# To Waffinity and Beyond:
## A Scalable Architecture for Incremental Parallelization of File System Code

Matthew Curtis-Maury, PhD
Vinay Devadas, PhD
Vania Fang
Aditya Kulkarni

NetApp, Inc

**NetApp**

# Background

- Data ONTAP is a storage operating system

- WAFL File System processes operations in the form of messages

- Competitive performance requires CPU scaling
  - WAFL is millions of lines of complicated code
  - A pure locking model is impractical
  - Many other techniques in the literature
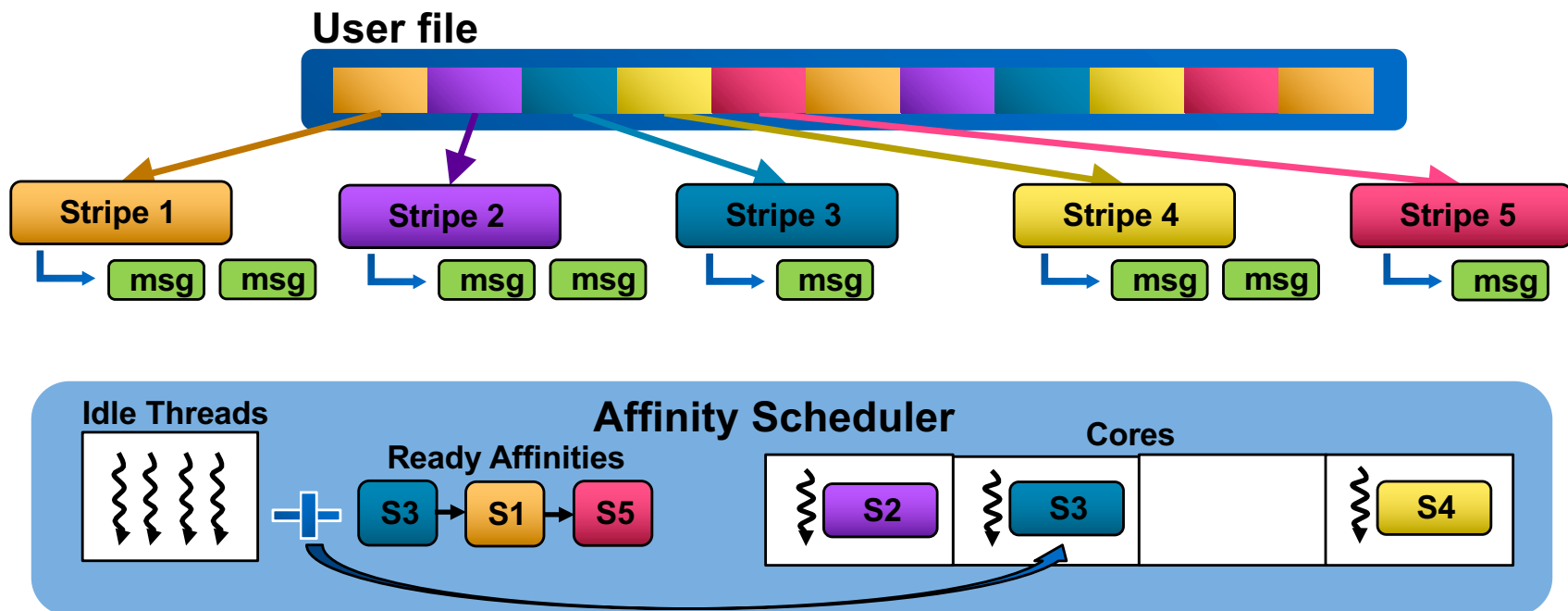    - Barrelfish, fos, Corey, Multikernel, …
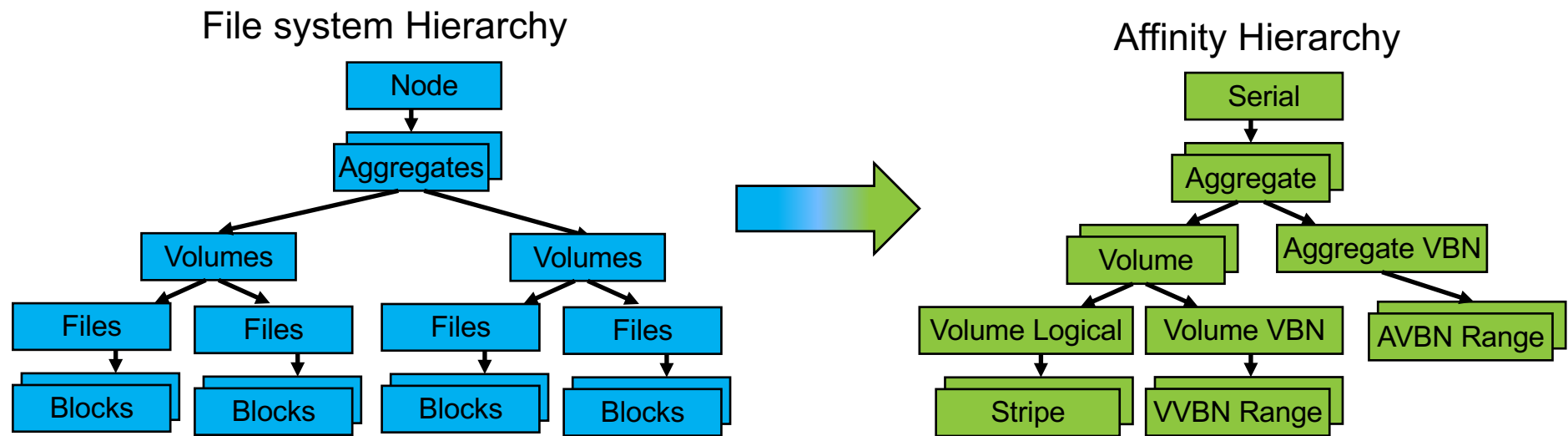
**NetApp**

