# Adapting RAID technology to large heterogeneous clusters

**Toni Cortes**

Barcelona Supercomputing Center

# Special thanks to

- **José Luis Gonzalez**
  - Universitat Politècnica de Catalunya
- **Prof. Dr. Jesús Labarta**
  - Universitat Politècnica de Catalunya
  - Barcelona Supercomputing Center

# Raid 0

■ **Data distribution**

– Round-robin on all discs

■ **Advantages**

– Highest bandwidth

– Highest capacity

■ **Disadvantage**

– Not tolerant to any failure

| 0 | 1 | 2 | 3 |
|---|---|---|---|
| 0 | 1 | 2 | 3 |
| 4 | 5 | 6 | 6 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

# Raid 1

- **Data distribution**
  - Two replicated RAID0
- **Writes**
  - Done on both copies
- **Reads**
  - From any of the copies
  - Possible optimizations
- **Advantages**
  - Fault tolerant
- **Disadvantages**
  - Less parallelism
  - Wasted space

| 0 | 1 | 2 | 3 |
|---|---|---|---|
| 0 | 1 | 0 | 1 |
| 2 | 3 | 2 | 3 |
| 4 | 5 | 4 | 5 |
| 6 | 7 | 6 | 7 |

Barcelona Supercomputing Center

Storage System Research Group

# Raid 4

- **Data distribution**
  - RAID0 plus a parity discs
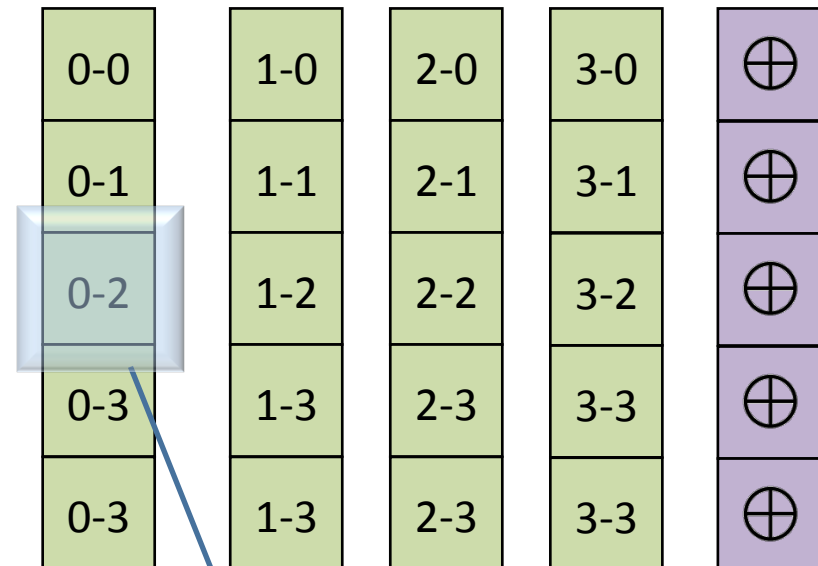  - Parity computed using XOR
- **Writes**
  - Modified data
  - Newly parity block
- **Advantages**
  - Allows one disc failure
- **Disadvantage**
  - No write parallelism
  - Always use parity disc

| | | | | |
|---|---|---|---|---|
| 0-0 | 1-0 | 2-0 | 3-0 | $\oplus$ |
| 0-1 | 1-1 | 2-1 | 3-1 | $\oplus$ |
| 0-2 | 1-2 | 2-2 | 3-2 | $\oplus$ |
| 0-3 | 1-3 | 2-3 | 3-3 | $\oplus$ |
| 0-3 | 1-3 | 2-3 | 3-3 | $\oplus$ |

Legend: **x-y** means block **x** from stripe **y**

# Terminology

- **Striping unit**
  - Block used to distribute data
- **Stripe**
  - Set of striping units that share parity computation
- **Parity block**
  - Block that keeps the "parity" of a stripe
    - Same size as a striping unit

# Raid5

- **Data distribution**
  - RAID4 plus interleaved parity
- **Advantages**
  - Good performance
  - Similar to RAID0
- **Disadvantage**
  - Only allows one disc failure
- **Slow small writes**
  - More later ☺

| | | | | |
|---|---|---|---|---|
| 0-0 | 1-0 | 2-0 | 3-0 | ⊕ |
| 0-1 | 1-1 | 2-1 | ⊕ | 3-1 |
| 0-2 | 1-2 | ⊕ | 2-2 | 3-2 |
| 0-3 | ⊕ | 1-3 | 2-3 | 3-3 |
| ⊕ | 0-3 | 1-3 | 2-3 | 3-3 |

# Parity Computation

- **Writing performance**
  - Full stripe ➜ Efficient
  - Small write ➜ Problem

- **"Small writes" implementation**
  - Read-write-modify
    - Better for "small" requests
  - Regenerate write
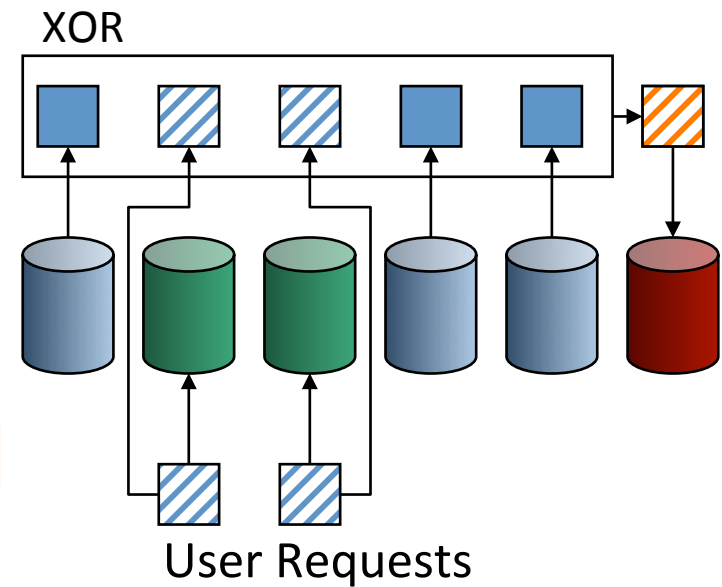    - Better for "large" requests

