

A GPU-Based Storage System

Matthew L. Curry

Sandia is a multiprogram laboratory operated by Sandia Corporation, a Lockheed Martin Company, for the United States Department of Energy's National Nuclear Security Administration under contract DE-AC04-94AL85000.





A RAID Overview

- RAID has served us well
- Several levels introduced initially
 - Some have been added, some have been removed
- Today, there are five major RAID types in use in reliable production systems
- Going forward, they will not be reliable enough.

Current Parity-Based RAID

RAID 5

Disk 0	Disk 1	Disk 2	Disk 3	Disk 4	Disk 5	Disk 6	Disk 7
Data Block 0	Data Block 1	Data Block 2	Data Block 3	Data Block 4	Data Block 5	Data Block 6	Parity Block 0
Data Block 7	Data Block 8	Data Block 9	Data Block 10	Data Block 11	Data Block 12	Parity Block 1	Data Block 13
Data Block 14	Data Block 15	Data Block 16	Data Block 17	Data Block 18	Parity Block 2	Data Block 19	Data Block 20
...

RAID 6

Disk 0	Disk 1	Disk 2	Disk 3	Disk 4	Disk 5	Disk 6	Disk 7
Data Block 0	Data Block 1	Data Block 2	Data Block 3	Data Block 4	Data Block 5	Parity Block 0	Parity Block 1
Data Block 6	Data Block 7	Data Block 8	Data Block 9	Data Block 10	Parity Block 2	Parity Block 3	Data Block 11
Data Block 12	Data Block 13	Data Block 14	Data Block 15	Parity Block 4	Parity Block 5	Data Block 16	Data Block 17

Hierarchical Parity-Based RAID

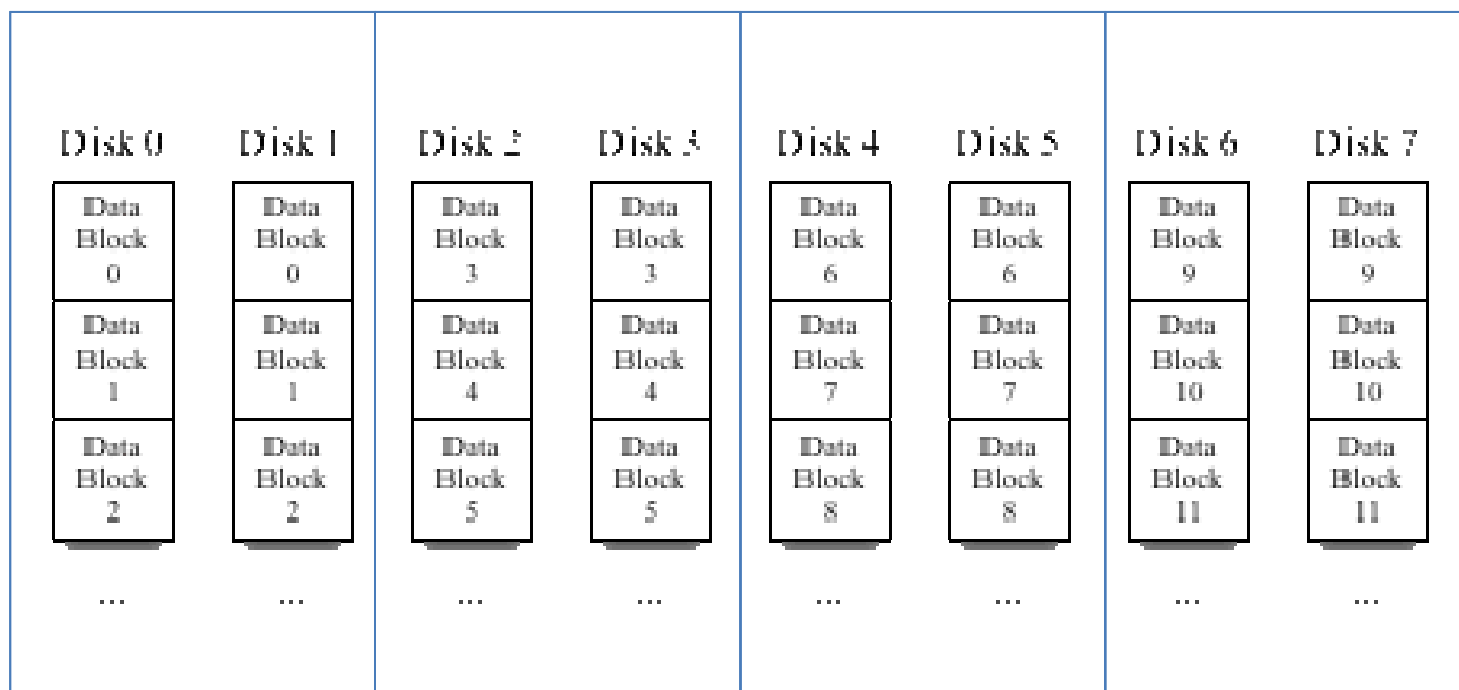
Disk 0	Disk 1	Disk 2	Disk 3	Disk 4	Disk 5	Disk 6	Disk 7
Data Block 0	Data Block 1	Data Block 2	Data Block 3	Data Block 4	Data Block 5	Parity Block 0	Parity Block 1
Data Block 12	Data Block 13	Data Block 14	Data Block 15	Data Block 16	Parity Block 4	Parity Block 5	Data Block 17
Data Block 24	Data Block 25	Data Block 26	Data Block 27	Parity Block 8	Parity Block 9	Data Block 28	Data Block 29
...

RAID 6+0

Disk 8	Disk 9	Disk 10	Disk 11	Disk 12	Disk 13	Disk 14	Disk 15
Data Block 6	Data Block 7	Data Block 8	Data Block 9	Data Block 10	Data Block 11	Parity Block 2	Parity Block 3
Data Block 18	Data Block 19	Data Block 20	Data Block 21	Data Block 22	Parity Block 6	Parity Block 7	Data Block 23
Data Block 30	Data Block 31	Data Block 32	Data Block 33	Parity Block 10	Parity Block 11	Data Block 34	Data Block 35
...