# WAFL Parallelization Overview

- In the beginning… WAFL processed all messages sequentially

- WAFL parallelism leverages data partitioning

- Set of techniques to allow incremental parallelization
  - Classical Waffinity – Partition user files into chunks
  - Hierarchical Waffinity – Partition many FS data structures
  - Hybrid Waffinity – Add locking *within* the data partition framework

- These techniques have been implemented in our production OS and deployed on >200K systems

**n NetApp**

# Classical Waffinity (2006)

- Partition user files into fixed-size chunks called *file stripes*
- Rotated over a set of message queues called *Stripe affinities*
- Affinity scheduler dynamically assigns affinities to threads
- Include a *Serial affinity* to process work outside of file stripes

**User file**

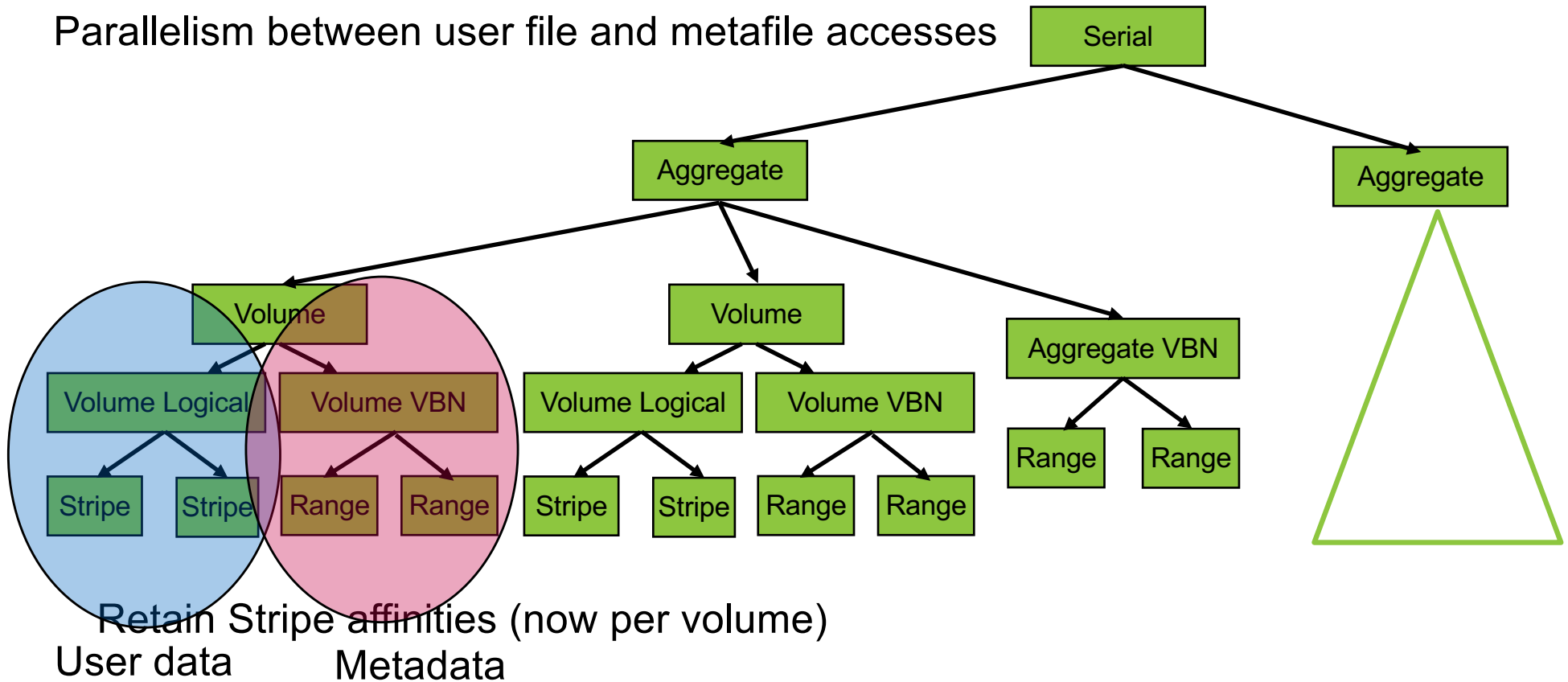| Stripe 1 | Stripe 2 | Stripe 3 | Stripe 4 | Stripe 5 |
|----------|----------|----------|----------|----------|
| msg msg | msg msg | msg | msg msg | msg |

**Affinity Scheduler**

**Idle Threads**

**Ready Affinities**  S3 → S1 → S5

**Cores**  S2   S3   S4

**NetApp**

# Hierarchical Waffinity (2011)

File system Hierarchy

```
Node
  ↓
Aggregates
  ↙        ↘
Volumes        Volumes
 ↓    ↘        ↙    ↘
Files  Files  Files  Files
 ↓      ↓      ↓      ↓
Blocks Blocks Blocks Blocks
```

Affinity Hierarchy

```
Serial
  ↓
Aggregate
  ↙        ↘
Volume        Aggregate VBN
 ↓    ↘            ↓
Volume Logical  Volume VBN   AVBN Range
 ↓              ↓
Stripe        VVBN Range
```