XOR

User Requests

# Parity Computation

■ **Writing performance**

  – Full stripe ➜ Efficient
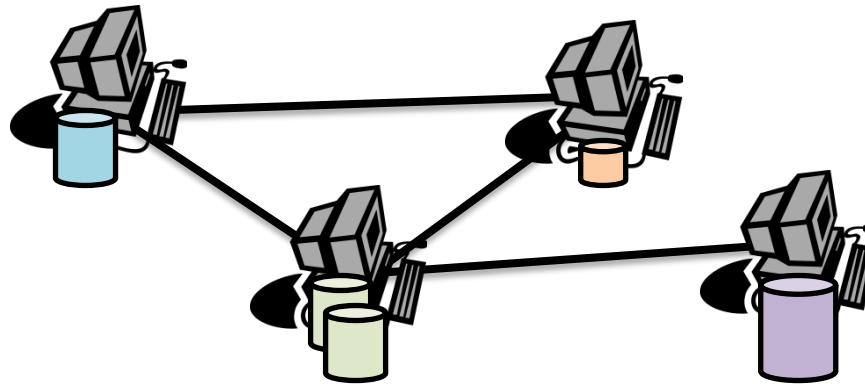  – Small write ➜ Problem

■ **"Small writes" implementation**

  – Read-write-modify
    • Better for "small" requests
  – Regenerate write
    • Better for "large" requests



XOR

User Requests

# Why is heterogeneity an issue?

- **Definition**
  - A heterogeneous set of disks is a set of disks with **different performance and capacity characteristics**



- **They are becoming a common configuration**
  - Replacing an old disk
  - Adding new disks
  - Cluster build from already existing (heterogeneous) components

# Traditional solution

- **Many systems just ignore it: all disks are treated as equal**
  - The usable size of all disks is like the smallest one
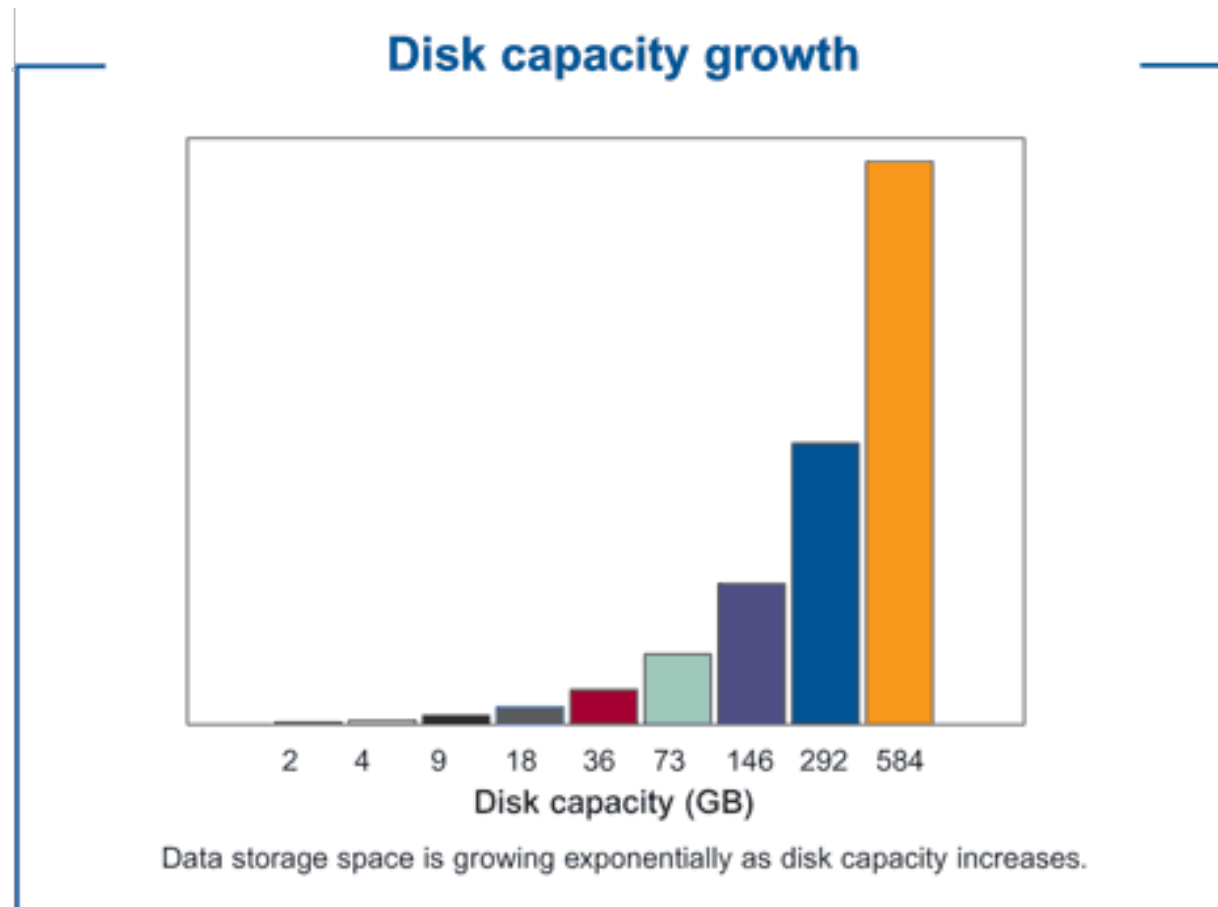  - The performance of all disks is assumed as the slowest one
- **Implications**
  - No performance gain is obtained
    - Except for some implicit side effect
  - Not all potential capacity gain is obtained
    - Some systems use the unused disk space to build a virtual disk
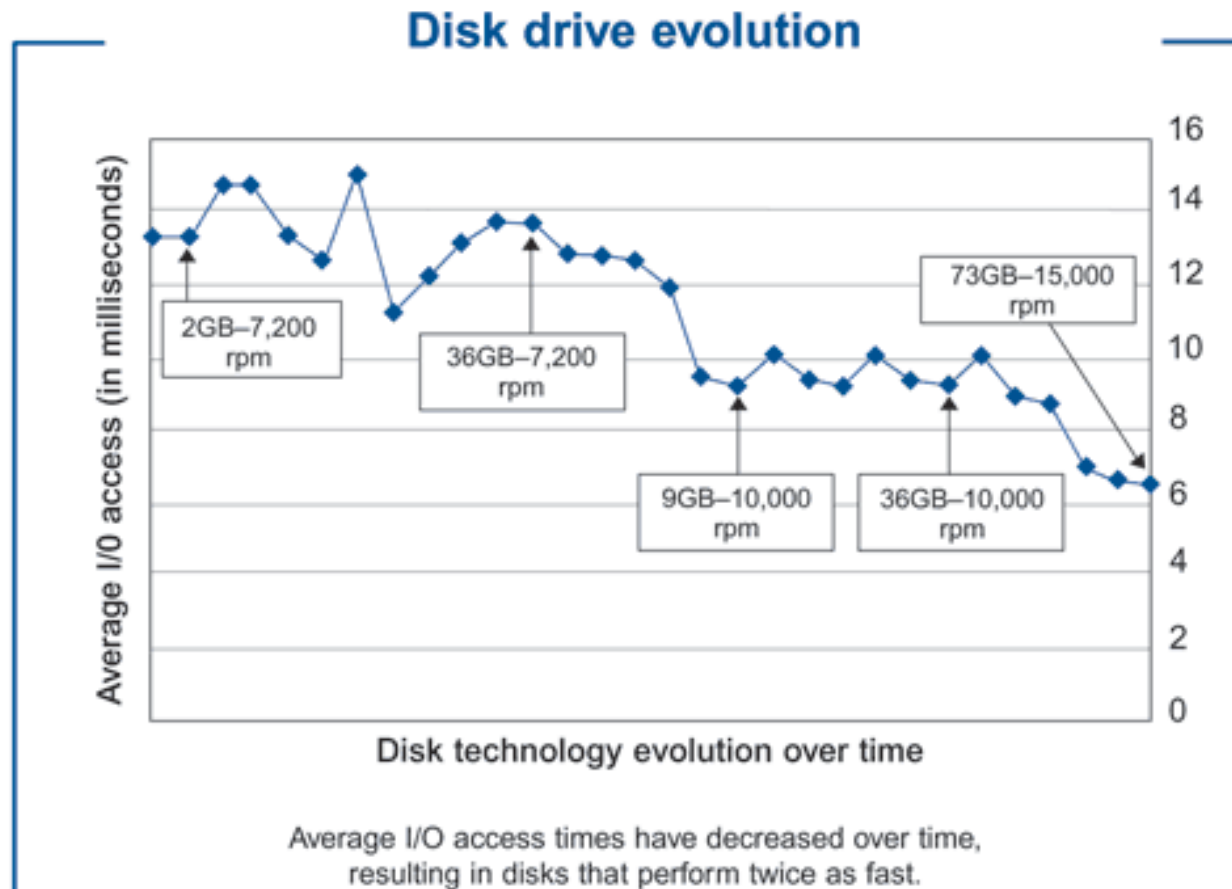- **Objective of this talk**
  - Show how to handle heterogeneous sets of disks
  - Show how to handle scalability
    - Specially how to grow storage systems with minimum overhead

