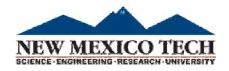
## Storage Systems (4)

Dr. Jun Zheng
CSE325 Principles of Operating
Systems
12/4/2019



# RAID 6: Recovering from 2 failures

- $\square$  Why > 1 failure recovery?
  - □ operator accidentally replaces the wrong disk during a failure
  - □ since disk bandwidth is growing more slowly than disk capacity, the MTT Repair a disk in a RAID system is increasing
    - $\Rightarrow$  increases the chances of a 2nd failure during repair since takes longer
  - □ reading much more data during reconstruction meant increasing the chance of an uncorrectable media failure, which would result in data loss



# RAID 6: Recovering from 2 failures

- $\square$  Network Appliance's row-diagonal parity or RAID-DP
- ☐ Like the standard RAID schemes, it uses redundant space based on parity calculation per stripe
- ☐ Since it is protecting against a double failure, it adds two check blocks per stripe of data.
  - $\square$  If p+1 disks total, p-1 disks have data; assume p = 5
  - ☐ Row parity disk is just like in RAID 4
  - ☐ Even parity across the other 4 data blocks in its stripe
- ☐ Each block of the diagonal parity disk contains the even parity of the blocks in the same diagonal

### Example p = 5

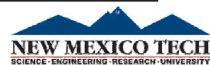
- ☐ Row diagonal parity starts by recovering one of the 4 blocks on the failed disk using diagonal parity
  - ☐ Since each diagonal misses one disk, and all diagonals miss a different disk, 2 diagonals are only missing 1 block
- □ Once the data for those blocks is recovered, then the standard RAID recovery scheme can be used to recover two more blocks in the standard RAID 4 stripes
- ☐ Process continues until two failed disks are restored

Data Disk o	Data Disk 1	Data Disk 2	Data Disk 3	Row Parity	Diagonal Parity
0	1	2	$\setminus$ 3	4	0
	2	3	4	0	
2	3	4	0	1	2
3	4	0	1	2	3



### **RAID Summary**

RAI	D level	Disk failures tolerated, check space overhead for 8 data disks	Pros	Cons	Company products
0	Nonredundant striped	0 failures, 0 check disks	No space overhead	No protection	Widely used
1	Mirrored	1 failure, 8 check disks	No parity calculation; fast recovery; small writes faster than higher RAIDs; fast reads	Highest check storage overhead	EMC, HP (Tandem), IBM
2	Memory-style ECC	1 failure, 4 check disks	Doesn't rely on failed disk to self-diagnose	~ Log 2 check storage overhead	Not used
3	Bit-interleaved parity	1 failure, 1 check disk	Low check overhead; high bandwidth for large reads or writes	No support for small, random reads or writes	Storage Concepts
4	Block-interleaved parity	1 failure, 1 check disk	Low check overhead; more bandwidth for small reads	Parity disk is small write bottleneck	Network Appliance
5	Block-interleaved distributed parity	1 failure, 1 check disk	Low check overhead; more bandwidth for small reads and writes	Small writes → 4 disk accesses	Widely used
6	Row-diagonal parity, EVEN-ODD	2 failures, 2 check disks	Protects against 2 disk failures	Small writes → 6 disk accesses; 2X check overhead	Network Appliance



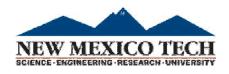
#### Question

Consider a RAID level 5 organization comprising five disks, with the parity for sets of four blocks on four disks stored on the fifth disk. How the blocks are accessed in order to perform the following?

- (1) A write of one block of data
- (2) A write of seven continuous blocks of data, assume the blocks begin at a four-block boundary

#### Answer

- (1) A write of one block of data requires the following: read of the parity block, read of the old data stored in the target block, computation of the new parity based on the differences between the new and old contents of the target block, and write of the parity block and the target block.
- (2) Assume that the seven contiguous blocks begin at a fourblock boundary. A write of seven contiguous blocks of data could be performed by writing the seven contiguous blocks, writing the parity block of the first four blocks, reading the eight block, computing the parity for the next set of four blocks and writing the corresponding parity block onto disk.



#### **In-class Work 10**

☐ Consider a 4-drive, 200 GB-per-drive RAID array. What is the available data storage capacity for each of the RAID levels, 0, 1+0, 3, 4, 5, and 6?

