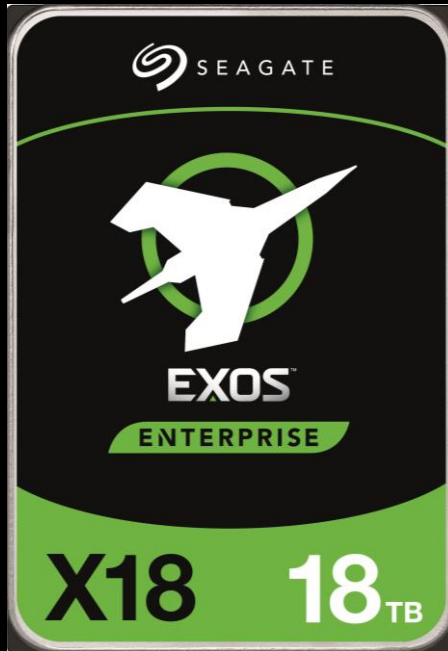




Protecting Data with Declustered Parity

John Bent, Seagate Senior Director



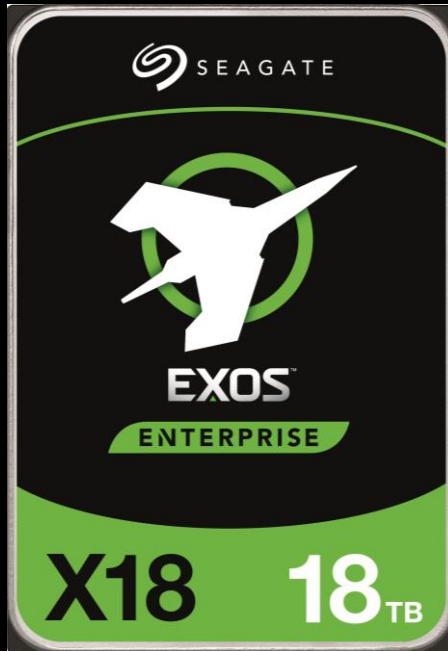


This is your drive



This is your precious





This is your drive

Store your precious
on your drive.
Happiness forever.



This is your precious



Happiness Forever?



Unfortunately, disks fail

Some of John's Favorite Papers about Storage Failure

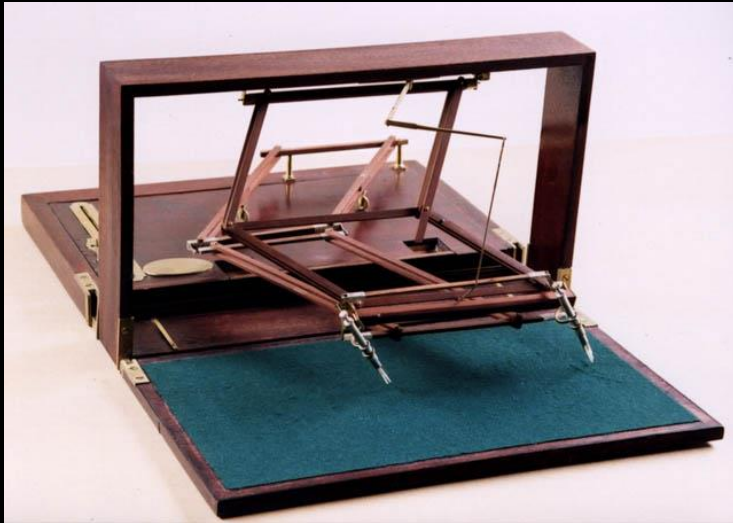
- Availability in Globally Distributed Storage Systems
 - <https://static.googleusercontent.com/media/research.google.com/en//pubs/archive/36737.pdf>
- An Analysis of Data Corruption in the Storage Stack
 - <https://www.cs.toronto.edu/~bianca/papers/fast08.pdf>
- Disk failures in the real world
 - <http://www.cs.toronto.edu/~bianca/papers/fast07.pdf>
- Mean time to meaningless
 - https://www.usenix.org/legacy/event/hotstorage10/tech/full_papers/Greenan.pdf

Also many good blog posts at backblaze.

This is Not Just a Modern Problem



As long as the world has stored data, the world has needed mechanisms to protect that data from the inevitable failure of the storage medium.