# Disk capacity



**Disk capacity growth**

Data storage space is growing exponentially as disk capacity increases.

*THE DATA STORAGE EVOLUTION. Has disk capacity outgrown its usefulness?*
*by Ron Yellin (Terada magazine 2006)*

# Disk performance



THE DATA STORAGE EVOLUTION. Has disk capacity outgrown its usefulness?
by Ron Yellin (Terada magazine 2006)

# Growth storage needs

- **Information point of view**

  – *Increase of 30% each year*

    • How much information 2003?
    Peter Lyman and Hal R. Varian
    School of Information Management and Systems
    University of California at Berkeley

- **Manufacturers point of view**

  – *Increase capacity 50% each year*

    • Drive manufacturers

    • THE DATA STORAGE EVOLUTION. Has disk capacity outgrown its usefulness?
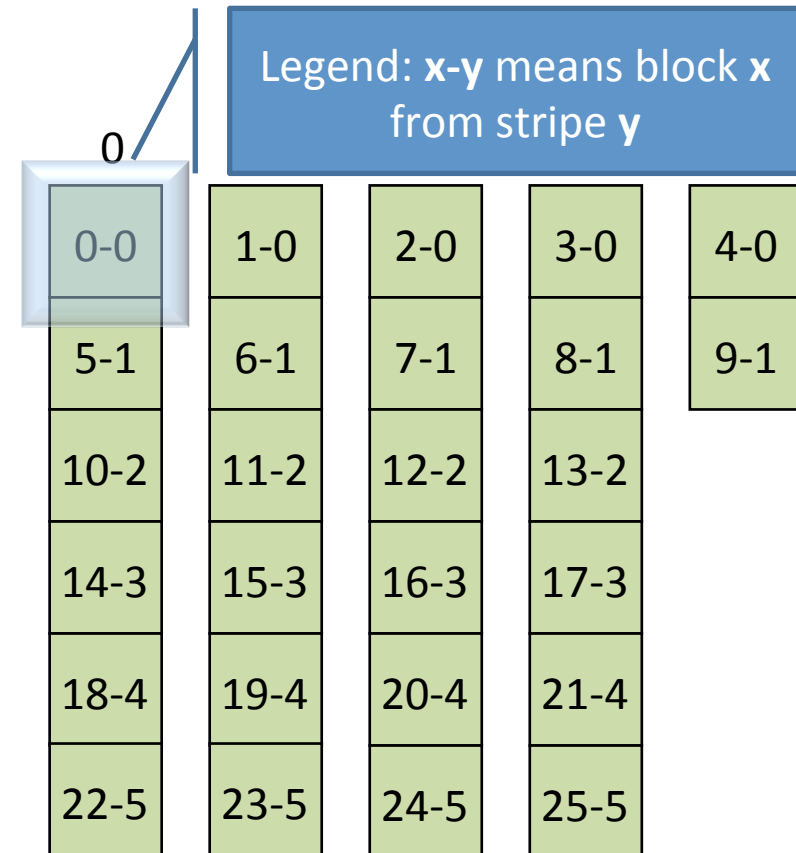    by Ron Yellin
    Terada magazine 2006

# AdaptRaid0: Intuitive idea

**Place more blocks on larger disks**
- Assumed to be faster
- Stripes with different sizes

**Problem**
- Variable parallelism
  - Long stripes in low @
  - Short stripes in high @

Legend: **x-y** means block **x** from stripe **y**

| 0 | | | | |
|---|---|---|---|---|
| 0-0 | 1-0 | 2-0 | 3-0 | 4-0 |
| 5-1 | 6-1 | 7-1 | 8-1 | 9-1 |
| 10-2 | 11-2 | 12-2 | 13-2 | |
| 14-3 | 15-3 | 16-3 | 17-3 | |
| 18-4 | 19-4 | 20-4 | 21-4 | |
| 22-5 | 23-5 | 24-5 | 25-5 | |

See T. Cortes and J. Labarta. Taking Advantage of Heterogeneity in Disk Arrays

# Reducing variance in parallelism

- **Use distribution like a repetition pattern**
  - All areas have short and long stripes
  - Hopefully, most large files will too

| 0 | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| 0-0 | 1-0 | 2-0 | 3-0 | 4-0 |
| 5-1 | 6-1 | 7-1 | 8-1 | 9-1 |
| 10-2 | 11-2 | 12-2 | 13-2 | 26-5 |
| 14-3 | 15-3 | 16-3 | 17-3 | 31-6 |
| 18-4 | 19-4 | 20-4 | 21-4 | |
| 22-5 | 23-5 | 24-5 | 25-5 | |
| 27-6 | 28-6 | 29-6 | 30-6 | |
| 32-7 | 33-7 | 34-7 | 35-7 | |
| 36-8 | 37-8 | 38-8 | 39-8 | |
| 40-9 | 41-9 | 42-9 | 43-9 | |

# AdaptRaid0 parameters

- **Utilization factor (UF)**
  - One factor per disk
    - Larger disks have more blocks?
    - Faster disks have more blocks?