- Hierarchical data partitioning to match hierarchical data
  - Particular shape fine-tuned for WAFL
  - Hierarchical permissions / exclusion
- Allows parallelization of work that used to run in Serial affinity
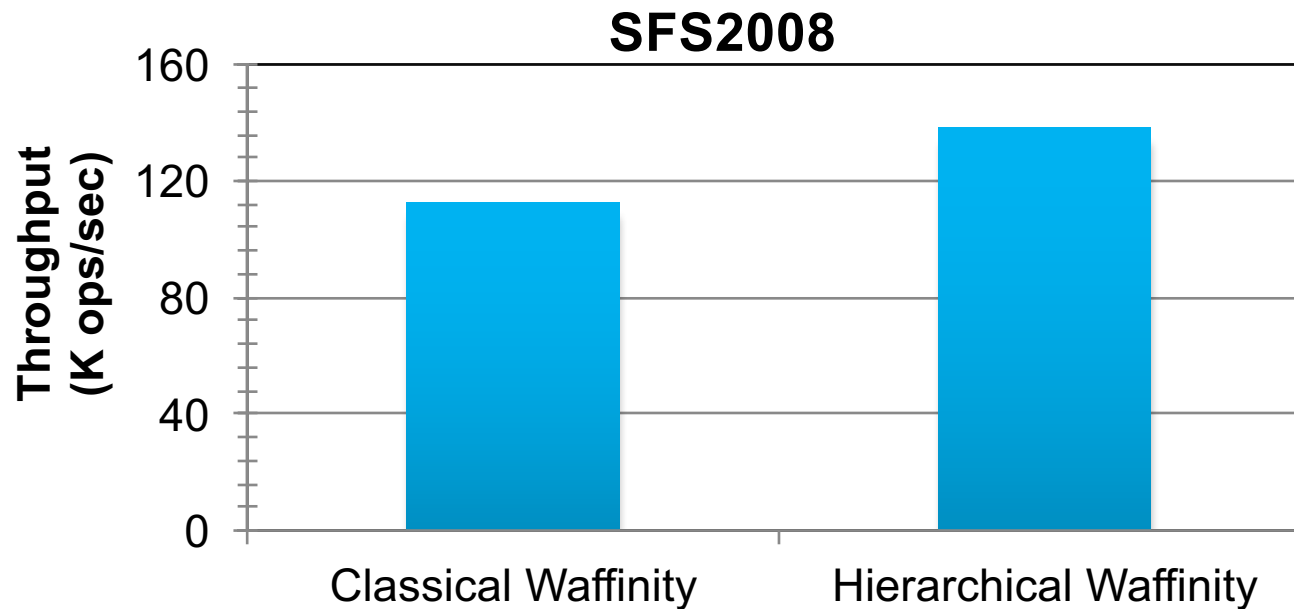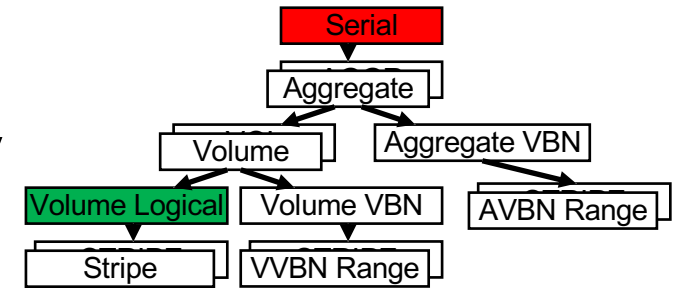- Friendly