# S$^2$-RAID: A New RAID Architecture for Fast Data Recovery

Jiguang Wan*, Jibin Wang*, Qing Yang+, and Changsheng Xie*

*Huazhong University of Science and Technology, China

+University of Rhode Island,USA

# Overview

- A reconstruction solution-S$^2$-RAID

  - Using parallel data layout to boost data construction

- Online reconstruction performance

  - Average user response time

  - Shorten reconstruction time by a factor of 3~6

    - Comparing with the traditional RAID

# **Outline**

- Reconstruction background

- Data layout strategy

- $S^2$-RAID prototype

- Evaluation results
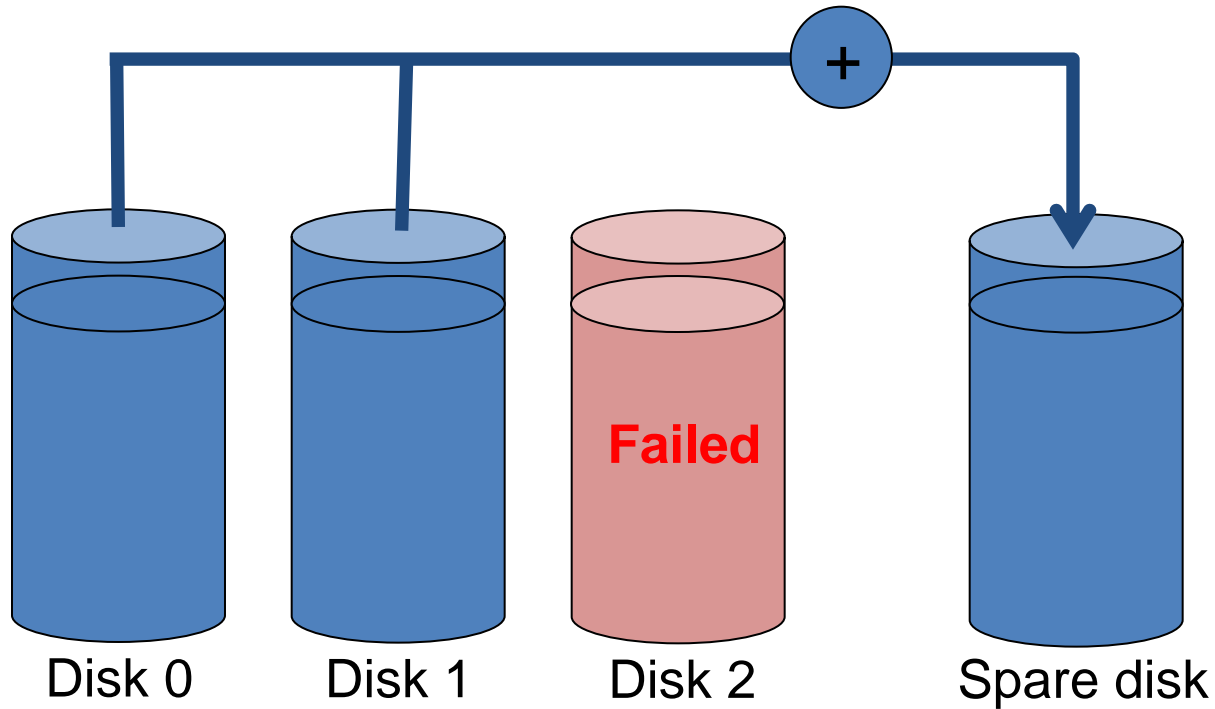
- Performance analyse

# Background

- High-capacity disk keep increasing.

- Offline reconstruction is result in service down time.

- Existing reconstruction solutions
  - Long reconstruction time and Average user response time

# S$^2$-RAID Idea

- ## Our goals
  - Reducing construction time sharply
  - Maximizing Parallel reconstruction
  - Minimizing the impact on front end performance.