- **Stripes in pattern (SIP)**
  - We define a pattern using the UF
    - Large patterns allow more requests with good disks
    - Small patterns allow a better distribution
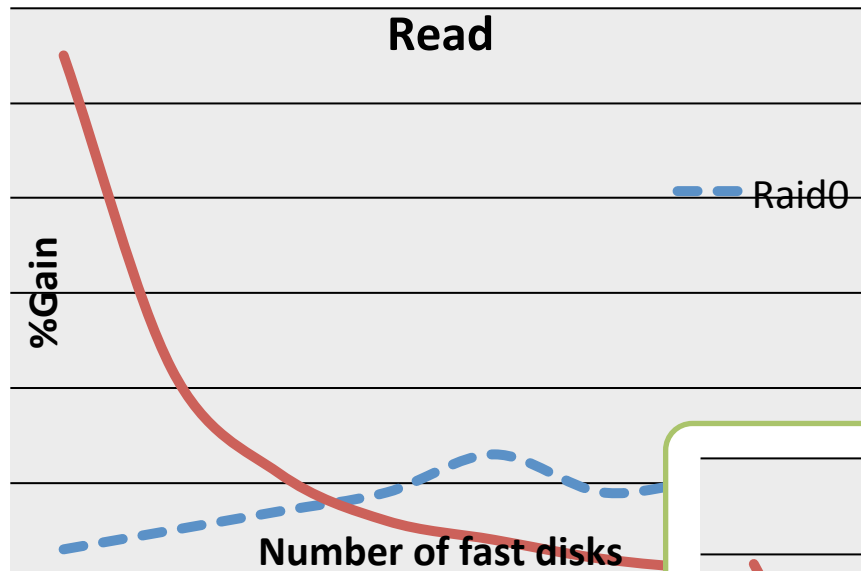
# Computing the location of a block

- **Formulas**
  - *Disk(B) =*
    **location**[B % Blks_in_a_pattern].disk
  - *Pos(B) =*
    **location**[B % Blks_in_a_pattern].pos +
    (B/Blks_in_a_pattern) *
    **Blks_per_disk_in_patern**[*Disk(B)*]

- **Metadata**
  - Location[Blks_in_a_pattern] // [SIP * DSKS]
    - int disk, pos
  - Blks_per_disk_in_patterns[DSKS]

# Performance



Using HP traces from 1999

# Parameter sensitivity

■ **Utilization factor (UF)**

    – Depends on what the administrator wants

■ **Strips in pattern (SIP)**

    – No big difference between the different values

    – The best option is SIP larger than DISKS*2

        • Measured experimentally

# AdaptRaid5: Intuitive idea

■ **Place more blocks on larger disks**
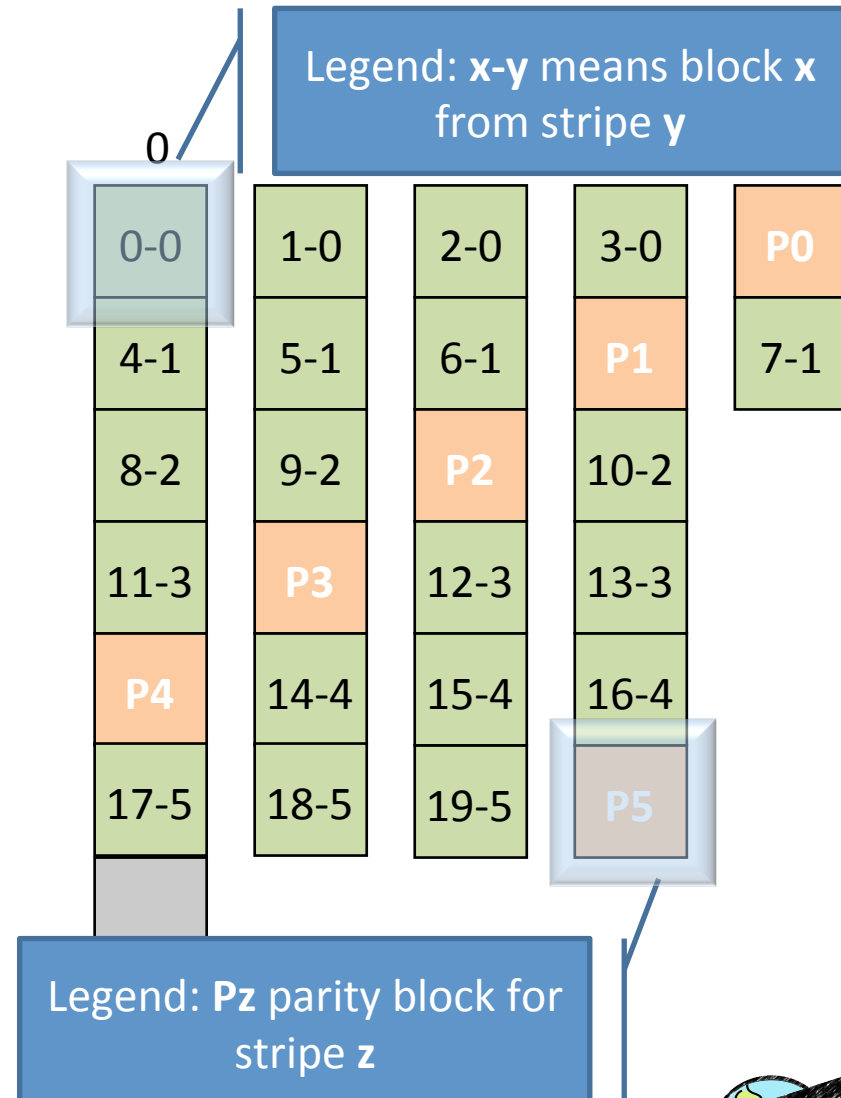  - Assumed to be faster
  - Stripes with different sizes

