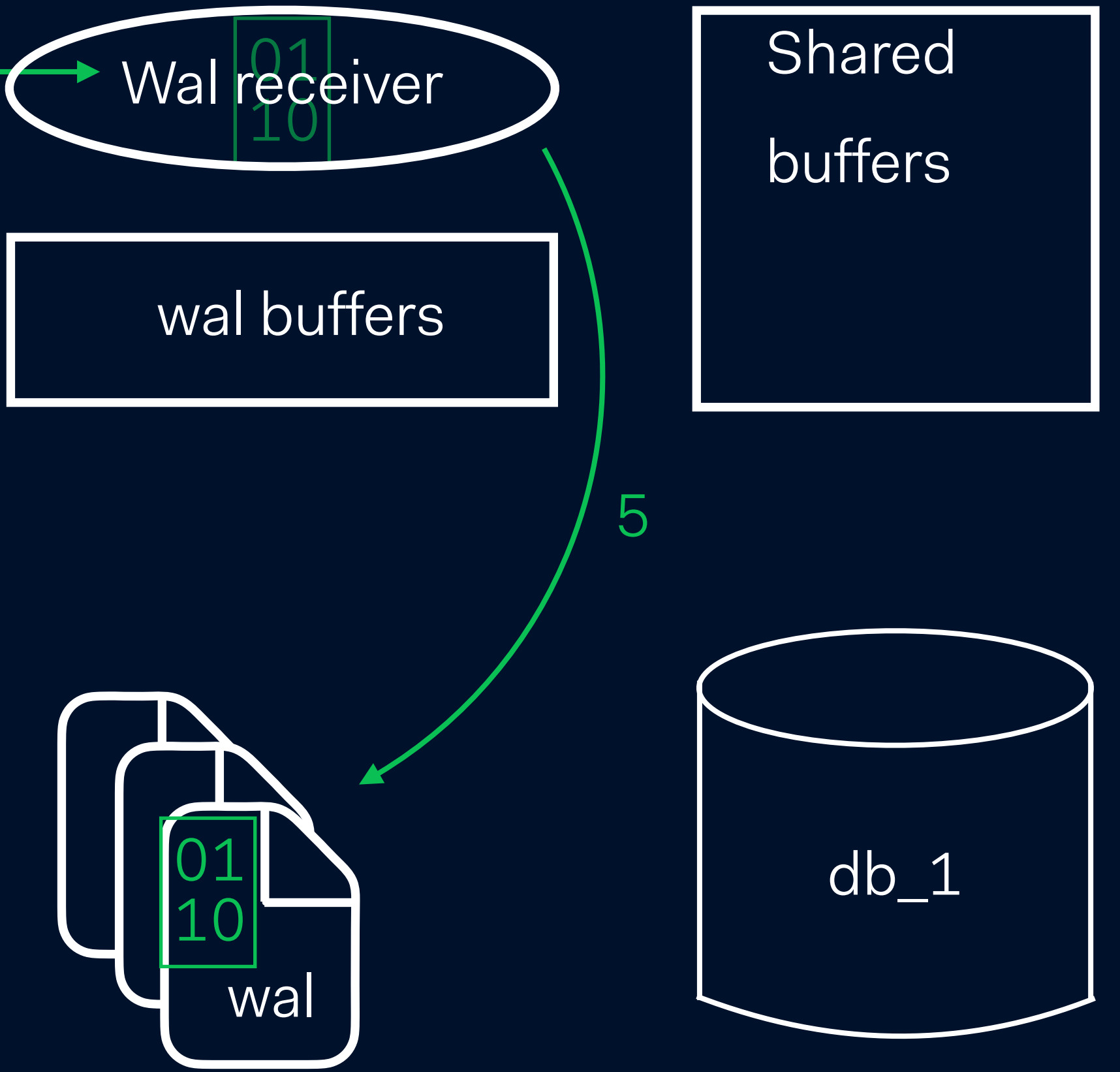


# Data flow

## Primary

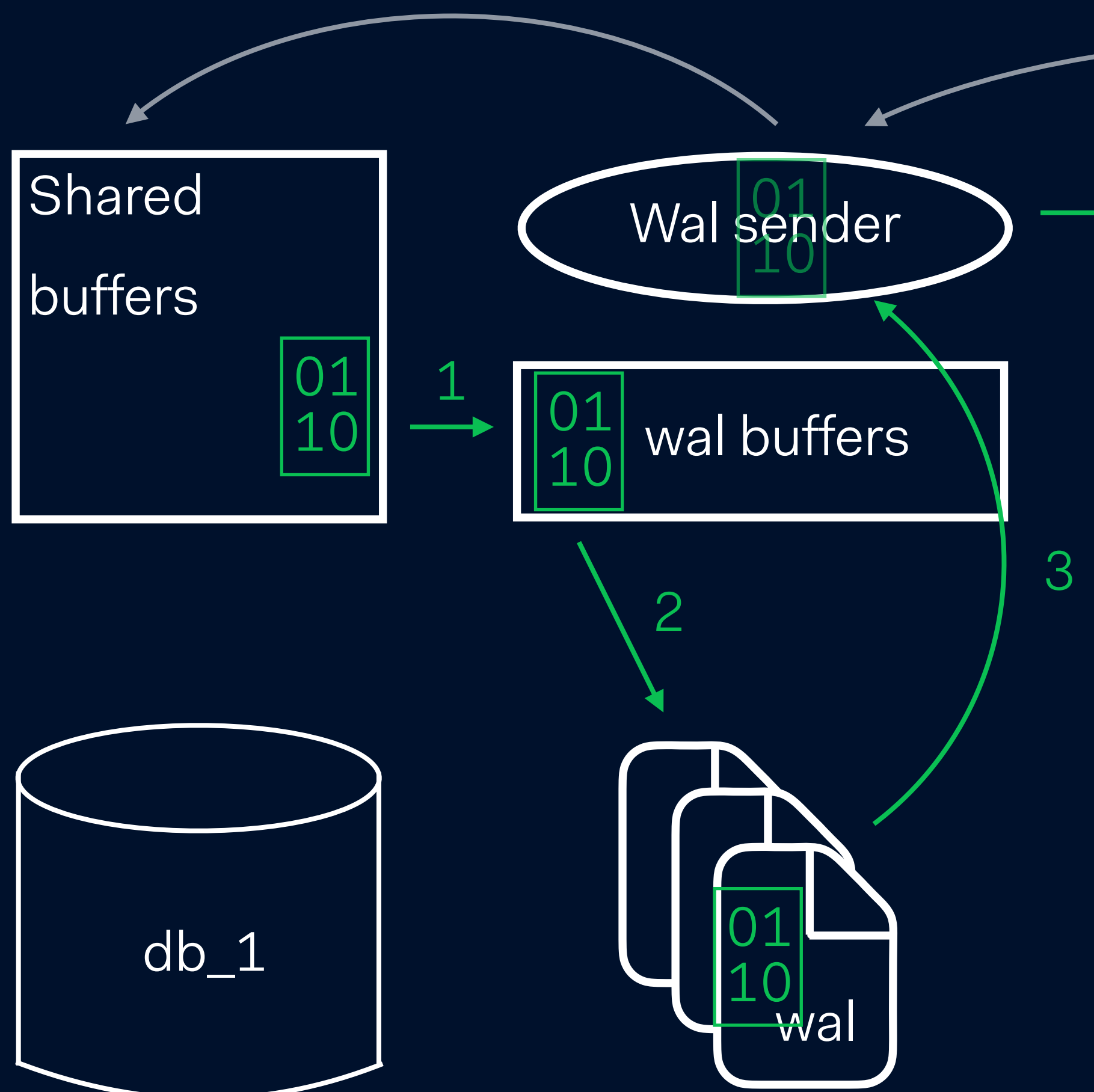


## Standby

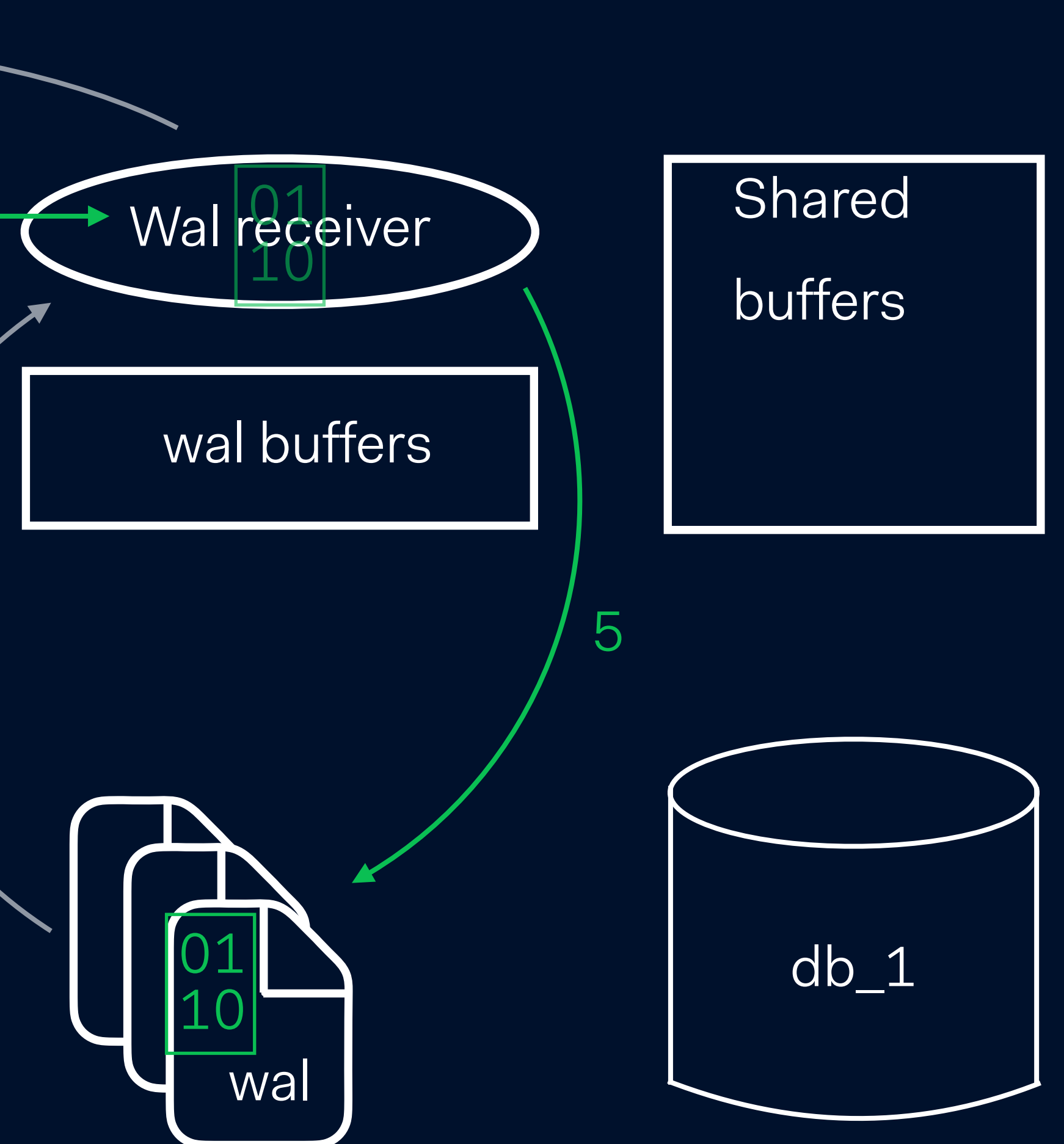


# Confirm

## Primary



## Standby



Ack

4

3

2

1

Ack

5

db\_1

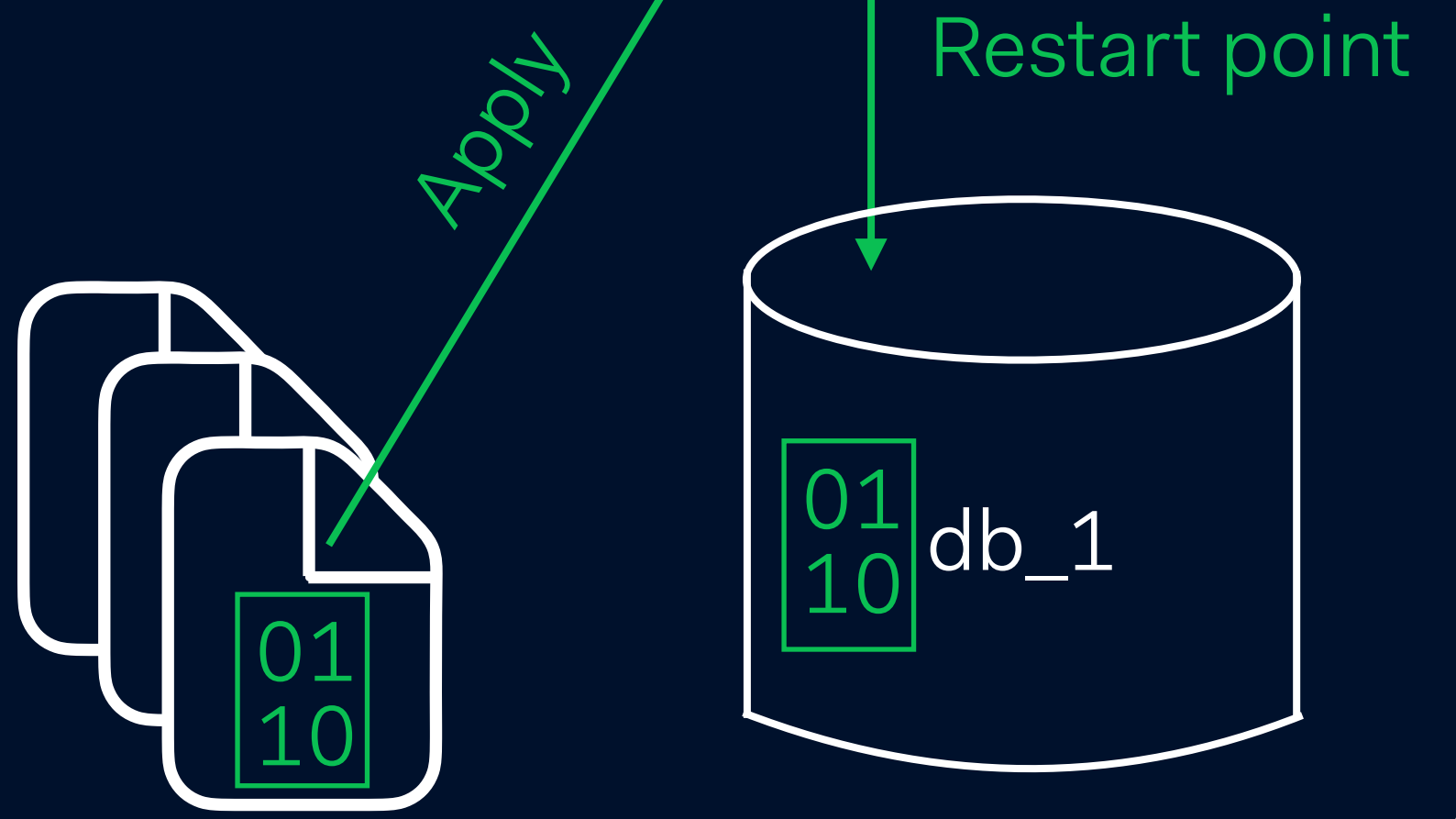
db\_1

# Restart point

## Primary



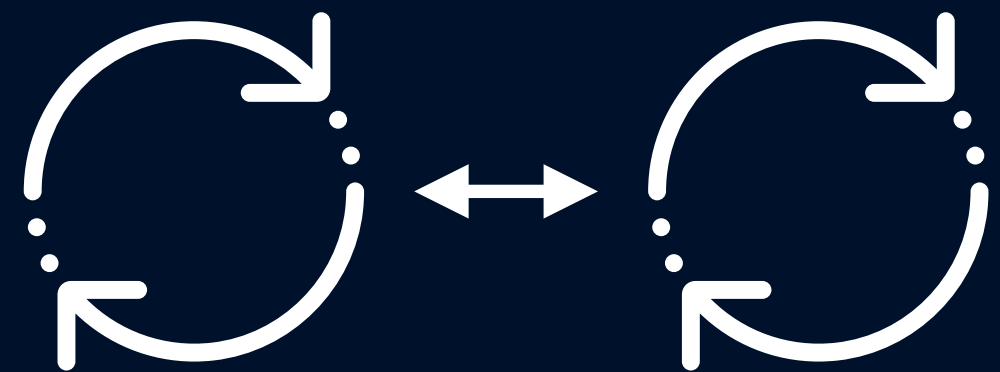
## Standby



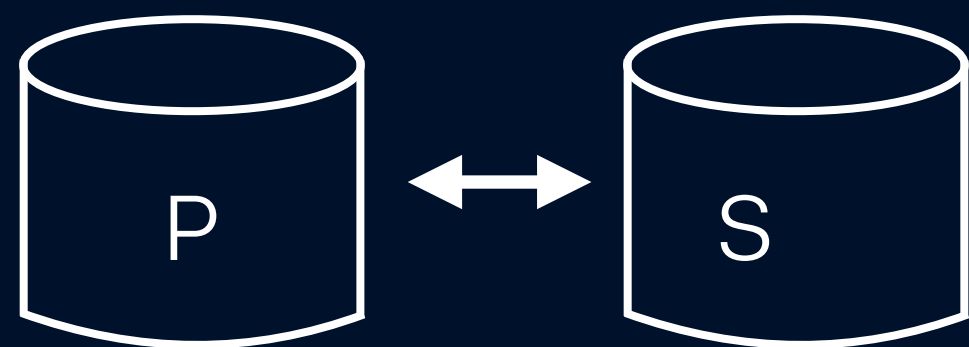
# Wait events



**Wait events**



**IPC wait events**



**Syncprep wait events**

# SyncRep

## Primary



Shared buffers