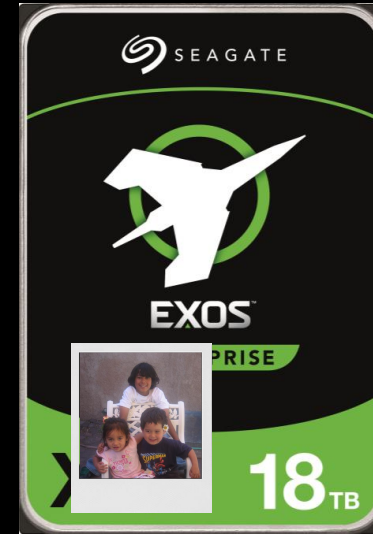
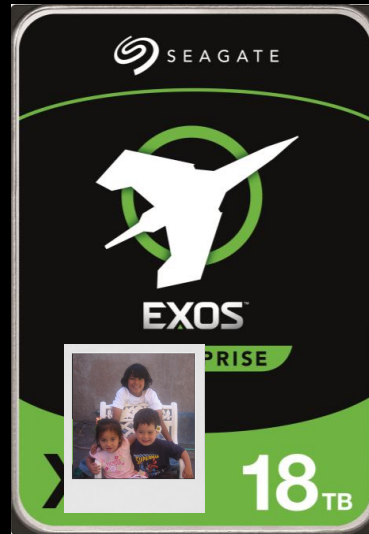
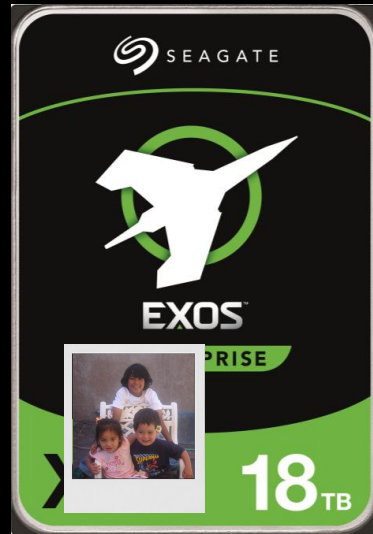


Polygraph: create a copy of your document as you write it.



Duplicating books before the invention of the printing press.

Replication Has Been Around For a Long Time



But high capacity overhead....

RAID: Modern Technique to Gain Durability without High Capacity Overhead

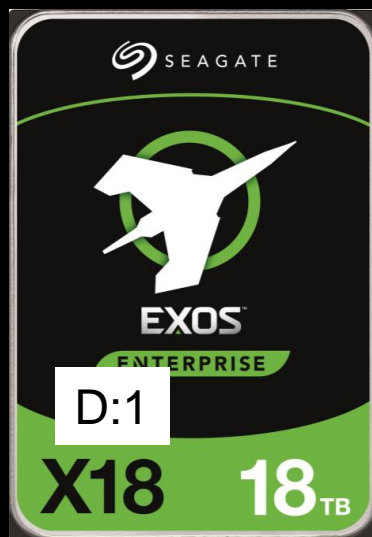
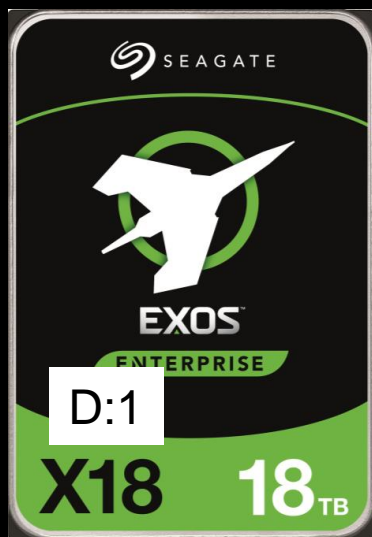


A case for redundant arrays of inexpensive disks (RAID)

[DA Patterson](#), [G Gibson](#), [RH Katz](#) - Proceedings of the 1988 ACM ..., 1988 - dl.acm.org
Increasing performance of CPUs and memories will be squandered if not matched by a similar performance increase in I/O. While the capacity of Single Large Expensive Disks (SLED) has grown rapidly, the performance improvement of SLED has been modest ...

