

Yandex

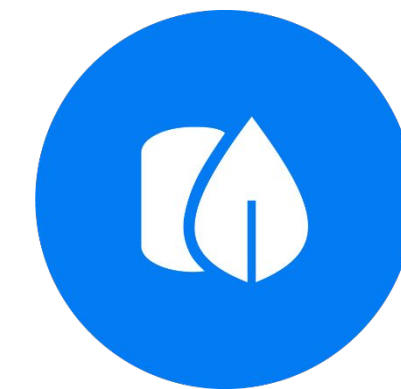


Unified Point-in-Time Recovery in the Cloud

Andrey Borodin, Team lead of opensource RDBMS development, Yandex.Cloud

Dmitry Smal, Team lead of Managed MySQL and SQL Server Development, Yandex.Cloud

- › Managed Service for PostgreSQL
- › Managed Service for MySQL
- › Managed Service for MongoDB
- › Many more DBs





What's interesting about Yandex.Cloud managed databases?

Yandex

PostgreSQL at Yandex.Cloud

Yandex.Mail

- › Hundreds of millions of users
- › 10^{12} rows, 10^6 queries per second, ~1 PB of data

Many others services use managed PostgreSQL

- › Taxi, carsharing, food delivery, self driving cars, etc.
- › Total of 3 million queries per second, 6 PB of space used

MySQL at Yandex.Cloud

Yandex.Direct

- › Online advertising network
- › Uses Percona build of MySQL at scale

Managed MySQL

- › 400+ TB as of 2021

Less is more

- Cloud providers usually sell computing resources.
- We aim to utilize fewer resources for the same workload.

Point-in-time recovery



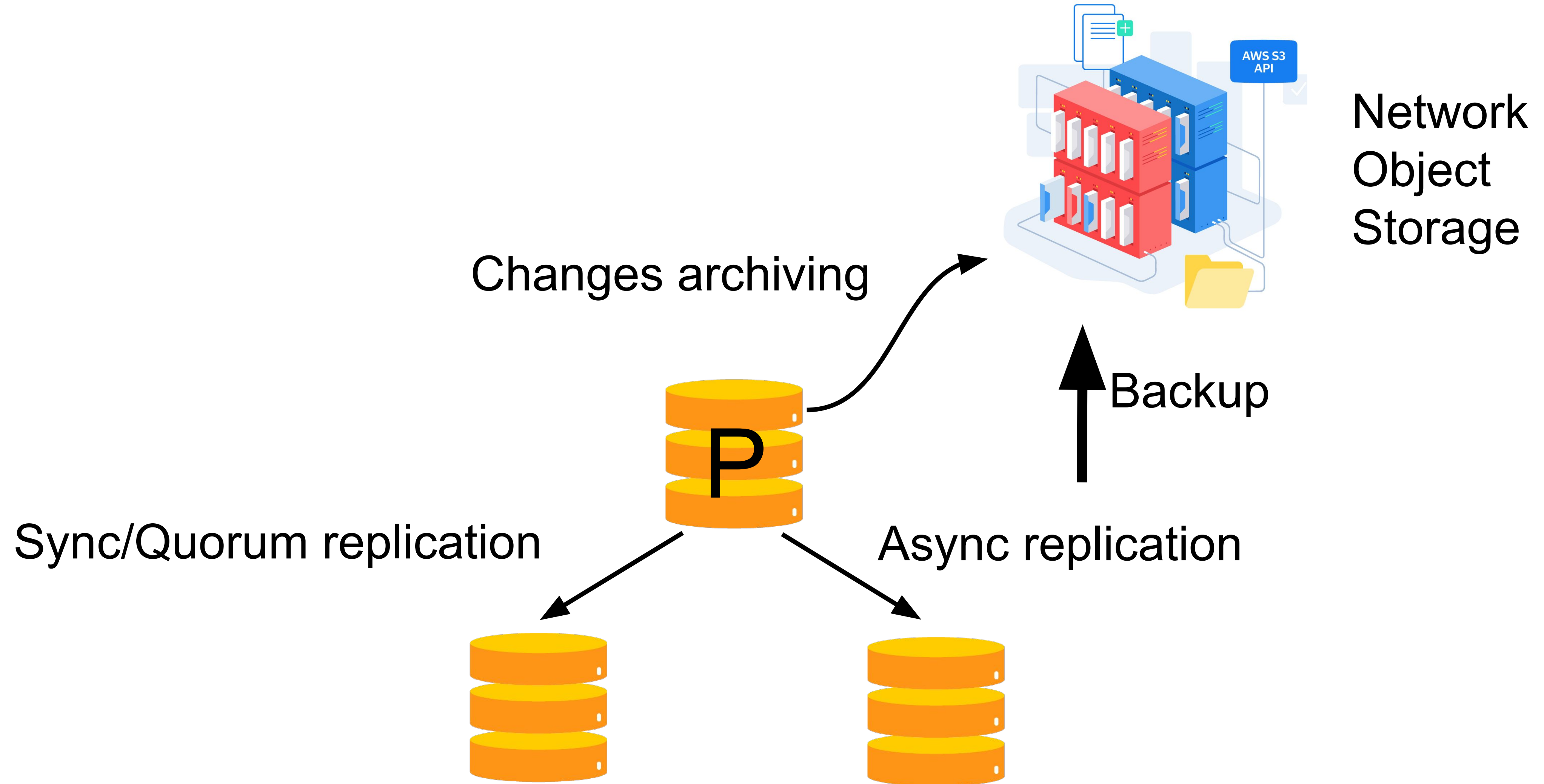
Backup + changes

- › Scalable
- › Reliable
- › Efficient
- › Fast

Scalability

- › Data: from 10 GB to 10 TB on a host
- › RAM: from 2 GB
- › Number of CPUs: from 0.05 to ~100
- › Async and parallel whenever possible
- › Don't spill anything on a local disk

HA cluster in the cloud



Resources

› Storage space

Resources

- › ~~Storage space~~
- › CPU
- › Net bandwidth
- › Disk IOPS

Reliability

- › Protection from human error via automation and safety checks
- › Prevention of data corruption
- › Consistency monitoring
- › Integration with other systems (HA tool)
- › Extensibility and unification of approaches
- › Encrypted data in storage

Fast recovery

- › OLAP

- From start to consistency point

- › OLTP standby

- To starting streaming replication

- › OLTP primary

- Until recovery target and accept of write queries

Unacceptable

- › Data locks

- Business can't wait

- › Data loss

- We call it a “database” after all



Barman
Backup and recovery
manager for PostgreSQL

pgProBackup



pgBackRest

WAL-G



Release v0.2.19 · wal-g/wal-g

← → ↺ 🏠 🔒 https://github.com/wal-g/wal-g/releases/tag/v0.2.19 ⋮ ⏴ ⭐ ⬇ 📄 👤 🔄 ☰


Latest release

v0.2.19
8dba4f4

Verified

Compare ▾

v0.2.19


 x4m released this on Nov 30, 2020


Notable changes in this release include:


1. Fixes for S3 and GCP storages [#656](#) [#756](#).
2. Add wal-show command to get information about wal storage folder.
3. Add wal-verify command. It checks the integrity of WAL history starting from the oldest backup available in current timeline history.
4. Add wal-receive command. You can use WAL-G as a replica running on another host to ensure RPO=0. This is beta functionality: API may change in the future.
5. Add reverse delta unpack for backup-fetch (`--reverse-unpack` flag).
6. Add redundant archives skipping for backup-fetch (`--skip-redundant-tars` flag, designed to work in pair with reverse delta unpack).
7. Add page checksum verification for backup-push (`--verify` flag).

You can find more about some of these new features in [Daniil Zakhlystov's post](#)

▼ Assets 4

 [wal-g.linux-amd64.tar.gz](#) 13.1 MB

 [wal-g.linux-amd64.tar.gz.sha256](#) 91 Bytes

 [Source code \(zip\)](#)


```
wal-g — -bash — 85x33
~/GoglandProjects/src/github.com/wal-g/wal-g/cmd/wal-g — -bash  ~/project/bin — psql postgres +
x4mmm-osx:wal-g x4mmm$
x4mmm-osx:wal-g x4mmm$
x4mmm-osx:wal-g x4mmm$
x4mmm-osx:wal-g x4mmm$ AWS_ENDPOINT=https://storage.yandexcloud.net AWS_ACCESS_KEY_ID
=wIRAxwOPLI3VrGwtYWL AWS_SECRET_ACCESS_KEY=ne[REDACTED]vsXX
WALE_S3_PREFIX=s3://wal-g-test/ ./wal-g backup-list
Path:
name                last_modified        wal_segment_backup_start
base_00000001000000000000000004 2019-02-02T18:39:30Z 00000001000000000000000004
x4mmm-osx:wal-g x4mmm$
x4mmm-osx:wal-g x4mmm$
```



```
x4mmm-osx:wal-g x4mmm$  
x4mmm-osx:wal-g x4mmm$ AWS_ENDPOINT=https://storage.yandexcloud.net AWS_ACCESS_KEY_ID=  
=wIRAxw0PLI3VrGwtYWL AWS_SECRET_ACCESS_KEY=neh7EEYANpqGS5GJEbm0ywhznxcIBukG3IamvsXX  
WALE_S3_PREFIX=s3://wal-g-test/ ./wal-g backup-push ~/DemoDb  
Path:  
INFO: 2019/02/02 21:56:42.509465 Doing full backup.  
WARNING: 2019/02/02 21:56:42.526434 It seems your archive_mode is not enabled. This w  
ill cause inconsistent backup. Please consider configuring WAL archiving.  
INFO: 2019/02/02 21:56:42.740377 Walking ...  
INFO: 2019/02/02 21:56:42.742571 Starting part 1 ...  
INFO: 2019/02/02 21:56:43.112485 Finished writing part 1.  
INFO: 2019/02/02 21:56:48.744337 Starting part 2 ...  
INFO: 2019/02/02 21:56:48.761395 /global/pg_control  
INFO: 2019/02/02 21:56:48.764006 Finished writing part 2.  
INFO: 2019/02/02 21:56:48.878931 Starting part 3 ...  
INFO: 2019/02/02 21:56:48.894990 backup_label  
INFO: 2019/02/02 21:56:48.895030 tablespace_map  
INFO: 2019/02/02 21:56:48.895056 Finished writing part 3.  
INFO: 2019/02/02 21:56:49.523658 Uploaded 3 compressed tar Files.  
x4mmm-osx:wal-g x4mmm$
```



```
# - Archiving -
```

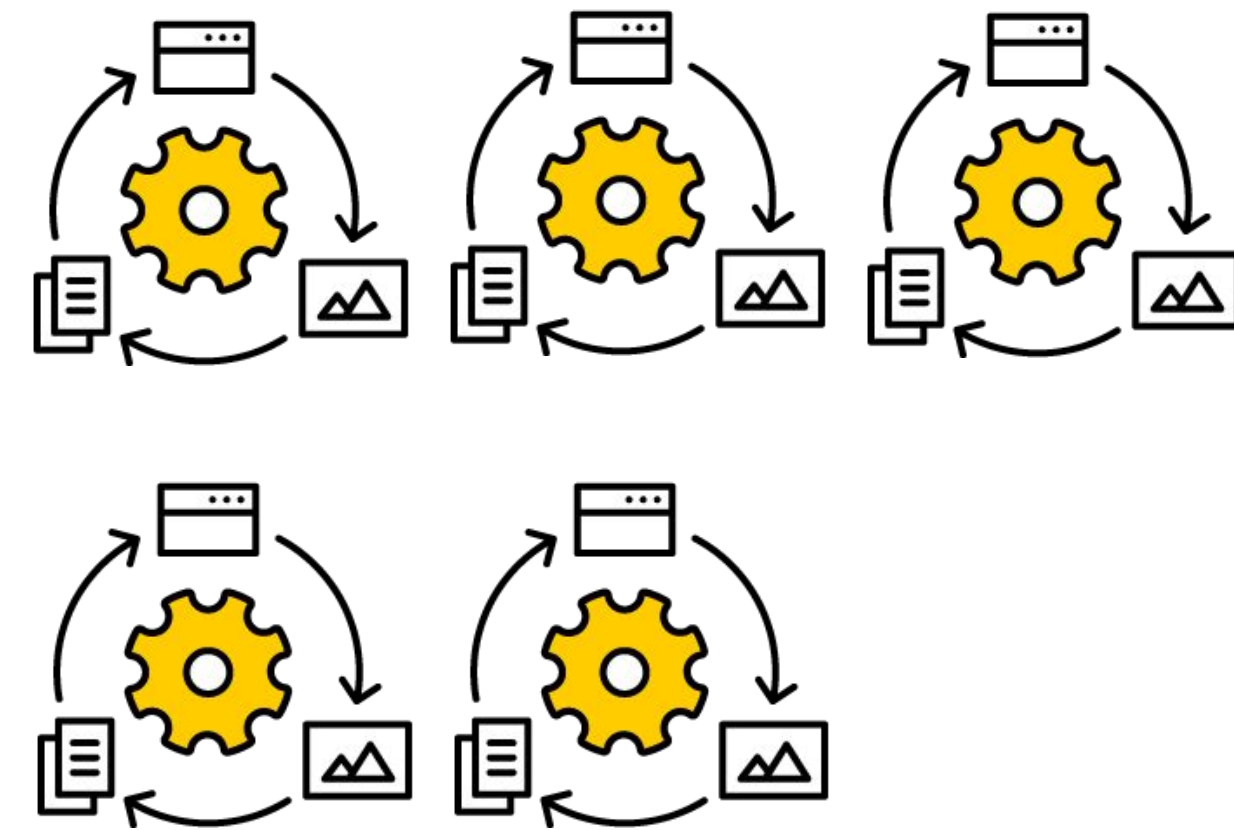
```
archive_mode = on
```

```
archive_command = '/usr/bin/envdir /etc/wal-g/envdir  
/usr/bin/timeout 600 /usr/bin/wal-g wal-push %p'
```