- ## S$^2$-RAID data layout
  - Parallel reconstruction model
  - Using "subRAID" concept
  - Each subRAID uses standard RAID
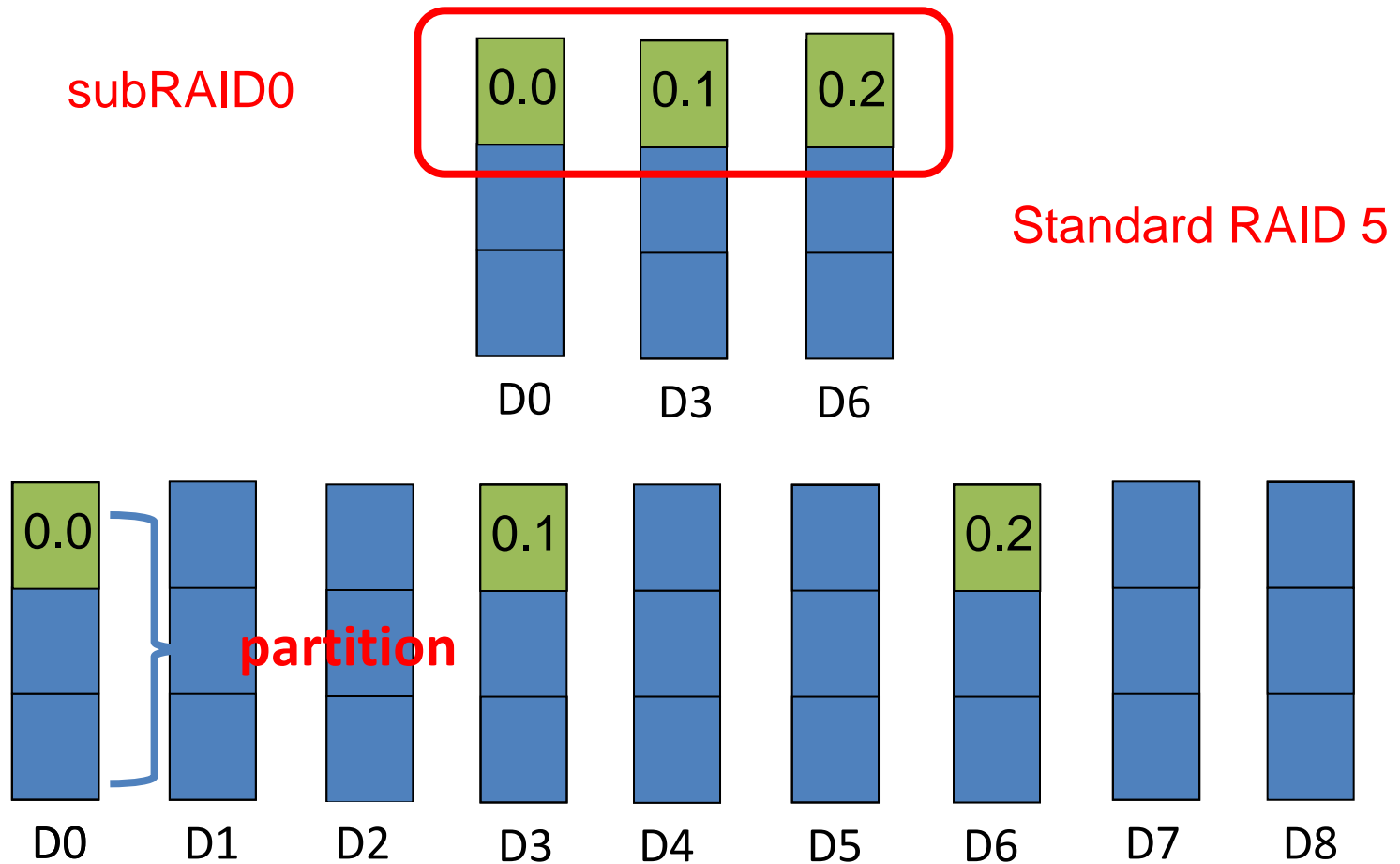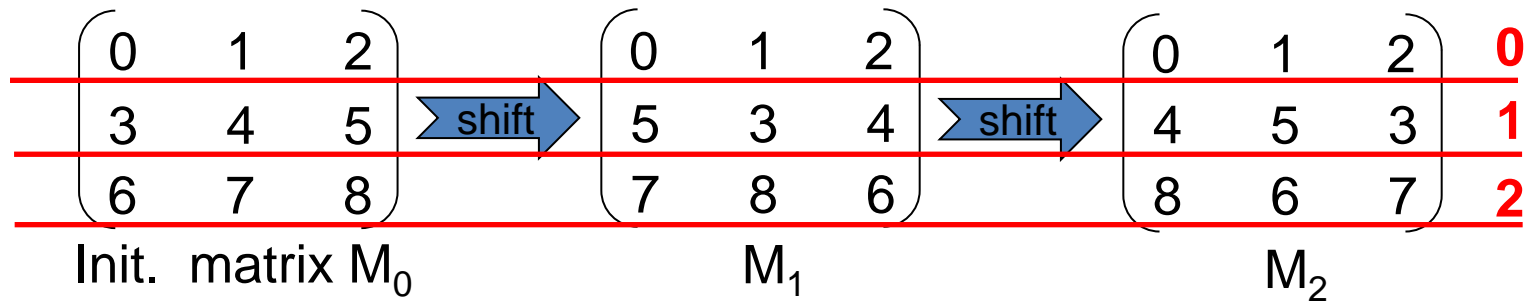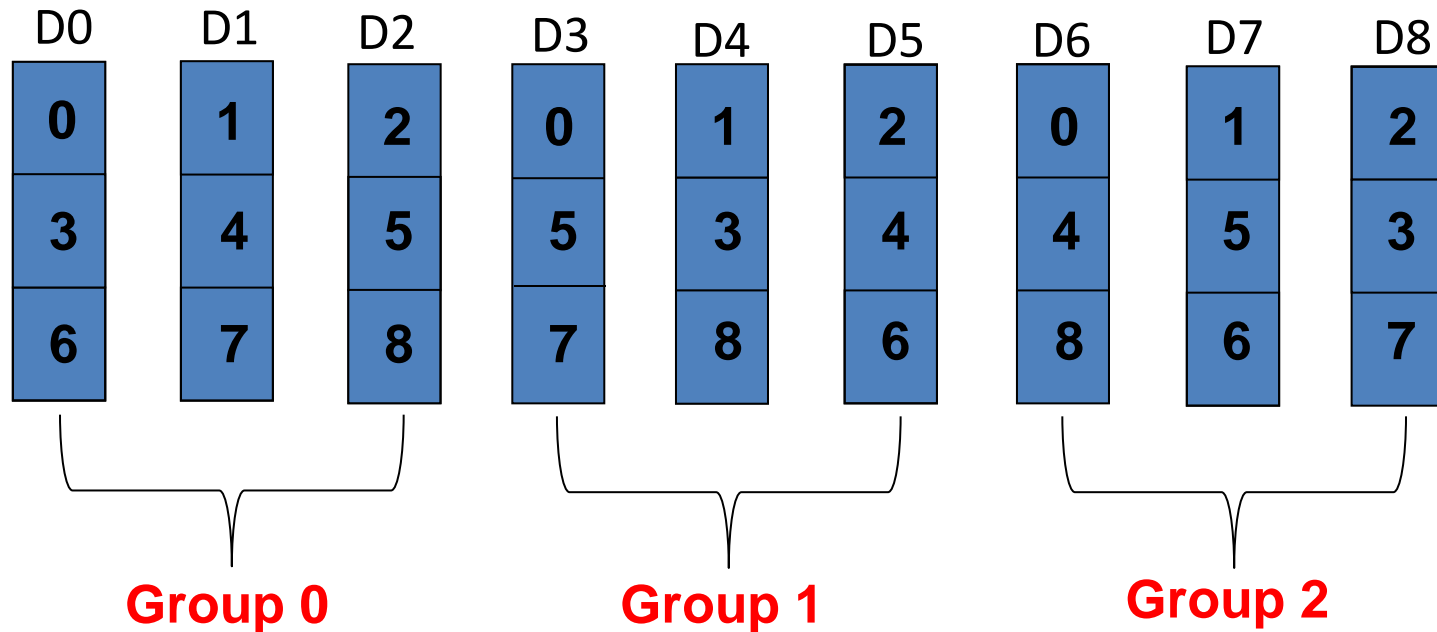
# Traditional RAID 5 reconstruction