■ **Singularities**
  - Last block disk 0 unused
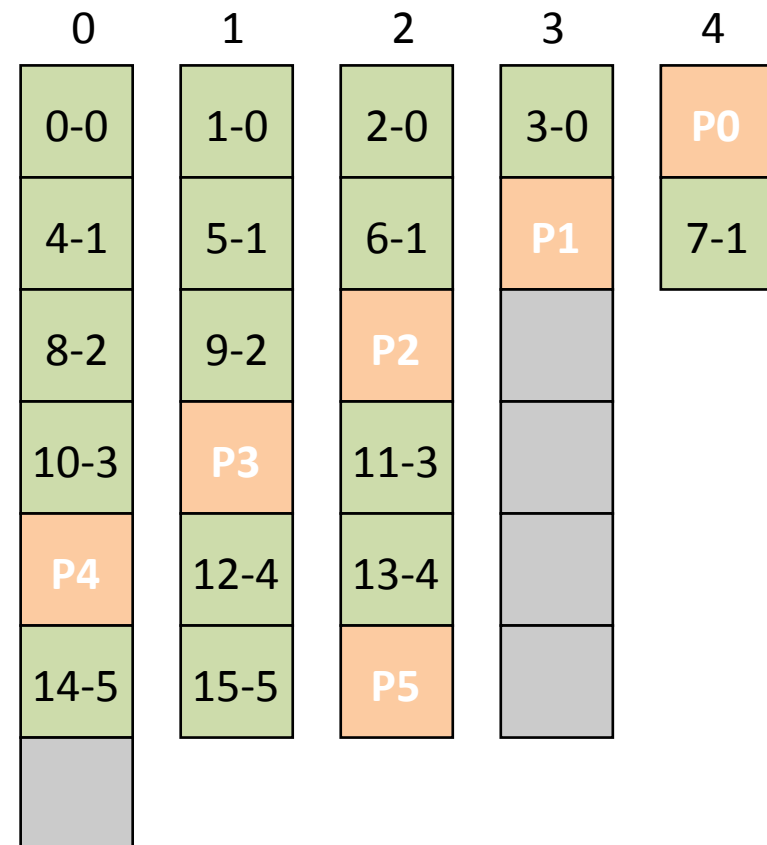    • Stripes need 2 blocks

■ **Problem**
  - Small-writes

Legend: **x-y** means block **x** from stripe **y**

| 0 | | | | |
|---|---|---|---|---|
| 0-0 | 1-0 | 2-0 | 3-0 | P0 |
| 4-1 | 5-1 | 6-1 | P1 | 7-1 |
| 8-2 | 9-2 | P2 | 10-2 | |
| 11-3 | P3 | 12-3 | 13-3 | |
| P4 | 14-4 | 15-4 | 16-4 | |
| 17-5 | 18-5 | 19-5 | P5 | |

Legend: **Pz** parity block for stripe **z**

See T. Cortes and J. Labarta. Taking Advantage of Heterogeneity in Disk Arrays

# Reducing the small-write problem

- **Define the number of data blocks per stripe**
  - Divisor of the number of data blocks in largest stripe
- **Problem**
  - Capacity wasted

| 0 | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| 0-0 | 1-0 | 2-0 | 3-0 | P0 |
| 4-1 | 5-1 | 6-1 | P1 | 7-1 |
| 8-2 | 9-2 | P2 | | |
| 10-3 | P3 | 11-3 | | |
| P4 | 12-4 | 13-4 | | |
| 14-5 | 15-5 | P5 | | |
| | | | | |

Barcelona Supercomputing Center

Storage System Research Group

# Increasing effective capacity

■ **Step 1**

– Use all disks leaving the unused block in RR way
  • Holland and Gibson 92
– Better load distribution
– No capacity is gained

■ **Step 2**

– Push blocks down "tetris-like"

■ **Problem**

– Variable parallelism

|  0  |  1  |  2  |  3  |  4  |
|-----|-----|-----|-----|-----|
| 0-0 | 1-0 | 2-0 | 3-0 | P0  |
| 4-1 | 5-1 | 6-1 | P1  | 7-1 |
|     | 8-2 | P2  | 9-2 |     |
| 10-3| P3  | 11-3| 13-4|     |
| P4  | 12-4|     | P5  |     |
| 14-5|     | 15-5|     |     |

# Increasing effective capacity

**Step 1**

- Use all disks leaving the unused block in RR way
- Better load distribution
- No capacity is gained

**Step 2**

- Push blocks down "tetris-like"
- Similar to parity declustering

**Problem**

- Variable parallelism
  - Long stripes in low @
  - Short stripes in high @

| 0 | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| 0-0 | 1-0 | 2-0 | 3-0 | P0 |
| 4-1 | 5-1 | 6-1 | P1 | 7-1 |
| 10-3 | 8-2 | P2 | 9-2 | |
| P4 | P3 | 11-3 | 13-4 | |
| 14-5 | 12-4 | 15-5 | P5 | |
| | 16-6 | P6 | 17-6 | |

# Reducing variance in parallelism

- **Use distribution like a repetition pattern**
- **Problem**
  - Minor problems with long/short stripes
  - Too detailed for this tutorial
    - Ask me if interested

| 0 | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| 0-0 | 1-0 | 2-0 | 3-0 | P0 |
| 4-1 | 5-1 | 6-1 | P1 | 7-1 |
| 10-3 | 8-2 | P2 | 9-2 | P7 |
| P4 | P3 | 11-3 | 13-4 | 23-8 |
| 14-5 | 12-4 | 15-5 | P5 | |
| 16-7 | 17-7 | 18-7 | 19-7 | |
| 20-8 | 21-8 | 22-8 | P8 | |
| 26-10 | 24-9 | P9 | 25-9 | |
| P11 | P10 | 27-10 | 28-11 | |
| 29-12 | 29-11 | 30-12 | P12 | |

# AdaptRaid5 parameters

■ **Utilization factor (UF)**

– One factor per disk

- Larger disks have more blocks?

- Faster disks have more blocks?

■ **Stripes in pattern (SIP)**

– We define a pattern using the UF

- Large patterns allow more requests with good disks

- Small patterns allow a better distribution

# Computing the location of a block

- **Formulas**
  - *S =*
    **stripe[**B % Blks_in_a_pattern**]** +
    (B / Blks_in_a_pattern) * SIP
  - *Disk(Parity of S) =*
    **parity[**S % SIP**].disk**
  - *Pos(Parity of S) =*
    **parity[**S % SIP**].pos** +
    (S/SIP) *
    **Blks_per_disk_in_patern[***Disk(Parity of S)***]**
- **Metadata**
  - Stripe[Blks_in_a_pattern] // [SIP * DSKS]
  - Parity[SIP]
    - int disk, pos