Normal backup

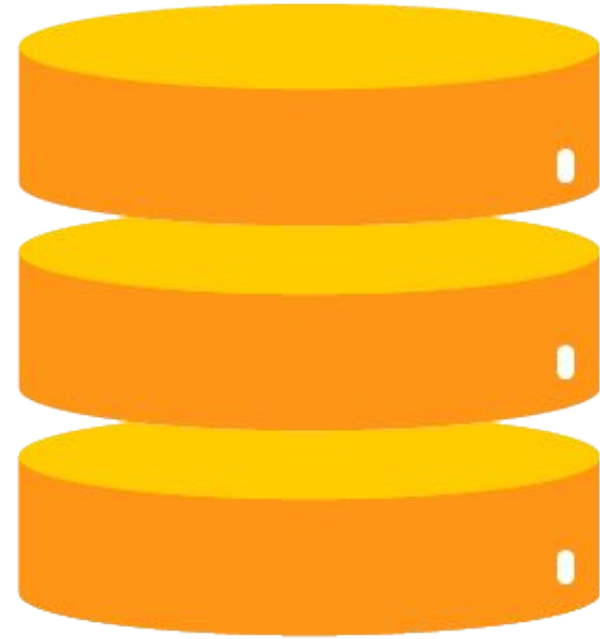


DB copy

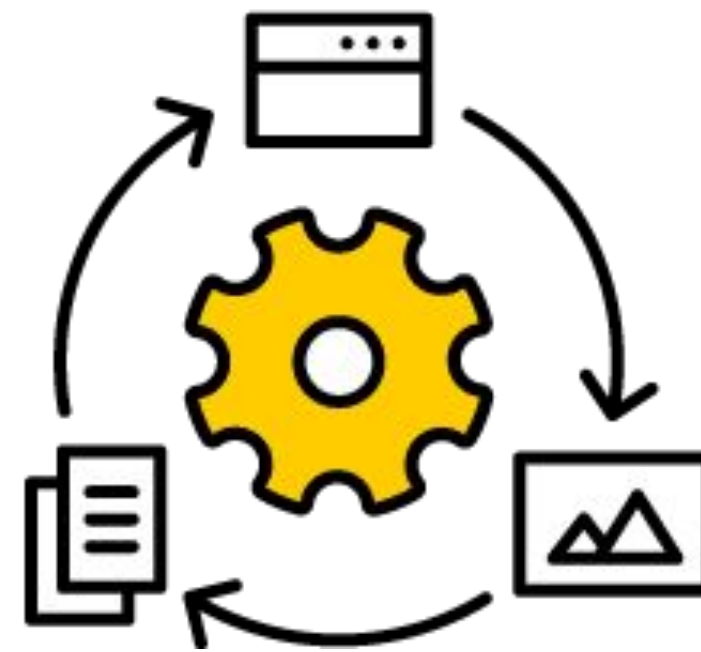


Changes (WAL)

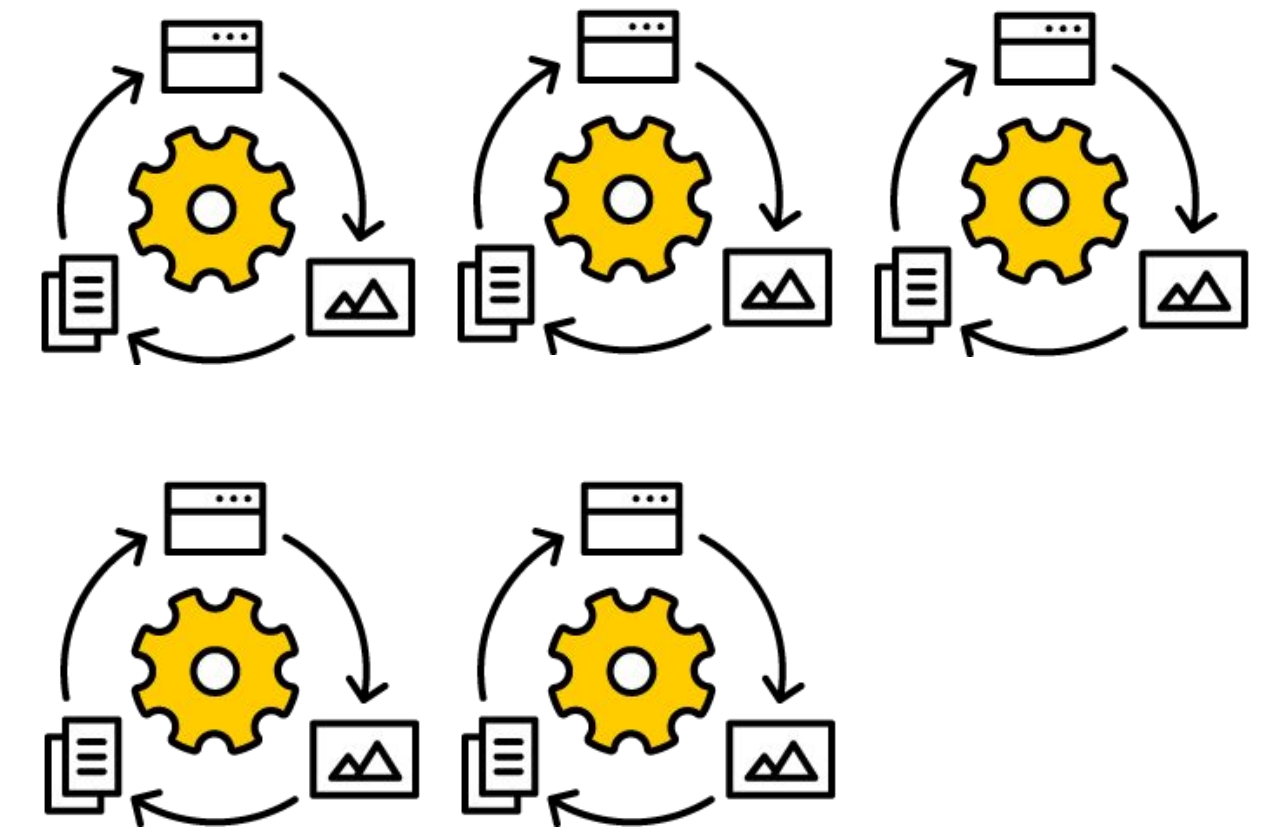
Delta backups



DB copy

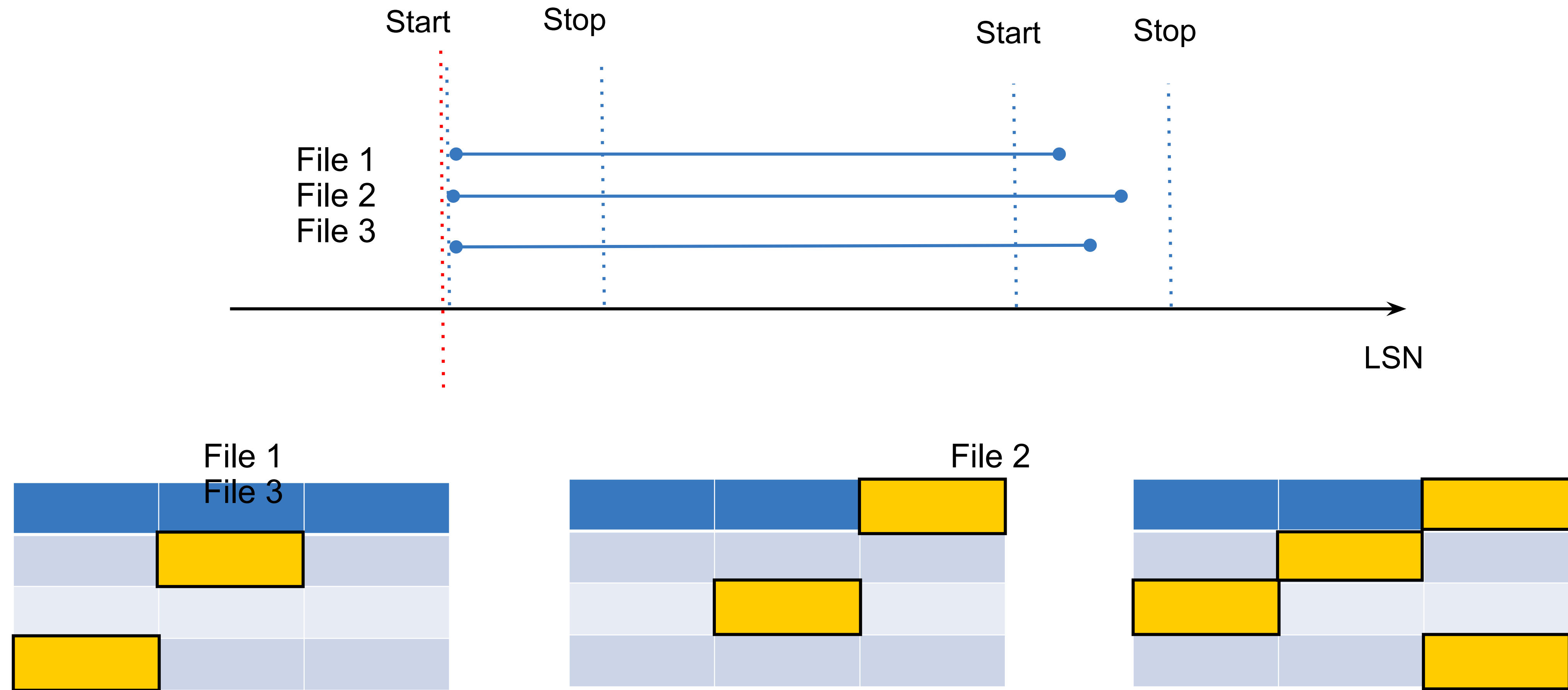


Delta copy

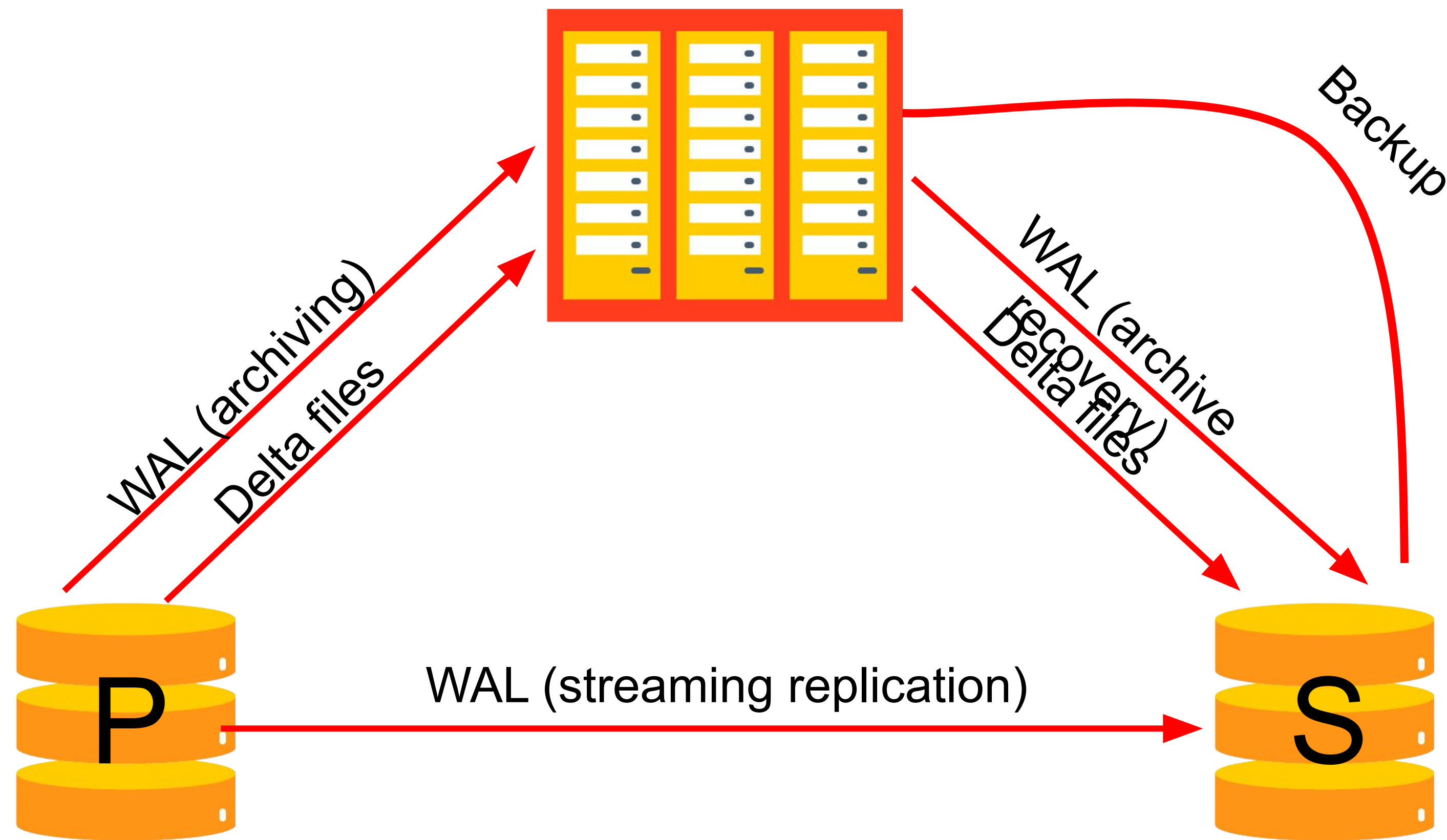


Changes (WAL)