Hierarchical Mirror-Based RAID

RAID 1+0

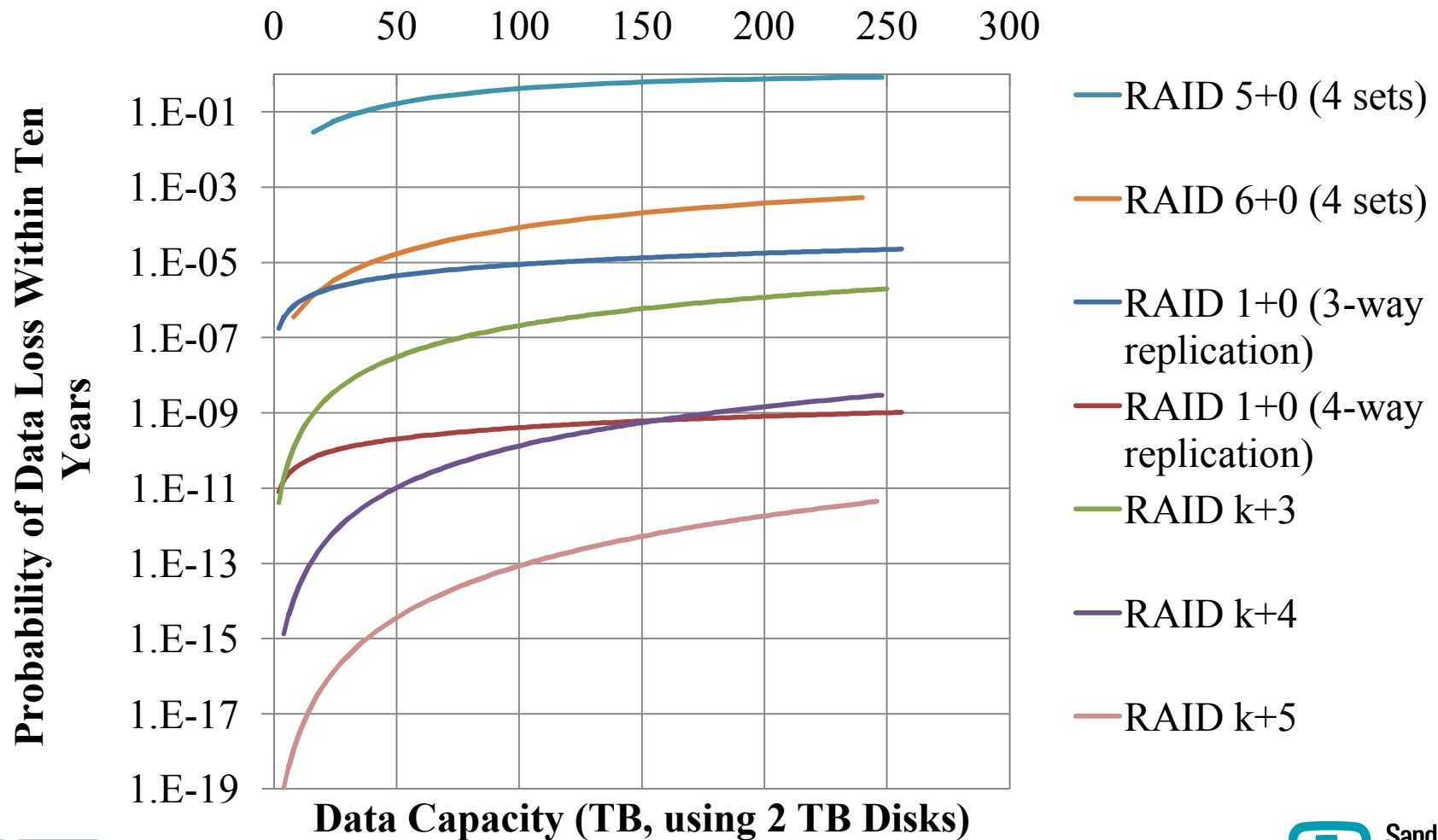




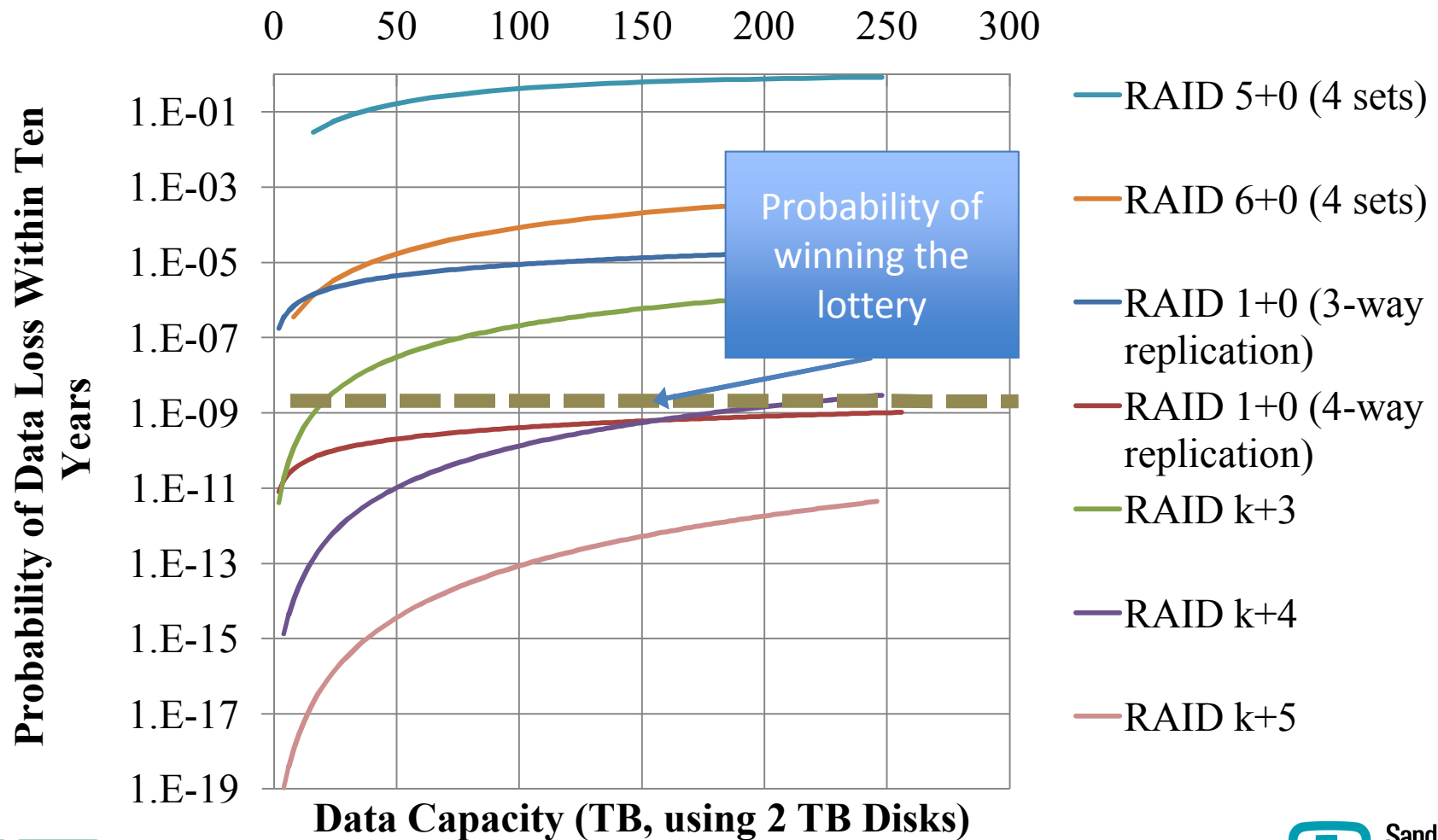
$k+m$ RAID, an Alternative

- The user can choose the amount of parity required for a RAID set
 - $k+1$ RAID = RAID 5
 - $k+2$ RAID = RAID 6
 - $k+3$ RAID = RAID TP (Esoteric)
 - $k+4$ RAID or better does not exist anywhere. (RAID 6+0 does not satisfy the requirements.)