# **Performance**



**Read**

%Gain

Number of fast disks

- - - Raid5
— OnlyFast

Using HP traces from 1999

**Write**

%Gain

Number of fast disks

- - - Raid0
— OnlyFast

# Parameter sensitivity

■ **Utilization factor (UF)**

    – Depends on what the administrator wants

■ **Strips in pattern (SIP)**

    – There is big difference between the different values

        • Especially for writes

            – Up to 5 times slower with small SIPs

            – More sensible for parity distribution

    – The best option is a SIP larger than DISKS*2

        • Measured experimentally

# Recovering a disk

- **Performed in two steps**
  - Recovery of lost data (like in RAID5)
    - During this step no other disk can fail (**vulnerability window**)
  - Reorganization to improve disk usage
    - During this step a disk failure would not be fatal

- **Vulnerability window comparison**
  - In heterogeneous arrays reduced up to 30% (depending on load)
    - Disk are used better ➔ reconstruction is faster
    - Some kind of parity declustering ➔ reconstruction faster

Barcelona
Supercomputing
Center

Storage
System
Research
Group

# Overhead for disk recovery



Using Server synthetic traces [Hsu03]
C = 250GB

Cumulative service time (min)

5F-4S

5F+4S (S->R)

Raid5 vulnerability window

Assimilation time (hours)

See J. L. Gonzalez, T. Cortes Evaluating the Effects of Upgrading Heterogeneous Disk Arrays

Barcelona Supercomputing Center

Storage System Research Group

# Scaling RAID architectures

■ **Using traditional RAID architecture does not scale**

  – Including AdaptRAID

■ **Adding news disk implies reorganizing the whole data**

  – Re-striping requires the movement of all data-blocks

  – Time $t_{striping}$ for re-layout grows linear in capacity:

$$t_{striping} = k * C_{old}$$

  where $k$ is a constant and $C_{old}$ is the already stored capacity

■ **Trend**

  – Newly integrated capacity $C_{new}$ is always smaller than $C_{old}$

Barcelona
Supercomputing
Center

Storage
System
Research
Group

# How expensive is re-striping?

■ **Assumptions**

– 36 GByte of data can be re-distributed in each hour

– 100 GByte of new capacity $C_{new}$ have to added

– Already existing capacity $C_{old}$ between 100 GByte and 1 PByte

# If re-striping is not the way …

- **Objective: find a way to**
  - Only migrate the needed amount of data
  - Continue having balanced load
  - Do not lose the deterministic behavior

# AdaptiveZ: intuitive idea

- **Divide the address space in zones**
  - Created dynamically each time new disks are added
  - Each zone has its own heterogeneous stripping "policy"
- **When new disks are added, only one zone is restriped**



See J. L. Gonzalez, T. Cortes. An Adaptive Data Block Placement based on Deterministic Zones (AdaptiveZ)

**Let's assume a first array**

- 3 disks 1Gbytes each
- Striping units ➜ 128KB
  - 8192 per disk
  - 24576 in total
- **Initially we have one zone**

# Adding 2 disks: naïve way

**Adding 2 disks**

– 2 disks 2Gbytes each

– Striping units ➔ 128KB

  • 16384 per disk

  • 32768 of new storage

  • 57344 final capacity

**Create a new zone with only two disks**

– No data movement

**Problem**

– NOT balanced

  • Only new data in new disks

– No increase in parallelism

# Adding 2 disks: load balancing

- **Create two new zones**
- **Zone1: old data**
  - Theoretical minimum movement to balance load
    - Assume new disks should have 2 times more SU
    - $C_{old}*(1-C_{old}/(C_{new} + C_{old}))$
      - 14043 SU
- **Zone2: new data**
  - Uses all disks
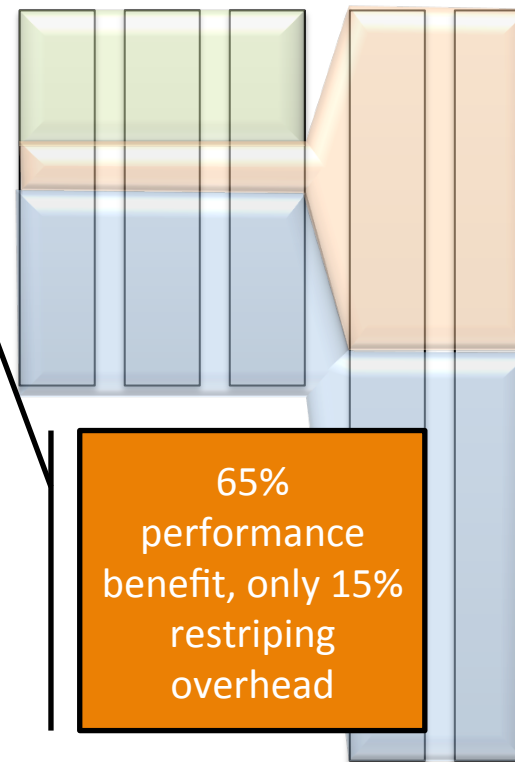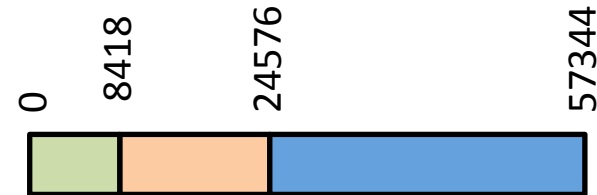  - Only for new data
- **Problem**
  - Old data looses parallelism

# Adding 2 disks: more parallelism

**Create two new zones**

**Zone1: old data**

- Increase 15% the theoretical minimum

  - *1.15 \* $C_{old}*(1-C_{old}/(C_{new} + C_{old}))$*
    - 16149 SU

- Restripe these blocks on all disks

  - Find the adequate UF for this zone to guaranty global disk UF

  - **Small compared to full restripe**
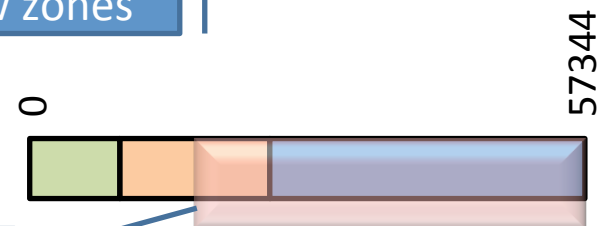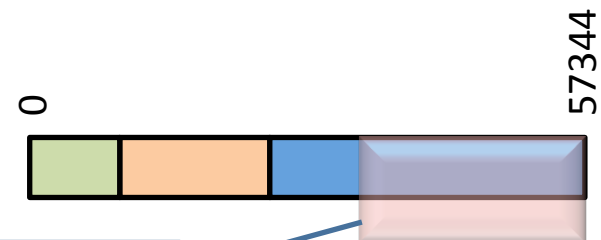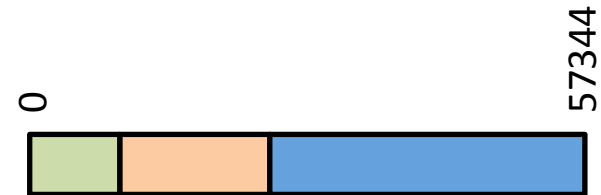
**Zone2: new data**