**Single reconstruction stream**          **long reconstruction time**

# S²-RAID data layout

# S²-RAID data layout structure
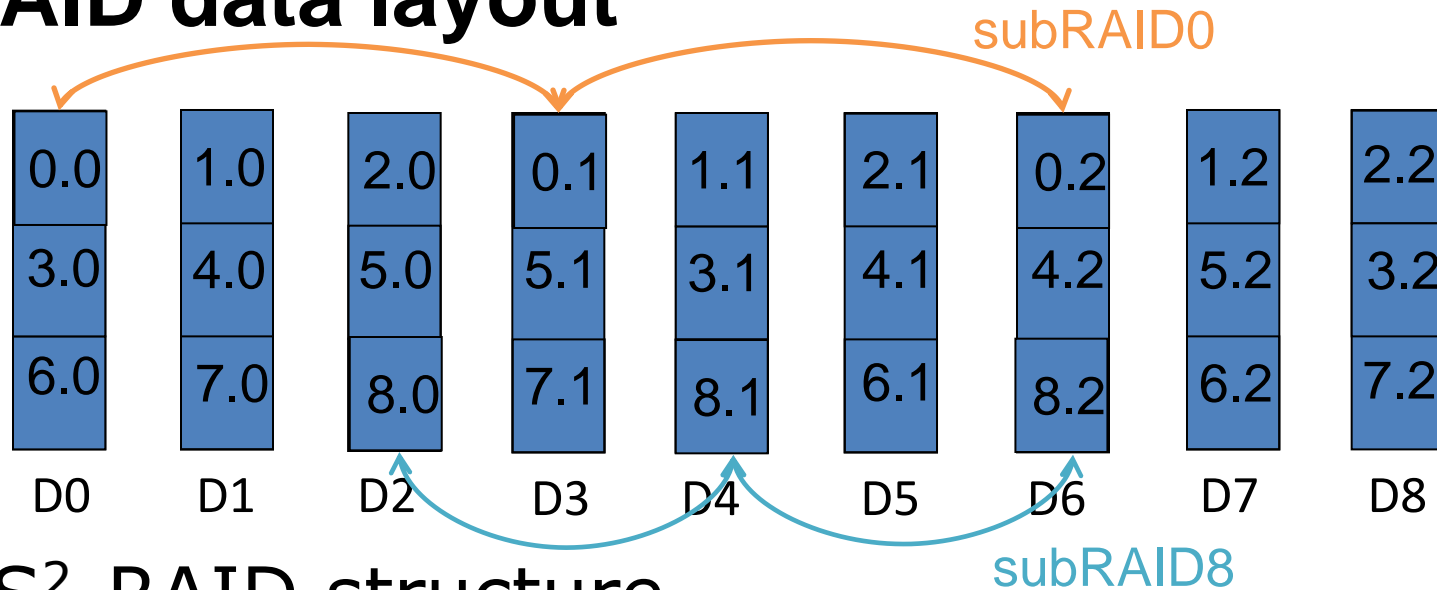
# S²-RAID data layout structure

$P_{i,j}$ :subRAID numbers of the $(j+1)^{th}$ partition on disks of $(i+1)^{th}$ group in the RAID
K: the partition number of the disk