★ 99 Cited by 4397 Related articles All 146 versions 99

RAID5: Simple XOR Logic to Compute P(arity) from D(ata)



Terminology

Chunk: the contiguous piece of data stored to a single device

Chunk size: you can guess this one

Data chunk: also this one

Parity chunk: also this one

Stripe: the collection of chunks that form a contiguous logical unit of data

Layout: the particular devices on which a stripe is stored

Pattern: the shared configuration of all stripes (e.g. 8+2 or 3-way replication)

Full-stripe write: an important concept for performance (not discussed here)

Partial-stripe write: a write that is smaller than a full-stripe write

Read-modify-write: the performance inhibition caused by partial-stripe write

Pool: for parity declustered (not yet explained), all devices on which a stripe could be stored

How Long Does Rebuild Take?

Find replacement drive

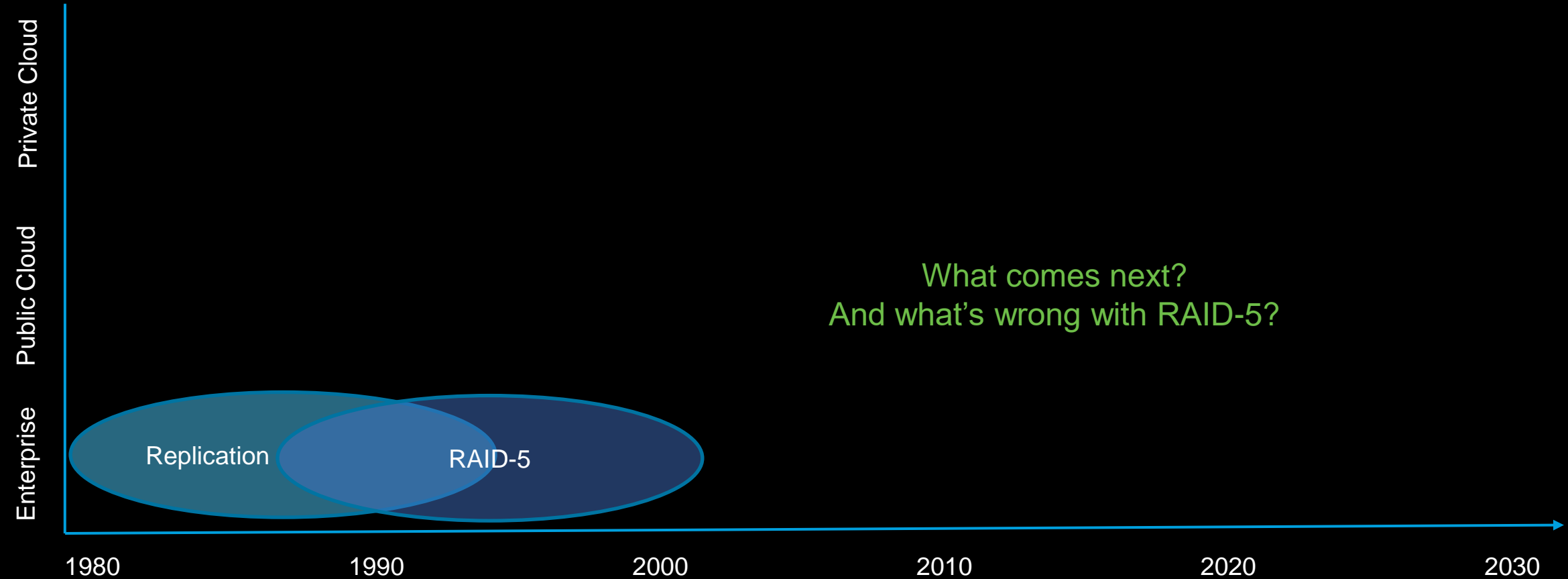
For each *chunk* that was on failed drive

