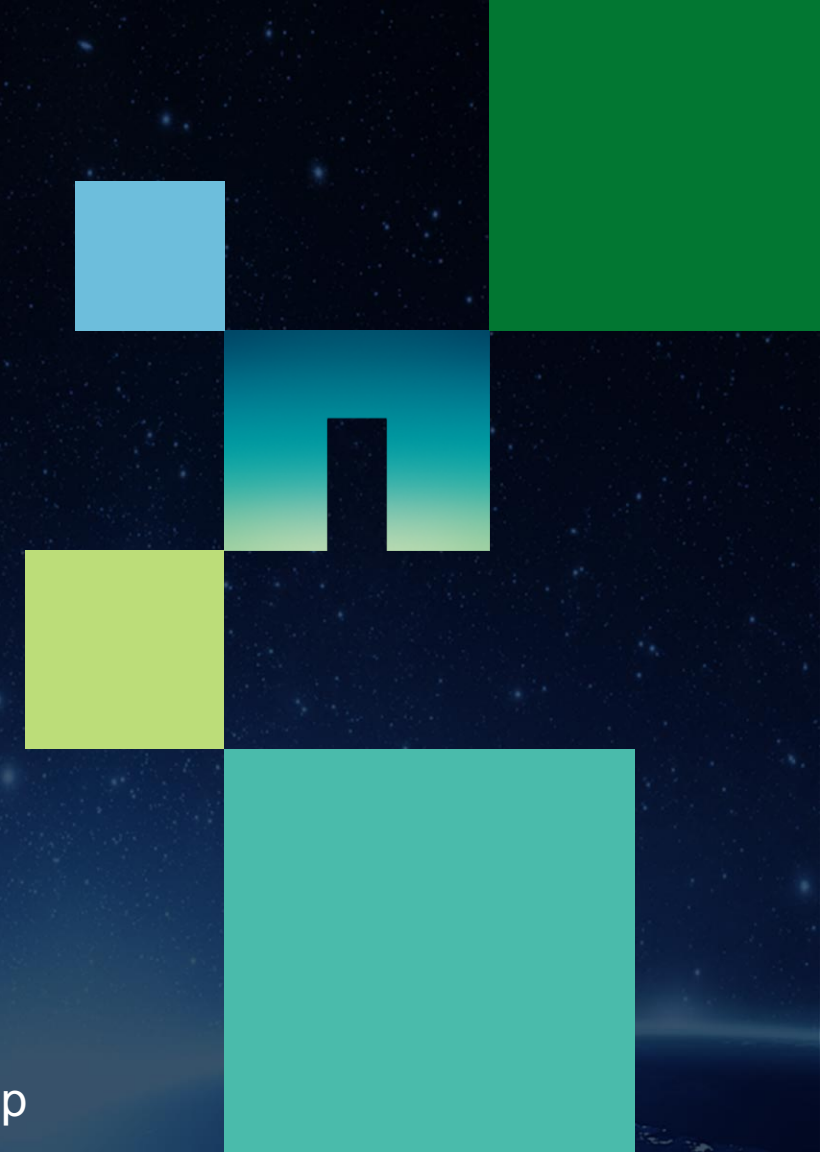# FlexGroup Volumes:

## A Distributed WAFL File System

Ram Kesavan, Google; *Jason Hennessey*, Richard Jernigan,
Peter Macko, Keith A. Smith, Daniel Tennant, and Bharadwaj V. R., NetApp

# 30 second overview

- Wanted an automatically balanced, simple performant distributed filesystem

- Reuse our existing filesystem technology
  - WAFL filesystem, ONTAP clusters
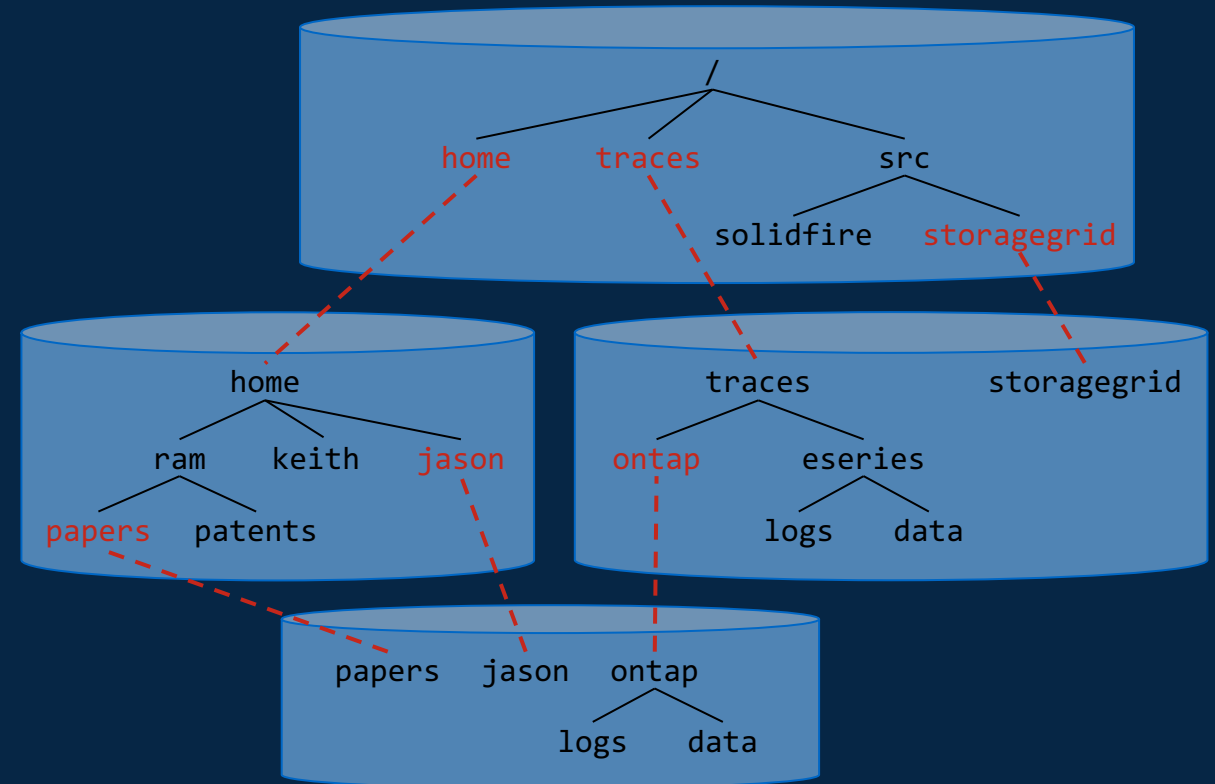  - Highly reliable local nodes