Wal sender



wal buffers

db\_1



## Standby

Wal receiver

Shared buffers

wal buffers



db\_1



# SyncRep

## Primary



Shared buffers

Wal sender

wal buffers

db\_1

wal

## Standby

Wal receiver

wal buffers

Shared buffers

db\_1

wal



# SyncRep

## Primary



Shared buffers

Wal sender

wal buffers

db\_1

wal

## Standby

Wal receiver

wal buffers

Shared buffers

db\_1

01  
10

wal



# SyncRep

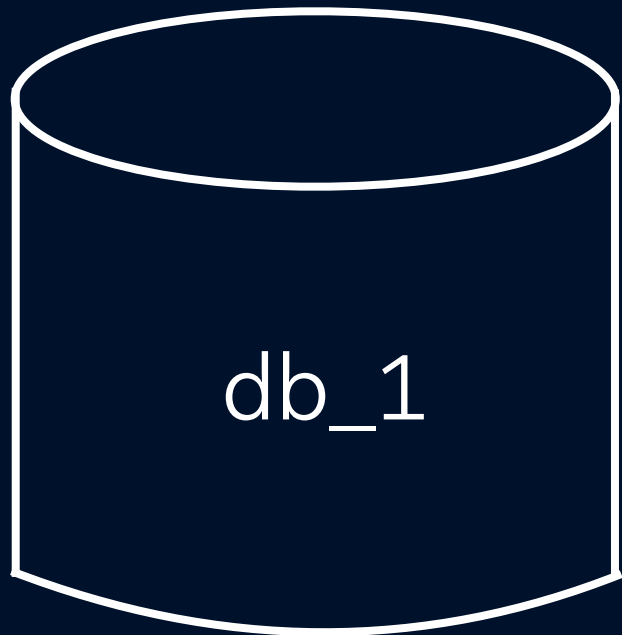