$k+m$ RAID Reliability, Standard Manufacturer Statistics



$k+m$ RAID Reliability, Standard Manufacturer Statistics





Why is $k+m$ RAID Necessary?

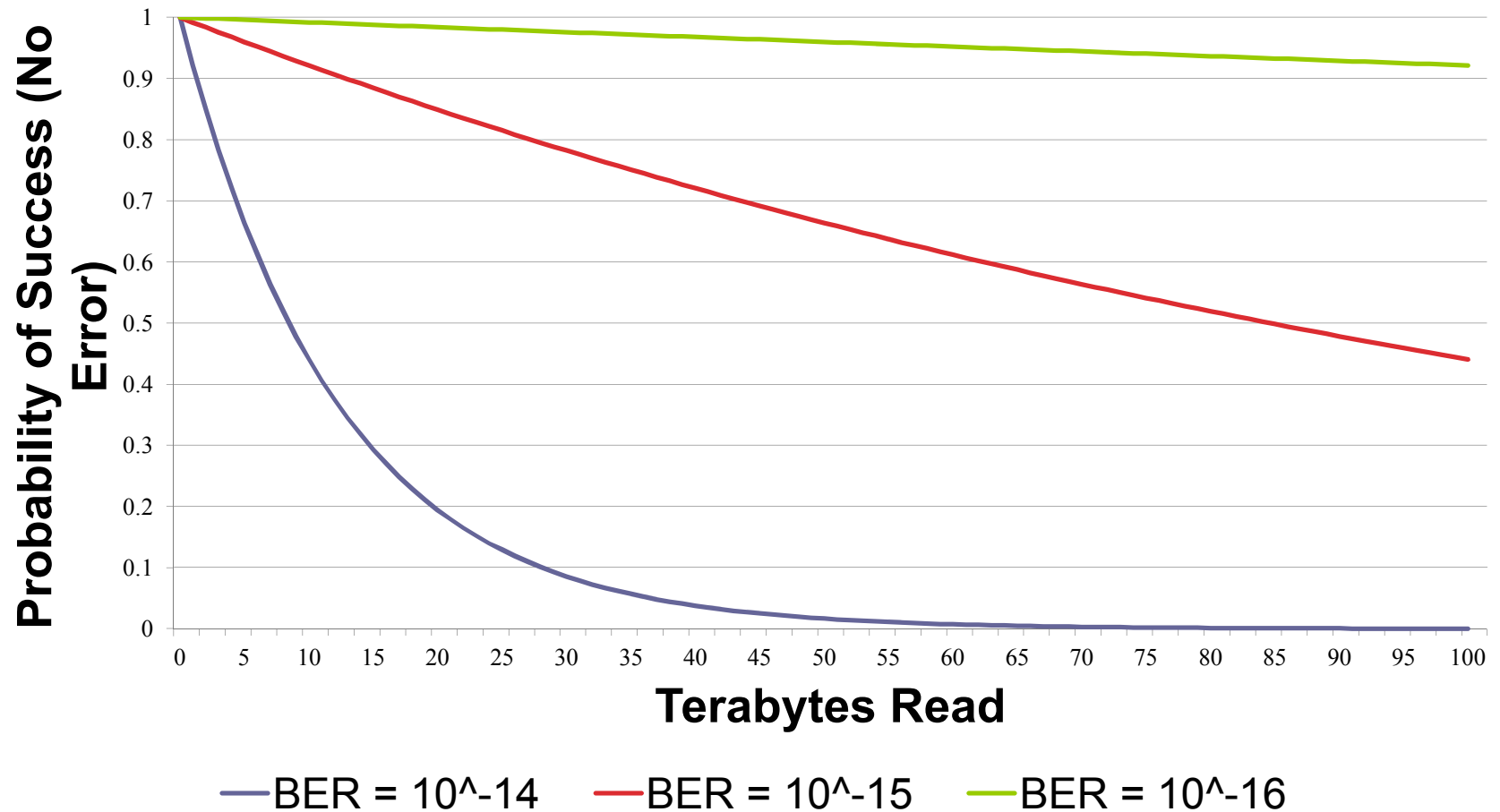
- Correlated failures
 - Manufacturer defects (batch correlated failures)
 - Environment and workload as factors
- Calculation methods of Mean Time To Failure
 - Disk age/failure correlation
 - Manufacturer-reported disk failure rates are more than 10x lower than observed
- Rebuild times are skyrocketing
 - Disk speeds are not increasing proportionally to disk capacity, and flash cannot save us yet



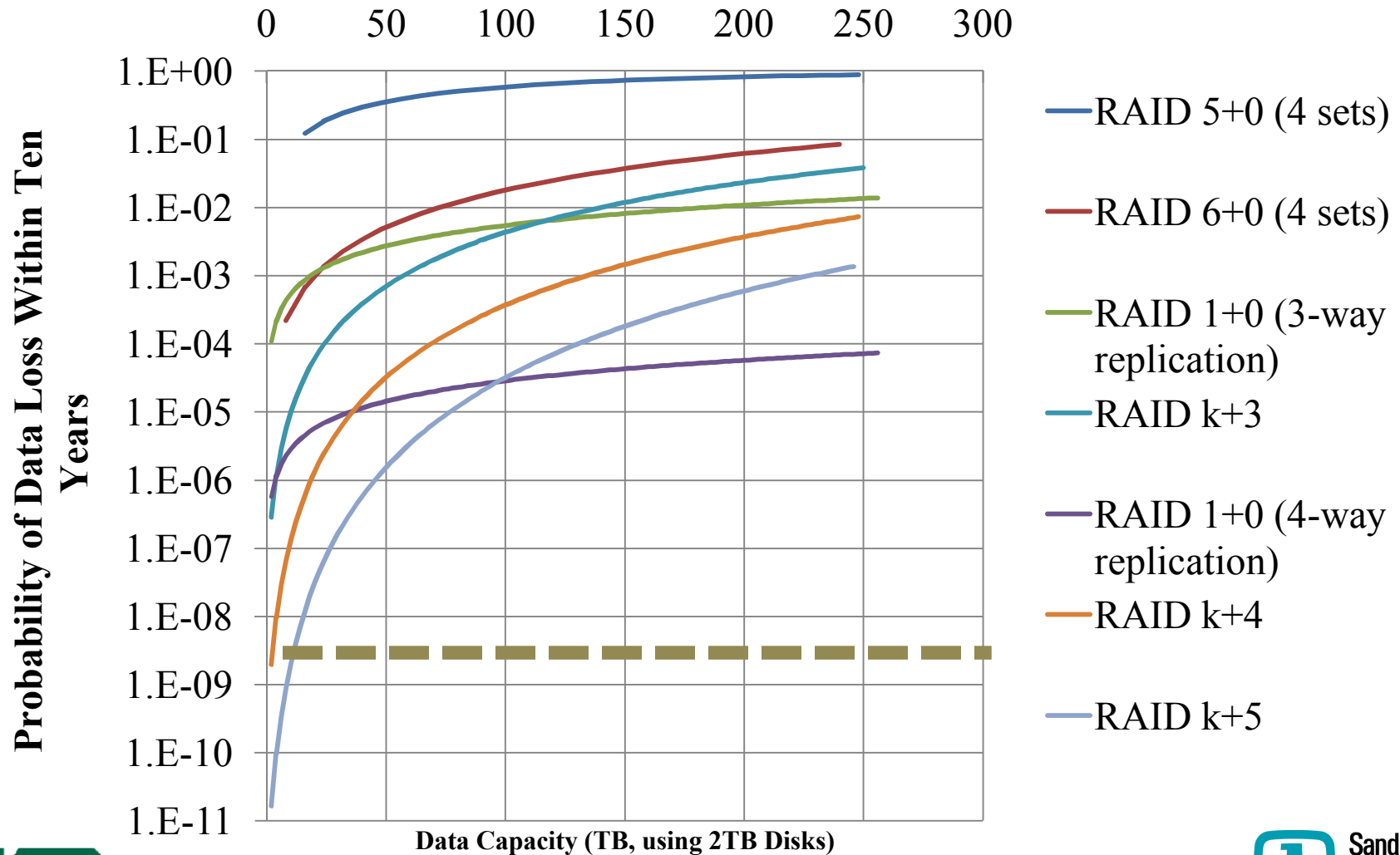
Why is $k+m$ RAID Necessary? (2)