to incremental parallelization

**NetApp**

# Hierarchical Waffinity – Data mappings

Parallelism between different volumes and aggregates
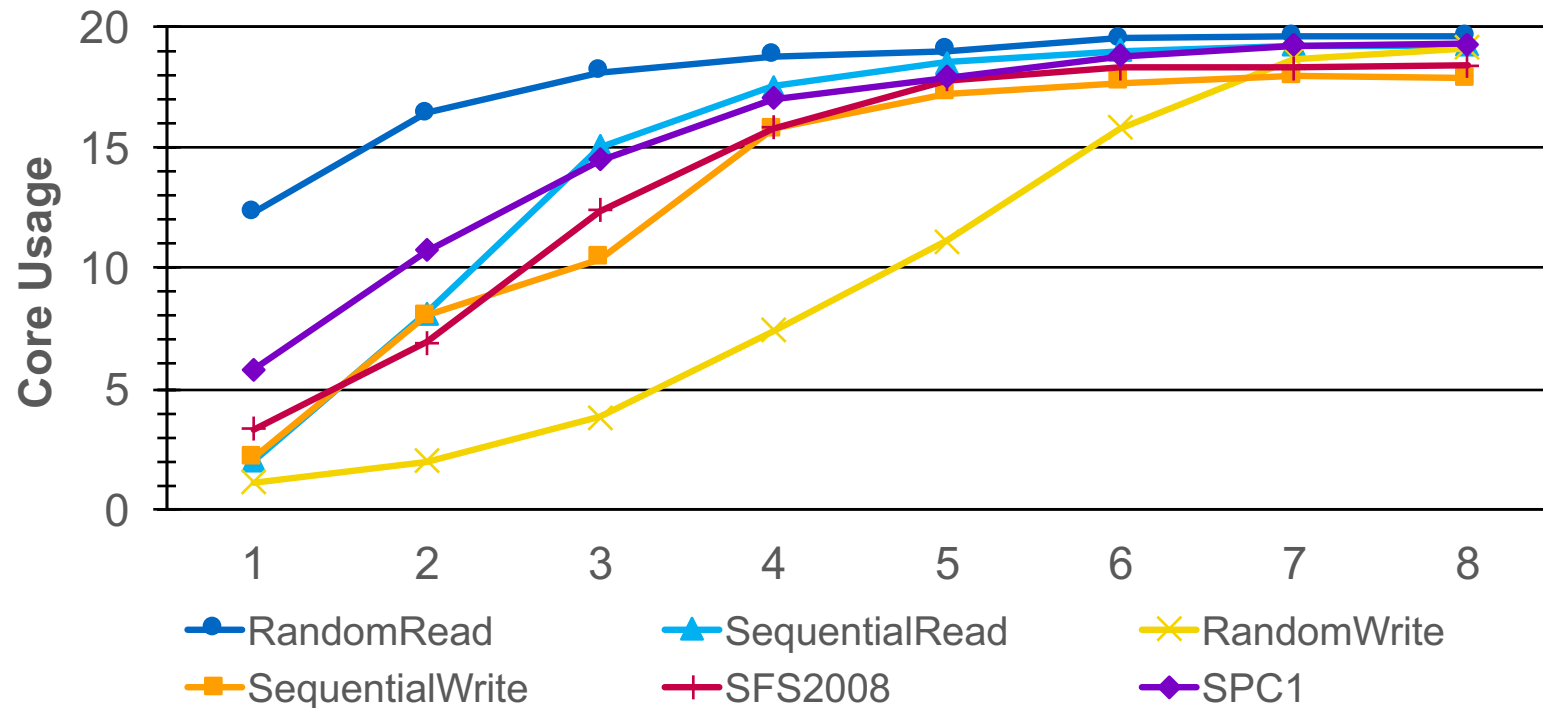
Parallelism between user file and metafile accesses