LSN-based deltas



Data flows in the system



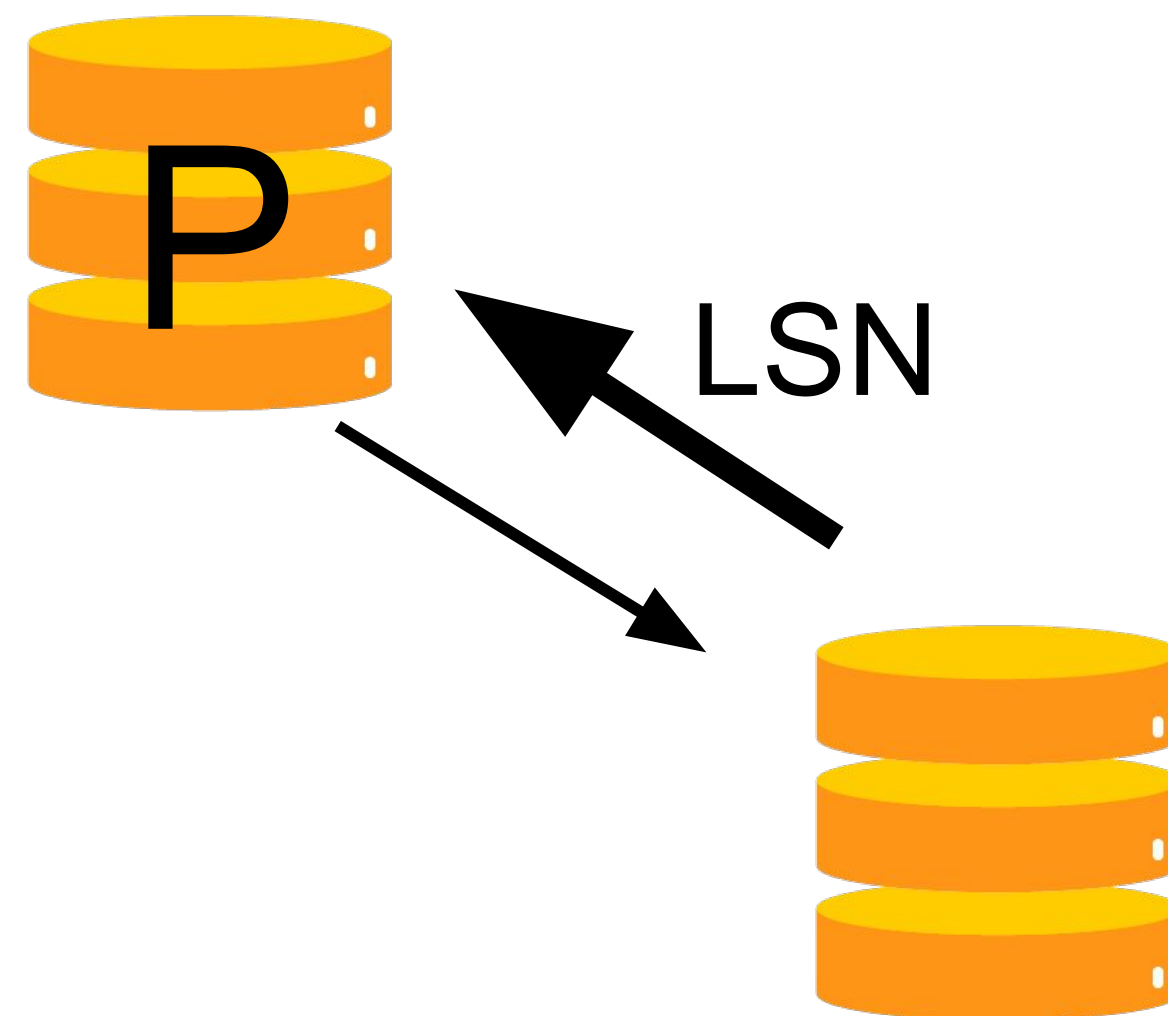
PG features



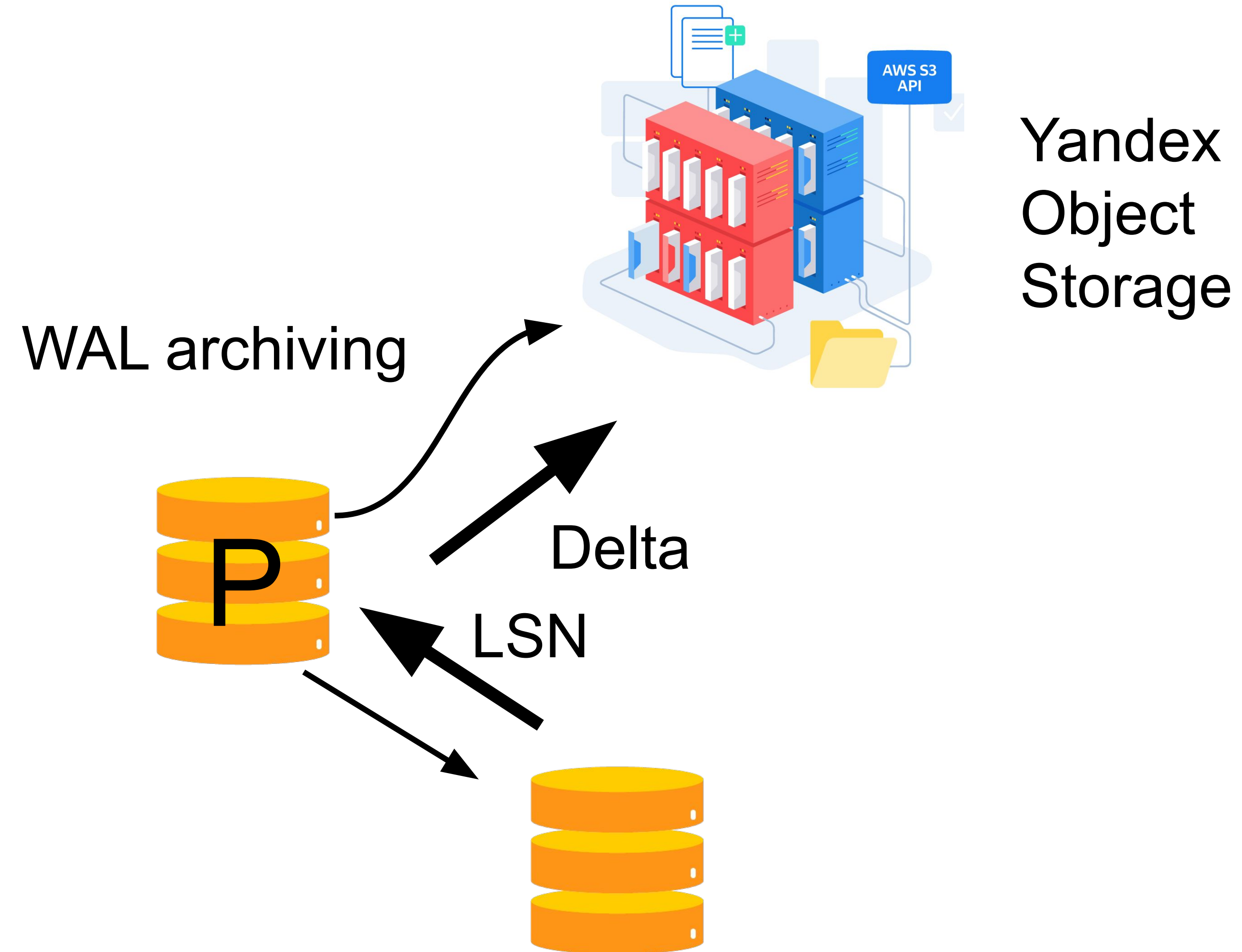
Catchup



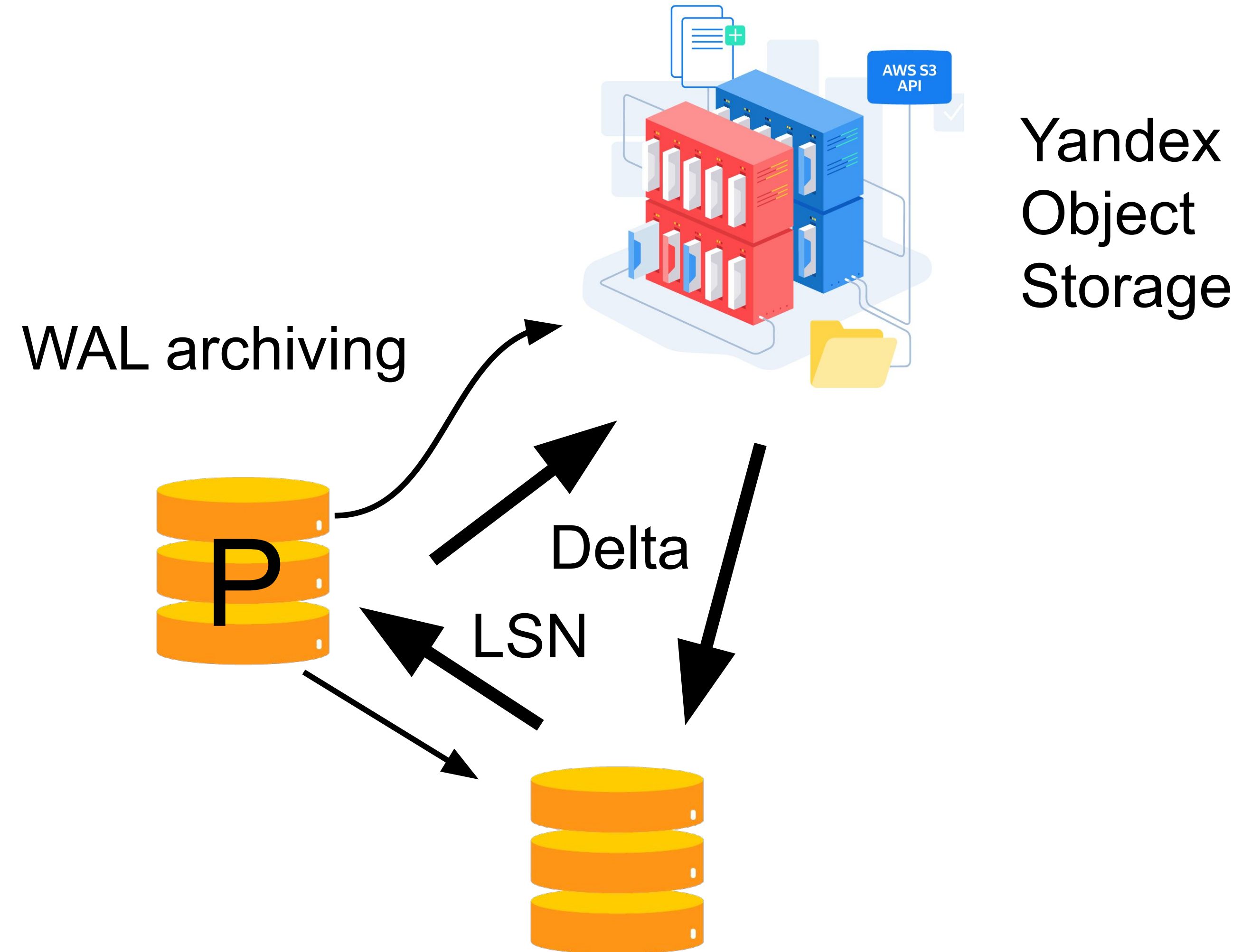
Yandex
Object
Storage



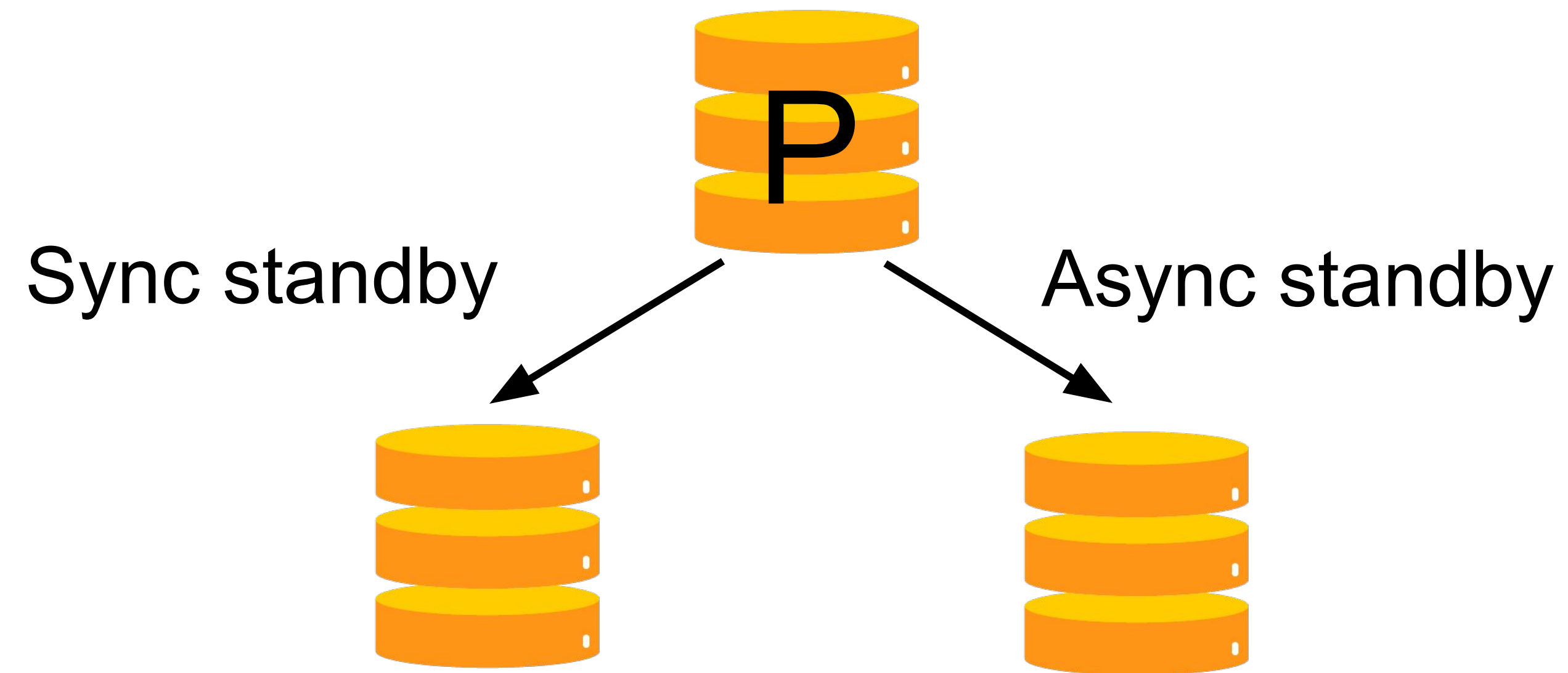
Catchup



Catchup



Backup-push quorum



Compress some bytes

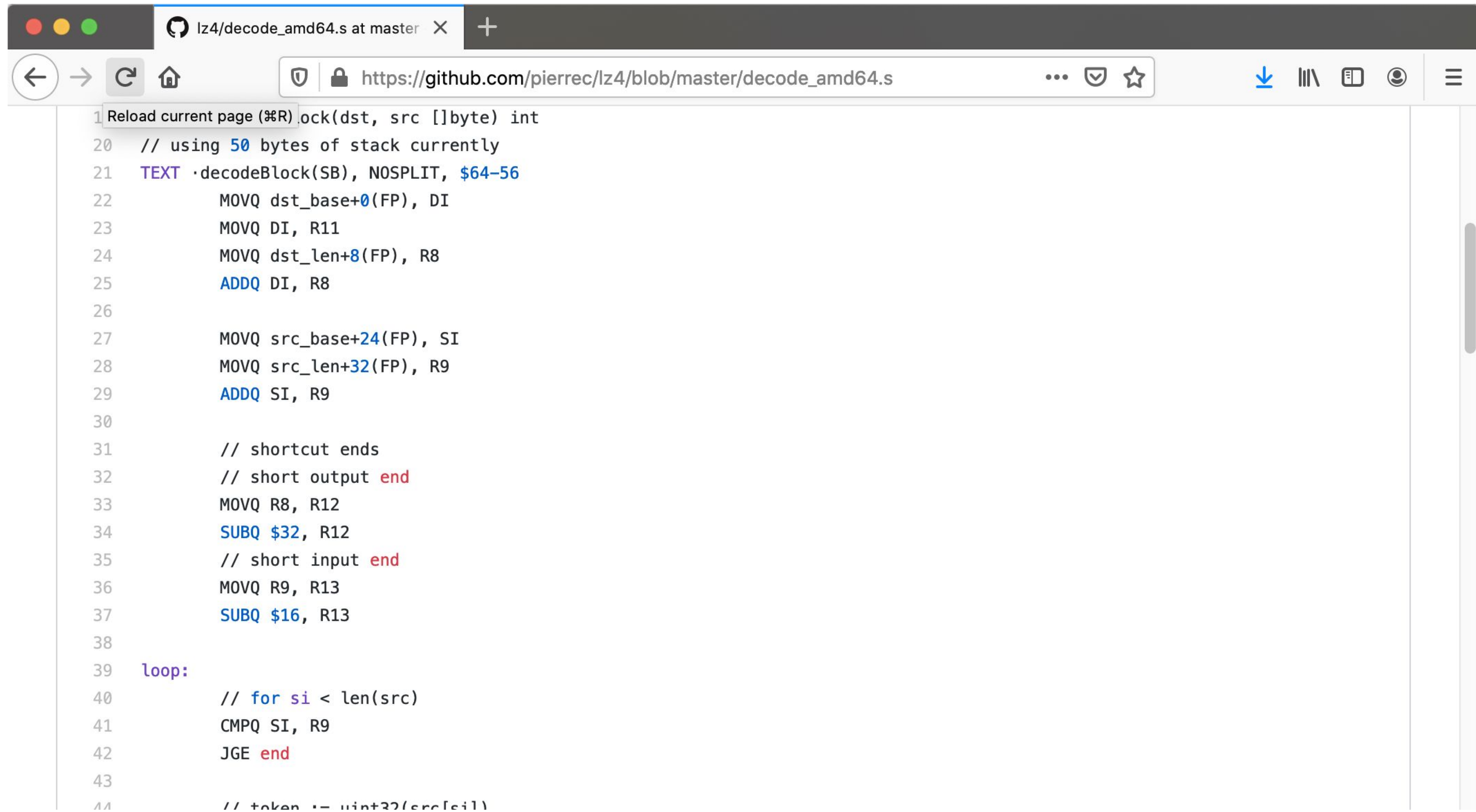


lz4

The screenshot shows the GitHub repository page for `pierrec/lz4`. The repository is titled "LZ4 compression and decompression in pure Go" and includes tags for `go`, `golang`, `lz4`, and `lz4-frame`. It has 137 commits, 4 branches, 0 packages, 34 releases, 17 contributors, and is licensed under BSD-3-Clause. The page features a navigation bar with links to Code, Issues (1), Pull requests (0), Actions, Projects (0), Wiki, Security, and Insights. A merge pull request #68 from `ivanov/patch-1` is highlighted, with the latest commit `cac5ed4` on Feb 20. Below the merge, a list of recent commits is shown:

<code>cmd/lz4c</code>	Writer: added concurrency support. Fixes #55	4 months ago
<code>fuzz</code>	Updated fuzz corpus 2.	10 months ago
<code>internal/xxh32</code>	<code>cmd/lz4c</code> : moved commands to local dir	10 months ago
<code>testdata</code>	<code>DecodeBlock</code> : handle case in <code>shortcut2</code> where the destination buffer is...	7 months ago
<code>.gitignore</code>	updated gitignore	11 months ago

lz4



```
1 Reload current page (%R) lock(dst, src []byte) int
20 // using 50 bytes of stack currently
21 TEXT ·decodeBlock(SB), NOSPLIT, $64-56
22     MOVQ dst_base+0(FP), DI
23     MOVQ DI, R11
24     MOVQ dst_len+8(FP), R8
