# Scaling Btrfs
# in an Enterprise

Motiejus Jakštys
m.jakstys.lt
@motiejus:jakstys.lt

Chronosphere

# This talk is about

- Btrfs cut our storage costs by 74%.
- Productionized Btrfs in Chronosphere and Google Cloud Platform.

# Our business & size (mid 2025)

- Observability SaaS, lots of time-series data.
- Customers so big we run thousands of VMs *per customer*.
- Multi-terabyte disks per VM.
  - ⇒ Some PiB in Google's Persistent Disks.
- All on ext4.

# Why Btrfs?

1.  Compression saves disk space.

2.  Compression is fully transparent.

3.  Checksums can save on CPU and application complexity.

# Disk compression demo

Initial tests showed ~65% disk savings.

# Btrfs at Facebook

By **Jonathan Corbet**
July 2, 2020

OSSNA

The Btrfs filesystem has had a long and sometimes turbulent history; LWN first wrote about it in 2007. It offers features not found in any other mainline

on a btrfs root partition recently as an experiment. What a pain.

# Hurdles

1. Confidence & reputation.
   a. Savings significant enough to worth trying.
2. ENOSPC.

   a. bg_reclaim_treshold + dynamic_reclaim.
3. Unsupported by Google — why?
   a. Didn't work? Or nobody just did it?
4. How will IO patterns affect performance?
   a. Will get to that 😈.
5. Unknown unknowns?

Surprises came up

# Reclaim causing massive IO on large deletes

Surprises came up
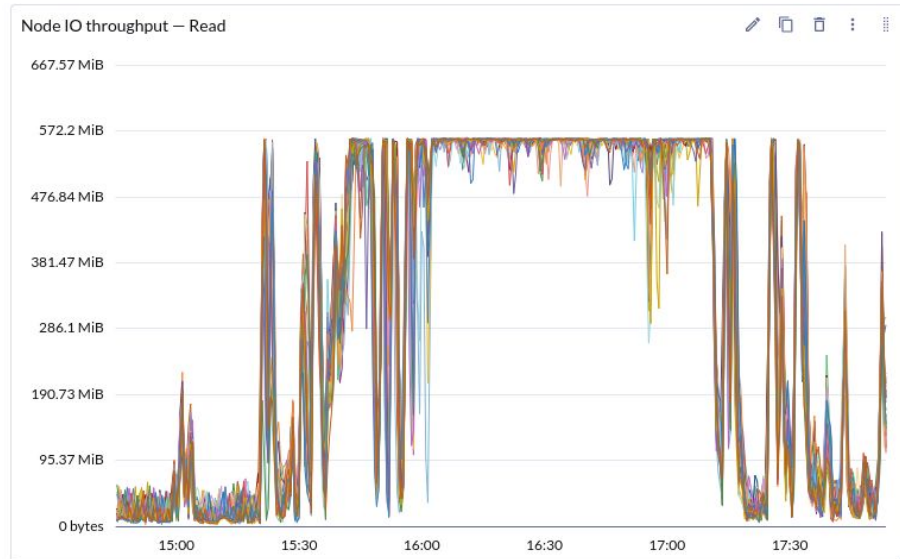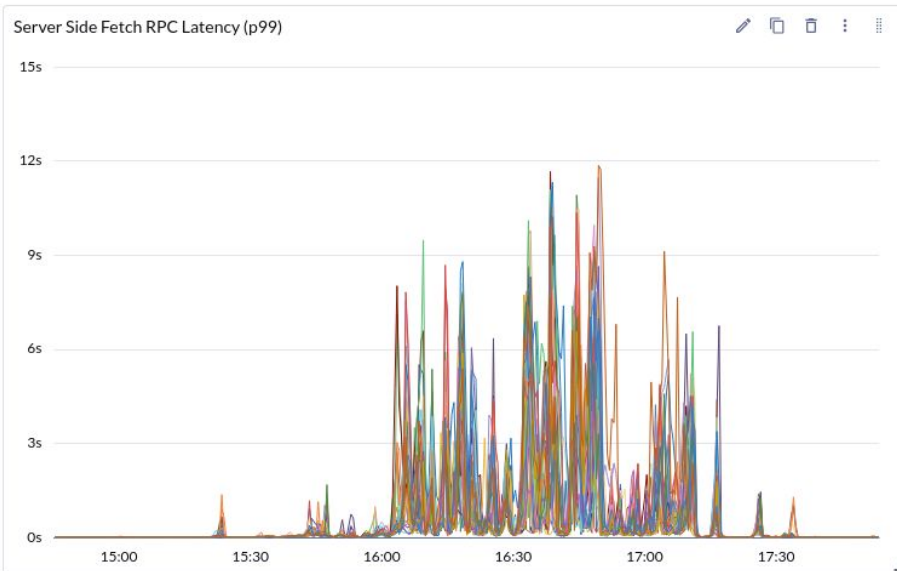
# Read Ahead!



Server Side Fetch RPC Latency (p99)

Node IO throughput — Read

# Read Ahead!

`/sys/block/<...>/queue/`**`read_ahead_kb`**

`128`

`/sys/fs/btrfs/<...>/bdi/`**`read_ahead_kb`**

`4096`

Can you tell the difference?

# Timeline

- 2024-06: Wrote a doc and redirected everyone there.
  - Every question ever asked was put to the doc.
- 2024-08: kick off infra work.
- 2025-01: first internal cluster.
- 2025-06: first customer cluster.
- 2025-08: mass migration.
- 2025-10: done.
- 2026-01: yours truly @ FOSDEM.

# Takeaways

- All metrics data @ Chronosphere is on Btrfs.
  - Our savings: 65% predicted, 74% actual.
- Btrfs is available upstream in GCP:
  - [(CSI) driver](#) and [Container Optimized OS (COS)](#).
  - Planning to go full-upstream in the next 1-2 months.
- You know one more company trusting Btrfs.

# Special thanks

- Boris Burkov, Josef Bacik from Meta.
  - Reclaim (defragment).
  - Spreading the word!
- All Btrfs maintainers.
- GCP Storage team for being Btrfs champions.

Motiejus Jakštys, Chronosphere
[m.jakstys.lt](m.jakstys.lt)

@motiejus:jakstys.lt

# Bonus slide: unknown unknowns

- Filesystem deadlock on significant memory pressure.
  - Fixed in mainline and backported to all stable kernels.
- Disk snapshots grew by 6-10x.
  - Replaced with file-based snapshots.