**NetApp**

# 30 second overview

- Wanted an automatically balanced, simple performant distributed filesystem

- Reuse our existing filesystem technology
  - WAFL filesystem, ONTAP clusters
  - Highly reliable local nodes

- FlexGroups: a distributed filesystem that seamlessly fuses WAFL volumes with automatic placement:
  - Remote links stich together filesystems from multiple nodes
  - Heuristics to keep them balanced

# Agenda Slide

1) Introduction:
   - Background, Problems, Requirements
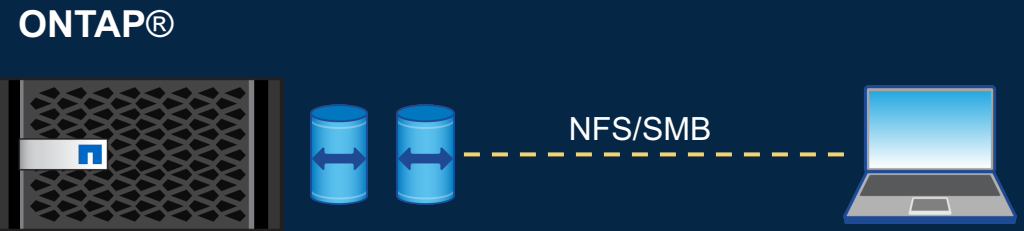
2) FlexGroups Design

3) Evaluation
   - Micro/macro benchmarks + customer experience

4) Conclusions + refs

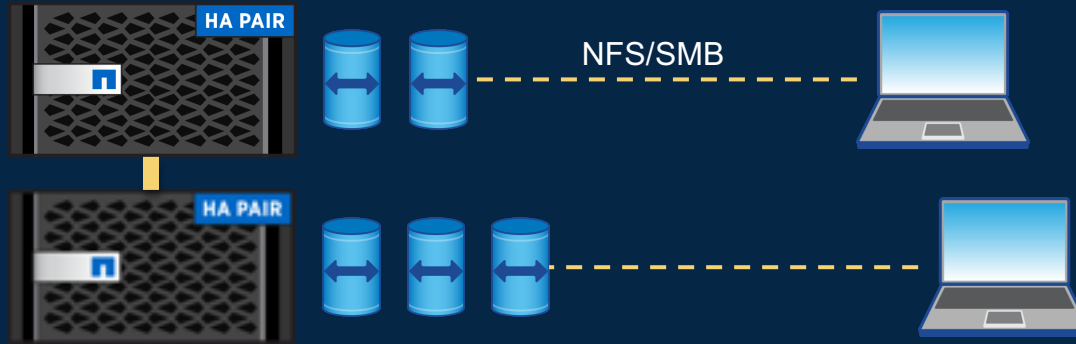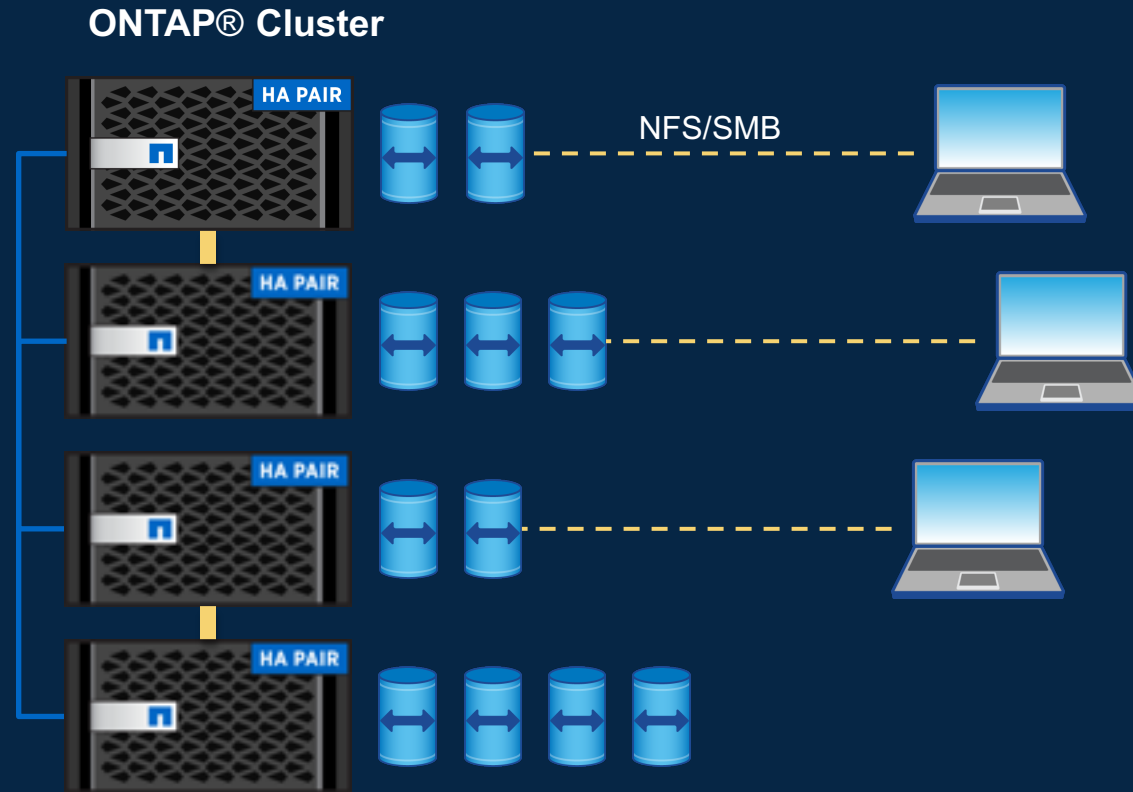**NetApp**

