

# zenith storage system

Heikki Linnakangas

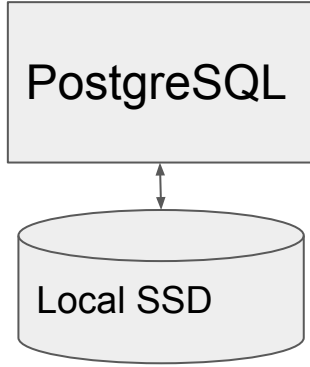
# What is zenith?

- New storage system for PostgreSQL
- Cloud service
- Startup

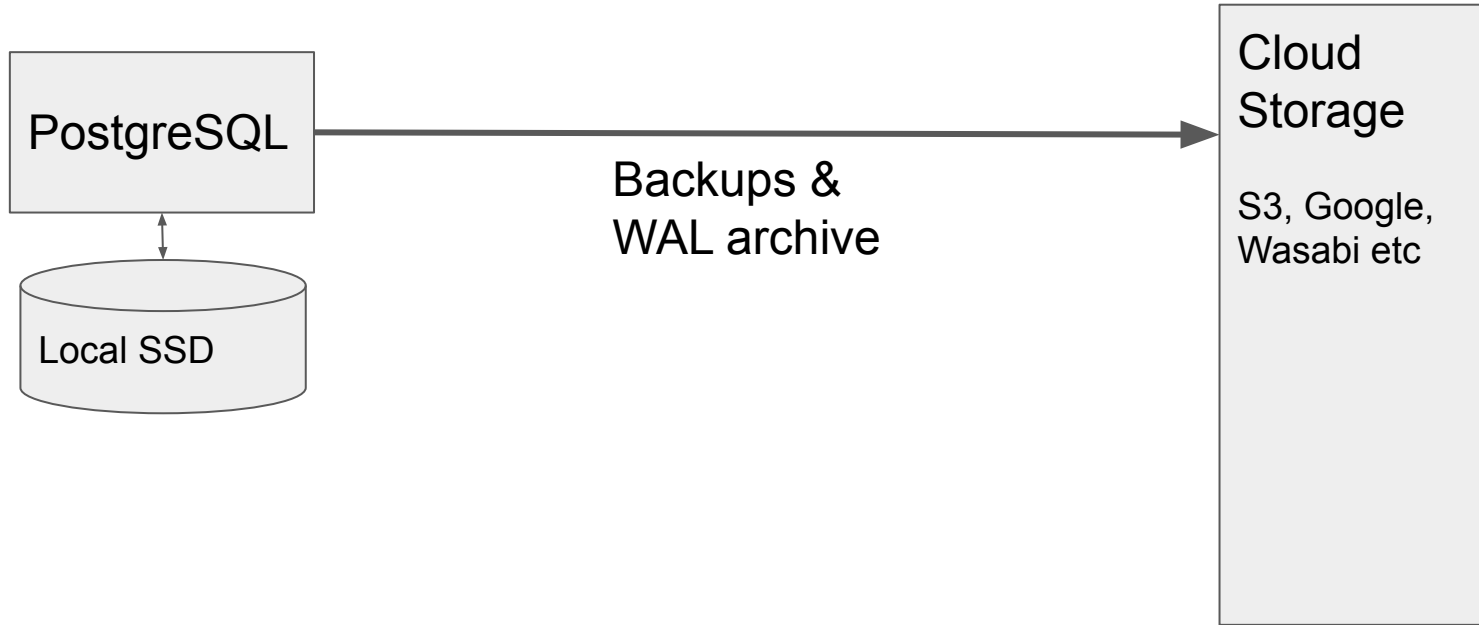
# Forewarning

- Work in progress
- Basic stuff works
- Architecture is sound
- But many of the features I'm describing have not actually been implemented yet.