$$m_0 = \begin{pmatrix} P_{0,0} \\ P_{0,1} \\ P_{0,2} \\ ... \\ P_{0,K-1} \end{pmatrix} \quad m_1 = \begin{pmatrix} P_{1,0} \\ P_{1,1} \\ P_{1,2} \\ ... \\ P_{1,K-1} \end{pmatrix} = \begin{pmatrix} SH_r^0(P_{0,0}) \\ SH_r^1(P_{0,1}) \\ SH_r^2(P_{0,2}) \\ ... \\ SH_r^{K-1}(P_{0,K-1}) \end{pmatrix} \quad m_i = \begin{pmatrix} P_{i,0} \\ P_{i,1} \\ P_{i,2} \\ ... \\ P_{i,K-1} \end{pmatrix} = \begin{pmatrix} SH_r^0(P_{i-1,0}) \\ SH_r^1(P_{i-1,1}) \\ SH_r^2(P_{i-1,2}) \\ ... \\ SH_r^{K-1}(P_{i-1,K-1}) \end{pmatrix}$$
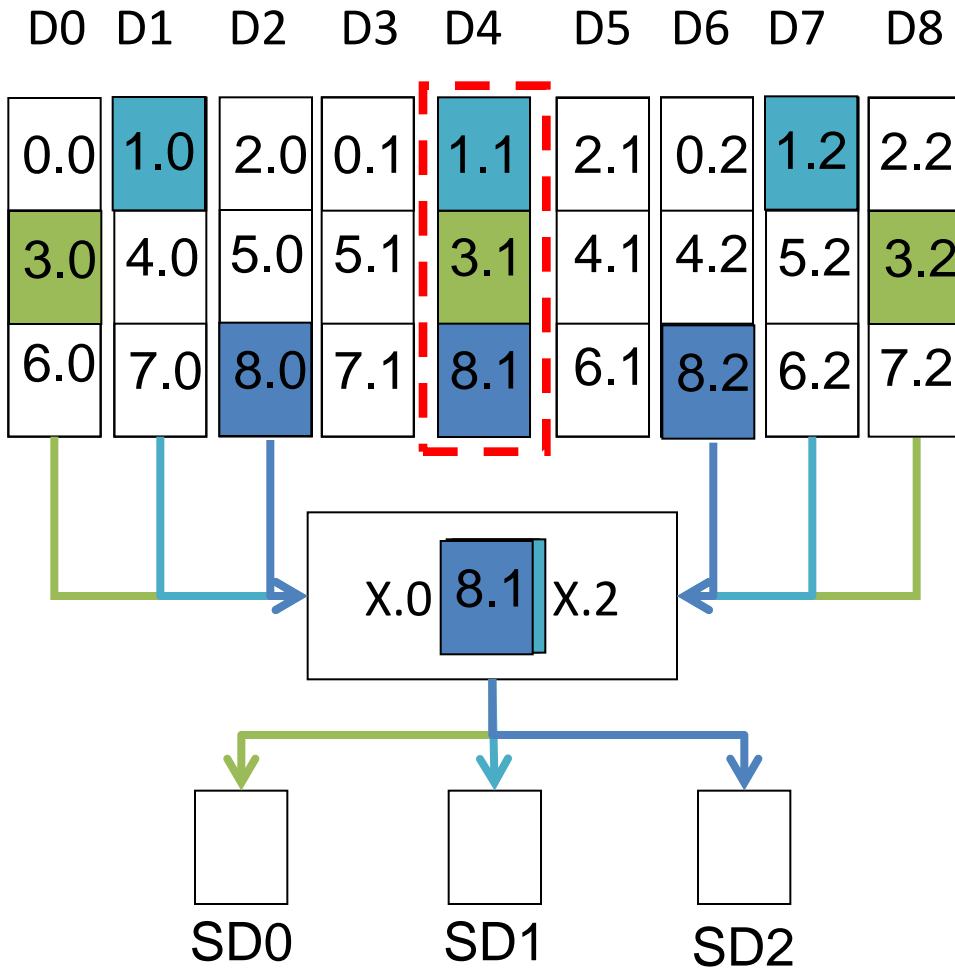
**Note:** the size of the group must be a prime number

8

# S²-RAID data layout



- S²-RAID structure

  - 9 disks

  - 9 subRAIDs

  - RAID type

  - RAID 5、RAID 10、RAID 6 etc.

# S²-RAID 5 reconstruction

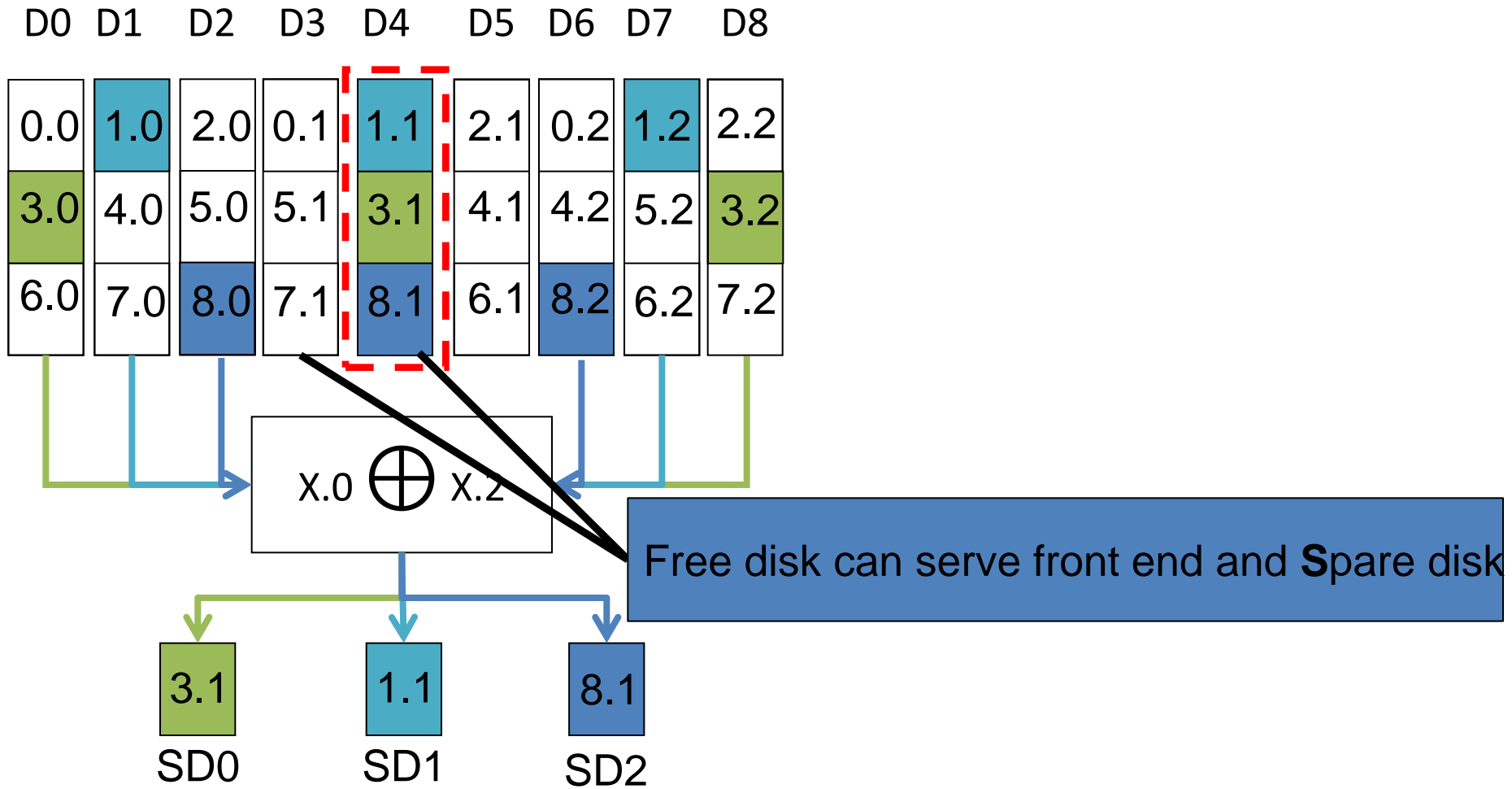| D0 | D1 | D2 | D3 | D4 | D5 | D6 | D7 | D8 |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| 0.0 | 1.0 | 2.0 | 0.1 | 1.1 | 2.1 | 0.2 | 1.2 | 2.2 |
| 3.0 | 4.0 | 5.0 | 5.1 | 3.1 | 4.1 | 4.2 | 5.2 | 3.2 |
| 6.0 | 7.0 | 8.0 | 7.1 | 8.1 | 6.1 | 8.2 | 6.2 | 7.2 |