Retain Stripe affinities (now per volume)

User data

Metadata

**NetApp**

# Classical vs. Hierarchical Waffinity



**SFS2008**



- SFS2008 contains metadata operations (Create, Remove, etc)
  - Classical Waffinity: Ran in Serial affinity (48% of wallclock time)
  - Hierarchical Waffinity allows the messages to run in Volume Logical
- **~3 additional cores** used translated into a **23% throughput increase**

**■ NetApp**

# Hierarchical Waffinity CPU Scaling



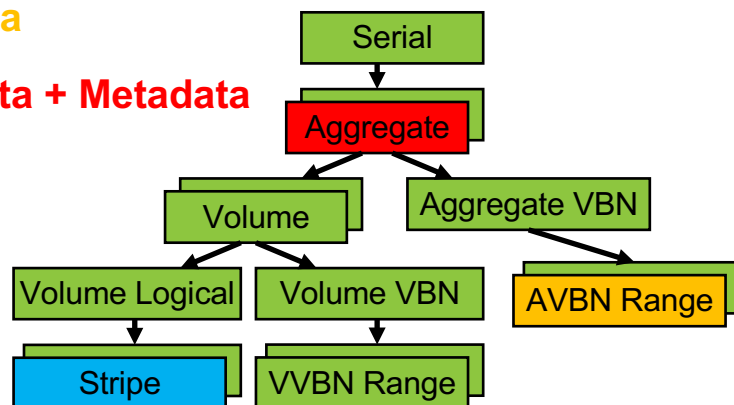- 95% average core occupancy across 6 key workloads

# Hybrid Waffinity (2016)

- ## Some important workloads access two different file blocks
  - Mappings optimized for traditional cases not well-suited here

- ## Hybrid Waffinity combines partitioning with fine-grained locking
  - **Particular blocks** are protected with locking from multiple affinities
  - Continues to allow **incremental** development