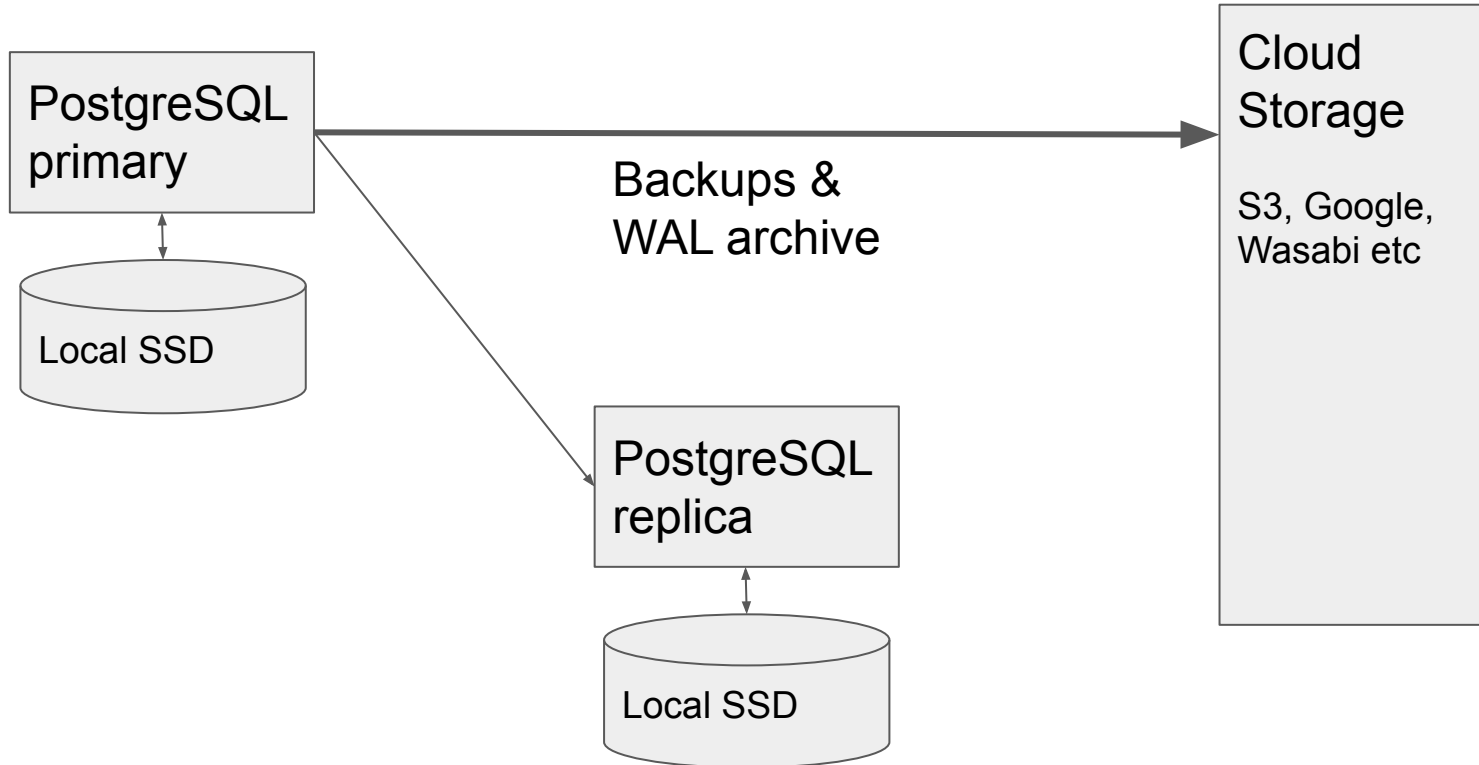
# Traditional PostgreSQL setup



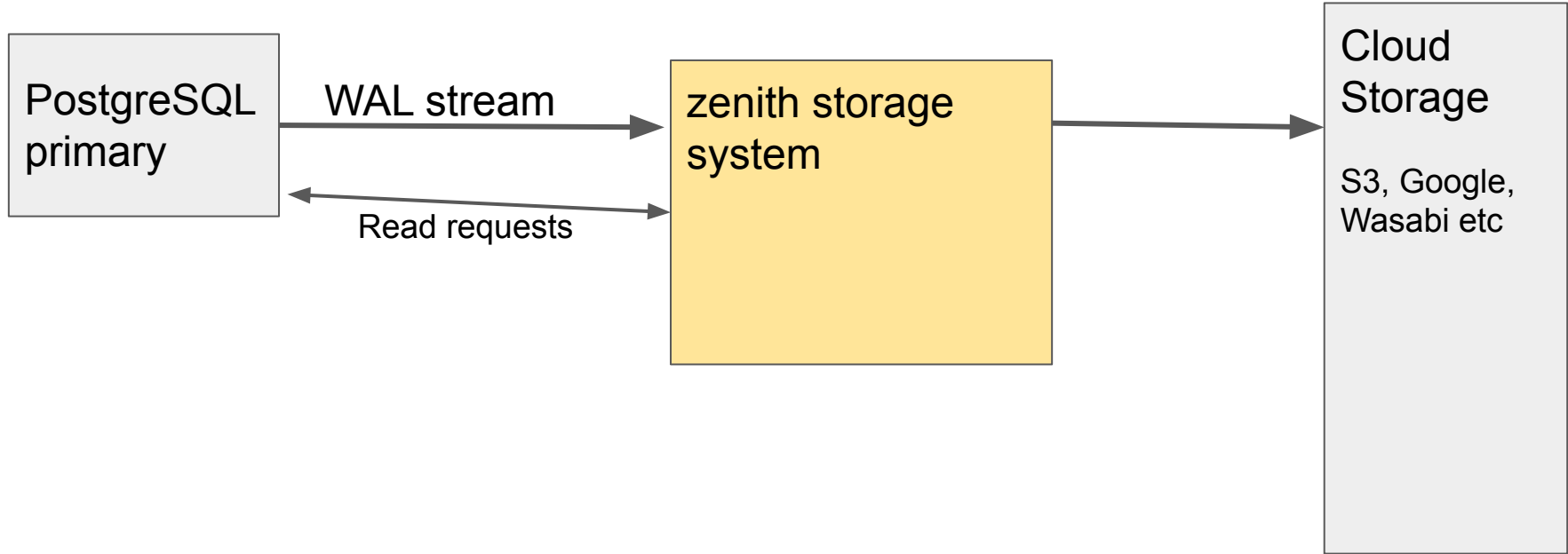
# Traditional PostgreSQL setup



# Traditional PostgreSQL setup



# Separate Compute and Storage



# Separating compute and storage

- Amazon Aurora
- PolarDB from Alibaba
- Microsoft Hyperscale
- Snowflake
- SingleStore



# zenith storage system

- Replaces local storage, backups and WAL archive
- Knows about PostgreSQL WAL and page format
- SAN on steroids

# No lengthy recovery needed

Traditional:

- Restore last backup, replay all WAL, can take hours

zenith:

- Start up **instantly**
- WAL is replayed on-demand, page at a time

# Setting up a read replica

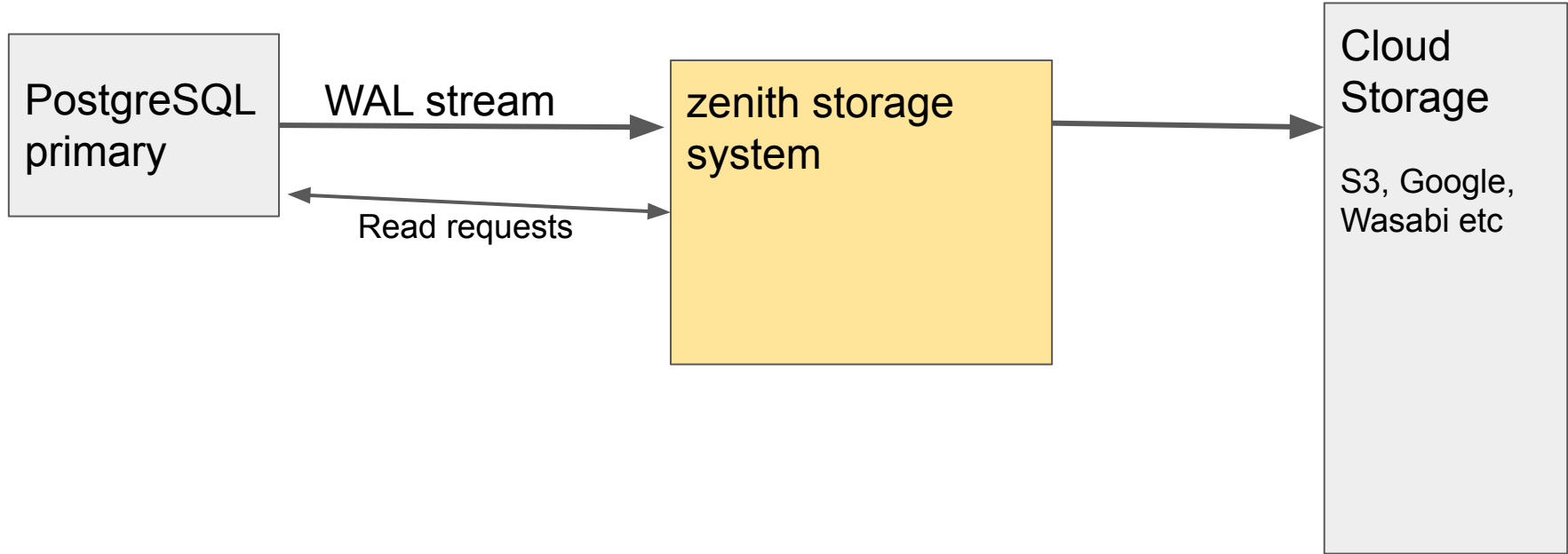