25     ADDQ DI, R8
26
27     MOVQ src_base+24(FP), SI
28     MOVQ src_len+32(FP), R9
29     ADDQ SI, R9
30
31     // shortcut ends
32     // short output end
33     MOVQ R8, R12
34     SUBQ $32, R12
35     // short input end
36     MOVQ R9, R13
37     SUBQ $16, R13
38
39 loop:
40     // for si < len(src)
41     CMPQ SI, R9
42     JGE end
43
44     // token := uint32(src[si])
```

lz4

The screenshot shows a web browser displaying a GitHub commit page for the repository `pierrec/lz4`. The commit message is "Fix short buffers in some specific cases". The commit was made by user `x4m` on March 26, 2019. The commit hash is `908ada6c279113e95a9a565b9f5fd8c22fb61263`. The commit shows 4 changed files with 4 additions and 0 deletions. The file `reader.go` is expanded, showing a diff where lines 161-163 were added. The diff shows a function `Read` that resets the uncompressed buffer and reads a block length.

Fix short buffers in some specific cases

master (#40) v3.2.1 v2.1.1

x4m committed on Mar 26, 2019 1 parent 062282e commit 908ada6c279113e95a9a565b9f5fd8c22fb61263

Showing 4 changed files with 4 additions and 0 deletions. Unified Split

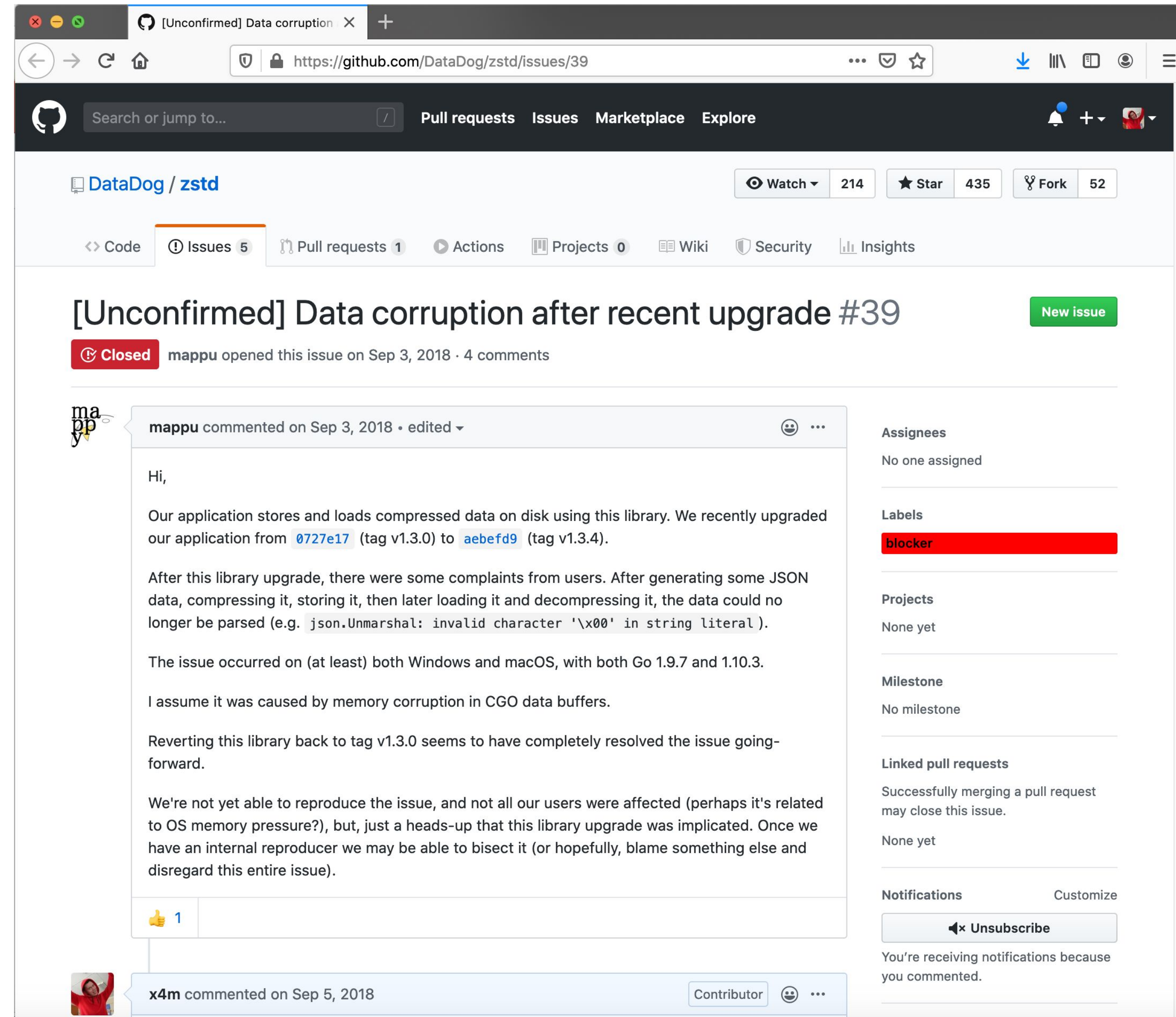
3 reader.go

```
@@ -158,6 +158,9 @@ func (z *Reader) Read(buf []byte) (int, error) {
158 158         if debugFlag {
159 159             debug("reading block from writer")
160 160         }
161 + // Reset uncompressed buffer
162 + z.data = z.zdata[:cap(z.zdata)][len(z.zdata):]
163 +
161 164 // Block length: 0 = end of frame, highest bit set: uncompressed.
162 165 bLen, err := z.readUint32()
163 166 if err != nil {
```