# Introduction

NetApp

# NetApp brief history

**ONTAP®**

NFS/SMB

**NetApp**

# NetApp brief history

**ONTAP® Cluster**

NFS/SMB

# NetApp brief history

ONTAP® Cluster

NFS/SMB

# NetApp brief history
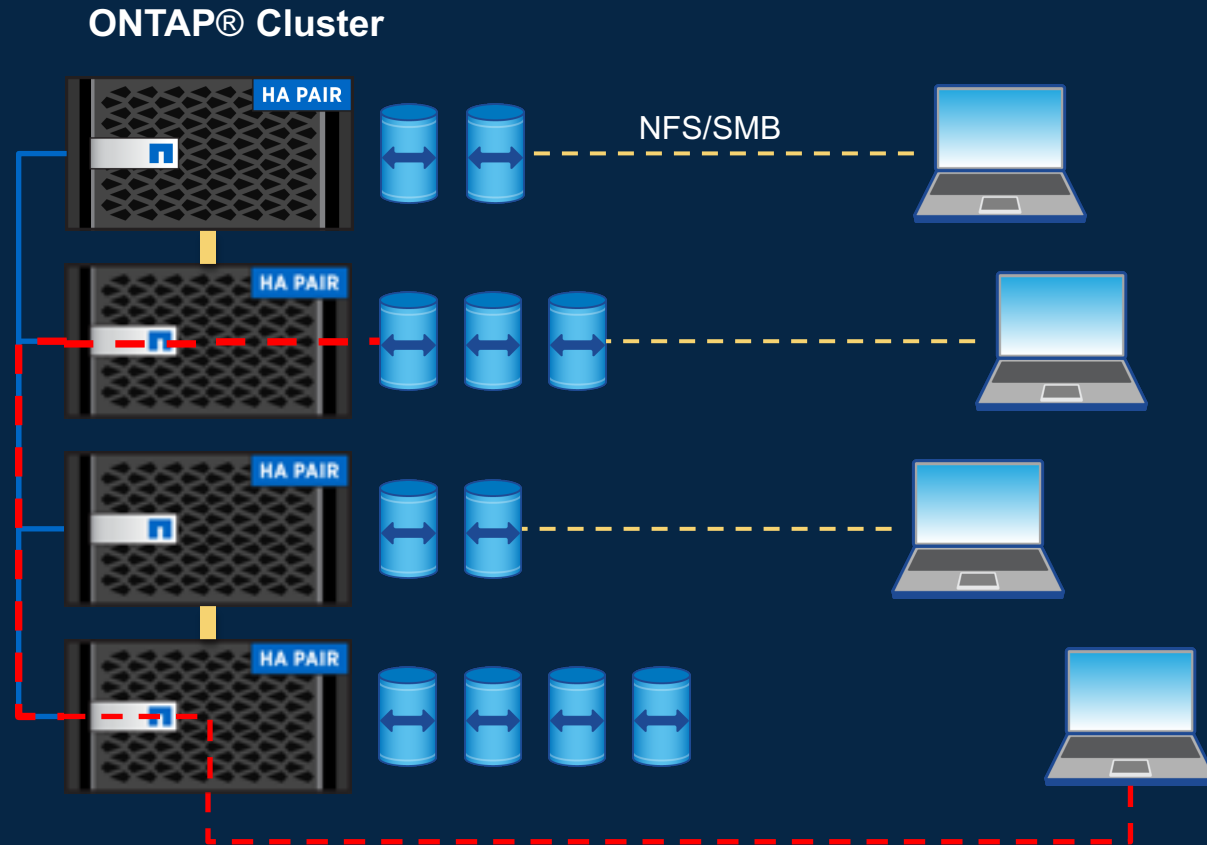


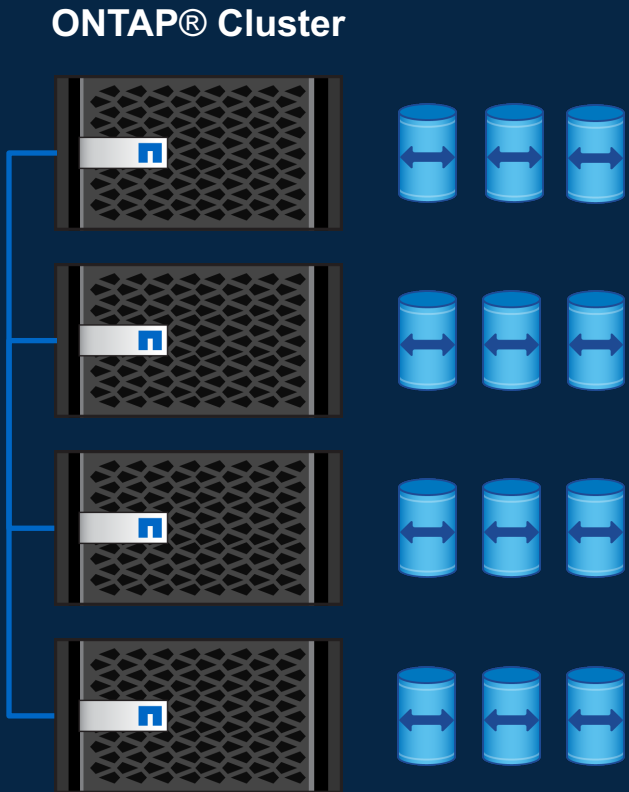ONTAP® Cluster

NFS/SMB

NetApp

# Problem

- Load balancing was manual

- Want to simplify/automate data placement

- Have fast, reliable, local technologies (like WAFL filesystem & ONTAP Clusters) with a lot of customer experience

**NetApp**

# Requirements

- Dynamic load balancing

- Ease of management

- Low latency

- Metadata scales with cluster