- Read other chunks from that chunk's stripe from the sibling drives in the stripe layout
- Use parity to reconstruct the lost unit
- Write the lost unit to the replacement drive

$$t(\text{rebuild}) \approx t(\text{write entire drive})$$

$$t(\text{write entire drive}) = \text{capacity} / \text{bandwidth}$$

The Evolution of Data Durability in Data Centers



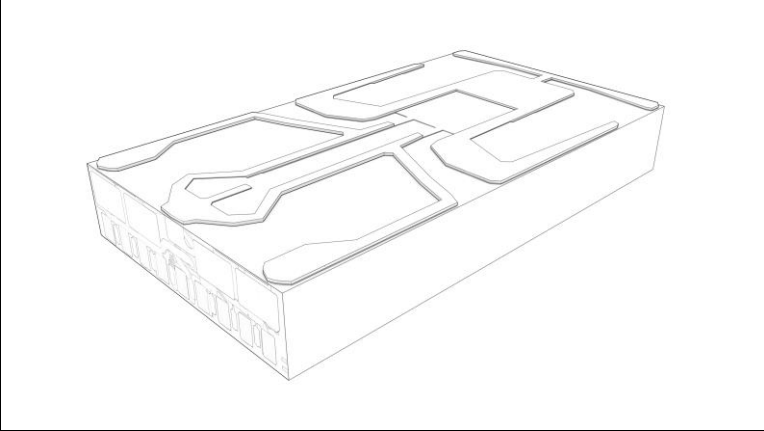
An Important Characteristic of Disks

$$\textit{growth}(\textit{bandwidth}) \approx \sqrt{\textit{growth}(\textit{capacity})}$$

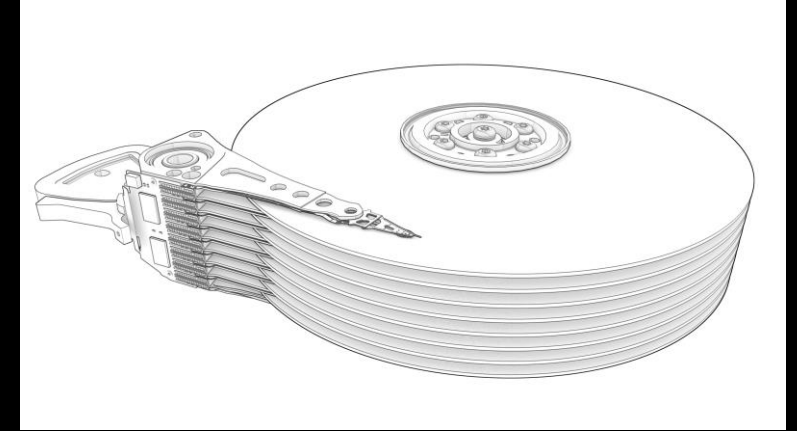
- alternately -

$$\textit{growth}(\textit{capacity}) \approx \textit{growth}(\textit{bandwidth})^2$$