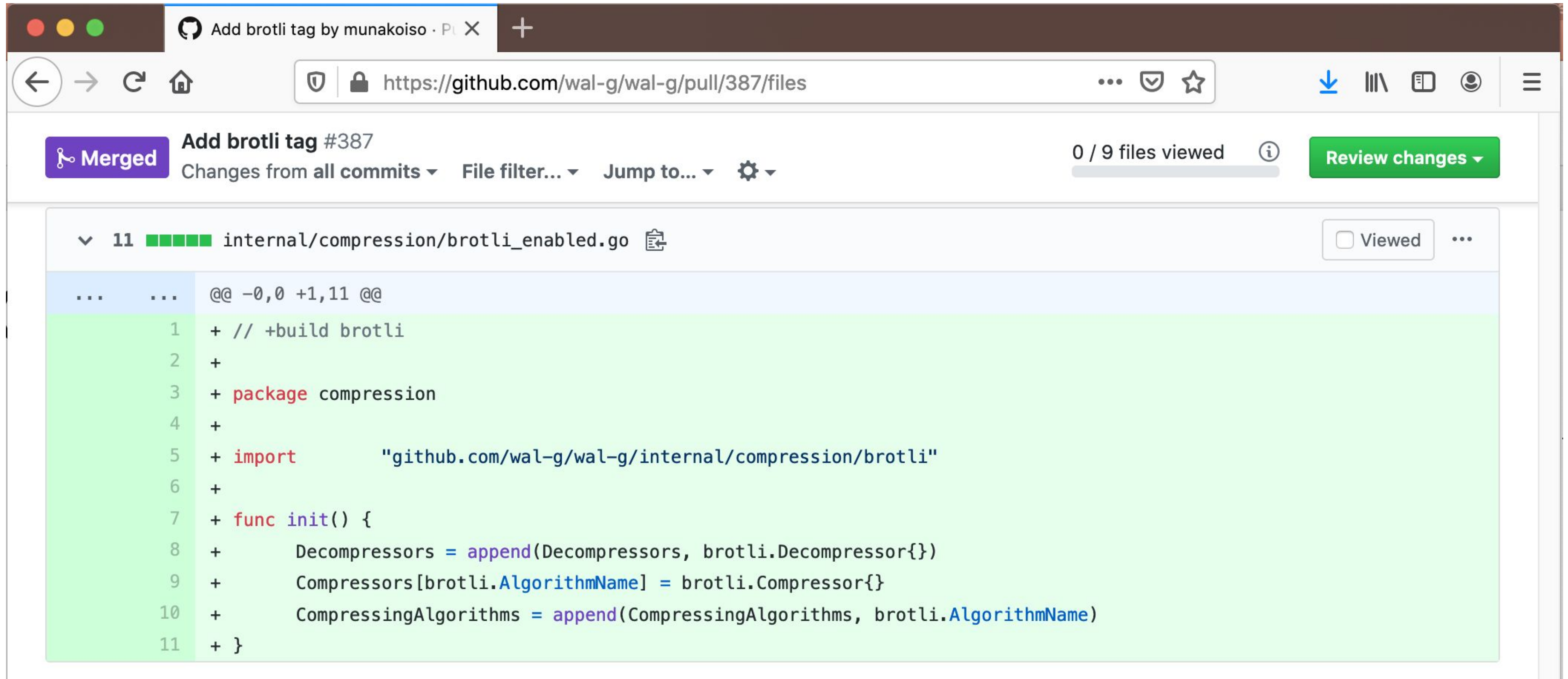

Izma

```
 1  [||||| 91.0%]  9  [||||| 91.1%] 17  [||||| 69.0%] 25  [||||| 68.2%]
 2  [||||| 89.5%] 10  [||||| 88.7%] 18  [||||| 69.5%] 26  [||||| 68.4%]
 3  [||||| 90.3%] 11  [||||| 86.4%] 19  [||||| 66.7%] 27  [||||| 68.4%]
 4  [||||| 88.0%] 12  [||||| 88.4%] 20  [||||| 64.7%] 28  [||||| 65.4%]
 5  [||||| 85.7%] 13  [||||| 86.1%] 21  [||||| 69.5%] 29  [||||| 65.4%]
 6  [||||| 85.1%] 14  [||||| 89.5%] 22  [||||| 69.3%] 30  [||||| 65.8%]
 7  [||||| 85.0%] 15  [||||| 84.9%] 23  [||||| 68.8%] 31  [||||| 64.5%]
 8  [||||| 82.4%] 16  [||||| 82.4%] 24  [||||| 72.1%] 32  [||||| 77.7%]
Mem [||||| 34.9G/252G] Tasks: 1371, 191 thr; 22 running
Swp [||||| 0K/0K]      Load average: 28.67 31.77 32.33
                        Uptime: 134 days(!), 19:58:41
```


ZStd



brothli



The screenshot shows a web browser window displaying a GitHub pull request. The browser's address bar shows the URL `https://github.com/wal-g/wal-g/pull/387/files`. The page title is "Add brotli tag by munakoiso · Pull Request". The pull request is titled "Add brotli tag #387" and is marked as "Merged". The interface shows "0 / 9 files viewed" and a "Review changes" button. The file being viewed is `internal/compression/brotli_enabled.go`. The diff shows the following changes:

```
...      ...      @@ -0,0 +1,11 @@
1      + // +build brotli
2      +
3      + package compression
4      +
5      + import      "github.com/wal-g/wal-g/internal/compression/brotli"
6      +
7      + func init() {
8      +     Decompressors = append(Decompressors, brotli.Decompressor{})
9      +     Compressors[brotli.AlgorithmName] = brotli.Compressor{}
10     +     CompressingAlgorithms = append(CompressingAlgorithms, brotli.AlgorithmName)
11     + }
```


OpenPGP

Go

CHANGES

DOCUMENTATION

BROWSE

Use correct default hashes and default ciphers when no preferences given

Updated Mar 23

Owner Gerrit Bot

Author Andrey Borodin

Assignee

Reviewers Filippo Valsorda

CC Adam Langley, Gobot Gobot

Repo / Branch crypto / master

Parent 891825f

Topic No topic

Strategy Cherry Pick

Hashtags

Code-Review No votes

Other labels

Run-TryBot No votes

TryBot-Result No votes

Use correct default hashes and default ciphers when no preferences given

Per comments in config.go

> // DefaultHash is the default hash function to be used.

> // If zero, SHA-256 is used.

> DefaultHash crypto.

> // DefaultCipher is the cipher to be used.

> // If zero, AES-128 is used.

> DefaultCipher CipherFunction

>

But instead we get RIPEMD160 and CAST5 if config is nil. Both are obsolete.

Fix <https://github.com/golang/go/issues/37646>

Change-Id: [I55f22580e1c3c00790cc0a9fad8c657fd68596b3](https://github.com/golang/go/pull/128)

GitHub-Last-Rev: 0791c42c2fa44d37386fd2f2514556b5b2ff217d

GitHub-Pull-Request: [golang/crypto#128](https://github.com/golang/go/pull/128)

Files

Checks

Findings

Base → Patchset 2

1ba5cf8

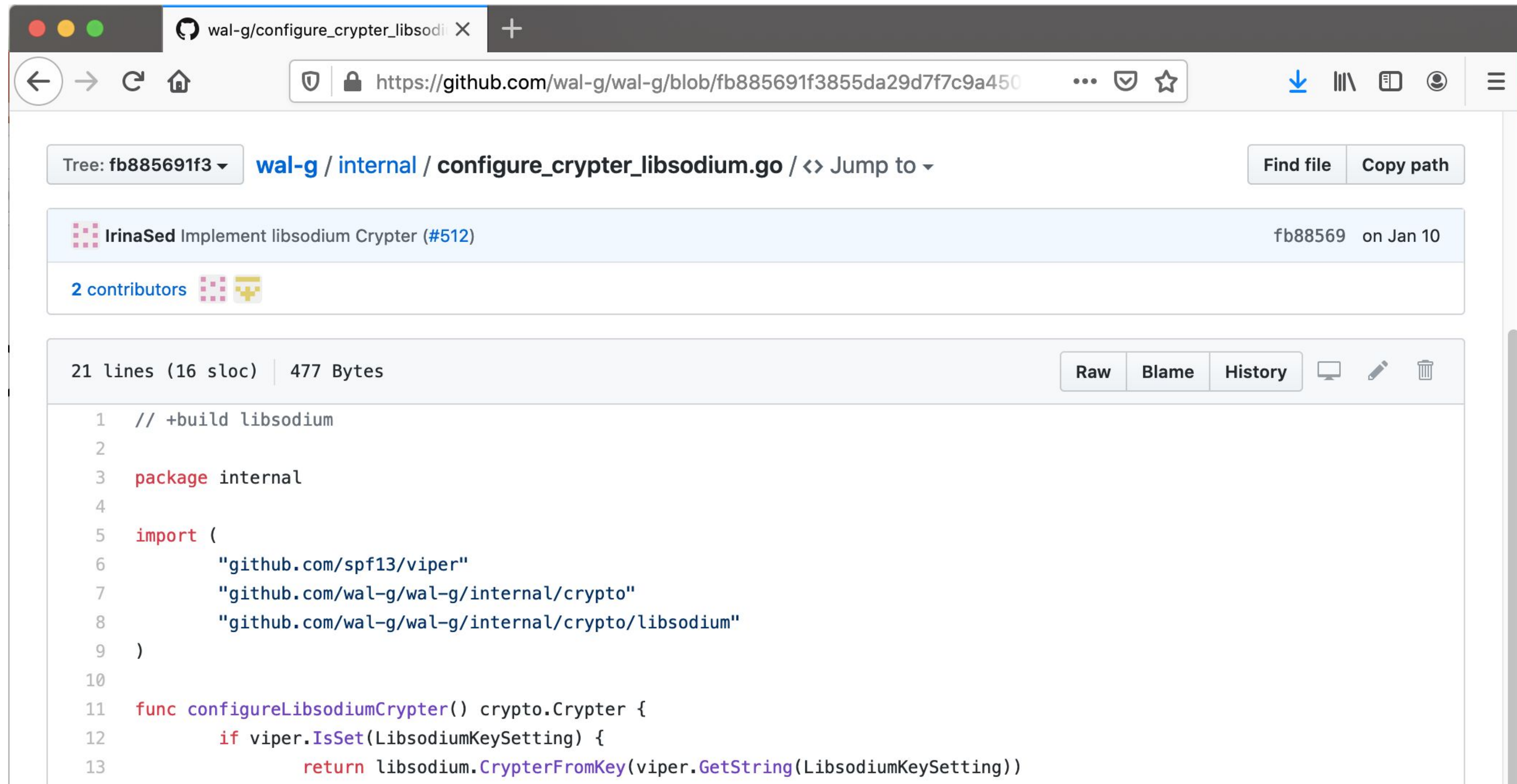
NO PATCHSET DESCRIPTION

DOWNLOAD


EXPAND ALL


File	Comments	Size	Delta
Commit message			
M openpgp/write.go			+3 -3


libsodium



Tree: fb885691f3 ▾ wal-g / internal / configure_crypter_libsodium.go / <> Jump to ▾ Find file Copy path

 IrinaSed Implement libsodium Crypter (#512) fb88569 on Jan 10

2 contributors 

21 lines (16 sloc) | 477 Bytes Raw Blame History 

```
1 // +build libsodium
2
3 package internal
4
5 import (
6     "github.com/spf13/viper"
7     "github.com/wal-g/wal-g/internal/crypto"
8     "github.com/wal-g/wal-g/internal/crypto/libsodium"
9 )
10
11 func configureLibsodiumCrypter() crypto.Crypter {
12     if viper.IsSet(LibsodiumKeySetting) {
13         return libsodium.CrypterFromKey(viper.GetString(LibsodiumKeySetting))
14     }
15 }
```


Push-based vs Pull-based executer



io.Reader vs io.Writer

```
type Reader interface {  
    Read(p []byte) (n int, err error)  
}
```

```
type Writer interface {  
    Write(p []byte) (n int, err error)  
}
```