- Data corruption/loss
 - Disk firmware bugs
 - Unrecoverable read errors
 - One sector in 100-1000 TB will be unreadable
 - Cabling, misdirected/torn writes
 - Also points to need for read verification
- Parity needs to be available to recover from these errors
 - RAID 6 double disk failures *do* happen

Unrecoverable Read Errors



RAID Reliability, Updated Statistics



Erasure Coding Computations

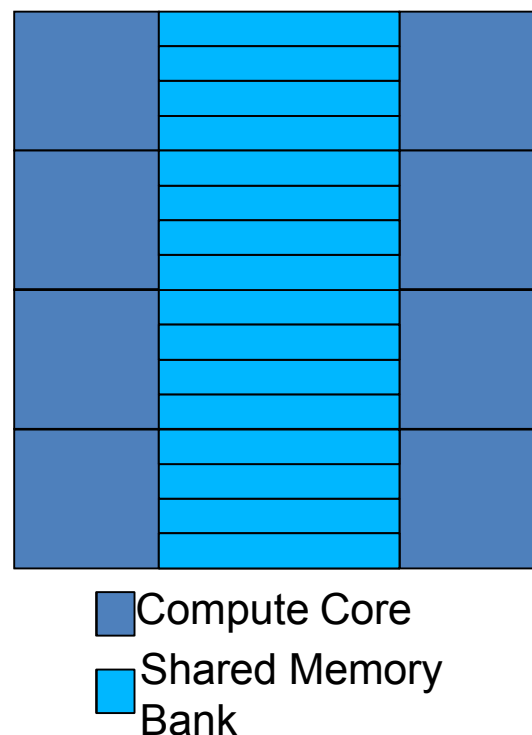
- $k+m$ RAID is necessary
- Reed-Solomon coding is required to implement it
 - $(k+m) \times k$ matrix
 - Finite field arithmetic
 - Table look ups abound

$$\begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ F_{1,0} & F_{1,1} & F_{1,2} & F_{1,3} \\ F_{2,0} & F_{2,1} & F_{2,2} & F_{2,3} \\ F_{3,0} & F_{3,1} & F_{3,2} & F_{3,3} \\ F_{4,0} & F_{4,1} & F_{4,2} & F_{4,3} \end{pmatrix} \begin{pmatrix} d_0 \\ d_1 \\ d_2 \\ d_3 \end{pmatrix} = \begin{pmatrix} c_0 \\ c_1 \\ c_2 \\ c_3 \\ c_4 \\ c_5 \\ c_6 \\ c_7 \end{pmatrix}$$

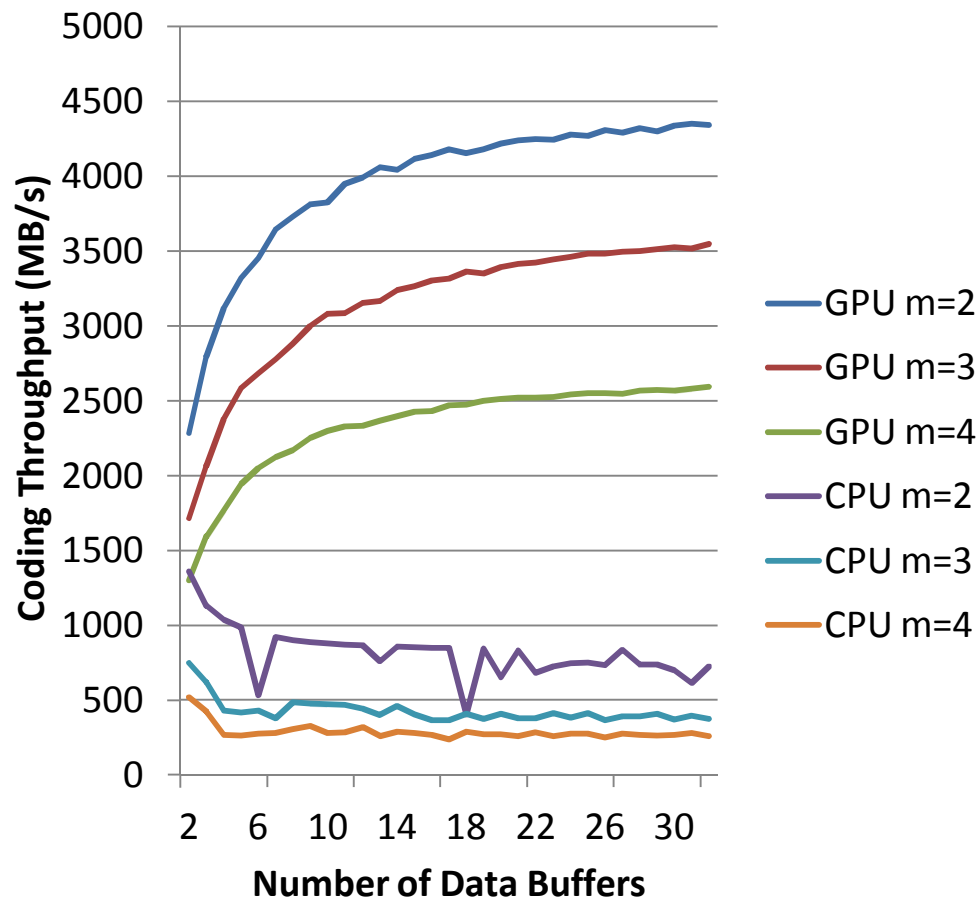
GPU Architecture Does This Well!