**■ NetApp**

# Solution: FlexGroups

**ONTAP® Cluster**

- FlexGroup internally composed of automatically distributed individual WAFL volumes

- Simple for:
  - **admins**: automatic load balance
  - **clients**: single volume

- FlexGroups are implemented with:
  - Data linked across volumes with remote link
  - Heuristics at ingest balance across nodes

# Solution: FlexGroups

**ONTAP® Cluster**



- FlexGroup internally composed of automatically distributed individual WAFL volumes

- Simple for:
  - **admins**: automatic load balance
  - **clients**: single volume

- FlexGroups are implemented with:
  - Data linked across volumes with remote link
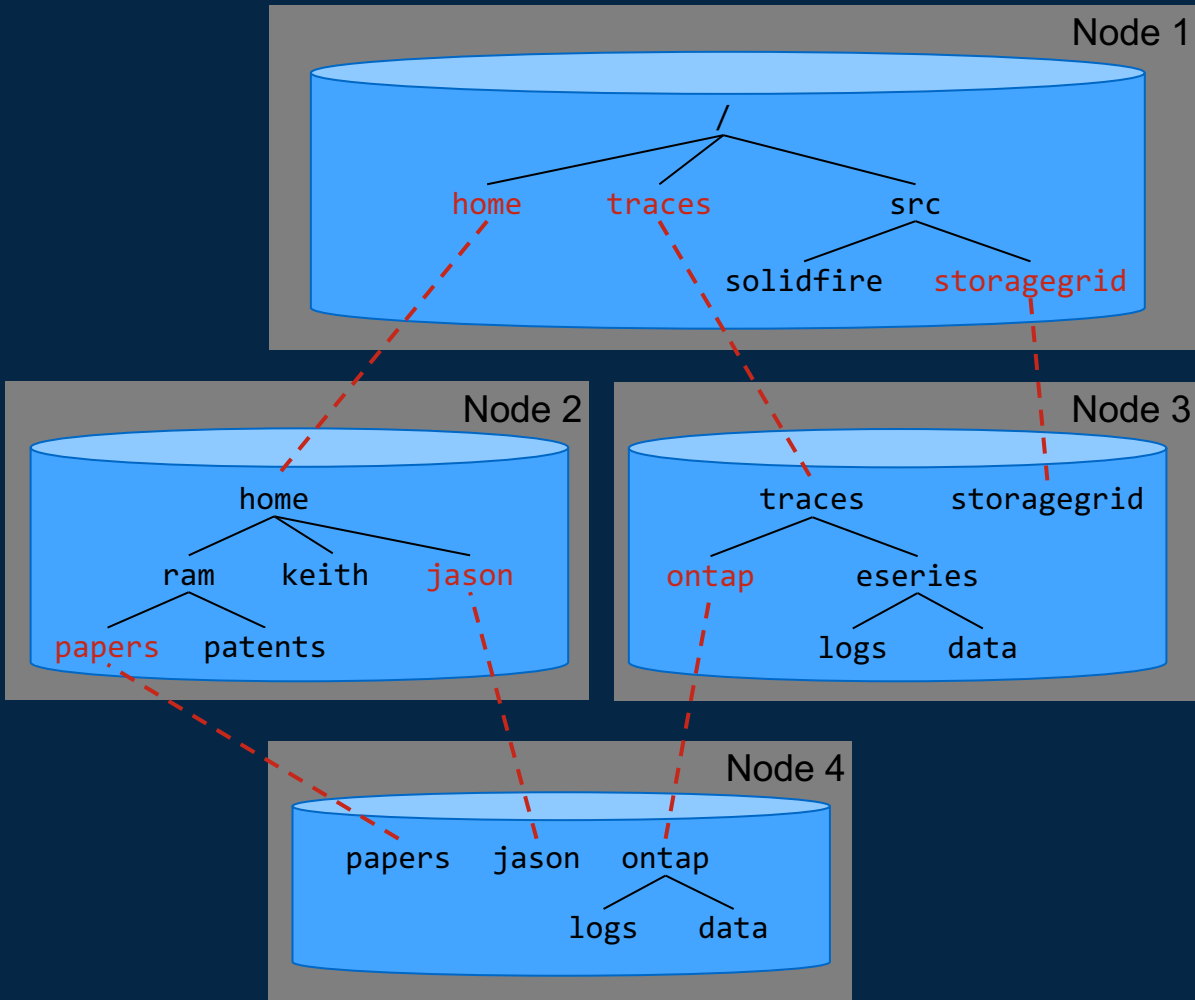  - Heuristics at ingest balance across nodes