X.0  8.1  X.2

SD0    SD1    SD2

D4 was divided into 3 partitions

Reconstruction speed!

No bottleneck in reconstruction

No operation conflict(write or read)

10

# S²-RAID 5 reconstruction

D0  D1  D2  D3  D4  D5  D6  D7  D8

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| 0.0 | 1.0 | 2.0 | 0.1 | 1.1 | 2.1 | 0.2 | 1.2 | 2.2 |
| 3.0 | 4.0 | 5.0 | 5.1 | 3.1 | 4.1 | 4.2 | 5.2 | 3.2 |
| 6.0 | 7.0 | 8.0 | 7.1 | 8.1 | 6.1 | 8.2 | 6.2 | 7.2 |

X.0 $\bigoplus$ X.2

Free disk can serve front end and **S**pare disk

3.1    1.1    8.1
SD0    SD1    SD2

# S²-RAID prototype structure



- S²-RAID prototype based on MD, are using the open source
- The *iSCSI target* module modifies the IET SCSI command handling and disk IO parts.
- The *Config* module provides RAID setup and configuration functions using mdadm commands to realize different S²-RAID subRAID functions.
- The *S²-RAID* module realizes the basic functions of RAID10 and RAID5 including RAID rebuilder based on MD.

12

# Experimental Setup

- Hardware of server and client
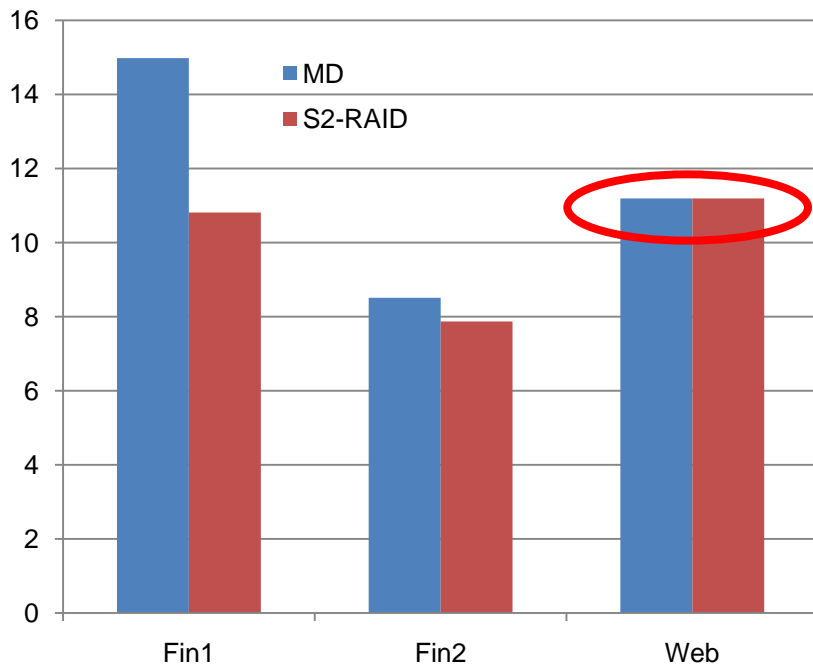- Evaluation tools of the storage server and client