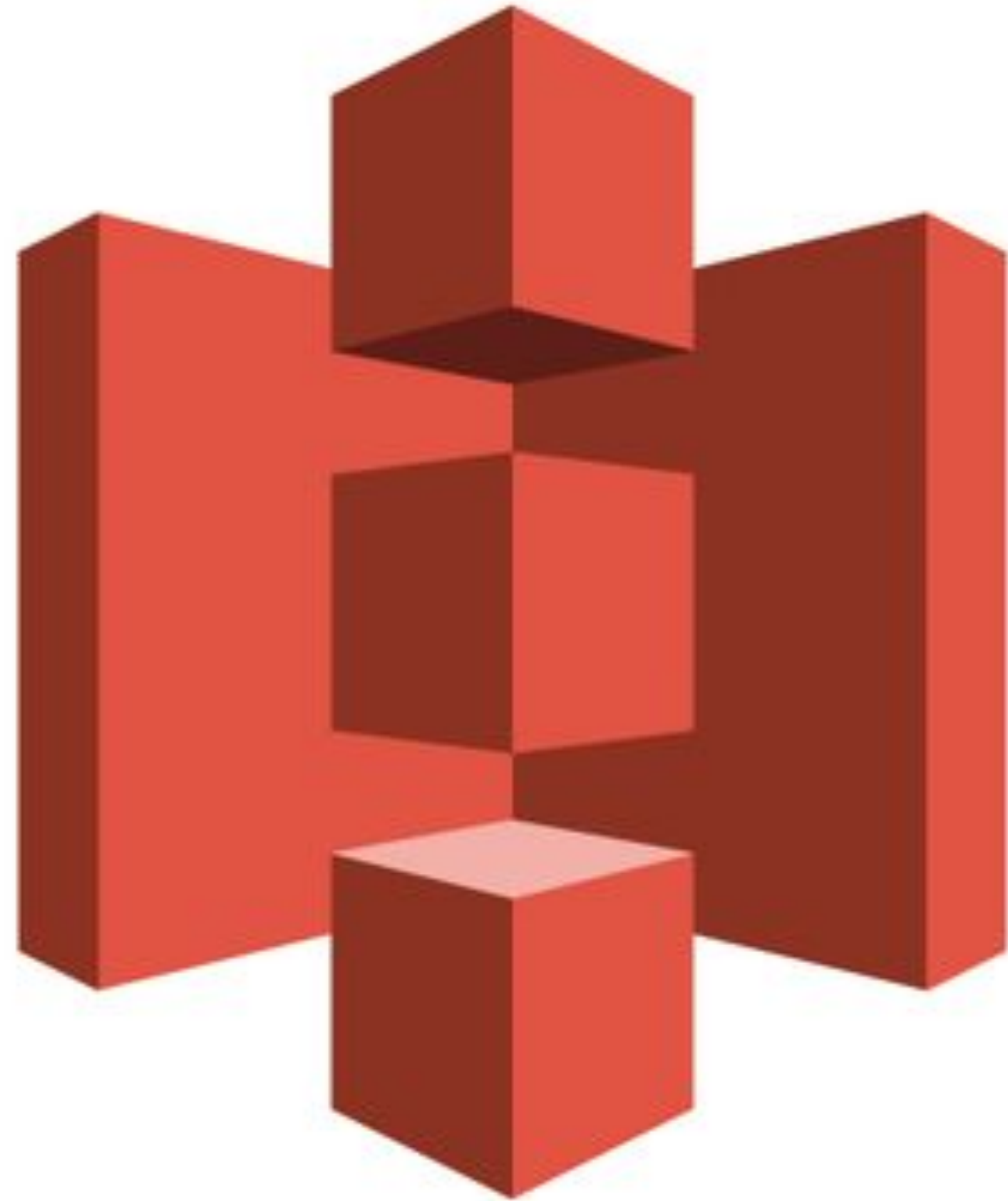
io.ReadFull

```
func ReadAtLeast(r Reader, buf []byte, min int) (n int, err error) {  
    if len(buf) < min : 0, ErrShortBuffer ↗  
    for n < min && err == nil {  
        var nn int  
        nn, err = r.Read(buf[n:])  
        n += nn  
    }  
    if n >= min {  
        → err = nil  
    } else if n > 0 && err == EOF {  
        err = ErrUnexpectedEOF  
    }  
    return  
}
```

Store some bytes



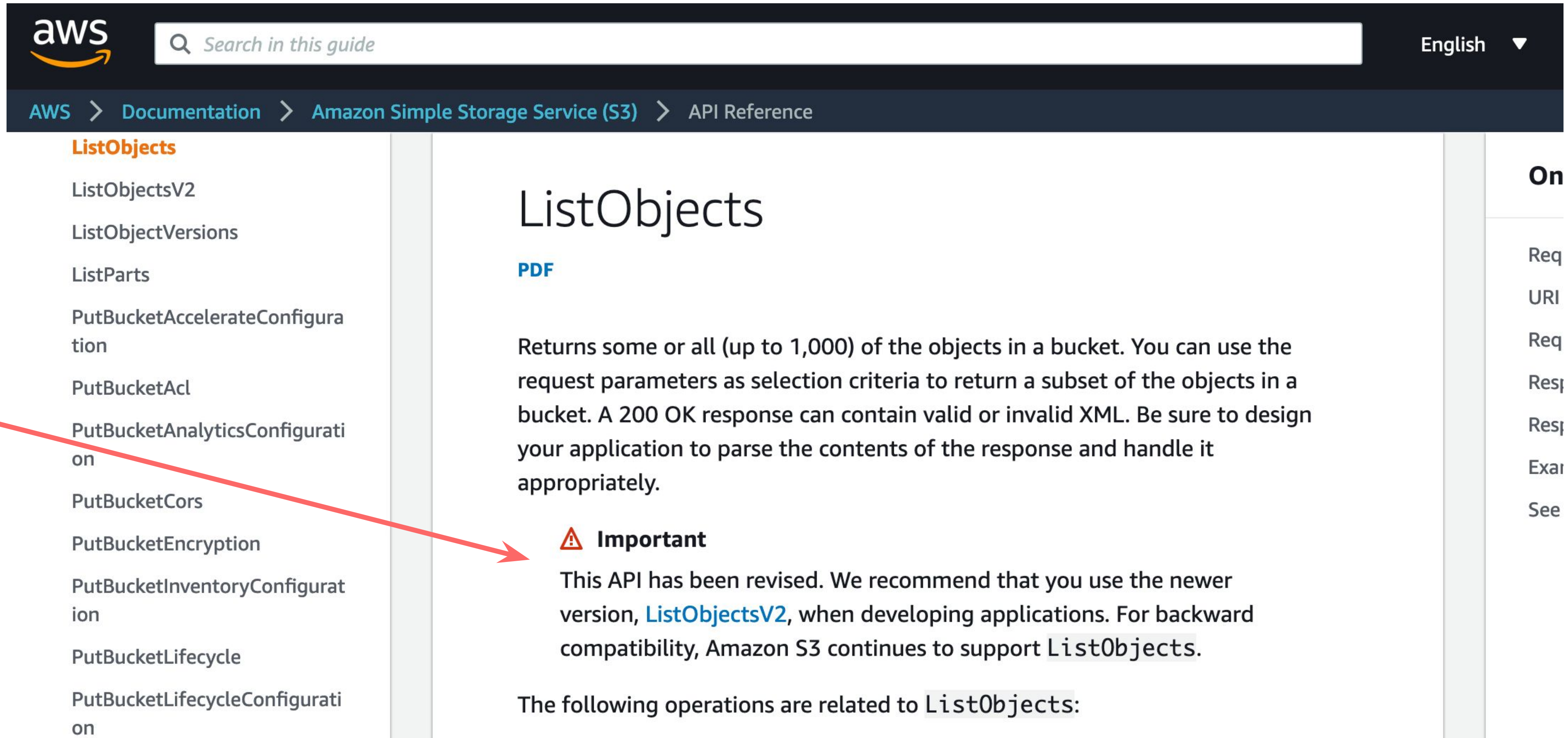
S3-based WAL-G



S3-based WAL-G

■ AWS has it's own cryptography

S3 ListObjectsV2()



The screenshot shows the AWS documentation page for the `ListObjects` API. The left sidebar contains a list of API actions, with `ListObjectsV2` highlighted. The main content area displays the `ListObjects` API details, including a description, a PDF link, and an important warning about the API being revised. The right sidebar shows a table of contents for the page.

aws Search in this guide English ▼

AWS > Documentation > Amazon Simple Storage Service (S3) > API Reference

ListObjects

- ListObjectsV2
- ListObjectVersions
- ListParts
- PutBucketAccelerateConfiguration
- PutBucketAcl
- PutBucketAnalyticsConfiguration
- PutBucketCors
- PutBucketEncryption
- PutBucketInventoryConfiguration
- PutBucketLifecycle
- PutBucketLifecycleConfiguration

ListObjects

[PDF](#)

Returns some or all (up to 1,000) of the objects in a bucket. You can use the request parameters as selection criteria to return a subset of the objects in a bucket. A 200 OK response can contain valid or invalid XML. Be sure to design your application to parse the contents of the response and handle it appropriately.

⚠ Important

This API has been revised. We recommend that you use the newer version, [ListObjectsV2](#), when developing applications. For backward compatibility, Amazon S3 continues to support `ListObjects`.

The following operations are related to `ListObjects`:

On

- Req
- URI
- Req
- Res
- Res
- Exa
- See

S3 infiltrated the whole codebase



Refactor storages

The screenshot shows the GitHub repository page for `wal-g/storages`. The repository is a fork of `Tinsane/storages`. The page displays the repository's name, description (which is empty), and various statistics: 34 commits, 4 branches, 0 packages, 0 releases, 8 contributors, and Apache-2.0 license. The 'Code' tab is selected, showing a list of files and their commit messages. The files listed are `azure`, `fs`, `gcs`, `memory`, `s3`, `storage`, and `swift`. The commit messages for these files are: 'use wal-g/tracelog instead of tinsane/tracelog' for `azure`, `fs`, `gcs`, and `swift`; 'Fix imports' for `memory`; and 'Remove s3 debug (#14)' for `s3`. The commit for `storage` is 'Add folder and object mocks'. The latest commit is `d4de6f5` on Jan 23.

wal-g / storages

forked from Tinsane/storages

Watch 0 Unstar 2 Fork 11

Code Pull requests 3 Actions Projects 0 Wiki Security Insights Settings

No description, website, or topics provided. Edit

Manage topics

34 commits 4 branches 0 packages 0 releases 8 contributors Apache-2.0

Branch: master New pull request Create new file Upload files Find file Clone or download

This branch is 27 commits ahead of Tinsane:master. Pull request Compare

perekalov and g0djan Remove s3 debug (#14) Latest commit d4de6f5 on Jan 23

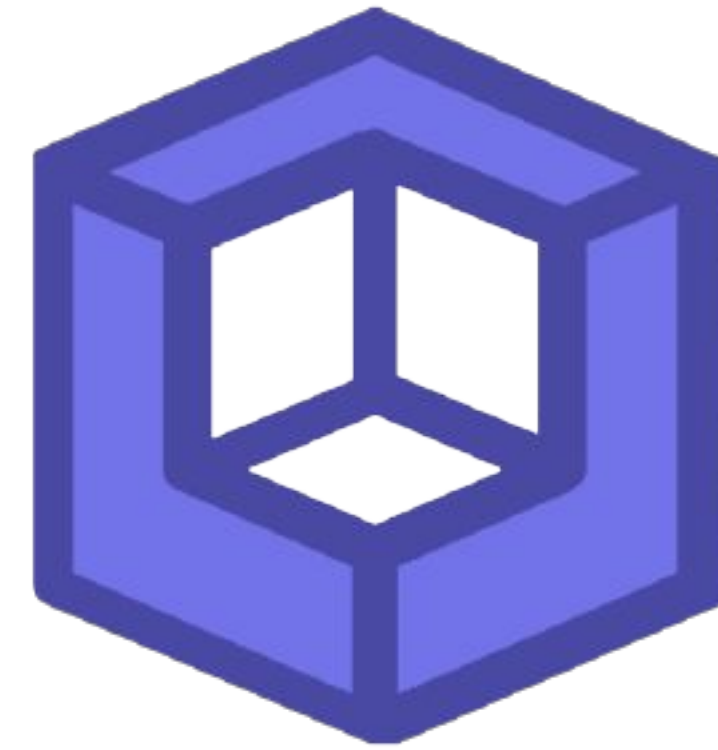
azure	use wal-g/tracelog instead of tinsane/tracelog	4 months ago
fs	use wal-g/tracelog instead of tinsane/tracelog	4 months ago
gcs	use wal-g/tracelog instead of tinsane/tracelog	4 months ago
memory	Fix imports	4 months ago
s3	Remove s3 debug (#14)	2 months ago
storage	Add folder and object mocks	2 months ago
swift	use wal-g/tracelog instead of tinsane/tracelog	4 months ago