**NetApp**

# Design of FlexGroups

**NetApp**

# Overview



- Remote links + Remote Access Layer (RAL)
- Heuristics

Legend:

 Member Volume

name Directory       name Remote Link

# Remote Access Layer (RAL)

- Problem: How to coordinate operations across nodes?

- Answer: Delegations!

- RAL Cache: remote inodes are cached in a local inode
  - Cached inodes stored persistently and crash consistently in local filesystem
  - Writeback cache with delegations
  - Takes advantage of existing failure filesystem recovery

**NetApp**

Remote Access Layer (RAL) Cache: New remote link to dir on Vol C
Step 1: create remote inode

# Remote Access Layer (RAL) Cache
## Step 2: delegate to Vol B



Vol A

Inode 1 "/"

.
..
dir1
file1

Inode 12

000101011

Vol B

Inode 2

.
..

Inode 25
cached

L2R: 25 → C:3

R2L: C:3 → 25

Vol C

Inode 3

L2R: 3 → B

Node 1

Node 2

Node 3

NetApp

# Remote Access Layer (RAL) Cache
# Step 3: Create remote link

# Remote Access Layer (RAL) Cache
## Step 4: Flush cached state back to origin

# Heuristics for data placement

- Problems: When to allocate remote? If so, where?

- Goal: just enough remote links to balance space & current load
  - Calculated on block usage, inode usage and recent ingest load

- Decisions made at ingest only

- Each node makes independent placement decisions

- More aggressive as volume space fills

**NetApp**

# Remote allocation

## Playing the percentages

- Probability of remote allocation relative to free space

- Simplified example

- Allocation proportional to free space

Total capacity: 40TB
Member capacity: 10 TB

Allocation likelihood

| 24% | 15% | 28% | 33% |
|-----|-----|-----|-----|
| 6.5TB Free | 4TB Free | 7.5TB Free | 9TB Free |

**Total free space: 27 TB**

**NetApp**

# Evaluation

**■ NetApp**

# Evaluation

- Goals:
  - Quantify Remote Access Layer overhead
  - How well do heuristics work?

- Micro-benchmarks (mdtest)

- Macro-benchmarks (SFS)

- Customer experiences: space balancing in the real world

**NetApp**

# Microbenchmarks

## NFS Operations

■ Always Local, No RAL  ■ FlexGroup  ■ Always Remote, No RAL



Normalized Latency

| Operation | Always Local, No RAL | FlexGroup | Always Remote, No RAL |
|---|---|---|---|
| ACCESS | 1.00 | 2.13 | 2.88 |
| GETATTR | 1.00 | 2.29 | 3.29 |
| READ-64KB | 1.00 | 1.45 | 1.85 |
| WRITE-64KB | 1.00 | 1.43 | 1.48 |
| LOOKUP | 1.00 | 4.25 | 3.00 |
| CREATE | 1.00 | 2.25 | 1.71 |
| MKDIR | 1.00 | 2.60 | 1.64 |
| REMOVE | 1.00 | 3.52 | 1.74 |
| RMDIR | 1.00 | 2.38 | 1.49 |

- RAL overhead when creating RO or RW caches

- No RAL overhead otherwise

■■ NetApp

# Microbenchmarks

## NFS Operations

■ Always Local, No RAL    ■ FlexGroup    ■ Always Remote, No RAL



| | ACCESS | GETATTR | READ-64KB | WRITE-64KB | LOOKUP | CREATE | MKDIR | REMOVE | RMDIR |
|---|---|---|---|---|---|---|---|---|---|
| Always Local, No RAL | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 |
| FlexGroup | 2.13 | 2.29 | 1.45 | 1.43 | 4.25 | 2.25 | 2.60 | 3.52 | 2.38 |
| Always Remote, No RAL | 2.88 | 3.29 | 1.85 | 1.48 | 3.00 | 1.71 | 1.64 | 1.74 | 1.49 |

Normalized Latency

No RAL (ACCESS, GETATTR)