**User Data**

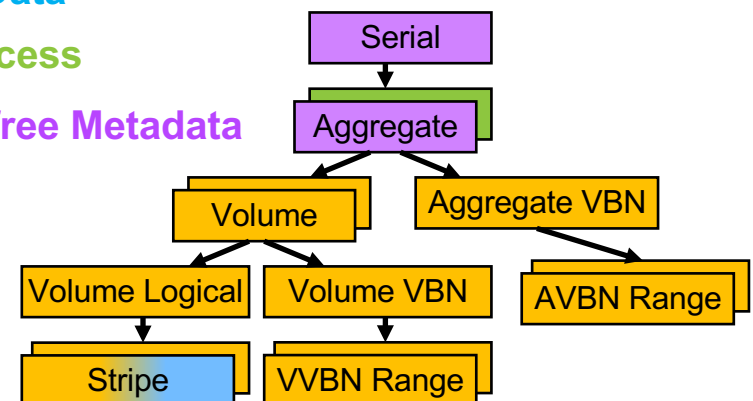**Metadata**

**User Data + Metadata**
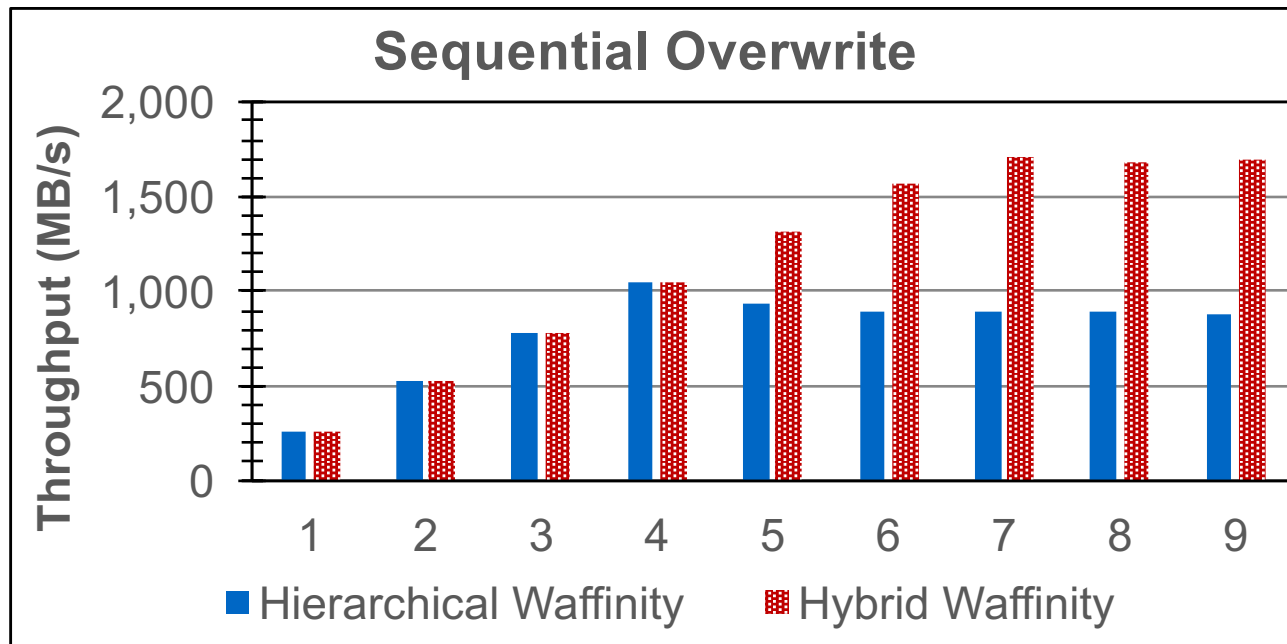
**Metadata with locking**

**User Data**

**No Access**

**Lock-free Metadata**

**NetApp**

# Hybrid vs. Hierarchical Waffinity



- Block free operations in Volume VBN for two metafile accesses
- Hybrid Waffinity parallelizes it further into VVBN Range
- **6 additional cores** translated into a **91% throughput increase**

**NetApp**

# Conclusion

- Developed a set of techniques to allow incremental parallelization of the WAFL file system
    - Focused on data partitioning
    - Selectively added in locking in a restricted way


- Provided insight into the internals of WAFL

**NetApp**

# Thank you.

**NetApp**

# History of Parallelism in ONTAP

- Data ONTAP was created for single-CPU systems of 1994

- Parallelism via "Coarse-grained Symmetric Multi-processing"
  - Each subsystem was assigned to a single-threaded *domain*
  - Minimal explicit locking required, message passing between domains
  - Scaled to 4 cores, but all of WAFL serialized

RAID    WAFL    Network    Storage    Protocol

...