## Primary



Shared buffers

Wal sender

wal buffers

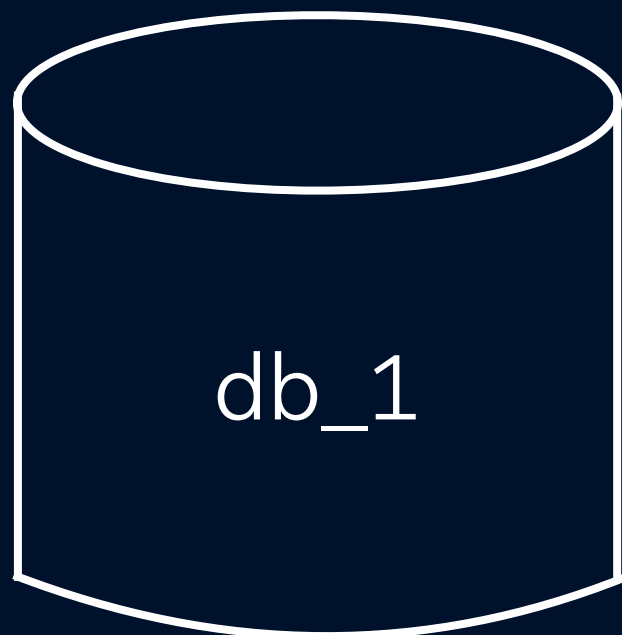


## Standby

Wal receiver

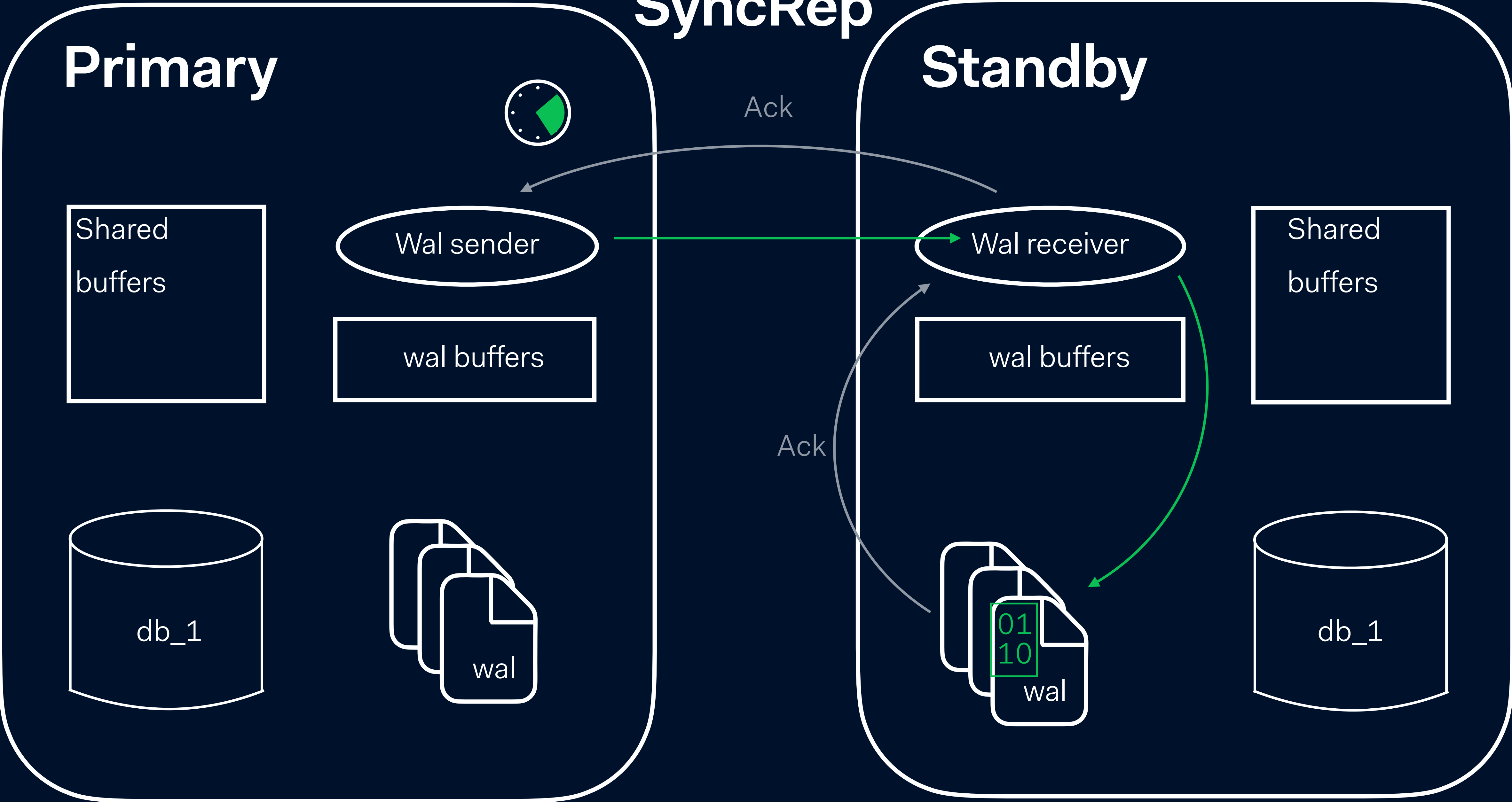
wal buffers

Shared buffers



Ack

Ack





**Problem**

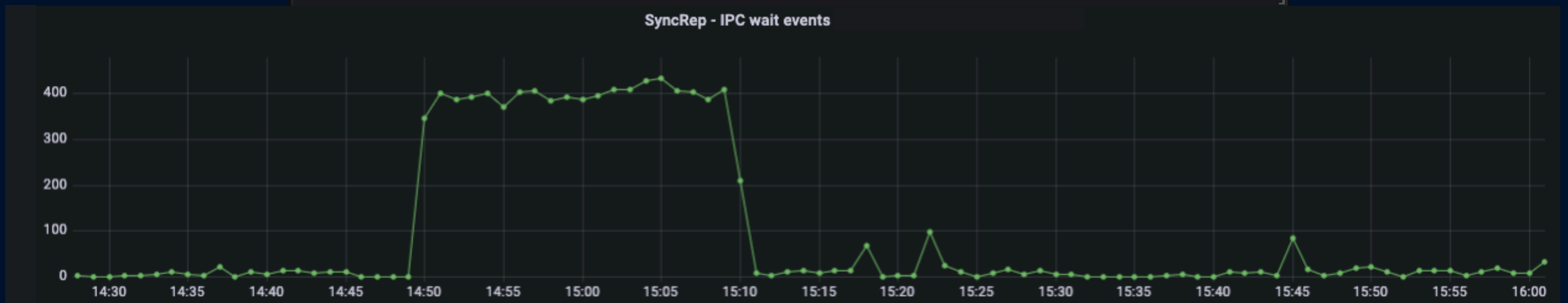


# Job is delayed

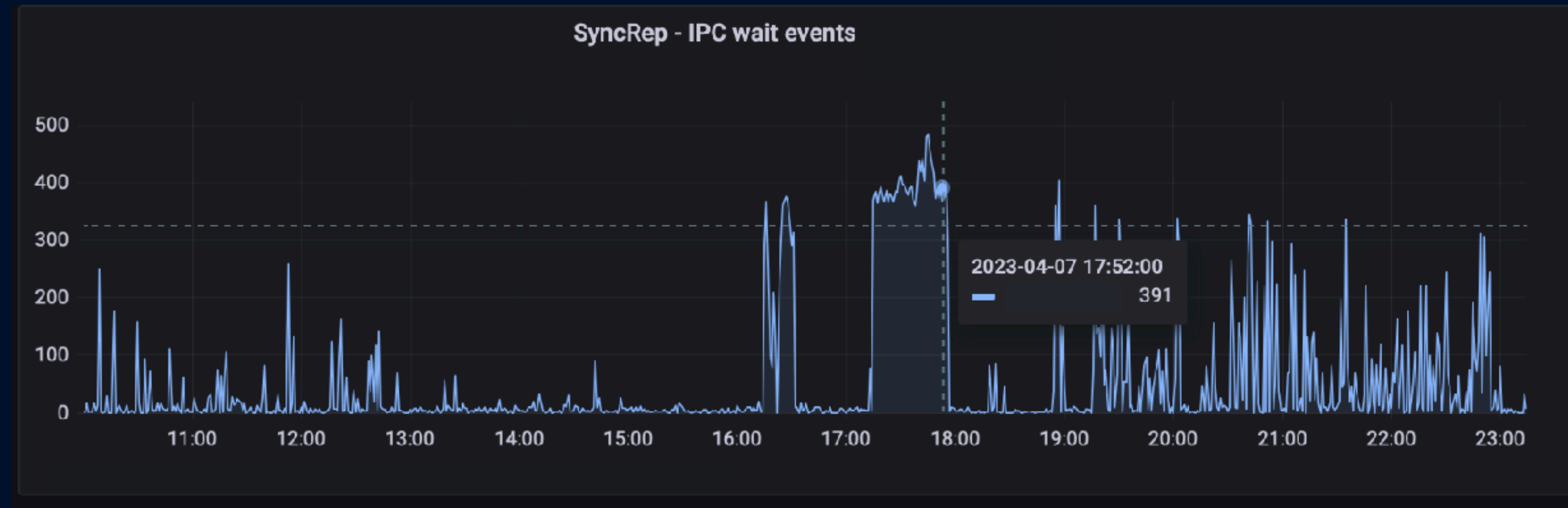




# Job is delayed

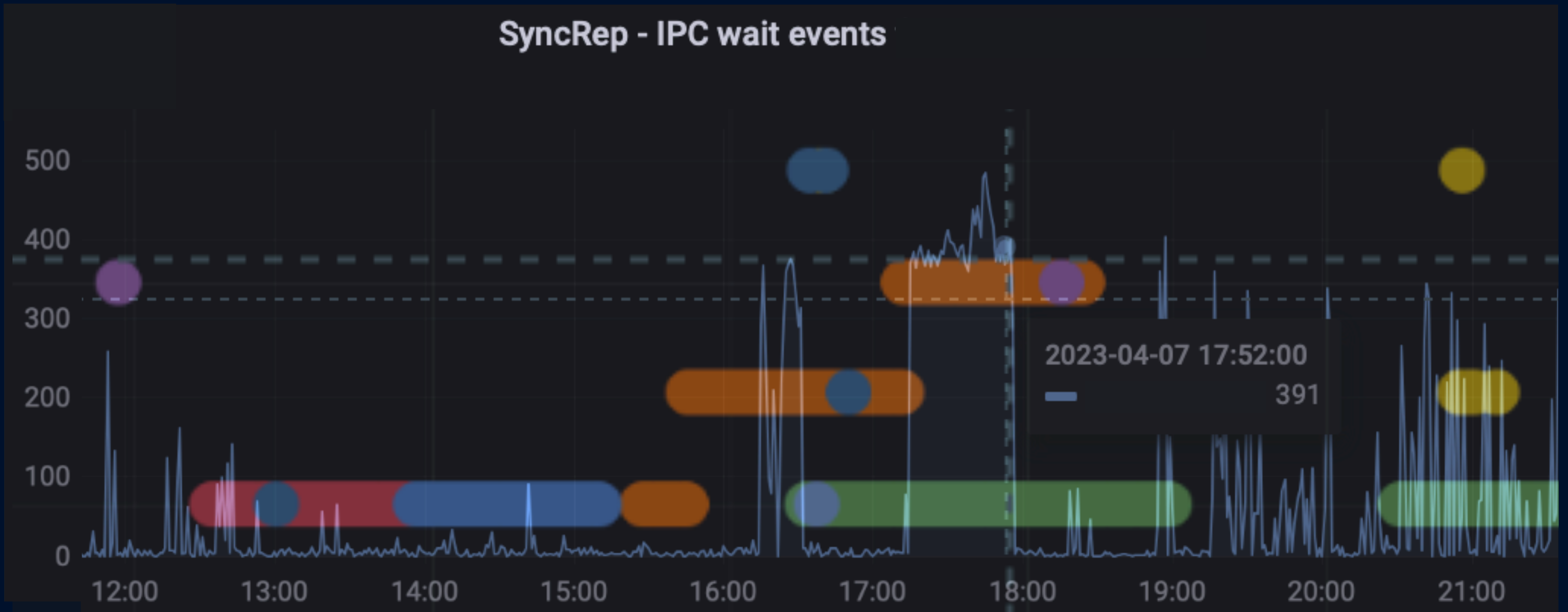


# Job is delayed



# Job is delayed

SyncRep - IPC wait events





**Mitigation:** make vacuum  
slower