Data Operations (No RAL) (READ-64KB, WRITE-64KB)
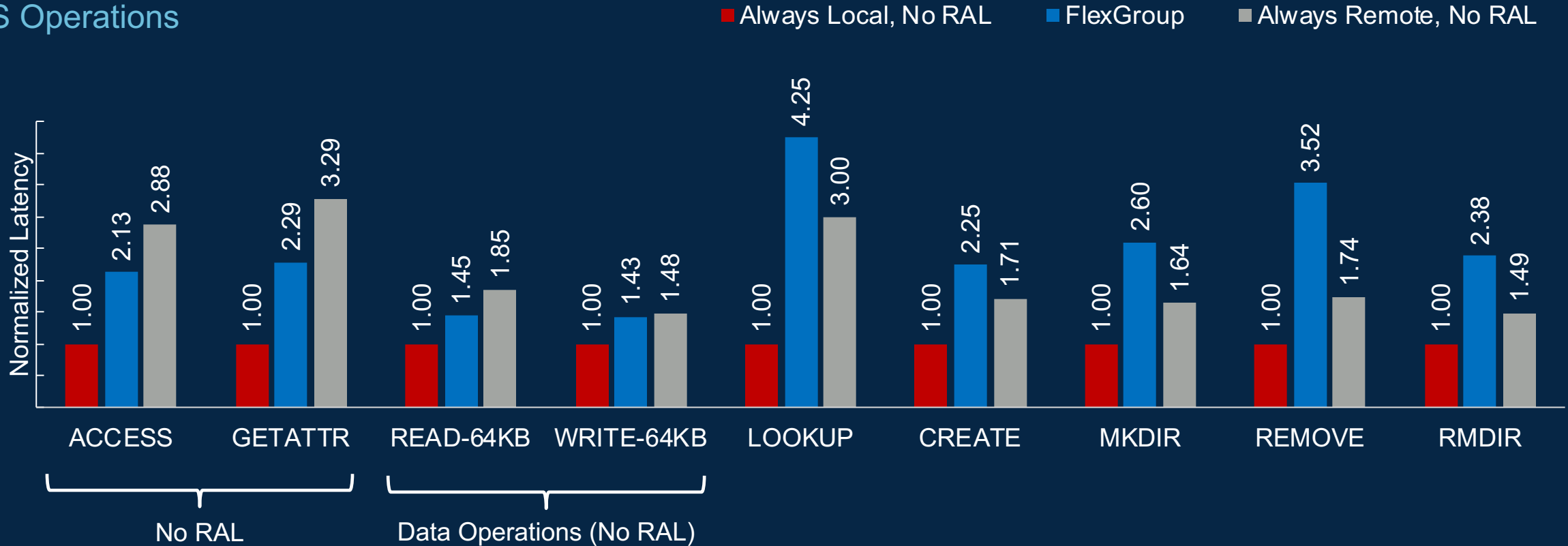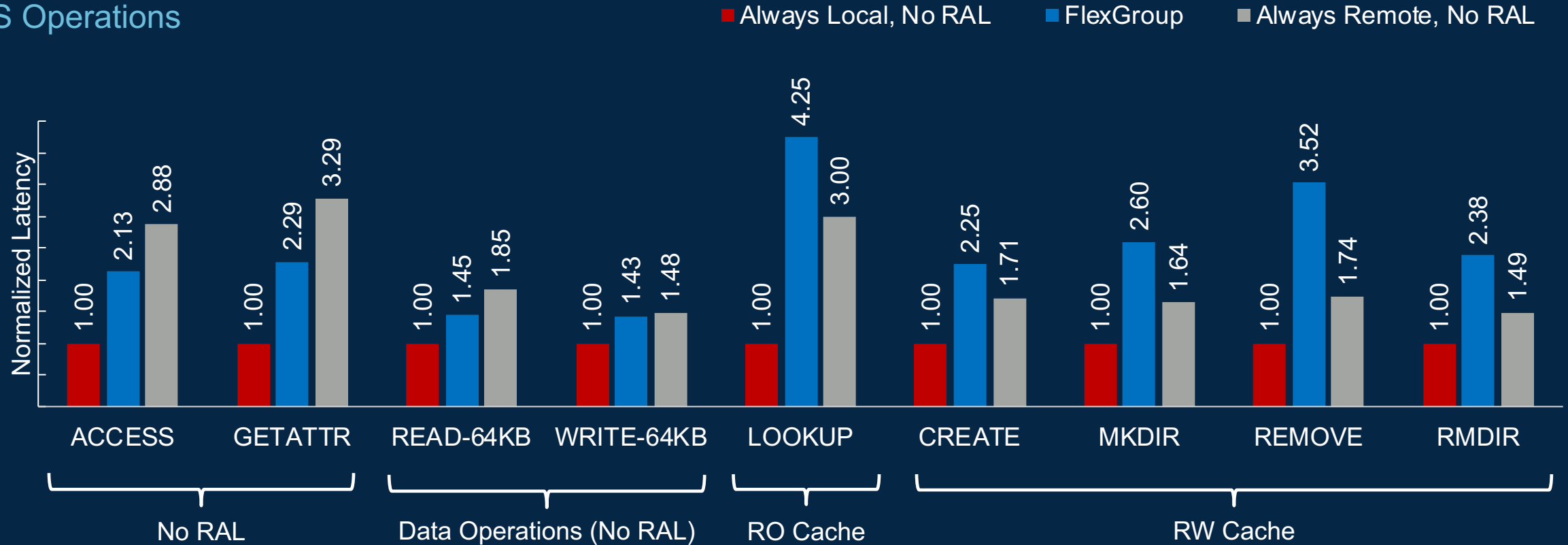
- RAL overhead when creating RO or RW caches