<span style="color:red">server</span> <span style="color:red">client</span>

| | |
|---|---|
| OS | Fedora Core 8.0 |
| blktrace | blktrace 1.0 |
| postmark | postmark 1.5 |
| TPC-C | TPC-C speculva-1.2.3 |
| postgresql | postgresql 8.1.15 |
| gnuplot | gnuplot 4.2.5 |
| TPC-W | TPC-W 1.5 |
| Jdk | jdk 1.5.0.06 |
| Tomcat | tomcat 5.5 |
| Mysql | mysql 5.0.45 |
| iscsi-initiator | iscsi-initiator-util-6.2.0.865 |

| | |
|---|---|
| OS | Fedora Core 8.0 |
| disks | 1 Seagate ST3160023AS, 160GB, 7200RPM. |
| Disks | 12 Seagate ST3160023AS, 160GB, 7200RPM |
| mainboard | SUPER X7DVL-L |
| Mainboard | GA-945GCMX-S2 |
| CPU | Intel(R) Xeon(R) CPU 5110 @ 1.60GHz |
| CPU | Intel(R) Celeron(R) CPU 2.80GHz |
| NIC | Tigon3 |
| NIC | Intel® PRO/1000 |
| Memory | 512MB DDR2 |
| HBA | Highpoint 2240 RAID, |

# S²-RAID 5 reconstruction performance

■Two evaluation parameters

- •Average User Response Time
- •Reconstruction Time

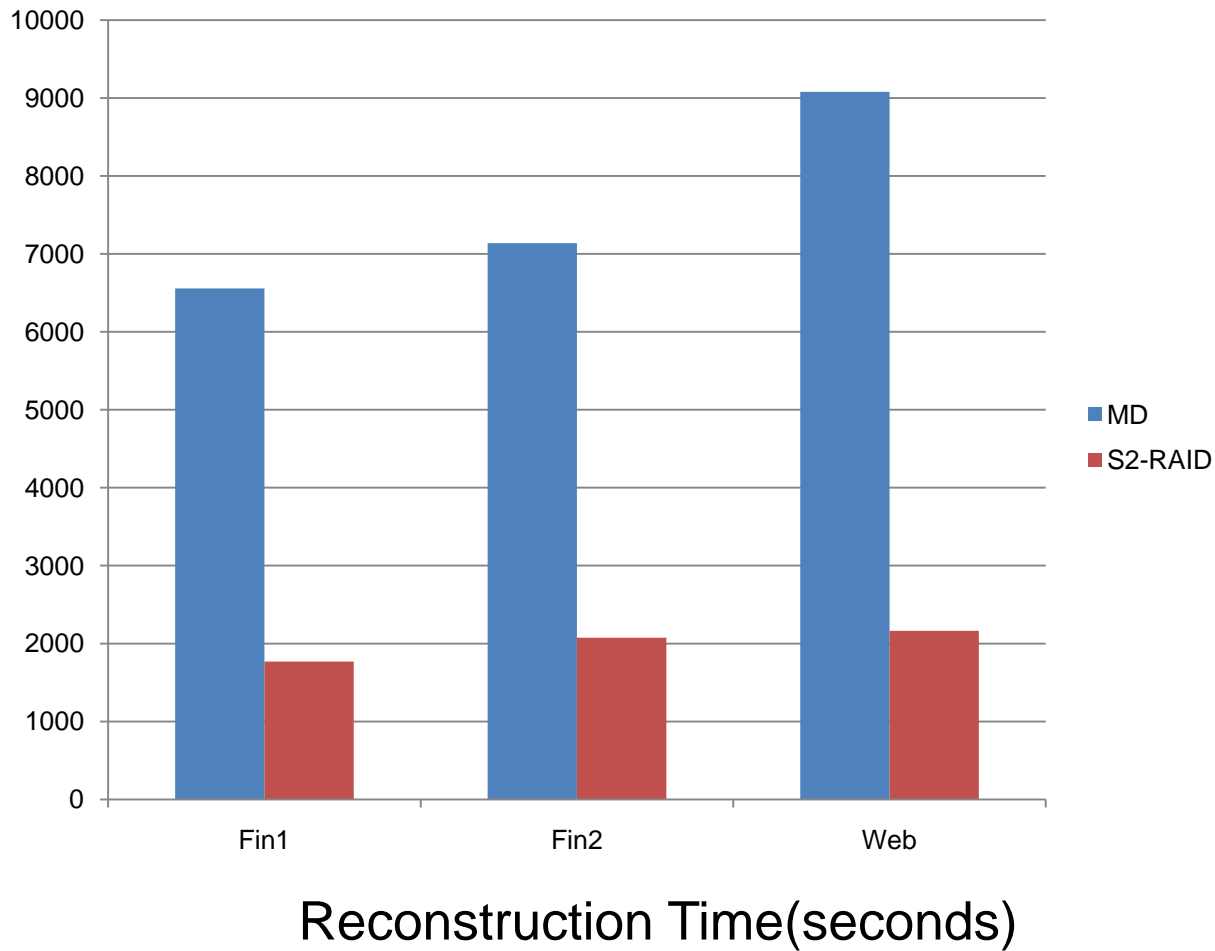Average User Response Time(ms)

SPC trace characteristics

| Trace File | Write Ratio | Ave Req Size: KB | Total Req |
|------------|-------------|------------------|-----------|
| Financial-1 | 76.84% | 3.38 | 5,334,987 |
| Financial-2 | 17.65% | 2.39 | 3,699,195 |
| Websearch | 0% | 15.07 | 4,579,809 |

# S²-RAID 5 reconstruction performance



Reconstruction Time(seconds)

# S²-RAID 5 Normal Performance



Average User Response Time(ms)

# S²-RAID 5 Degraded Performance



Average User Response Time(ms)