- Each GPU contains some number of Streaming Multiprocessors (SMs)
- Banks allow parallel non-conflicting accesses
- GPU RAID simulation results: 1.82 look-ups per cycle per SM on average
- GeForce GTX 285 has 30 SMs, so can satisfy 55 look-ups per cycle for device

CUDA 200-Series
Streaming
Multiprocessor



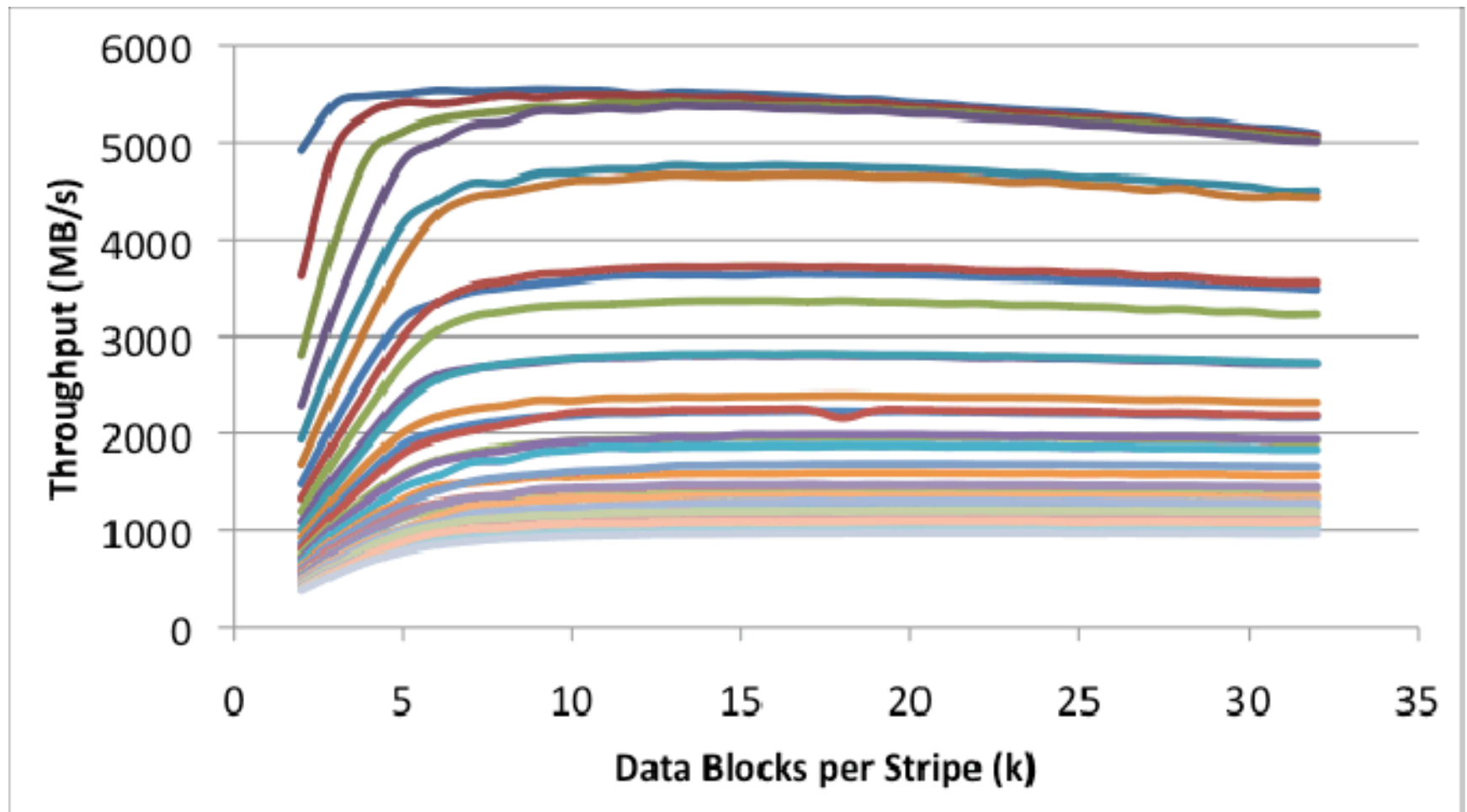
GPU-Based Reed-Solomon



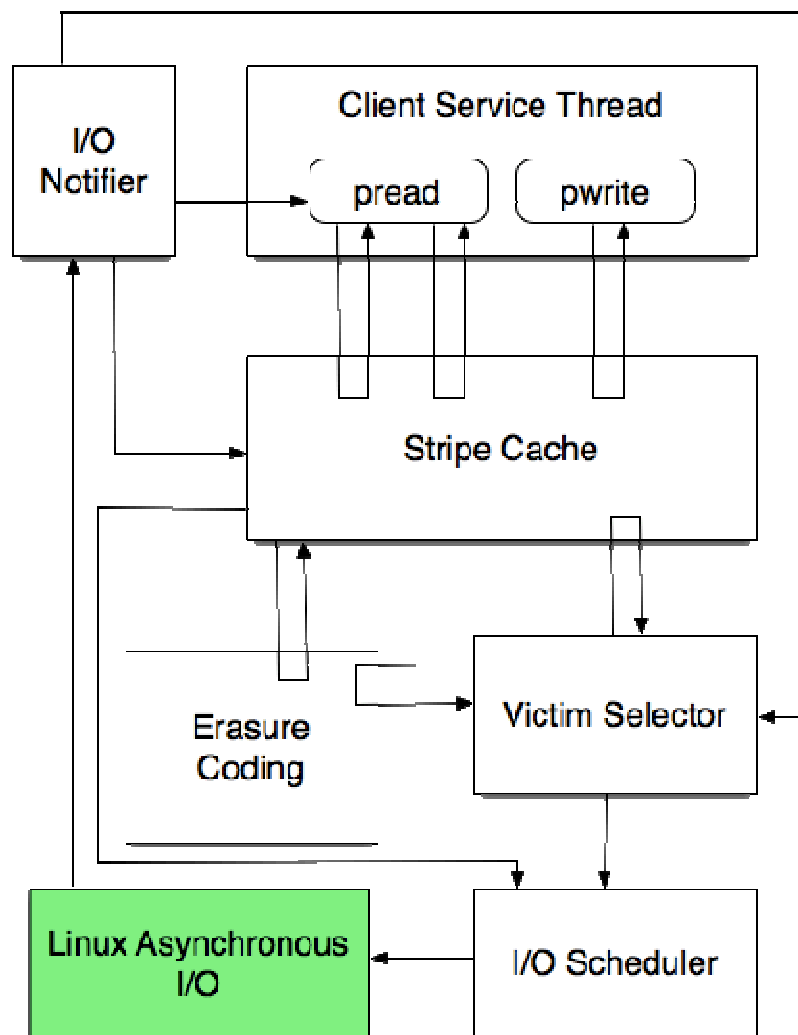
- GPU capable of providing high bandwidth at 6x-10x CPU speeds
- New memory layout and matrix generation algorithm for Reed-Solomon yields equivalent write and degraded read performance

k+m Coding on GeForce GTX 480

(Each line represents a value of m , $3 \leq m \leq 32$)