# SyncRep

## Primary



Shared buffers

Wal sender

wal buffers

db\_1

wal

## Standby

Wal receiver

wal buffers

Shared buffers

db\_1

wal





# Problem description



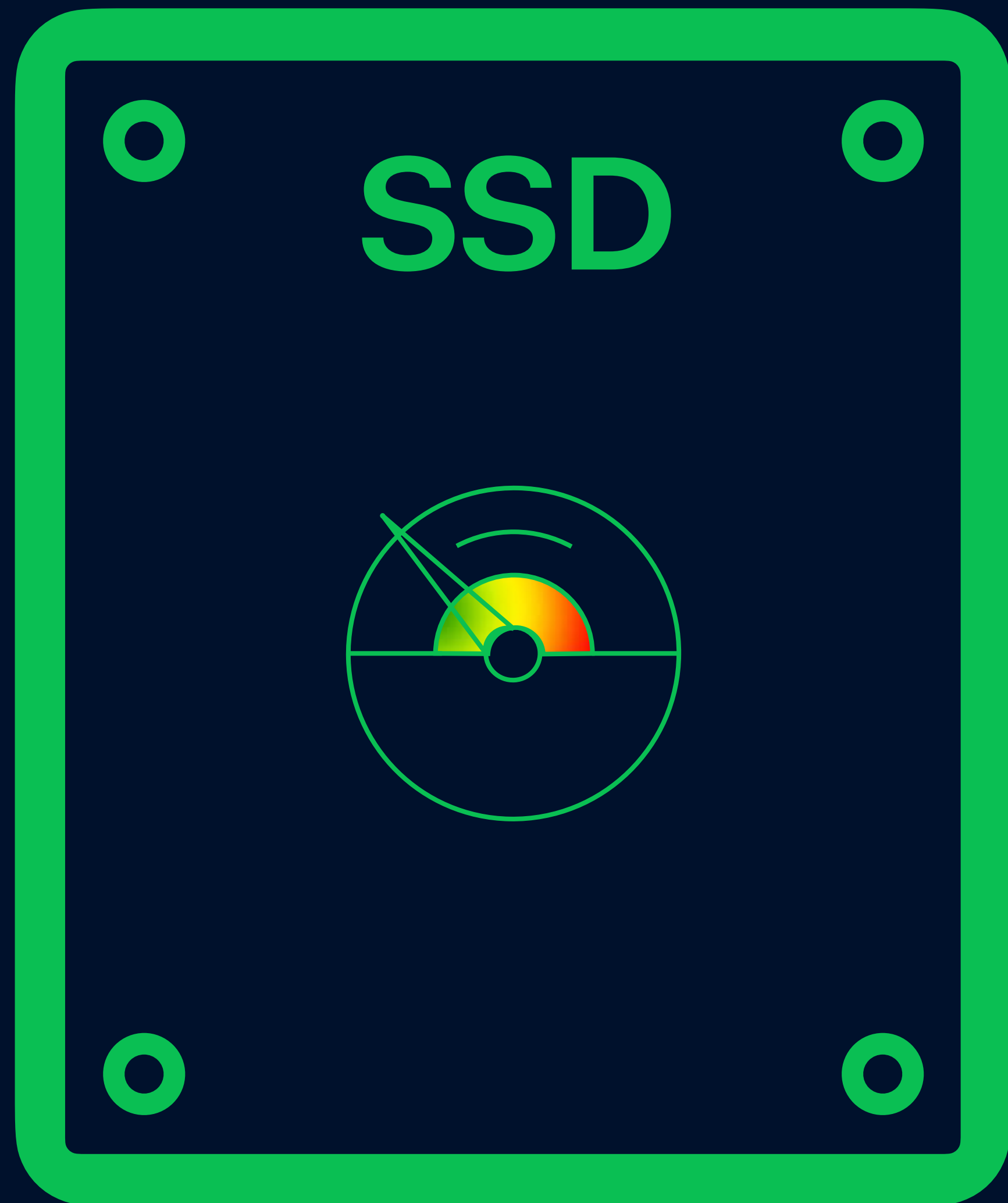
- We have significant job delays related to SyncRep
- Strong correlation with vacuum on large tables
- Increase in SyncRep directly after checkpoint starts
- While SyncRep piles up, transaction rate goes down
- Terminate vacuum helps
- Delay is equal for both standbys