## Traditional:

- Restore last backup, replay all WAL, can take hours
- Needs a fully copy of the data

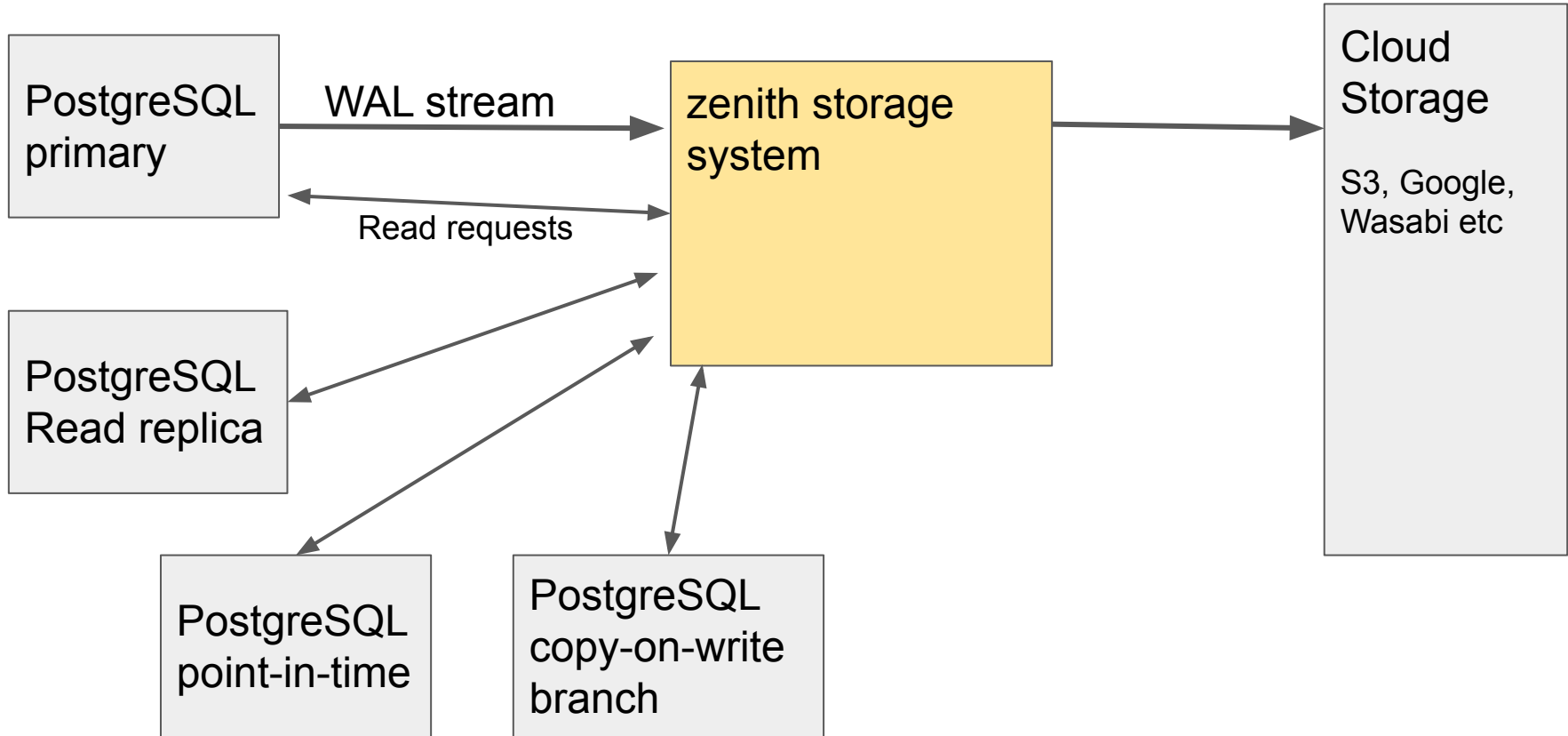
## Zenith:

- Start up **instantly**
- Lightweight, shared storage

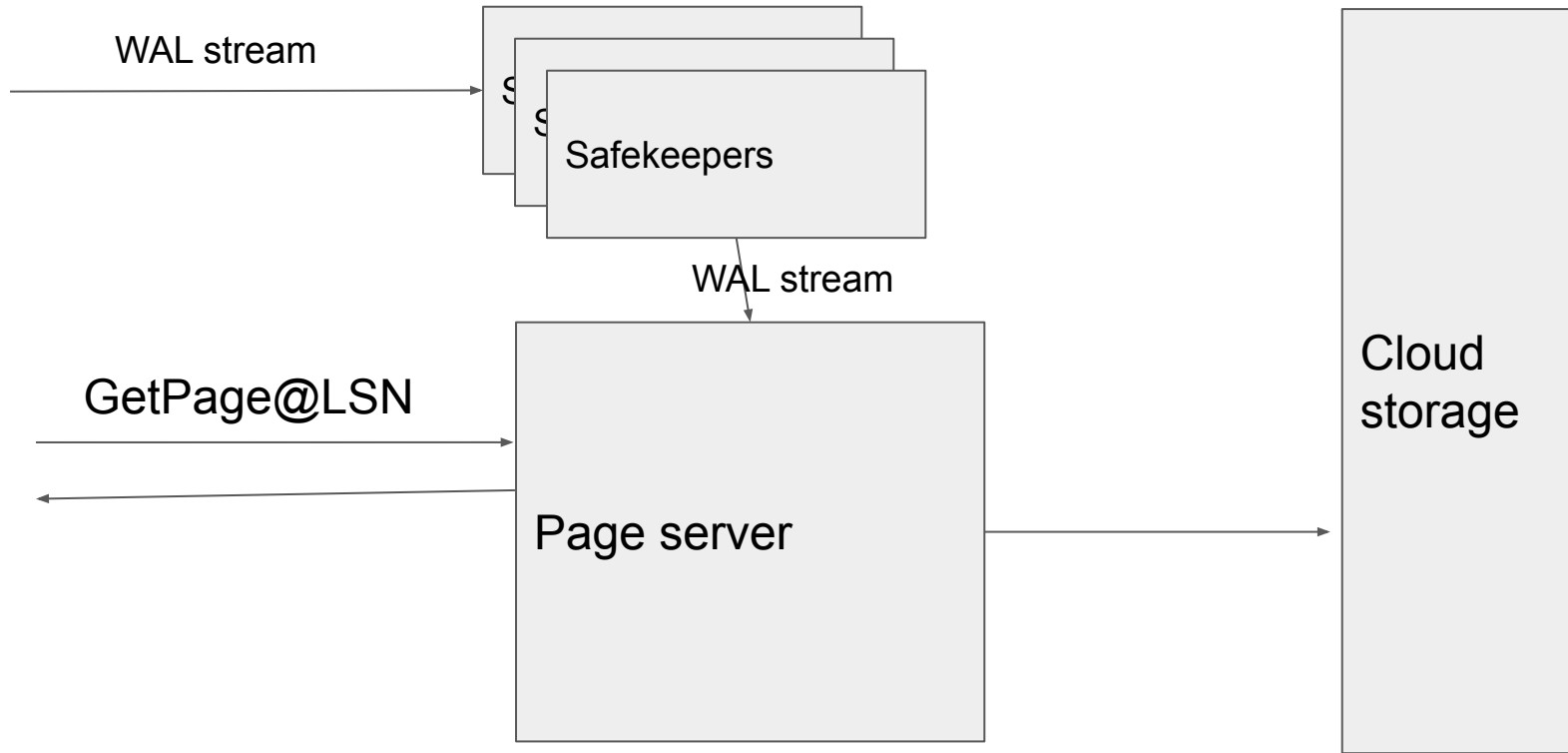
# Separate Compute and Storage



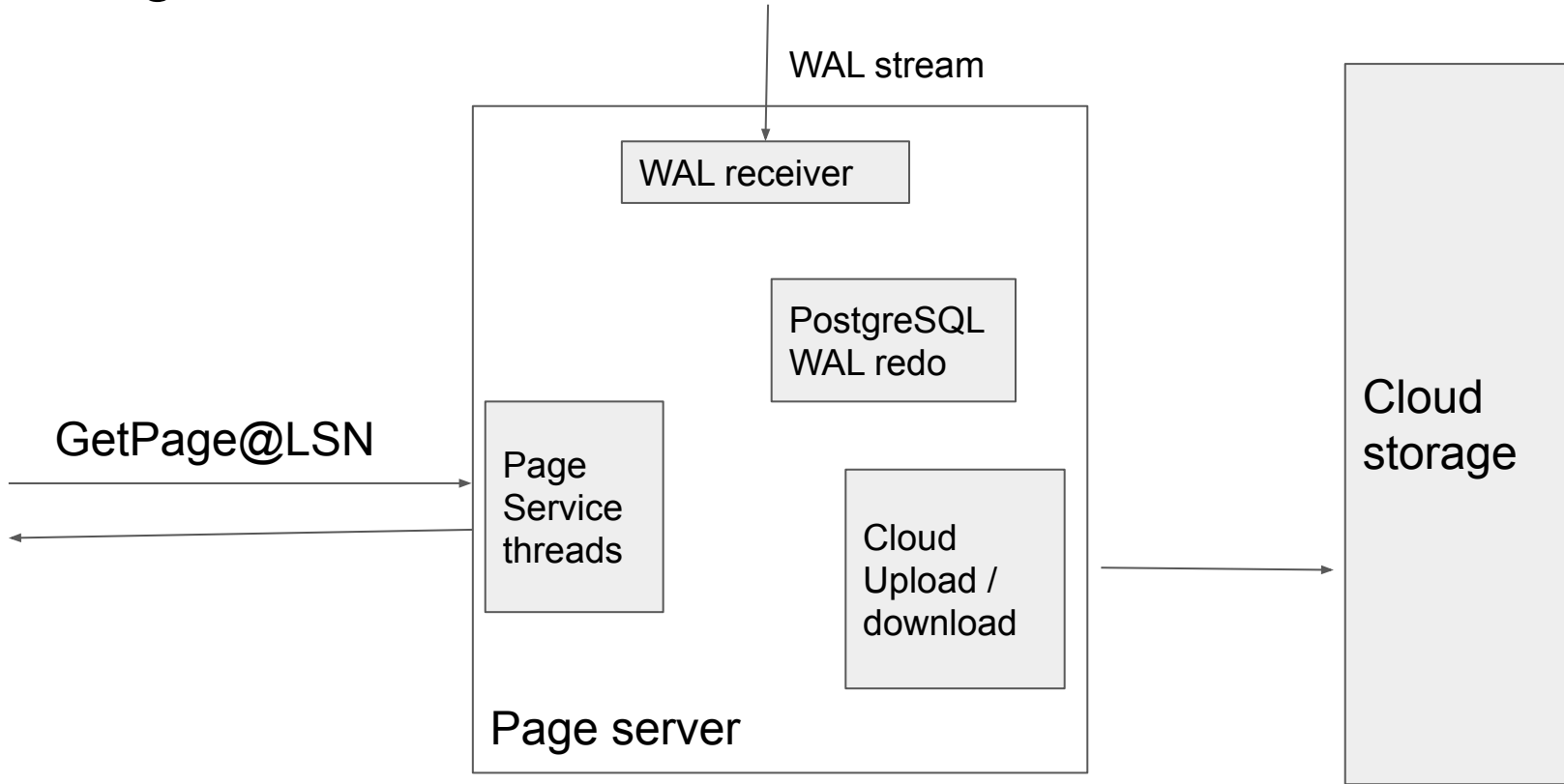
# Separate Compute and Storage



# Closer look



# Page server internals



# Page server

- Heart of the storage system
- Receives WAL
- Parses it
- Writes it in a pre-processed format to local disk
- Uploads it to cloud storage



# Storage format

- Immutable files
- Similar to an LSM tree
- Compaction and garbage collection in the background
- Heuristics on when to store WAL records, and when to materialize an image of a page by replaying the WAL records

# GetPage@LSN

- Find the last image of the requested page,  $\leq$  requested LSN, and any WAL records up to the LSN
- Replay the WAL records to reconstruct the requested version of the page

# Branching

- Like git branch, but for data
- Copy-on-Write
- CI/CD workflow:
  1. Create a branch of production database
  2. Run tests against the branch
  3. Drop the branch

Current status

# Current Status

- All the code is on github
- Apache 2 license
- Storage code is written in Rust
- Passes PostgreSQL regression tests

# TODOs

- Read replicas that follow primary
- Sharding/scaling the storage
- Performance
- On-demand download from cloud storage