- No RAL overhead otherwise

NetApp

# Microbenchmarks

## NFS Operations



- RAL overhead when creating RO or RW caches
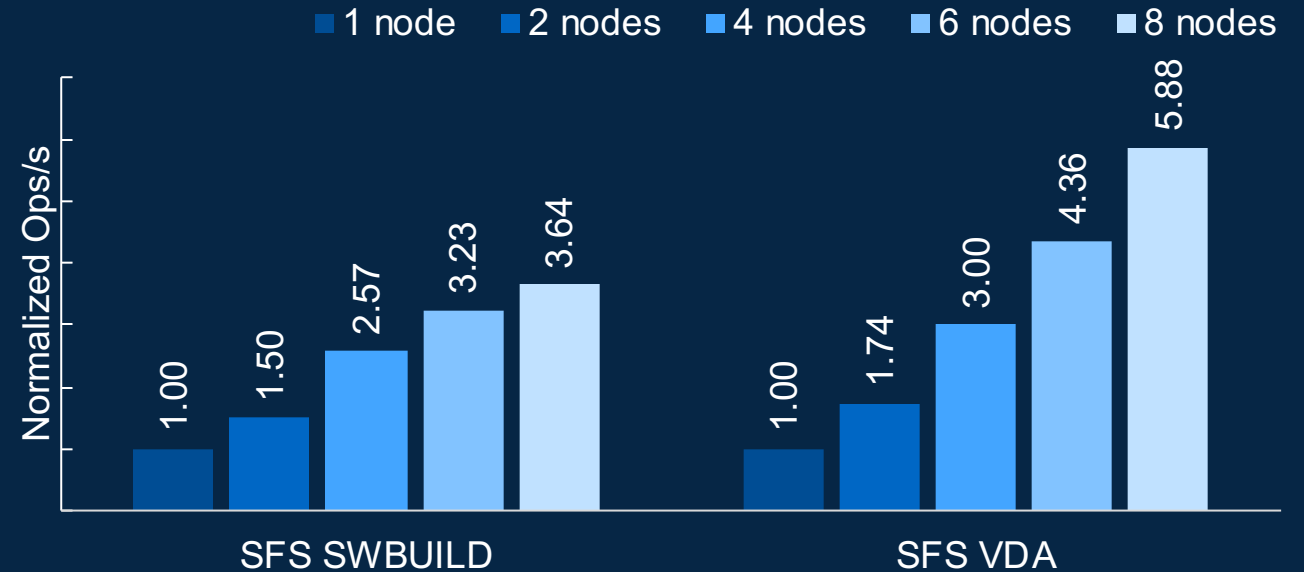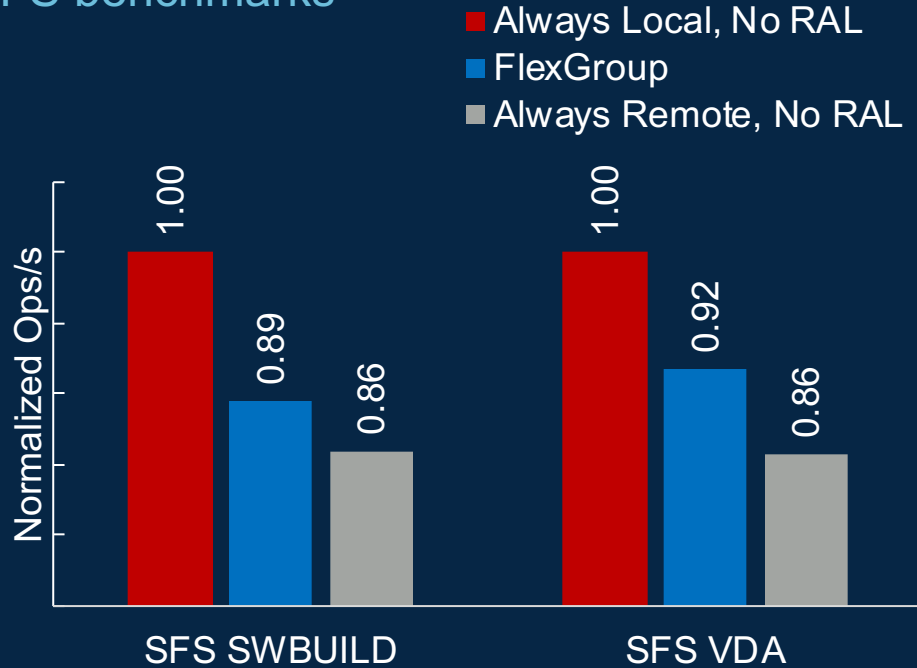
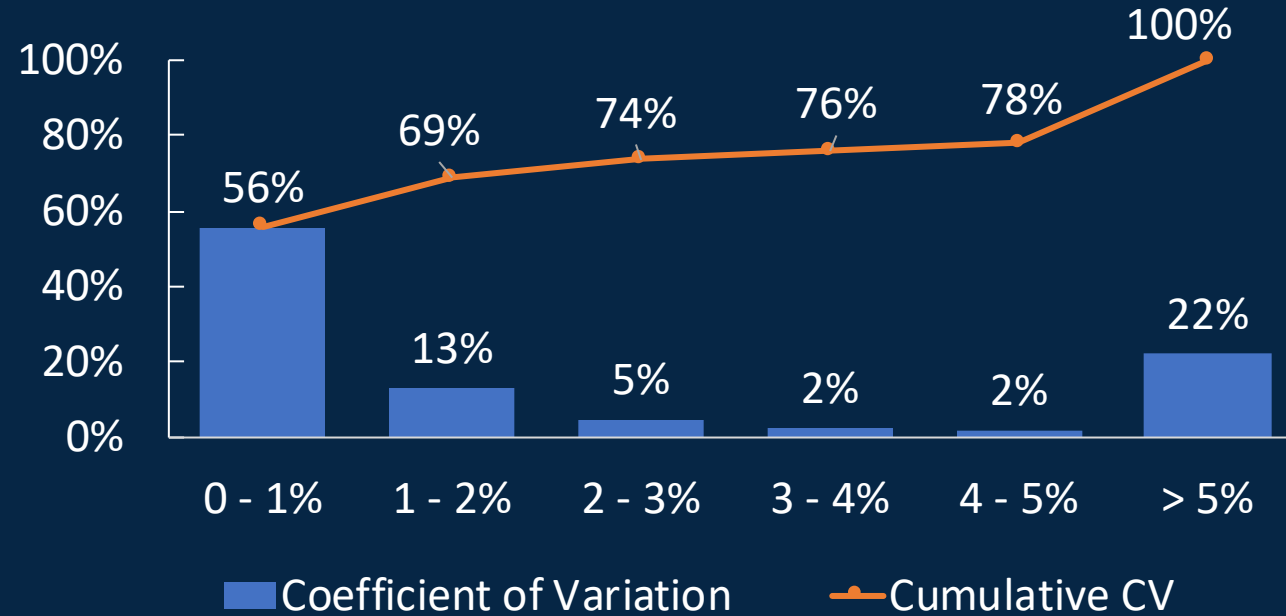- No RAL overhead otherwise