# **Problem investigation**

Lets **blame** the firewall/  
**network**

Lets **blame** the  
**storage**





# SyncRep

## Primary



Shared buffers

Wal sender

wal buffers

db\_1

wal

## Standby

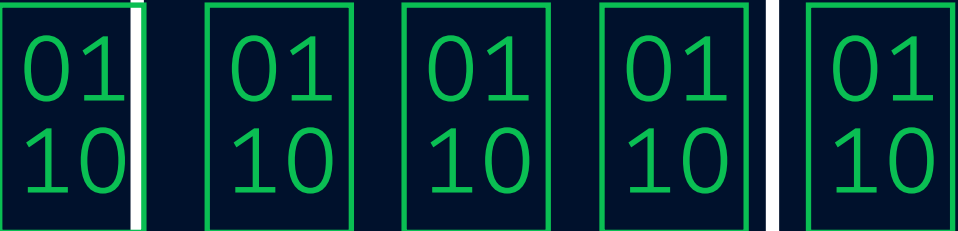
Wal receiver

wal buffers

Shared buffers

db\_1

wal



# SyncRep

## Primary



Shared buffers

Wal sender

wal buffers

db\_1

wal

01 01 01 01 01  
10 10 10 10 10

## Standby

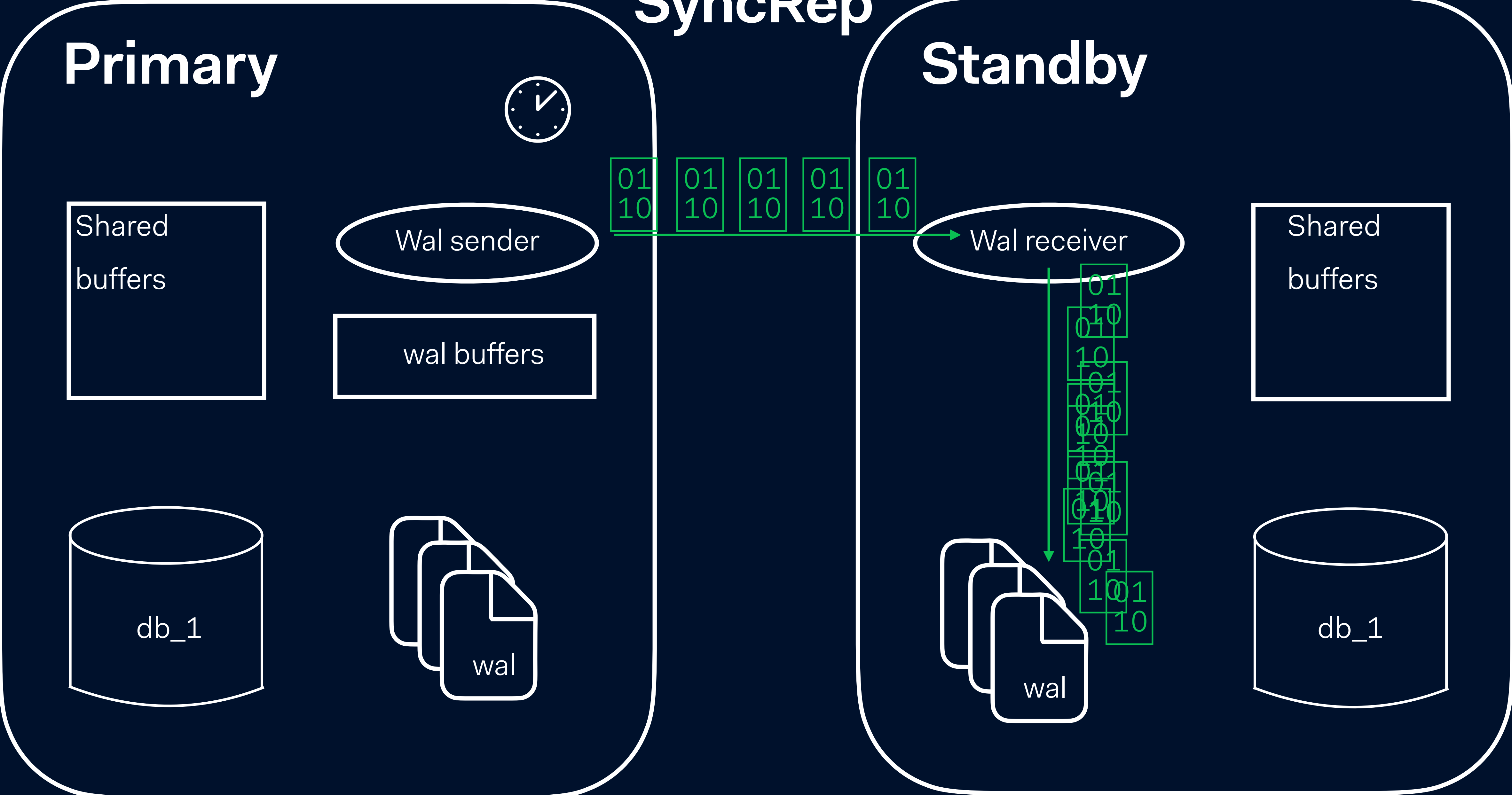
Wal receiver

Shared buffers

db\_1

wal

01  
010  
10  
01  
010  
10  
01  
010  
10  
01  
10  
01  
10



# If everything is sequential

- Tcp
- Wal Isn
- OS writes
- Commits

## What is the theoretical limit?

“**Standby** can’t keep up  
with the amount of **wal**  
writes”

# Lets start some debug

- SAR report
- IOStat
- SOSReport
- ss
- perf
- ioping
- pidstat
- Netstat
- iotop

# Lets start some debug

- SAR report
- IOStat
- SOSReport
- ss
- perf
- ioping
- pidstat
- netstat
- iotop

# SS

“A **socket** is one **endpoint** of a **two-way communication** link between two programs running on the **network**”

# Socket Statistics

<b>Round Trip Time</b>	average round trip time with mean deviation
<b>Recv-Q</b>	receive buffer, application needs to retrieve it
<b>Send-Q</b>	send buffer; not sent or sent but not ACKed
<b>lastsnd</b>	How long since last package sent