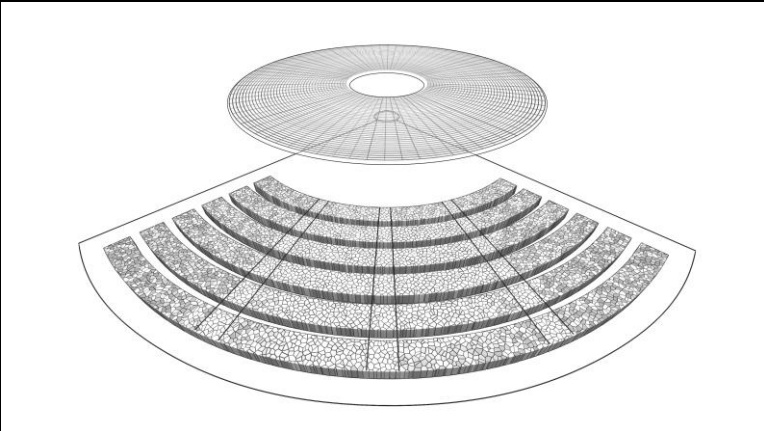
How to Increase Disk Capacity



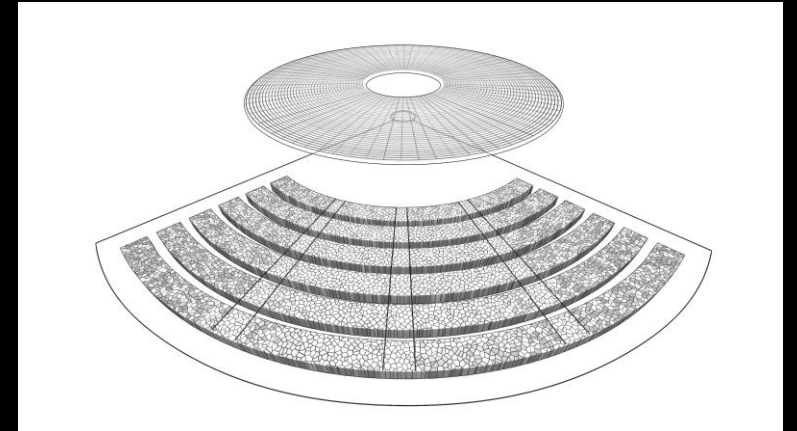
Increase the Form Factor



Insert more Platters



More Tracks per Surface

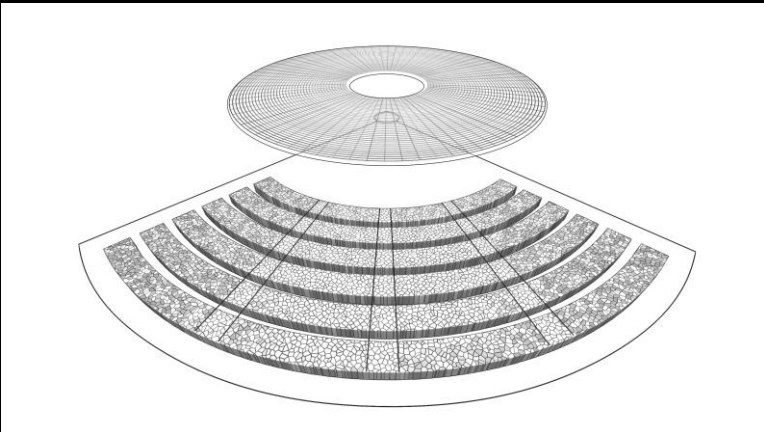


More Bits Per Track

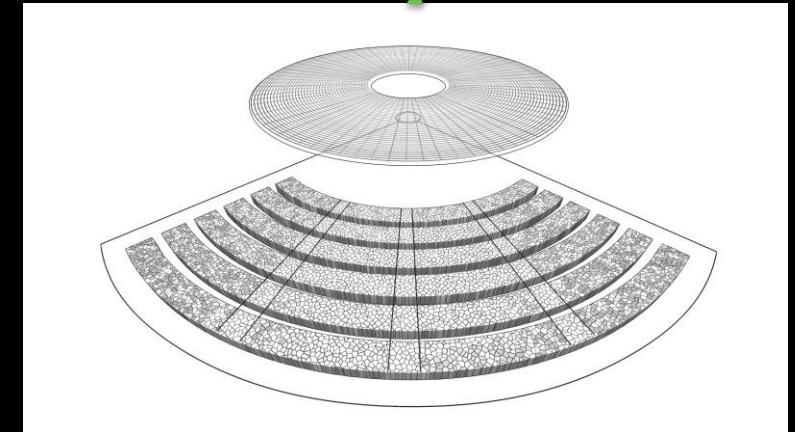
How to Increase Disk Bandwidth



Spin Faster



Access More Tracks
Simultaneously



More Bits Per Track

SAME!

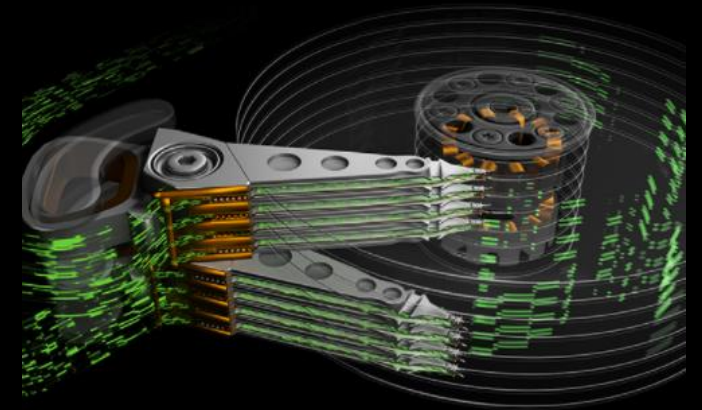
$$growth(bandwidth) \approx \sqrt{growth(capacity)}$$

How to Increase Capacity

- ~~Increase the Form Factor~~
- ~~Insert more platters~~
- More bits per track
- More tracks per surface

How to Increase Bandwidth

- ~~Spin the surface faster~~
- ~~Access more tracks simultaneously~~
- More bits per track