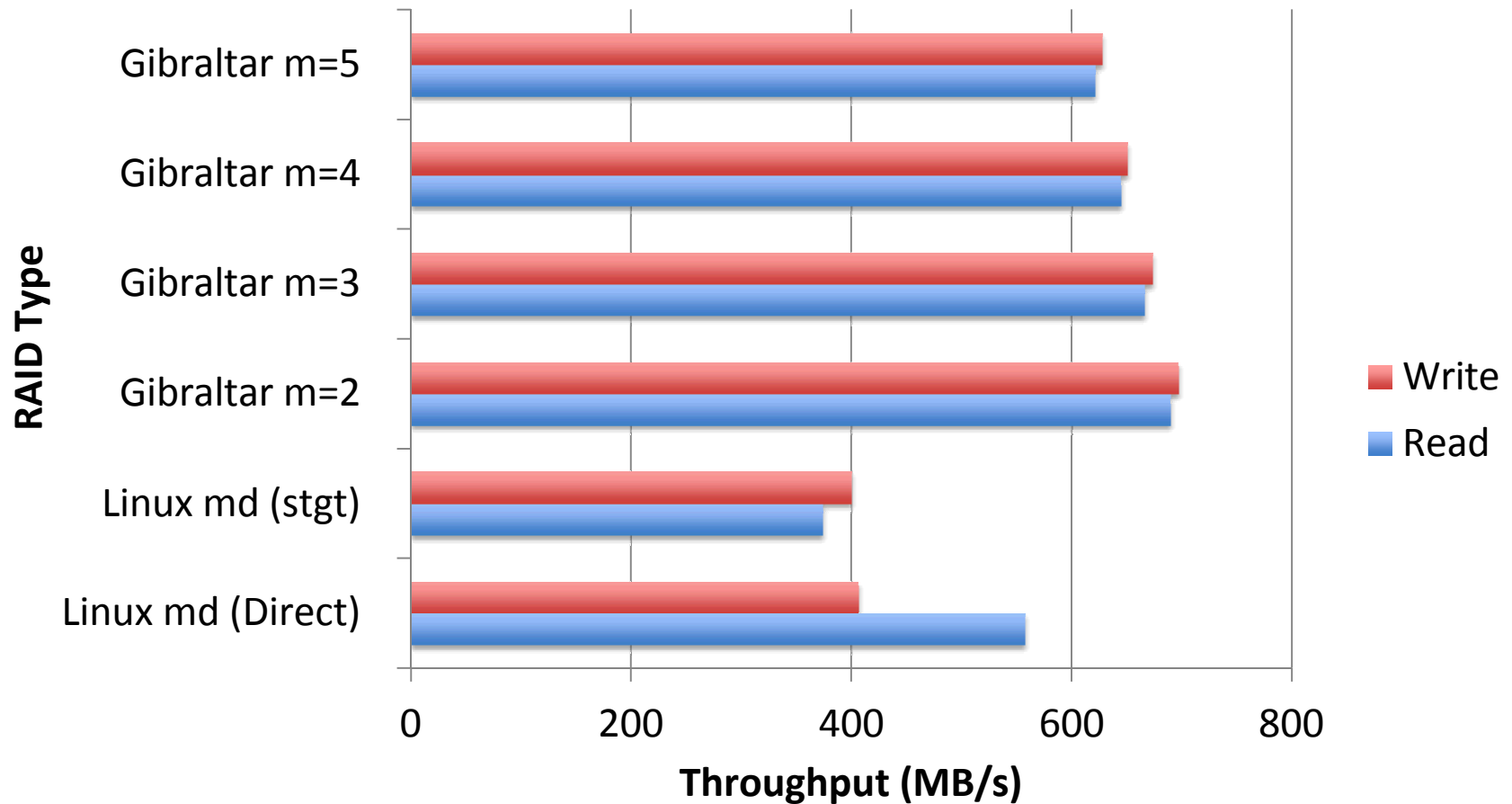


A User Space RAID Architecture

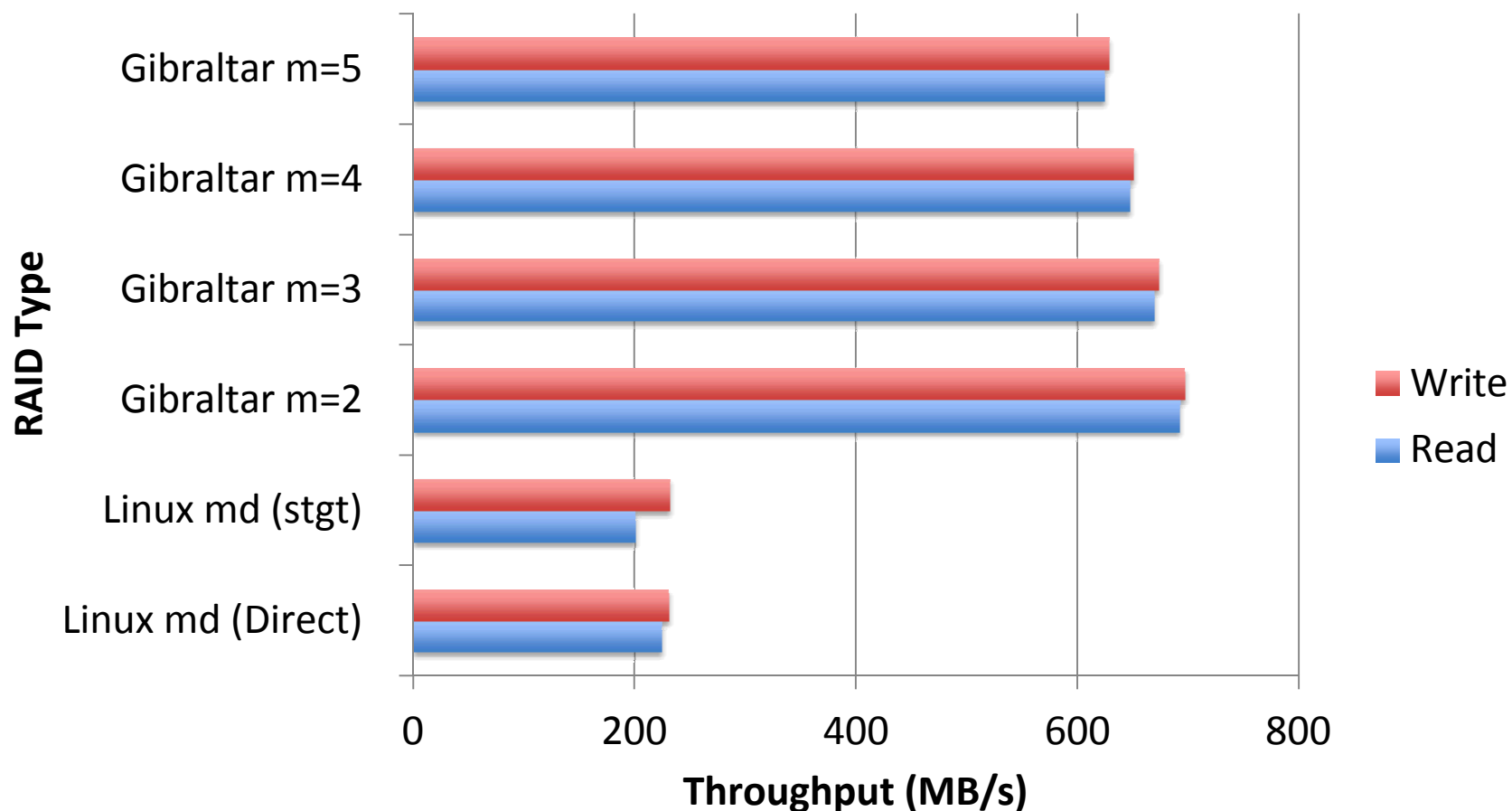


- GPU computing facilities are inaccessible from kernel space
- Provide I/O stack components in user space, accessible via iSCSI