GCP

- › No multipart uploads
- › A lot of problems with retries

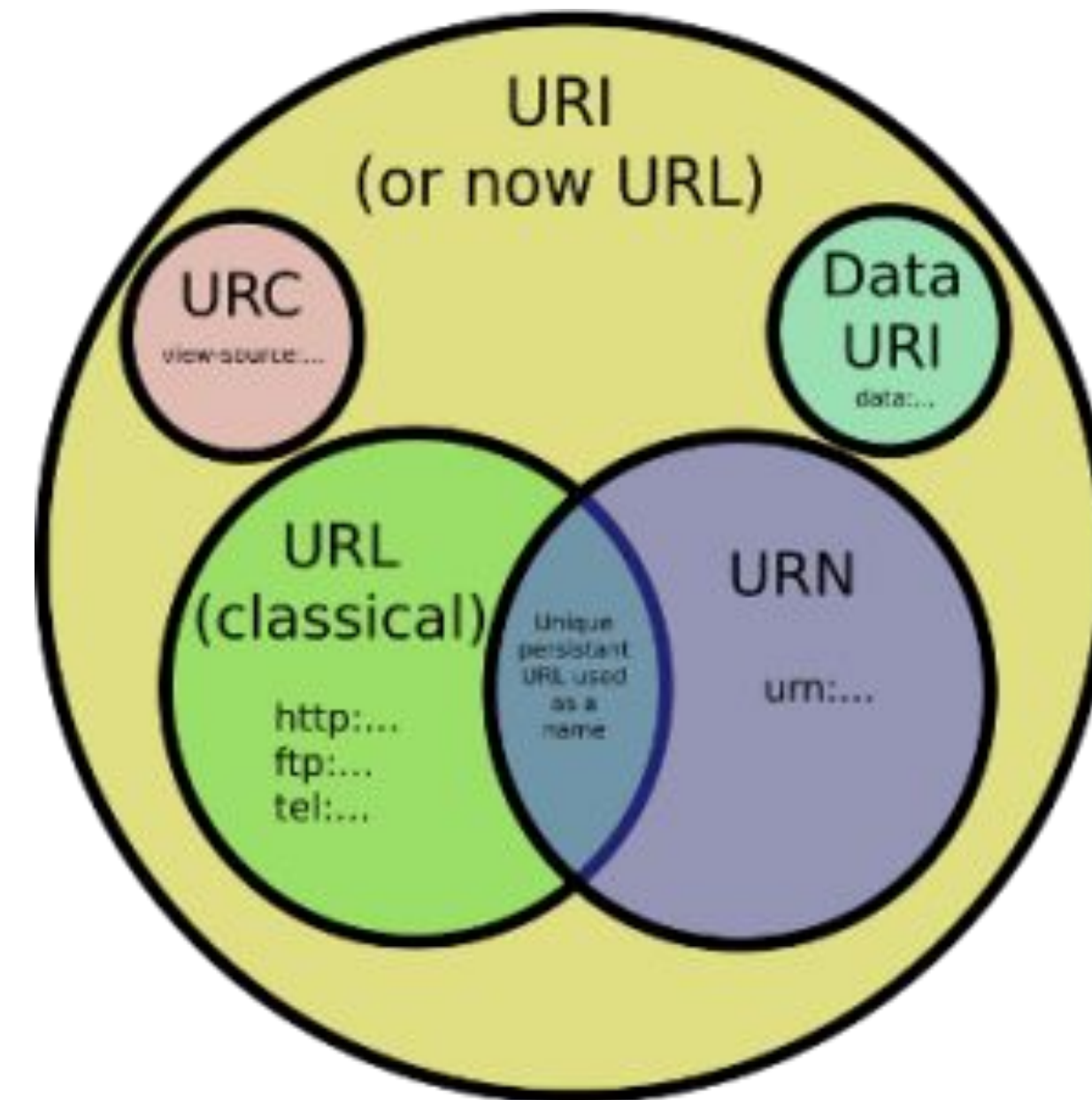
Azure

 Contribution from KubeDB



FS

Problems with Windows build



SWIFT

■ No problems, does anyone use WAL-G + SWIFT?

SSH/SCP

- › Because we need to back up S3 too

WAL-G for other databases



WAL-G for MySQL®



What to archive?

Data

Transactions log



Differences from PostgreSQL

- Storage engines
- Redo log is circular
- No archive command
- PITR is based on binlog, not redo



Data archiving

xtrabackup

mysqldump



What is binlog?

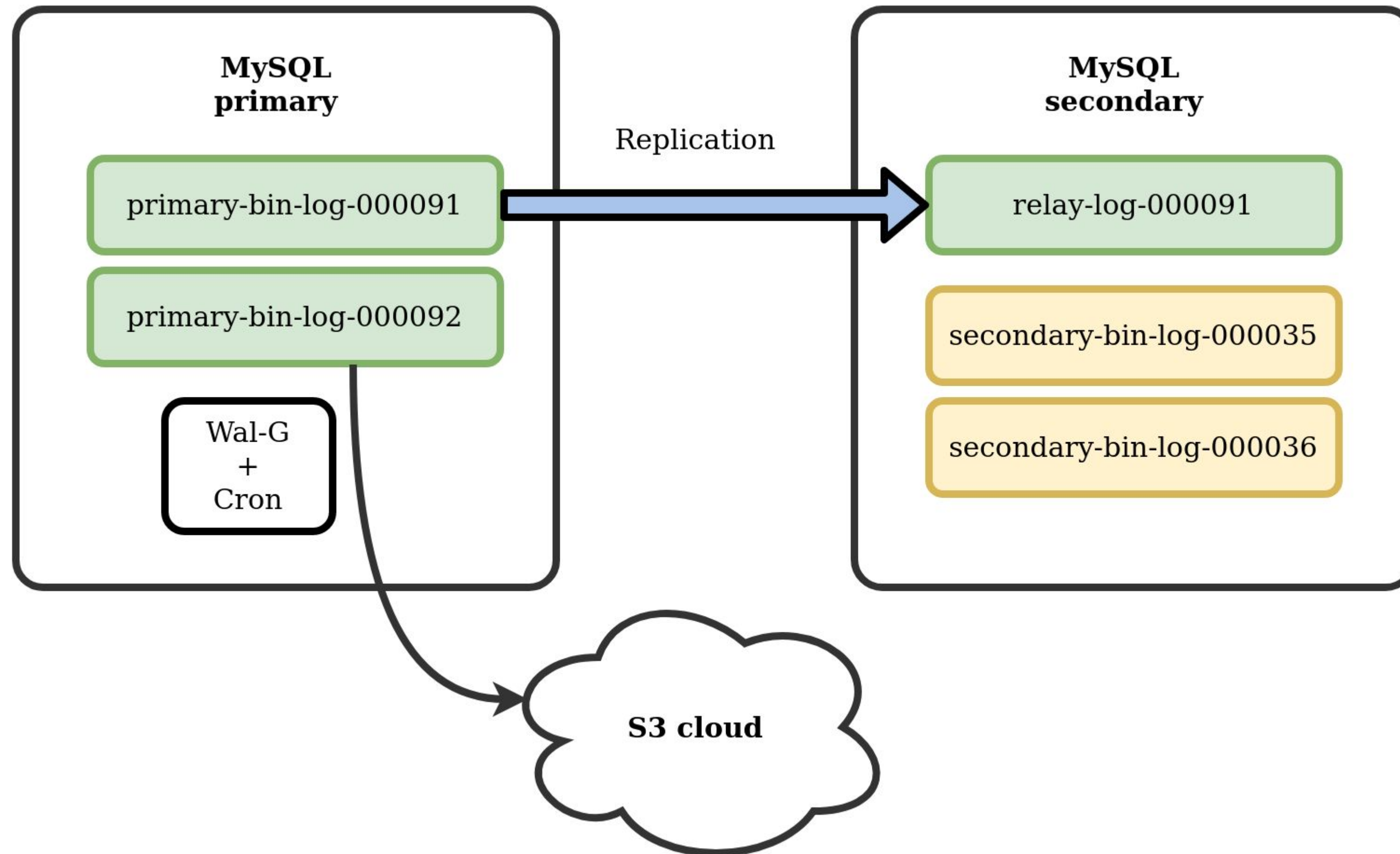
- Transactions (rows or queries)
- Timestamps
- GTID - unique ID

```
$ mysqlbinlog -v /var/lib/mysql/binlog.001242
....
SET @@SESSION.GTID_NEXT=
'63bebac7-c424-11e9-bb82-15ed63279c98:31214803';

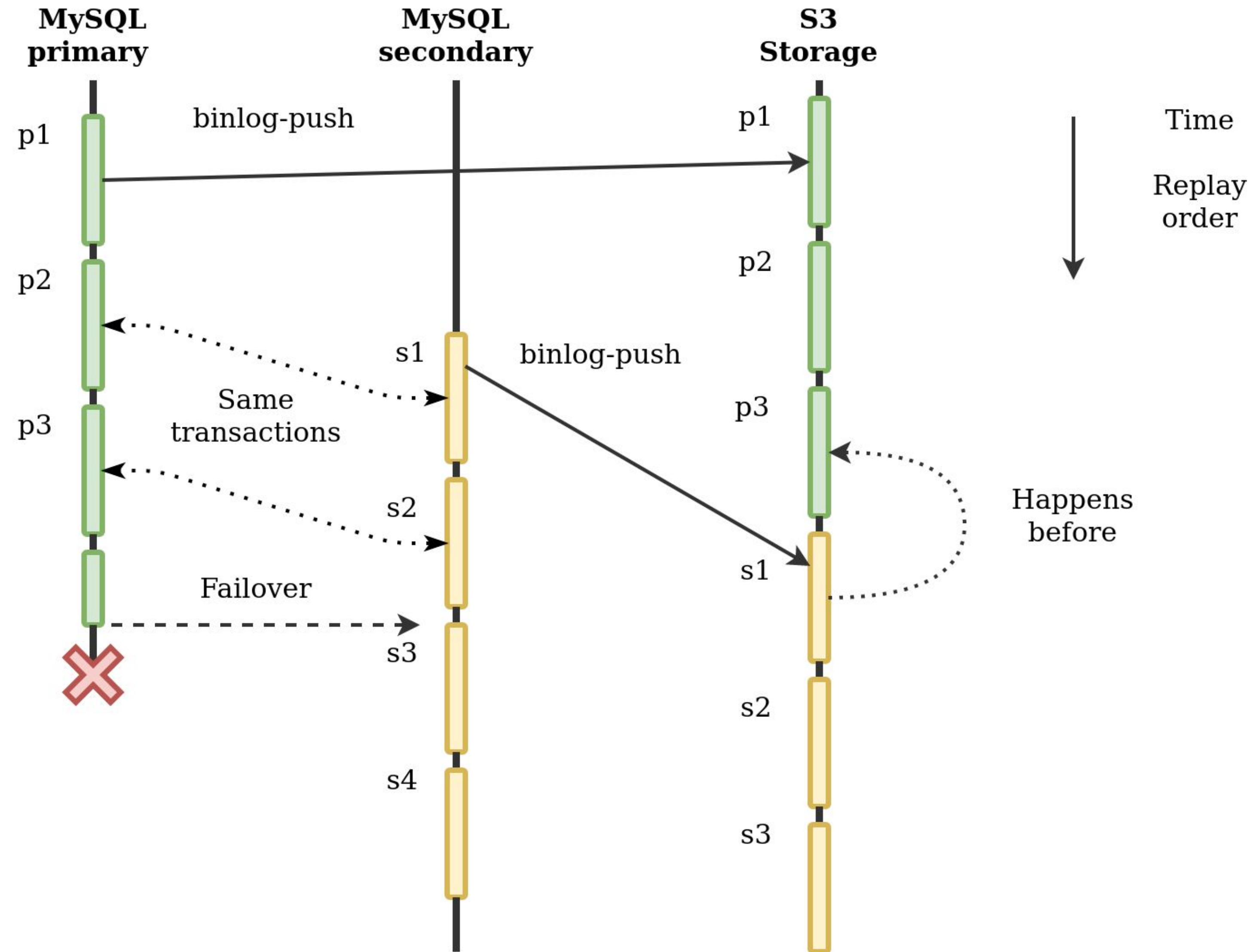
TIMESTAMP=1619172002;

BEGIN
BINLOG '
opqCYBMCAAAAOgAAALoxHgAAAG0AAAAAAAEABW15c3
FsAAxtZGJfcmVwbF9tb24AAgMRAQMC+ajWYw==
opqCYB8CAAAAOgAAAPQxHgAAAG0AAAAAAAEAAgAC///8
AQAAAGCCmqEaNvwBAAAAYIKaohpA3Bx6WA==
'
....
```


Binlog archivation



Binlog backup and replay order



Speed up replay

- Run MySQL on high port (3308)
- Disable all possible logging
- Don't forget to put everything
back into production!

```
double_write_buffer = OFF
```

```
innodb_flush_log_at_trx_commit = 0
```

```
sync_binlog = 1000
```

```
event-scheduler = OFF
```