MACH·2™

$$growth(bandwidth) \approx \sqrt{growth(capacity)}$$

How to Increase Capacity

- ~~Increase the Form Factor~~
- ~~Insert more platters~~
- More bits per track
- More tracks per surface

How to Increase Bandwidth

- ~~Spin the surface faster~~
- ~~Access more tracks simultaneously~~
- More bits per track

$$bandwidth(disk_{n+1}) \approx bandwidth(disk_n) * growth(bits \ per \ track)$$

$$capacity(disk_{n+1}) \approx capacity(disk_n) * growth(bits \ per \ track) * growth(tracks \ per \ surface)$$

$$growth(tracks \ per \ surface) \approx growth(bits \ per \ track)$$

$$capacity(disk_{n+1}) \approx capacity(disk_n) * growth(bits \ per \ track)^2$$

$$growth(bandwidth) \approx \sqrt{growth(capacity)}$$

$$t(\textit{write entire drive}) = \textit{capacity} / \textit{bandwidth}$$

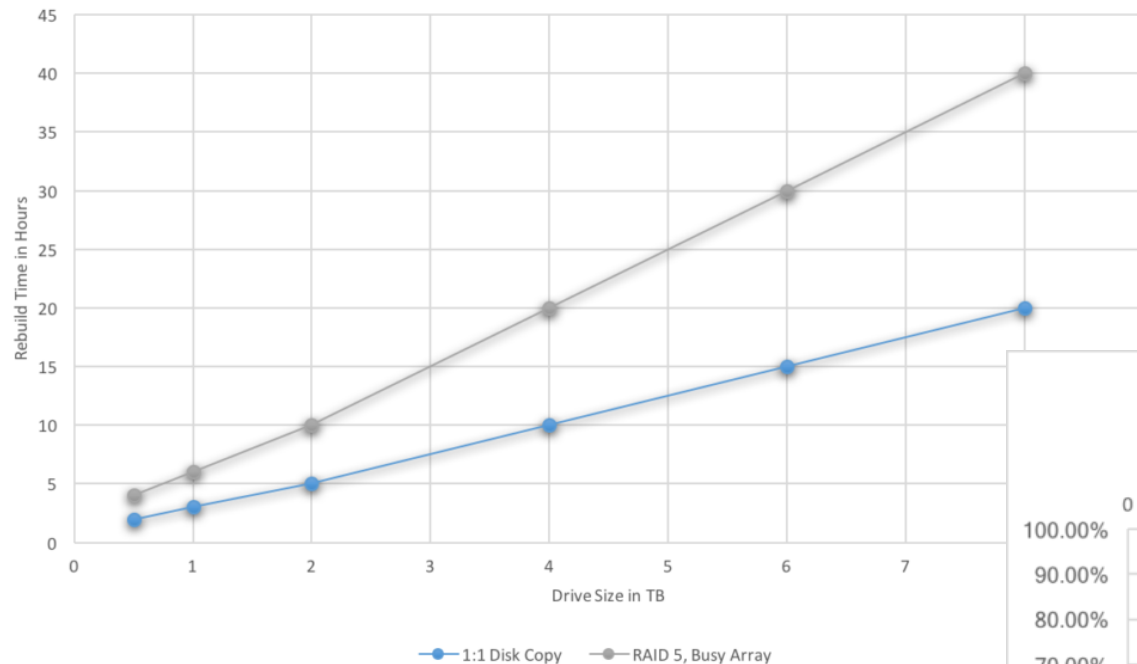
PLUS

$$\textit{growth}(\textit{bandwidth}) \approx \sqrt{\textit{growth}(\textit{capacity})}$$