- Garbage collection

# PostgreSQL modifications

- Goal is to run unmodified PostgreSQL on zenith
  - Not quite there yet
- Smgr API changes
- Handling non-relation data like pg\_xact
- Tracking LSN of last evicted page in buffer manager
- Unlogged tables

# Hosted Cloud service

- Serverless
- Web UI
- API for programmatic access
- Free tier

Try it: `psql -h start.zenith.tech`

Invite code: **nordic-2022**

Zenith console — Mozilla Firefox

File Edit View History Bookmarks Tools Help

Zenith console x +

Jump to... Documentation Website hinnaka

### Free tier during the Technical Preview

Zenith cloud service is free over the course of Technical Preview in 2022. The following limits apply:

- compute up to 1vCPU/256 mb
- up to 10GB storage
- 3 clusters per user

Clusters (3) + New Cluster (3/3)

Name	Region	Created at	Platform	Status
main	US East (N. Virginia)	Mar 21, 2022 11:03	Serverless	Idle
broken-mud-043209	US East (N. Virginia)	Mar 21, 2022 11:03	Serverless	Idle
bitter-salad-2982	US East (N. Virginia)	Dec 15, 2021 07:53	Serverless	Idle



Thank you!

Q & A

<https://github.com/zenithdb/zenith/>

Try it: `psql -h start.zenith.tech`

Invite code: **nordic-2022**

Feedback: [heikki.linnakangas@iki.fi](mailto:heikki.linnakangas@iki.fi)