# Macrobenchmarks

## SFS benchmarks

**Legend (left chart):**
- ■ Always Local, No RAL (red)
- ■ FlexGroup (blue)
- ■ Always Remote, No RAL (gray)

**Left chart — Normalized Ops/s**

| | Always Local, No RAL | FlexGroup | Always Remote, No RAL |
|---|---|---|---|
| SFS SWBUILD | 1.00 | 0.89 | 0.86 |
| SFS VDA | 1.00 | 0.92 | 0.86 |

**Legend (right chart):** ■ 1 node ■ 2 nodes ■ 4 nodes ■ 6 nodes ■ 8 nodes

**Right chart — Normalized Ops/s**

| | 1 node | 2 nodes | 4 nodes | 6 nodes | 8 nodes |
|---|---|---|---|---|---|
| SFS SWBUILD | 1.00 | 1.50 | 2.57 | 3.23 | 3.64 |
| SFS VDA | 1.00 | 1.74 | 3.00 | 4.36 | 5.88 |

- ▪ Effective balancing of load, minimal RAL overhead

- ▪ Scalability by cluster size

**NetApp**

# Customer experiences
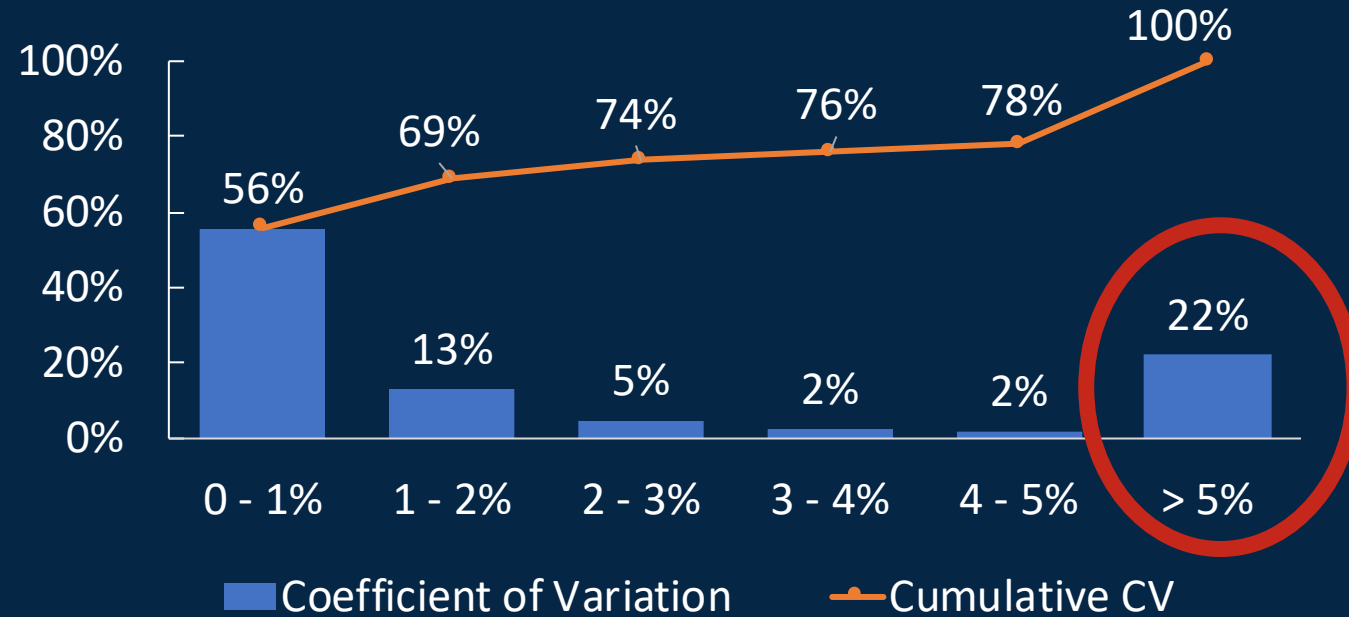
Achieved balancing

$$COV = \frac{stdev}{mean}$$



- 1000s of FlexGroups, 100s of PB deployed in the field

- Data well-balanced for most customers, with two exceptions:
  - Adding new member volumes to increase capacity
  - FlexGroups with small numbers of very large files

  NetApp

# Customer experiences

## Achieved balancing

$$COV = \frac{stdev}{mean}$$



- 1000s of FlexGroups, 100s of PB deployed in the field

- Data well-balanced for most customers, with two exceptions:
  - Adding new member volumes to increase capacity
  - FlexGroups with small numbers of very large files

 **■ NetApp**

# Conclusions

- Created distributed filesystem from existing technologies

- A reliable local filesystem changes design

- More papers on ONTAP/WAFL: https://atg.netapp.com/?tag=ontap

**■ NetApp**