17

# Other Benchmark Performance (MD vs S²-RAID)

■ Average User Response Time          ■ Reconstruction Time

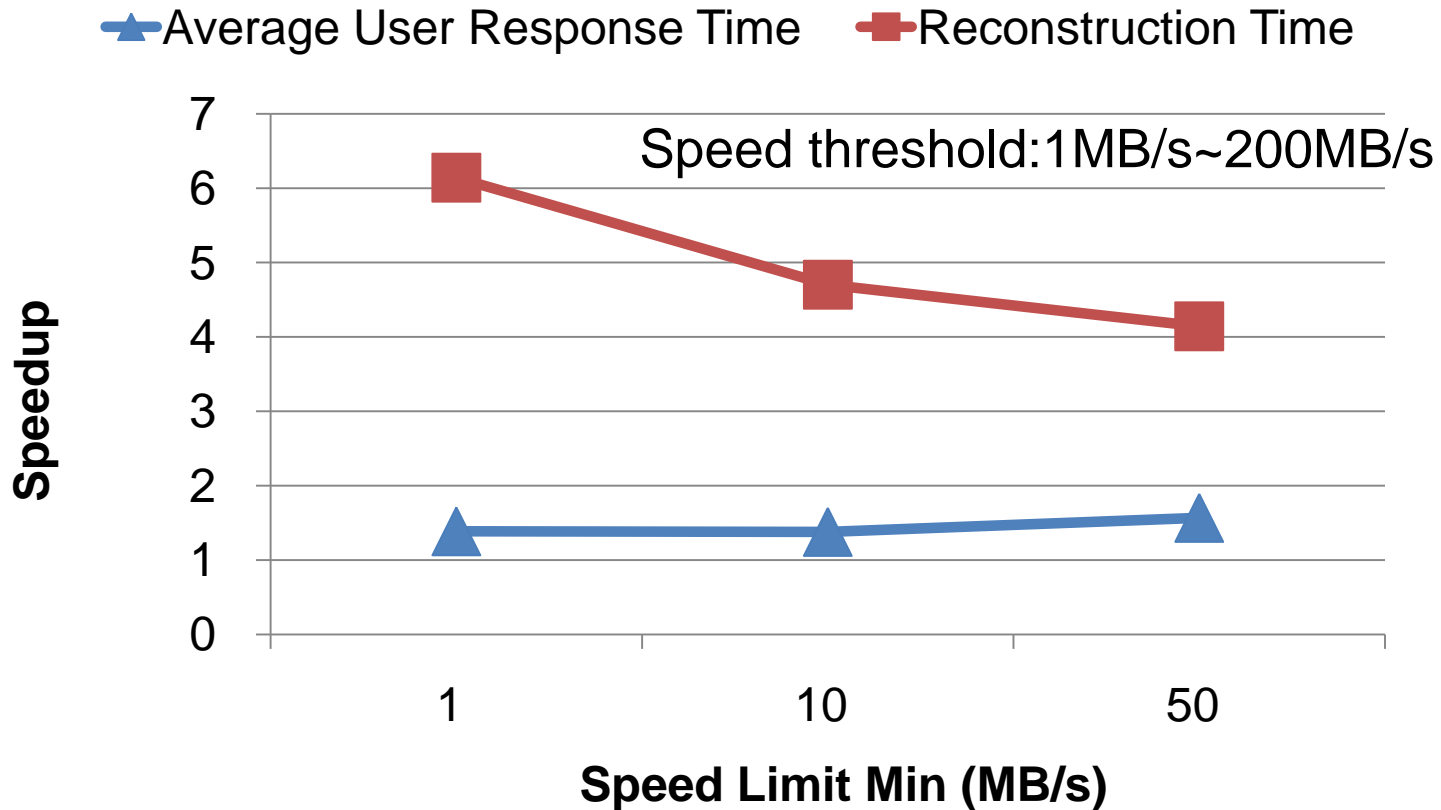TPCC:20 warehouses with 10 terminals per warehouse interval of 120 minutes

TPCW:150 emulated browsers

Postmark:20,000 files of size 4KB to 500KB and to perform 100,000 transactions
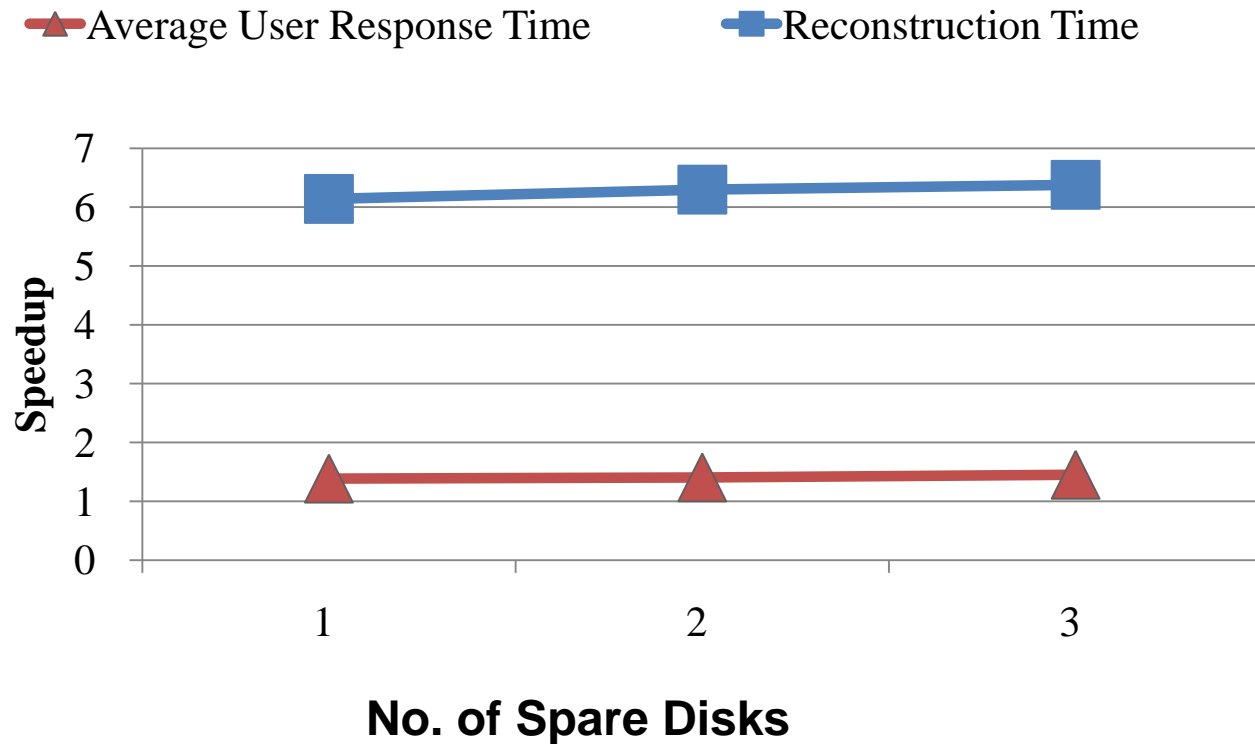
18