# Socket Statistics

rtt:5.06/0.428 send 254.1Mbps

rtt:2.657/2.413 send 457.8Mbps

# Socket Statistics

skmem:(

r6142208,

97.6% rec buffer full

rb6291456,

t0,

tb332800,

f161536,

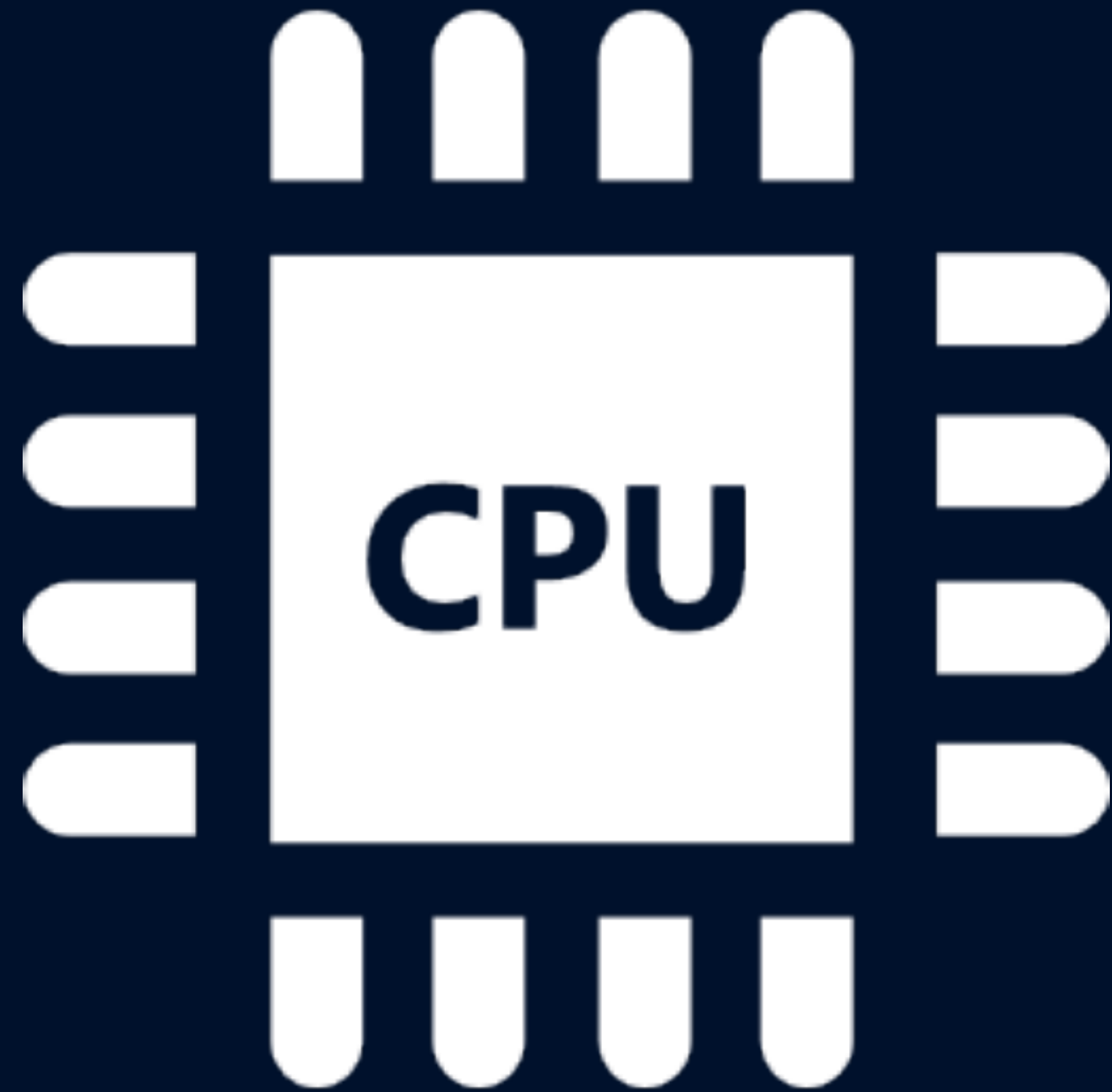
w0,

o0,

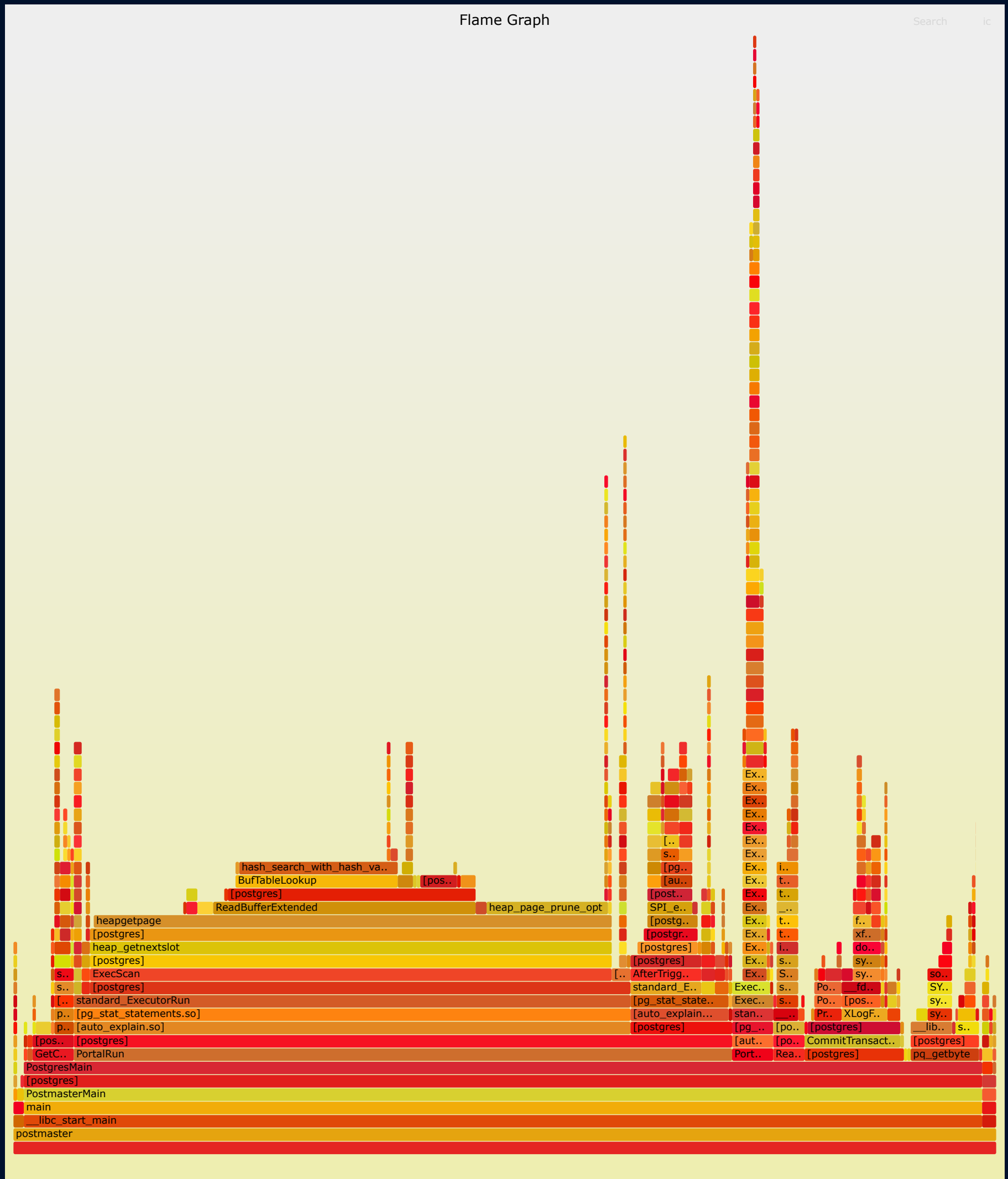
bl0,

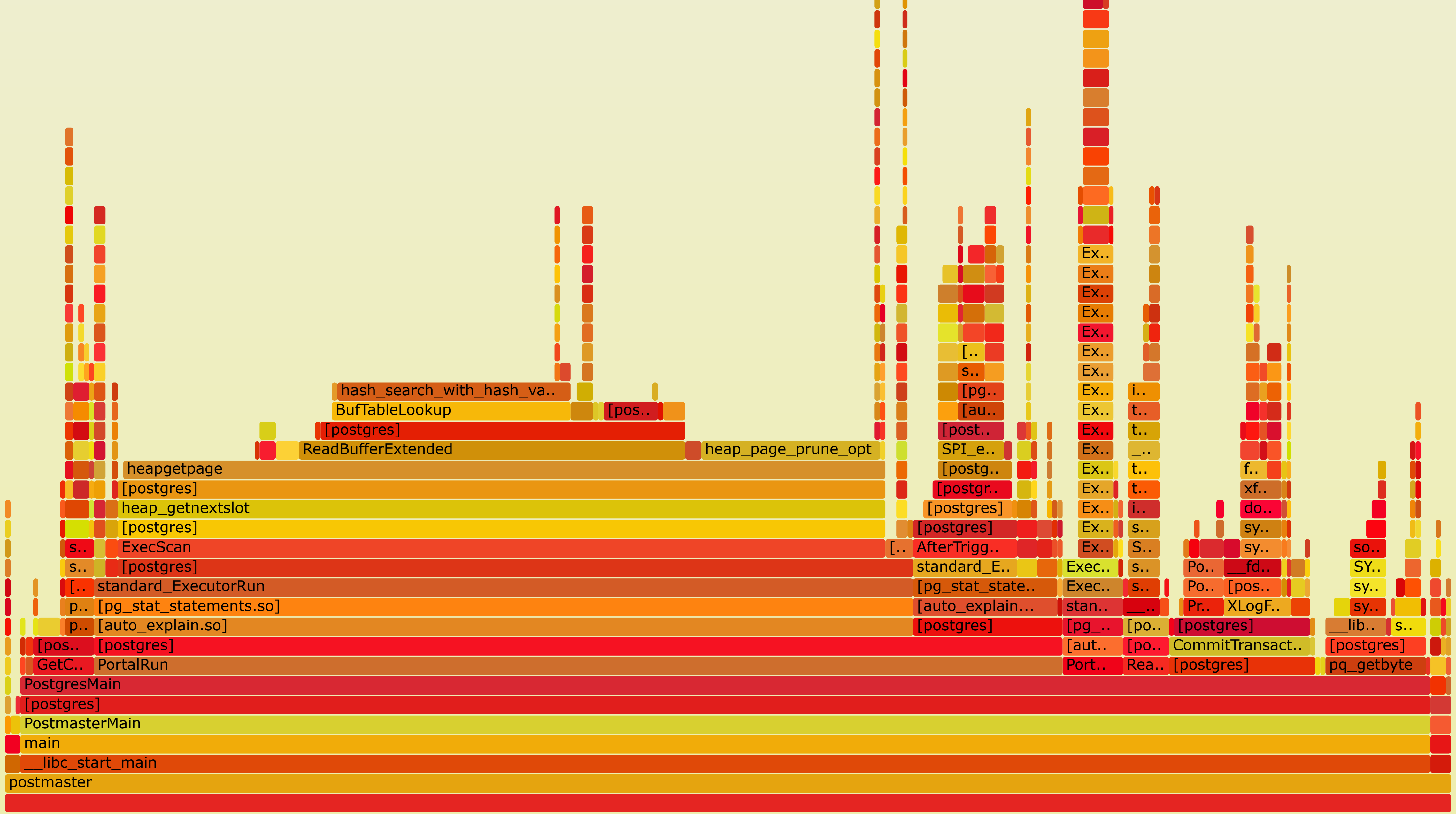
d580)