0  8418  24576  57344

65% performance benefit, only 15% restriping overhead

# Adding 2 more disks

- **Create two new zones**
- **Zone1: old data**
  - Size = *1.15 \* $C_{old}*(1-C_{old}/(C_{new} + C_{old}))$*
    - Depending on size
      - We create 2 zones
      - We merge zones
  - Restripe these blocks on all disks
- **Zone2: new data**

57344

0                                              57344

Creating 2 new zones

0                                              57344

Merging zones

# AdaptiveZ parameters

- **AdaptRaid parameters are used**
  - One set per zone
- **% increase on the size to restripe could be changed**
  - Our evaluation shows that 15% is a good tradeoff
  - and suggest no modifications

Barcelona
Supercomputing
Center

Storage
System
Research
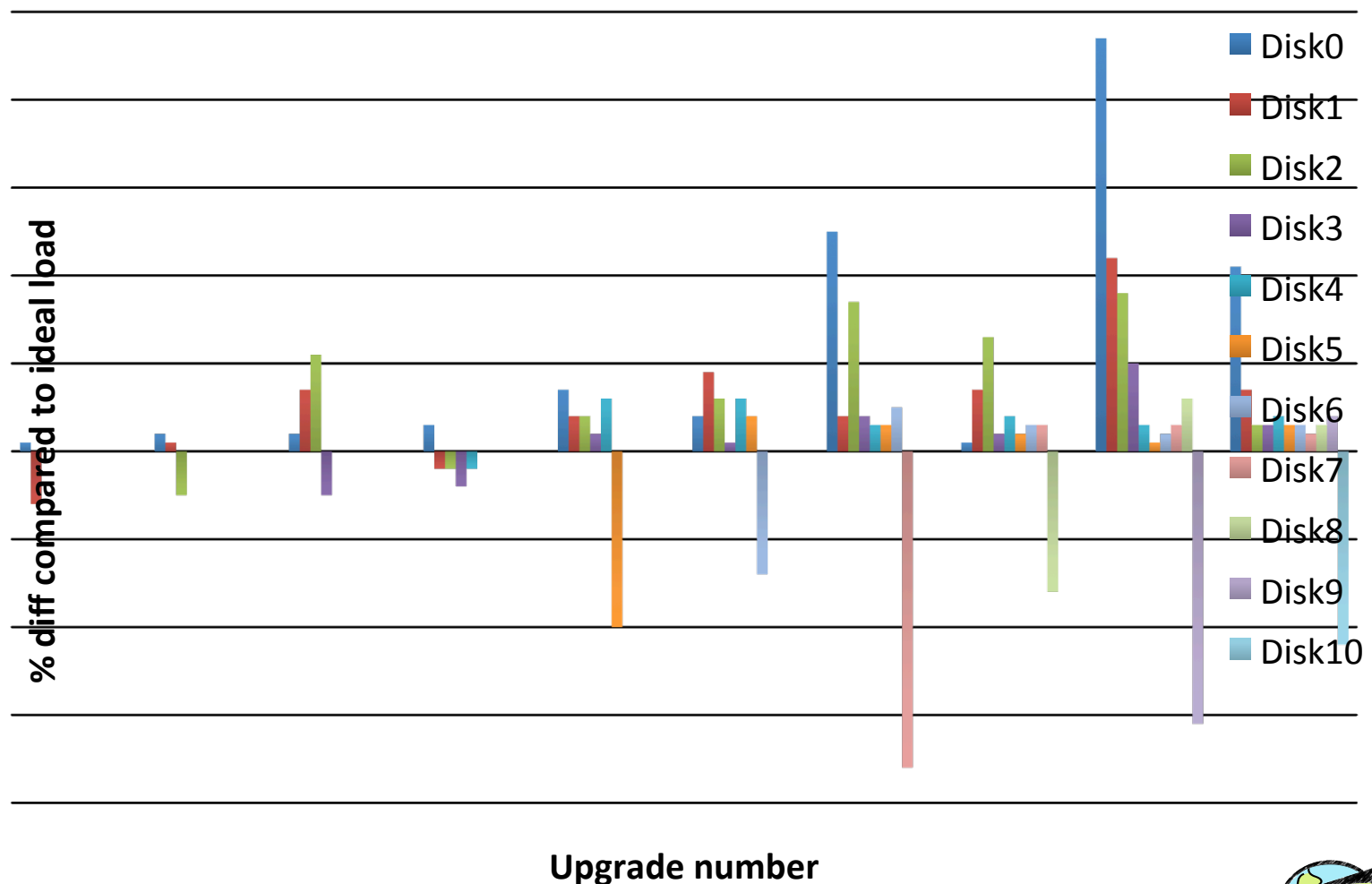Group

# Computing the location of a block

- **Formulas**
  - *Zone  =*
    - *Search in zone tree*
  - Then use AdaptRaid mechanisms and metadata
- **Metadata**
  - Zone // Tree of AdaptRaid patterns
    - // At most 2 new zones per upgrade
    - // If upgrades every 6 months:
    - //      40 zones in 10 years
    - //      6 levels in the tree