Complete scenario

Back up

- Daily:

wal-g backup-push
- Each minute:

wal-g binlog-push

Restore

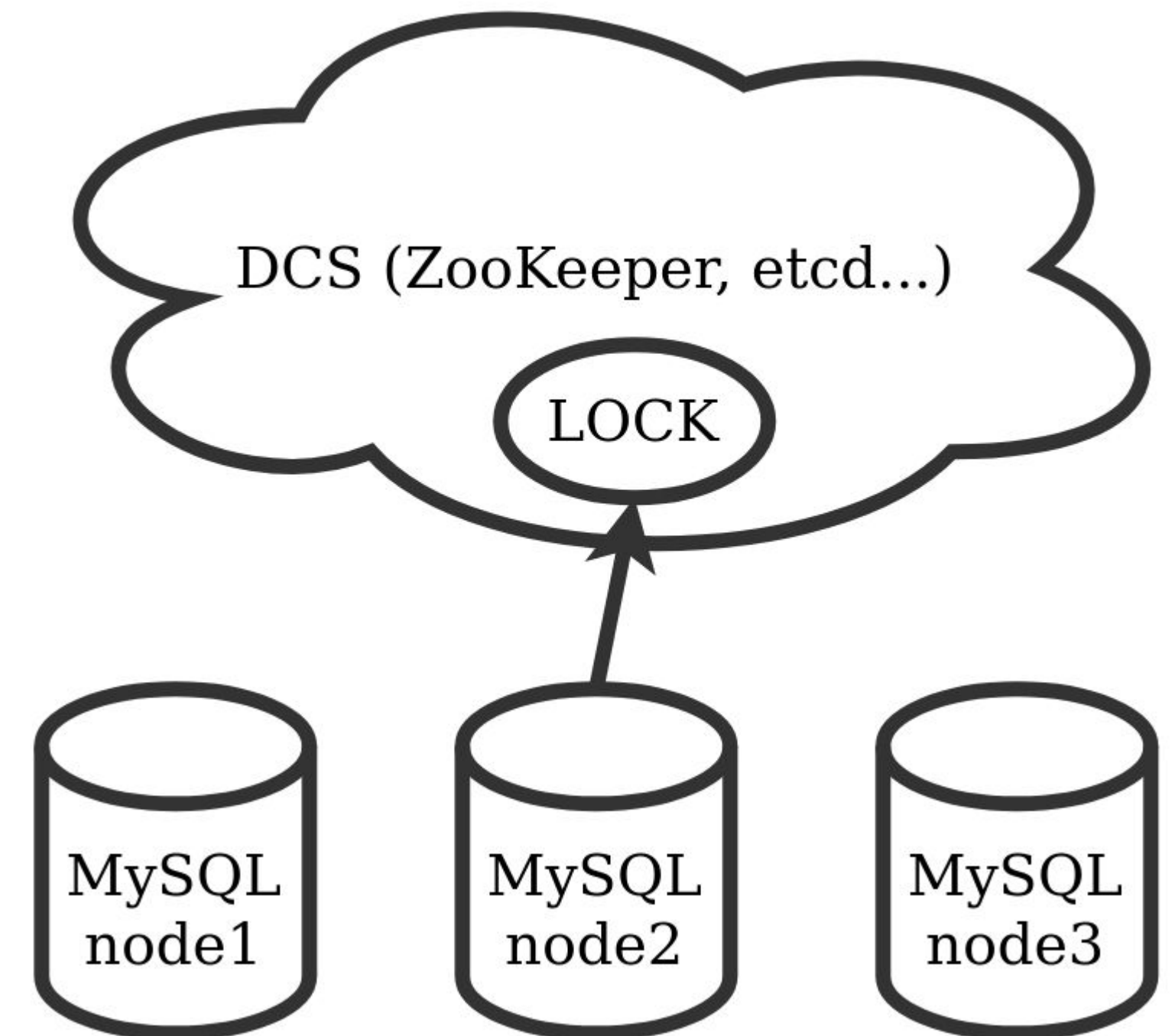
- wal-g backup-fetch LATEST
- Start mysql on high port
- wal-g binlog-replay

--since LATEST

--until 2021-05-12T12:00:00Z
- Restart MySQL with production config

Out of scope

- How to choose which node to back up?
- How to not upload binlog from a crashed master?
- How to handle binlog name collisions?



Does it work at all?

> 1k
clusters

> 1.5k
nodes

*In Yandex

WAL-G for Microsoft SQLServer™



SQLServer backup features

- Database backups
- Log backups
- Incremental backups
- Point-in-time recovery



Backup methods

- DISK
- TAPE
- VIRTUAL DEVICE
- URL (Azure)



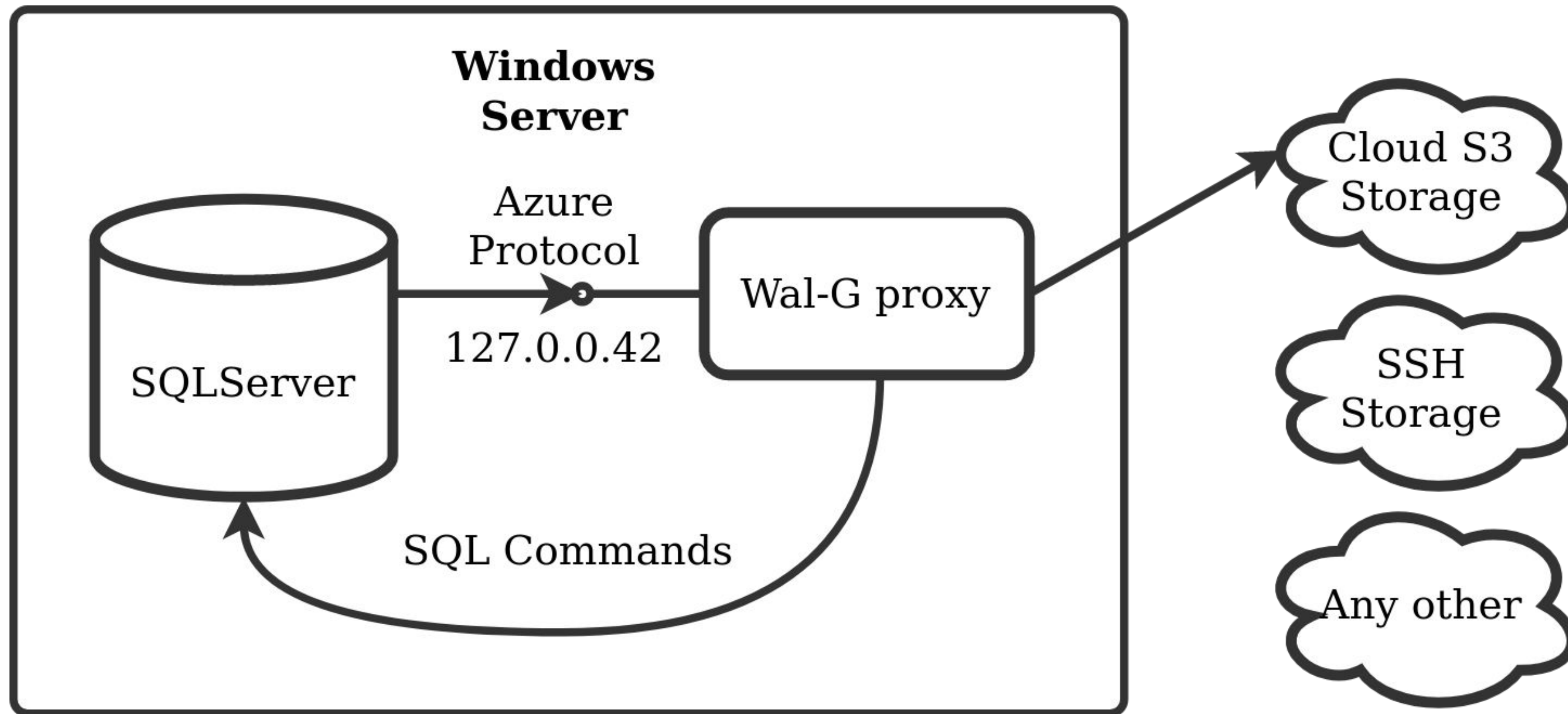
Backup to DISK

- Additional disk space required
- Additional IOPs required
- Still need to upload to storage

OR

- Additional network disk

WAL-G approach: emulate Azure



Features and hacks

1. Need domain name
2. Need HTTPS
3. Only default HTTPS port - 443
4. Azure credentials required

```
$ Add-Content -Path  
-Value '127.0.0.42 backup.local'  
-Path 'C:\Windows\System32\Drivers\etc\hosts'  
  
$ Import-Certificate  
-CertStoreLocation cert:\LocalMachine\Root \  
-FilePath 'C:\backup.local.cert.pem'  
  
$ Invoke-Sqlcmd "  
CREATE CREDENTIAL  
[https://backup.local/basebackups_005]  
WITH DENTITY='SHARED ACCESS SIGNATURE',  
SECRET = 'does_not_matter'  
"
```


Case1: flexibility

1. Run WAL-G in background

```
wal-g --config c:\walg.yaml proxy
```

2. Now you have your own private Azure

```
BACKUP DATABASE [db1]  
TO URL = 'https://walg.local/somepath/weekly/20210501'
```

Case 2: PG-like CLI

wal-g backup-push

wal-g backup-list

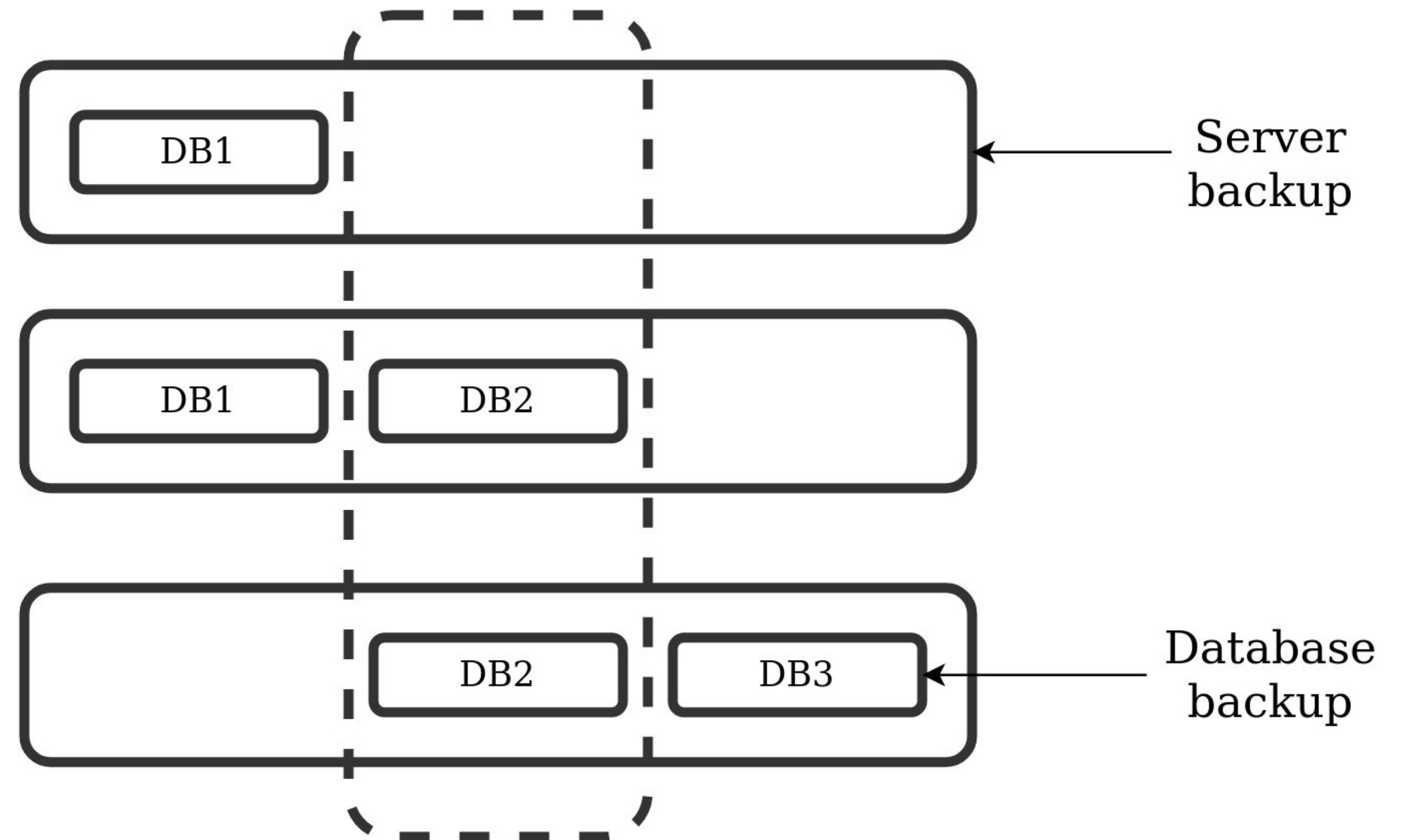
wal-g log-push

wal-g backup-restore LATEST

wal-g log-restore --since LATEST --until TS

Multiple databases

- backup-push
- backup-push -d db2
- log-restore
- log-restore -d db1,db2

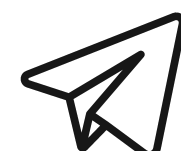


Questions

Andrey Borodin



x4mmm@yandex-team.ru

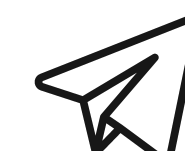


x4mmm

Dmitry Smal



mialinx@yandex-team.ru



mialinx