# Perf



# Perf flame-graph





# SS & Perf output

TIME	TCP_RECVQ	TCP_RBUFPC	TCP_RTT	TCP_SENDRATE	TOPSYSCALL
21:00:32	899	0.0%	0.695/0.058	50.0Mbps	13.74%--sys_fdatasync
21:02:04	0	0.0%	0.877/0.19	39.6Mbps	13.43%--sys_fdatasync
21:03:37	20896	0.7%	0.823/0.214	42.2Mbps	12.80%--sys_fdatasync
21:05:09	0	0.0%	0.765/0.143	45.4Mbps	14.50%--sys_fdatasync
21:06:41	0	0.0%	0.777/0.083	44.7Mbps	13.53%--sys_fdatasync
21:08:13	0	0.0%	0.804/0.087	43.2Mbps	14.01%--sys_fdatasync
21:09:45	139	0.0%	0.808/0.133	43.0Mbps	13.38%--sys_fdatasync
21:11:19	0	0.0%	0.79/0.11	44.0Mbps	11.99%--sys_fdatasync
21:12:51	0	0.0%	0.702/0.109	49.5Mbps	13.87%--sys_fdatasync
21:14:23	1211	0.0%	0.732/0.105	47.5Mbps	12.21%--sys_fdatasync
21:15:53	125633	4.0%	0.86/0.121	40.4Mbps	36.60%--sys_pwrite64
21:17:24	1379822	44.6%	0.835/0.175	41.6Mbps	34.97%--sys_pwrite64
21:18:55	683034	21.9%	0.791/0.122	43.9Mbps	35.22%--sys_pwrite64
21:20:25	1790854	57.3%	0.673/0.082	51.6Mbps	34.73%--sys_pwrite64
21:21:56	3288552	99.5%	1.361/0.406	25.5Mbps	34.99%--sys_pwrite64
21:23:26	716453	23.0%	0.679/0.16	51.2Mbps	36.18%--sys_pwrite64
21:24:57	500137	16.1%	0.777/0.134	44.7Mbps	34.08%--sys_pwrite64
21:26:28	365872	11.4%	0.81/0.21	42.9Mbps	25.82%--sys_pwrite64
21:27:59	4131	0.1%	0.628/0.118	55.3Mbps	26.19%--sys_pwrite64
21:29:30	838693	26.7%	0.674/0.062	51.6Mbps	27.37%--sys_pwrite64
21:31:00	115372	3.6%	0.846/0.173	41.1Mbps	35.10%--sys_pwrite64
21:32:31	139553	4.4%	0.719/0.164	48.3Mbps	35.50%--sys_pwrite64
21:34:03	0	0.0%	0.81/0.182	42.9Mbps	13.52%--sys_fdatasync
21:35:35	6803	0.2%	0.637/0.086	54.6Mbps	14.31%--sys_fdatasync

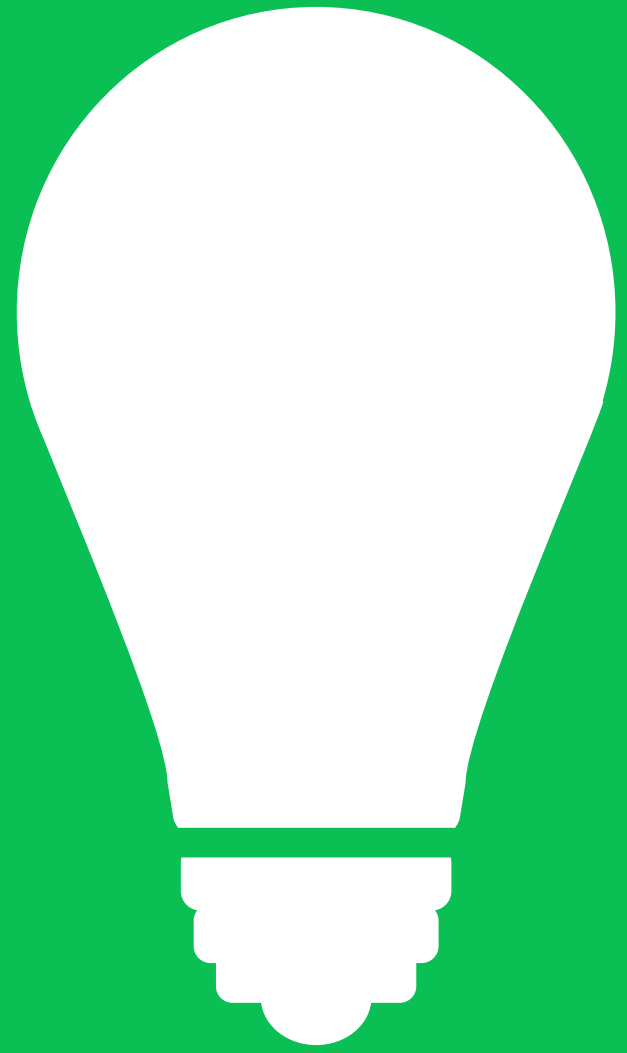


TIME	TCP_RECVQ	TCP_RBUFPCT	TCP_RTT	TCP_SENDRATE	TOPSYSCALL
21:05:09	0	0.0%	0.765/0.143	45.4Mbps	14.50%--sys_fdatasync
21:06:41	0	0.0%	0.777/0.083	44.7Mbps	13.53%--sys_fdatasync
21:08:13	0	0.0%	0.804/0.087	43.2Mbps	14.01%--sys_fdatasync
21:09:45	139	0.0%	0.808/0.133	43.0Mbps	13.38%--sys_fdatasync
21:11:19	0	0.0%	0.79/0.11	44.0Mbps	11.99%--sys_fdatasync
21:12:51	0	0.0%	0.702/0.109	49.5Mbps	13.87%--sys_fdatasync
21:14:23	1211	0.0%	0.732/0.105	47.5Mbps	12.21%--sys_fdatasync
21:15:53	125633	4.0%	0.86/0.121	40.4Mbps	36.60%--sys_pwrite64
21:17:24	1379822	44.6%	0.835/0.175	41.6Mbps	34.97%--sys_pwrite64
21:18:55	683034	21.9%	0.791/0.122	43.9Mbps	35.22%--sys_pwrite64
21:20:25	1790854	57.3%	0.673/0.082	51.6Mbps	34.73%--sys_pwrite64
21:21:56	3288552	99.5%	1.361/0.406	25.5Mbps	34.99%--sys_pwrite64
21:23:26	716453	23.0%	0.679/0.16	51.2Mbps	36.18%--sys_pwrite64
21:24:57	500137	16.1%	0.777/0.134	44.7Mbps	34.08%--sys_pwrite64
21:26:28	365872	11.4%	0.81/0.21	42.9Mbps	25.82%--sys_pwrite64
21:27:59	4131	0.1%	0.628/0.118	55.3Mbps	26.19%--sys_pwrite64
21:29:30	838693	26.7%	0.674/0.062	51.6Mbps	27.37%--sys_pwrite64
21:31:00	115372	3.6%	0.846/0.173	41.1Mbps	35.10%--sys_pwrite64
21:32:31	139553	4.4%	0.719/0.164	48.3Mbps	35.50%--sys_pwrite64
21:34:03	0	0.0%	0.81/0.182	42.9Mbps	13.52%--sys_fdatasync
21:35:35	6803	0.2%	0.637/0.086	54.6Mbps	14.31%--sys_fdatasync



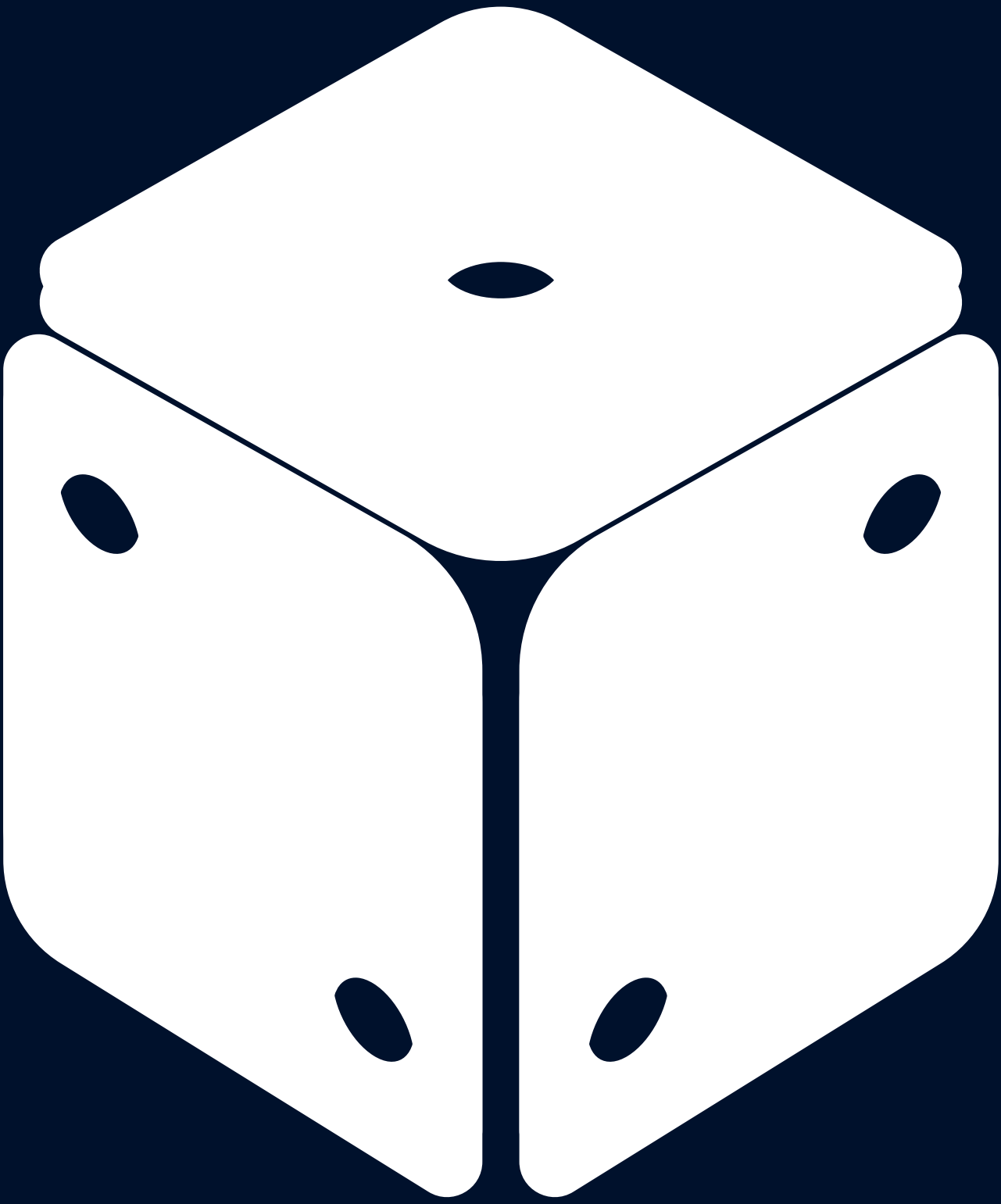
“On **high load** caused by vacuum, restart point and activity, the standby **filesystem can't** keep up but the **network** and the **disk can**”