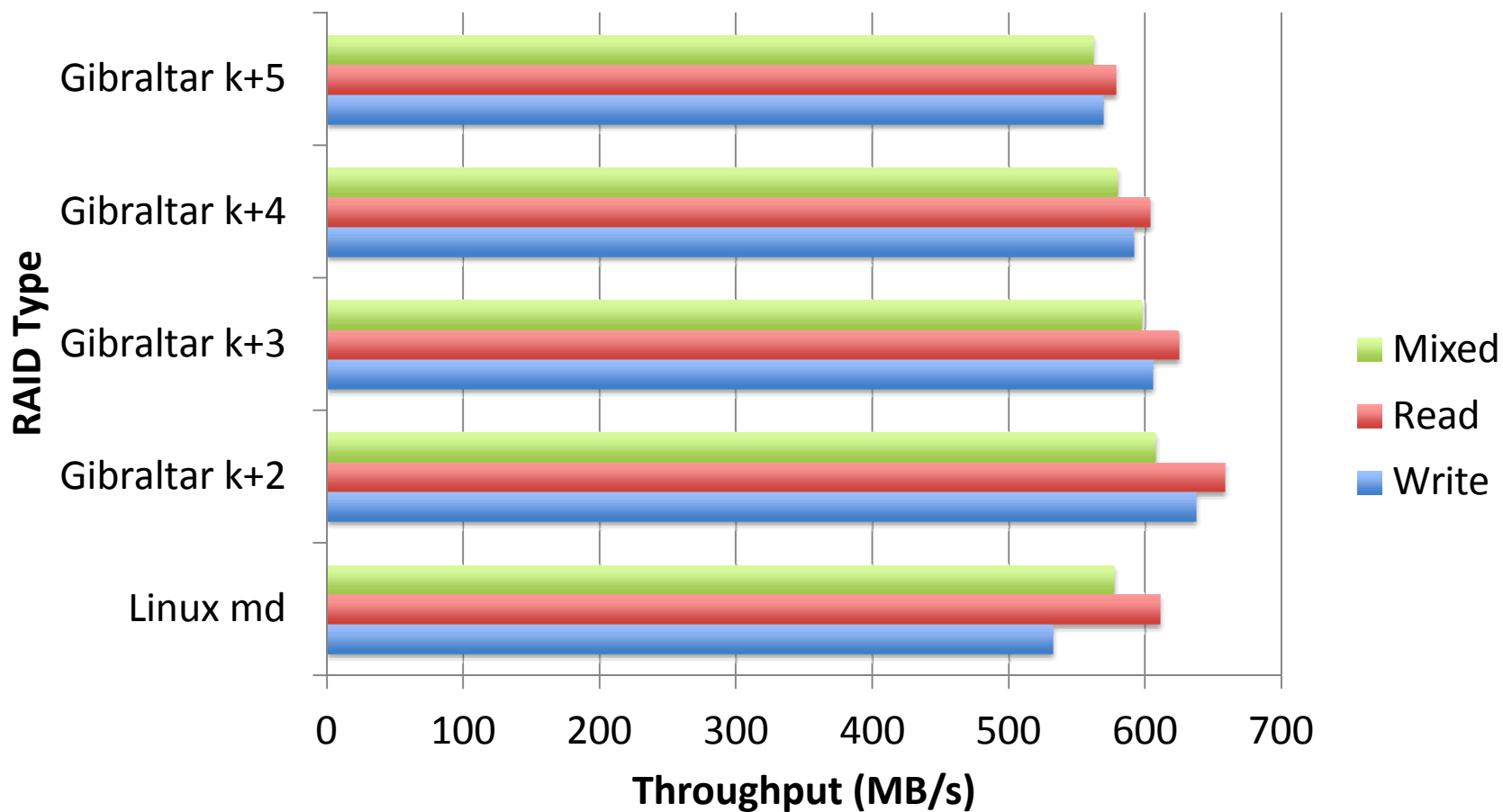
Direct Attached Storage



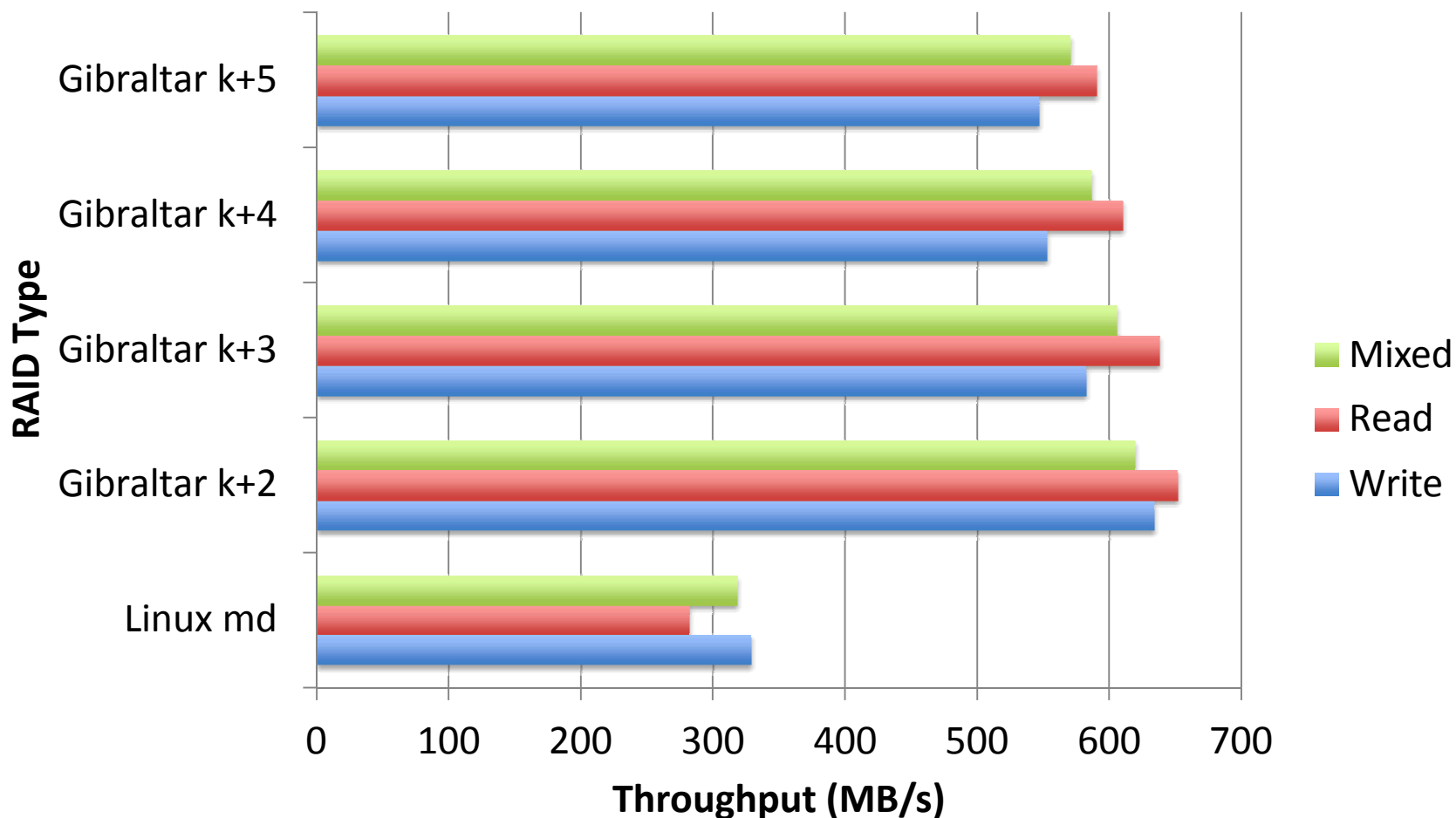
Direct Attached Storage, Degraded



Network Attached Storage, Four Clients



Network Attached Storage, Four Clients, Degraded





Combine RAID with Compression

- Before or After Reed-Solomon coding is performed, it is desirable to perform other operations as well
- Initial investigations into GPU compression with Robert Cloud have yielded insight into integration
- Huffman coding can be performed at 2.6 GB/s on a GPU



Combine RAID with Other Processing

- AES, allowing encryption on disk
 - “Design of a Parallel AES for Graphics Hardware using the CUDA framework,” Biago et al., IPDPS 2009
- SHA-1, for block-level deduplication or content-based addressing
 - “A Brief Implementation Analysis of SHA-1 on FPGAs, GPUs, and Cell Processors,” 2009 International Conference on Engineering and Computation



Fail in Place

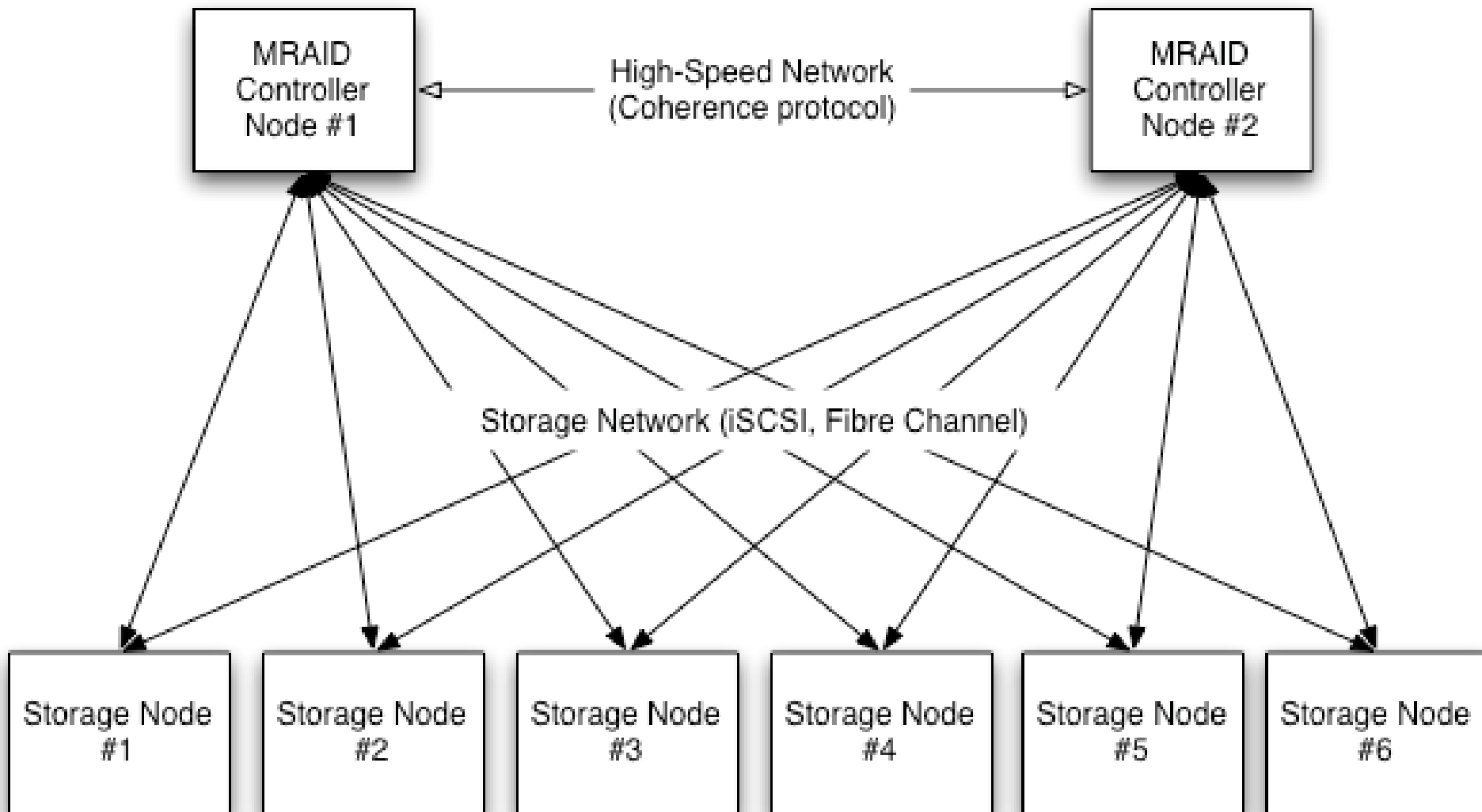
- Storage bricks in remote or harsh locations
- Degraded mode operation presents risks that affect hardware and software RAID
 - Lowered performance
 - Increased risk of multiple failures
 - Hot spare induced window of vulnerability
- High-parity RAID can mitigate the risks of degraded mode operation



Using MRAID for Data Center Resiliency

- MRAID is similar to RAID
 - Disks → Storage servers
- Multiple controllers
 - Multiple paths of access
- Node failures
 - Power supplies
 - Disks
 - Network

MRAID





Distributed Data Storage

- Current technologies currently use 3 x data replication, which is logically the same as RAID 1+0 with three sets
 - Google FS
 - Amazon S3
- User-managed distributed data storage can use high-speed coding
 - Multiple administrative domains
 - Low bandwidth overhead
 - Higher reliability



Conclusions

- Introduced a high-parity RAID variant that is impractical in software today
- Prototyped a GPU-based library and RAID controller that can provide high parity
 - High speed
 - Low cost
 - High flexibility
- Detailed future major research areas
 - Processing in storage stack
 - Embedded/low maintenance storage
 - MRAID