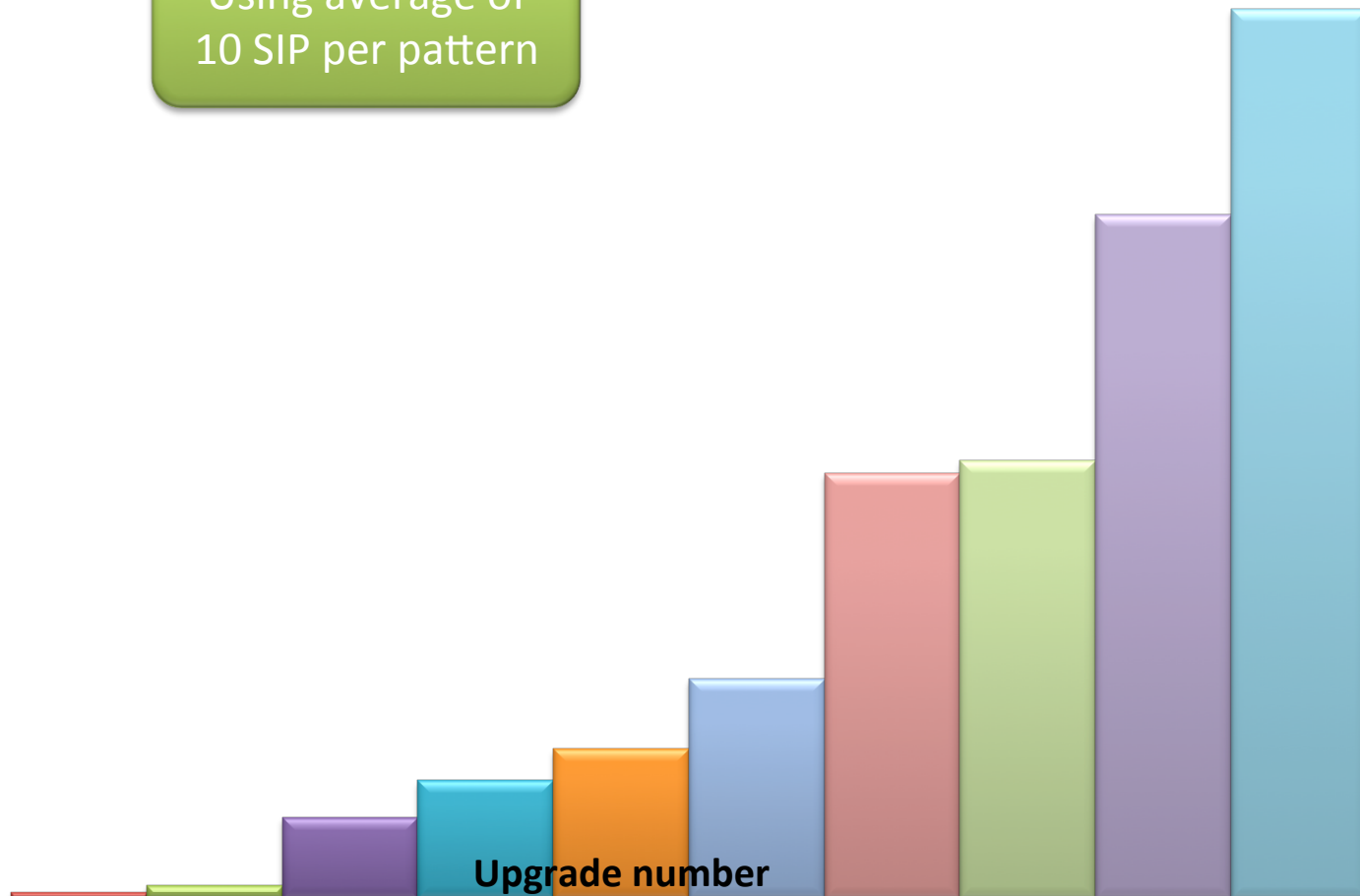# Load balance after 10 upgrades

# Metadata size after 10 upgrades
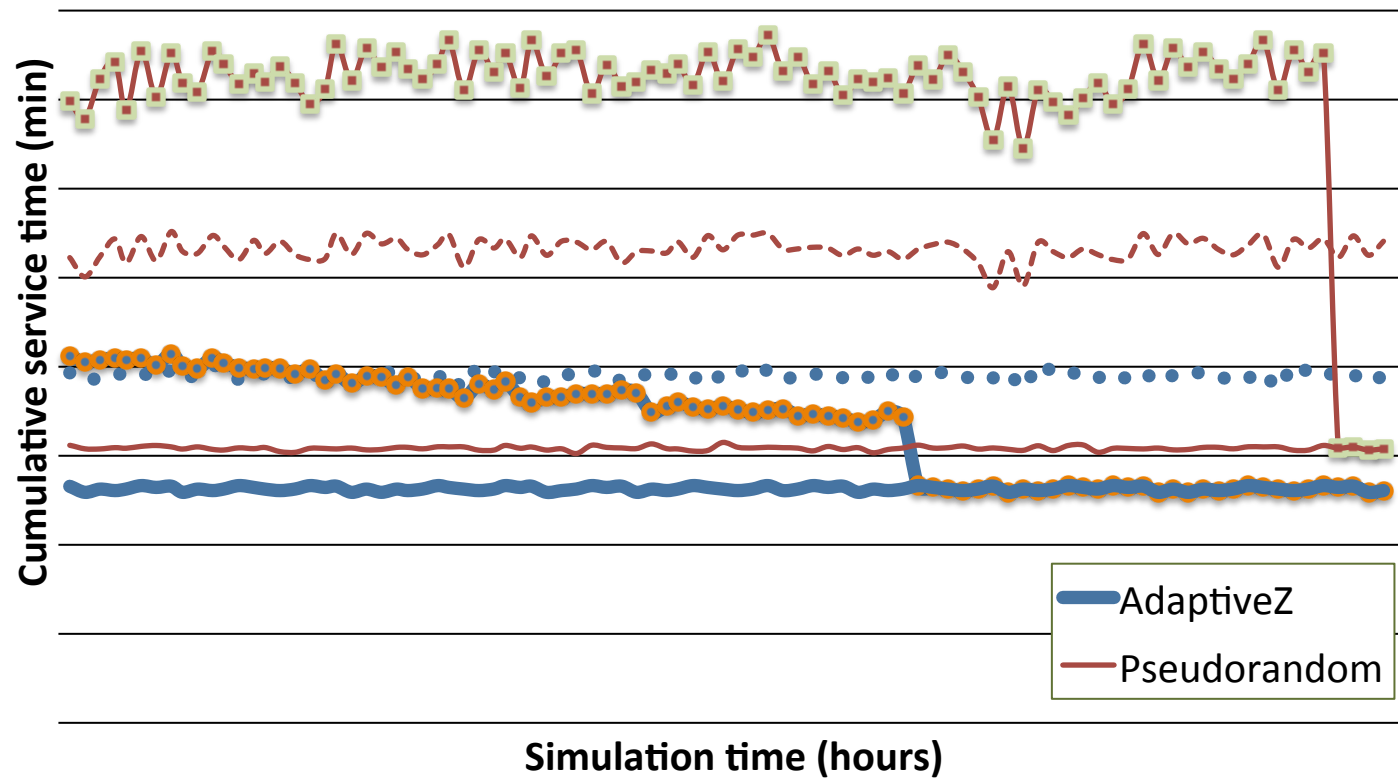
Using average of 10 SIP per pattern

Kbytes

Upgrade number

# Performance results

5 ms arrival time, 300 ms OFF periods, 70% reads
Size = poison 8K, Sequentiality 35% [Zhang 2004]



Cumulative service time (min)

Simulation time (hours)

AdaptiveZ

Pseudorandom

# Conclusions

- **Heterogeneous storage systems can be handled**
  - We can take advantage of the heterogeneity
- **They can be scalable**

- **There is another way to solve the same problem**
  - Randomization

Storage System Research Group