**Solution**





# Solution

`wal_recycle = off`

# Solution

```
wal_init_zero = off
```

# Solution

**Pin** process wal writer **single**  
**cpu**

```
taskset -pc 5 `pgrep -f walreceiver`
```



# Lessons learned



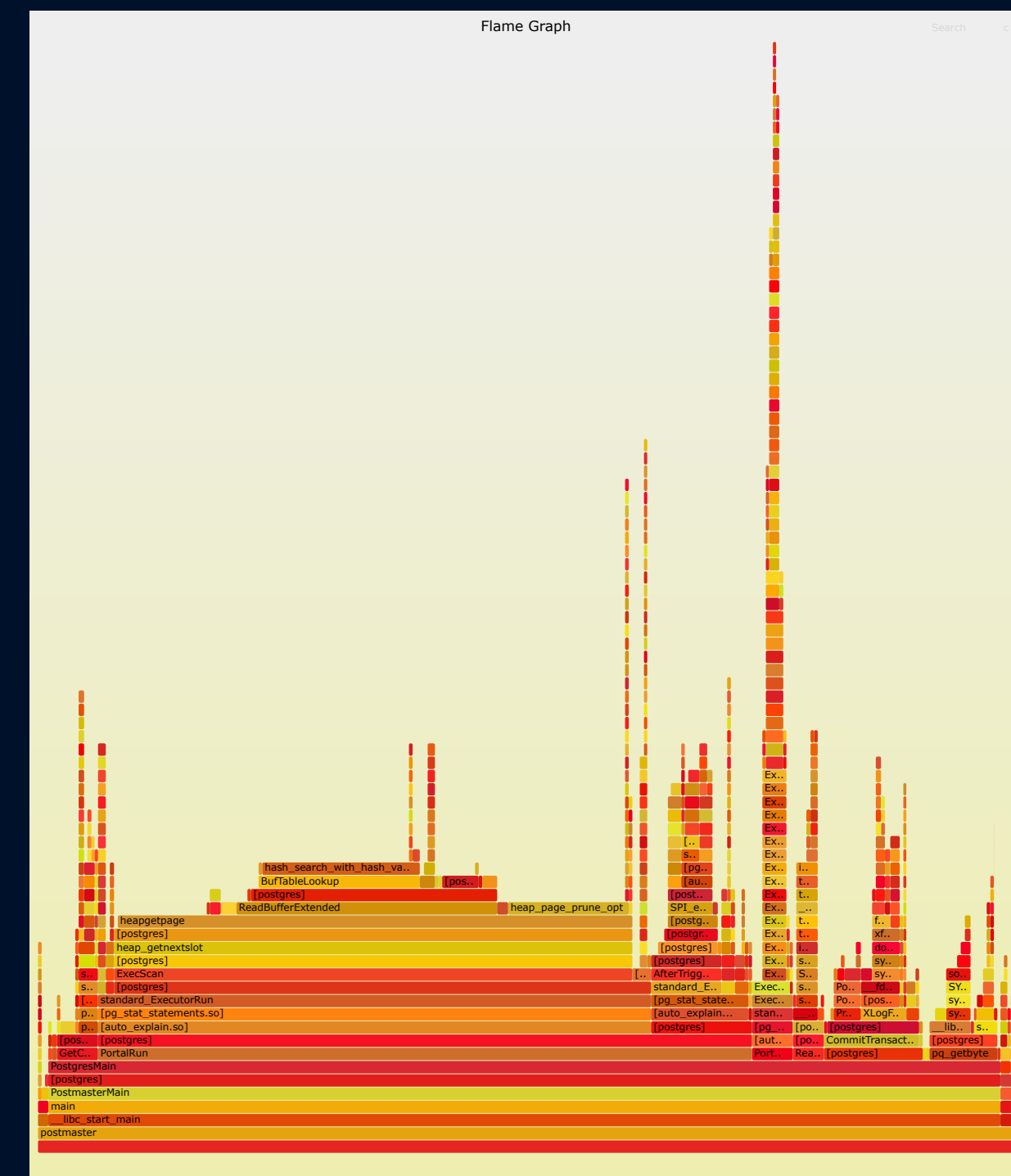
# Lessons learned



# Lessons learned

## SS & Perf

0	0.0%	0.777/0.085
0	0.0%	0.804/0.087
139	0.0%	0.808/0.133
0	0.0%	0.79/0.11
0	0.0%	0.702/0.109
1211	0.0%	0.732/0.105
125633	4.0%	0.86/0.121
1379822	44.6%	0.835/0.175
683034	21.9%	0.791/0.122
1790854	57.3%	0.673/0.082
3288552	99.5%	1.361/0.406
716453	23.0%	0.679/0.16
500137	16.1%	0.777/0.134
365872	11.4%	0.81/0.21
4131	0.1%	0.628/0.118
838693	26.7%	0.674/0.062
115372	3.6%	0.846/0.173
139553	4.4%	0.719/0.164
0	0.0%	0.81/0.182
6803	0.2%	0.637/0.086





# Lessons learned





**Special  
thanks to**



# Adyen

**Dmitry Fomin**



**Milen Blagojevic**





# EDB

**Jakub Wartak**



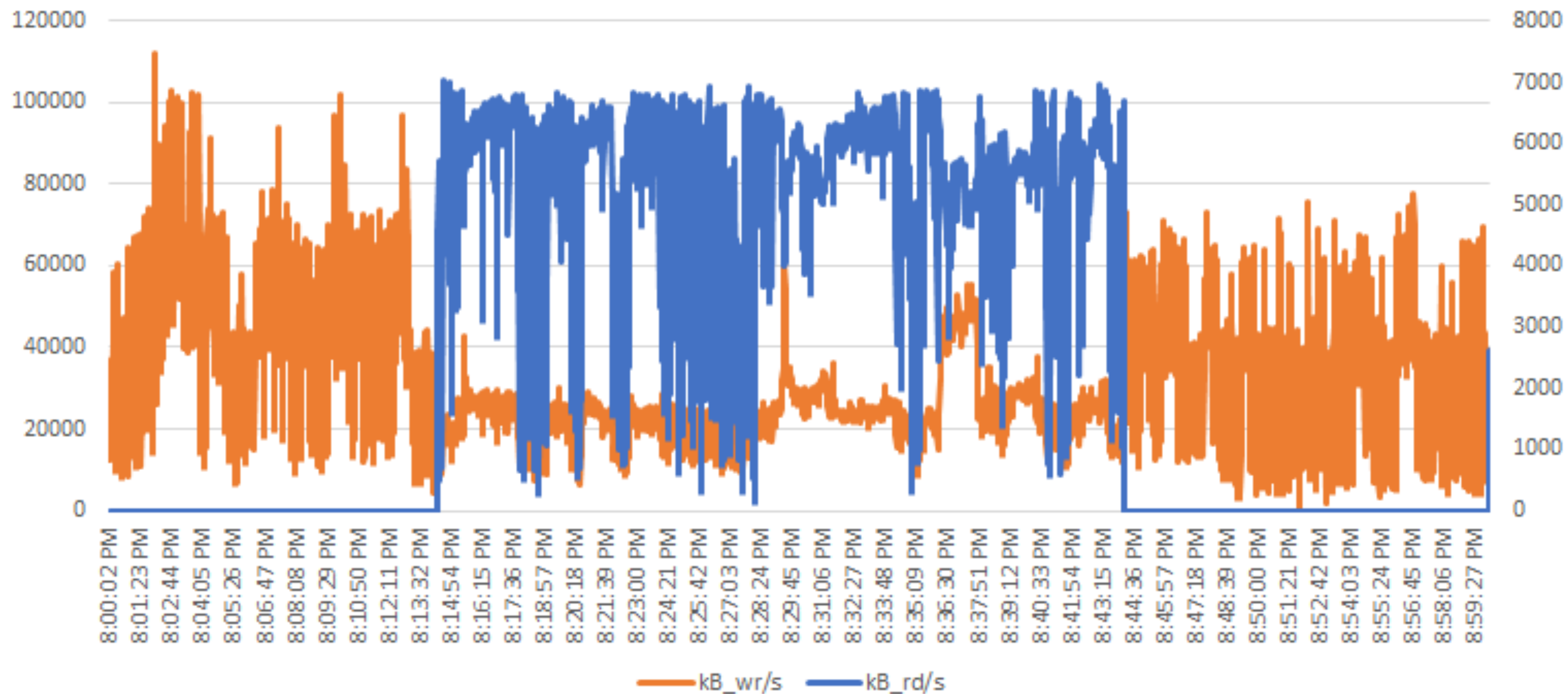
**Tomas Vondra**



**Álvaro Herrera  
Muñoz**



walreceiver pidstat -d (logical IO)



**Questions**