CPUs

**∩ NetApp**

# Example Scheduler State



**Affinity Hierarchy**

**Hierarchical Scheduler**

Idle threads

Runnable affinities

Running affinity    Runnable affinity    Excluded affinity

- Scheduler keeps FIFO list of *runnable* affinities
- Threads call into Affinity scheduler for work
- Work in *coarse* affinities starves the system of runnable affinities

**NetApp**

# Example Affinity Mappings

```
                          ┌─────────────────┐
                          │     Volume      │
                          │   Remove: V     │
                          │   W: A, MD      │
                          └────────┬────────┘
                   ┌───────────────┴────────────────┐
                   ▼                                 ▼
        ┌─────────────────┐               ┌─────────────────┐
        │ Volume Logical  │               │   Volume VBN    │
        │   Remove: A     │               │   Create: MD    │
        │  W: A[100..200] │               │  W: MD[10..20]  │
        └────────┬────────┘               └────────┬────────┘
        ┌────────┼────────────┐              ┌──────┴──────┐
        ▼        ▼            ▼              ▼             ▼
   ┌─────────┐┌─────────┐┌─────────┐   ┌─────────┐  ┌─────────┐
   │ Stripe0 ││ Stripe1 ││ Stripe2 │   │  VVR0   │  │  VVR1   │
   │W: A[100]││W: A[200]││W: B[100]│   │W: MD[20]│  │W: MD[10]│
   │R: B[200]││R: A[200]││R: A[300]│   │R: MD[20]│  │R: MD[10]│
   └─────────┘└─────────┘└─────────┘   └─────────┘  └─────────┘
```
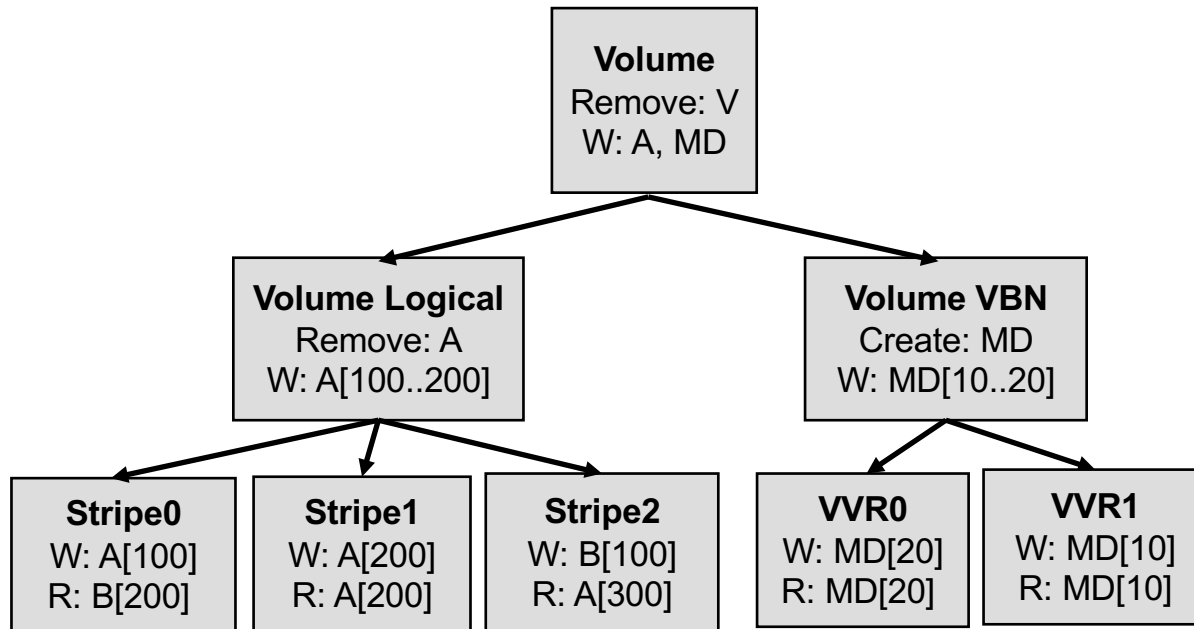
# Development Experiences

- Hierarchical Waffinity
  - Parallelization occurs at the message granularity, changed O(hundreds) LoC
  - Only parallelize critical messages, in common paths, and to a suitable affinity
  - Infrastructure required 22k LoC

- Hybrid Waffinity
  - Infrastructure for each access mode was ~3k LoC
  - Using Eject and Insert is easy, fewer than 20 lines per message optimized
  - Write involves updating and restructuring message handler -> 2k LoC
  - Now applying to Inodes with modest code changes

**■ NetApp**