# Sensitivity Parameters for Reconstruction

- Some sensitivity parameters

  Reconstruction speed bandwidth

  I/O request block size

  Number of spare disk(additional disk not system disk)

# Reconstruction speed bandwidth
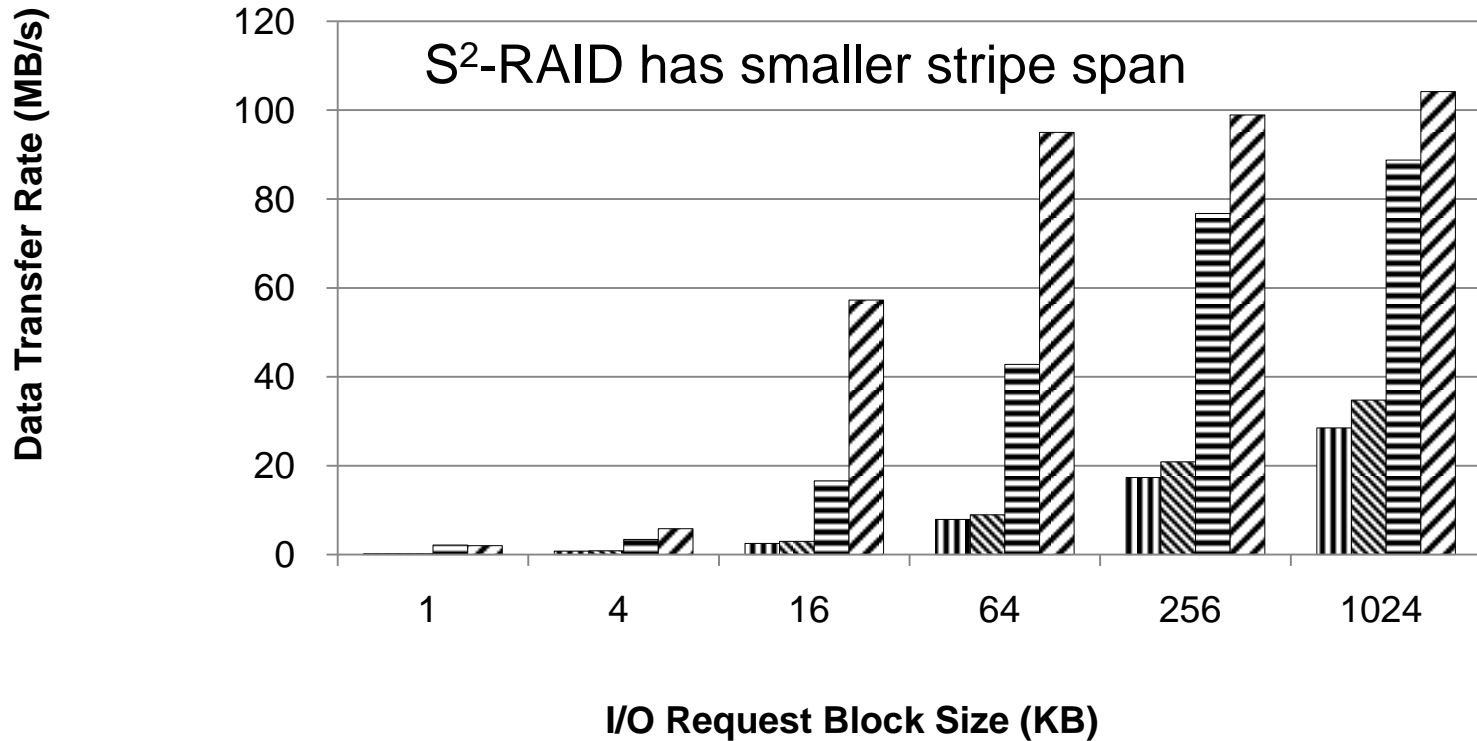


The result is based on Financial-1 traces

# Number of spare disk



Speed bandwidth and No. of spare disk is insensitive to $s^2$-RAID

# I/O request block size



RAID5-random    S2-RAID5-random    RAID5-sequence    S2-RAID5-sequence

$S^2$-RAID has smaller stripe span

Data Transfer Rate (MB/s)

I/O Request Block Size (KB)

# Conclusion

- A parallel reconstruction data layout
- Implement the $s^2$-RAID prototype and evaluation of this structure

- $S^2$-Raid reduces the reconstruction time greatly.
- User response time of $S^2$-Raid is comparable to that of MD.

- Optimization?
  - Embedding existing rebuilding process (distributed sparing)------Reduce the number of disks
  - Tolerate the mulit-disk failures.

23

# Thank you for your attention!

**Questions?**