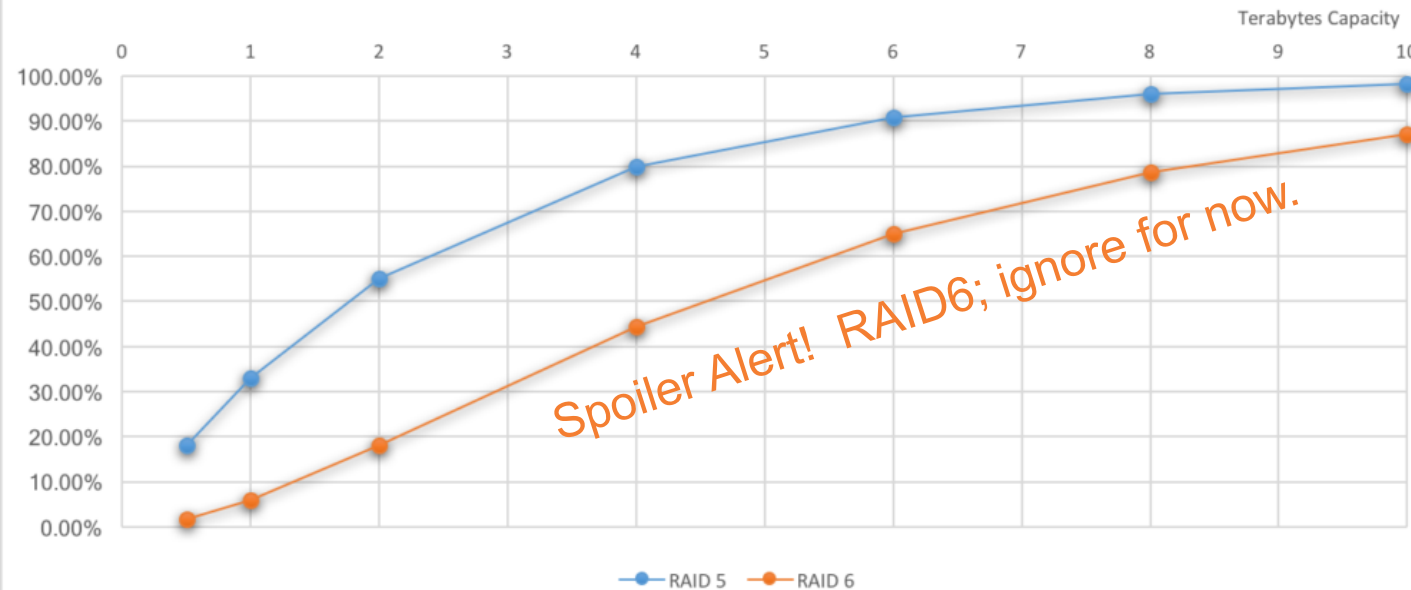
EQUALS

$$\textit{Probability}(\textit{DataLoss}, \textit{RAID5}) \Rightarrow 1$$

Time to Rebuild (Hours)



Statistical Probability of Failure During Rebuild

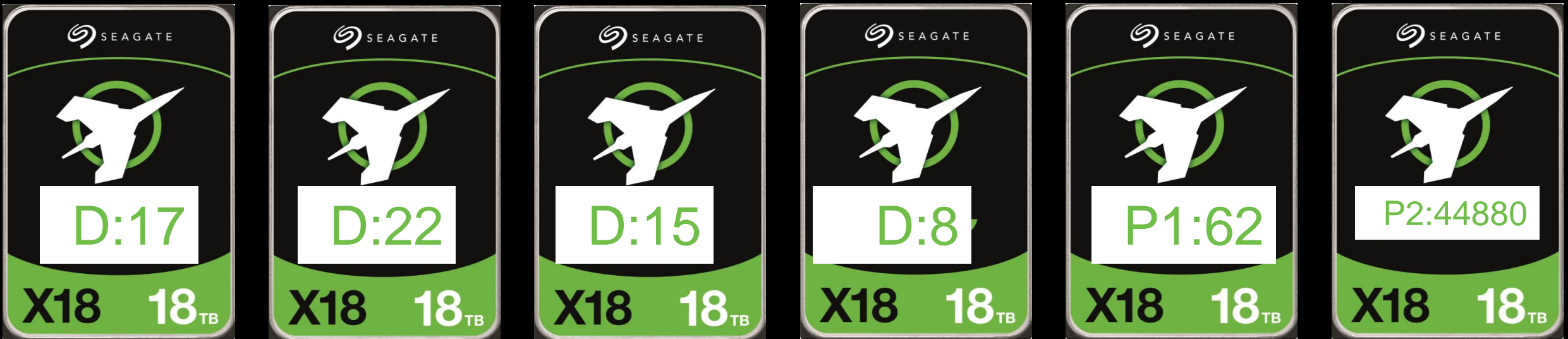


What Comes After RAID5? RAID6!



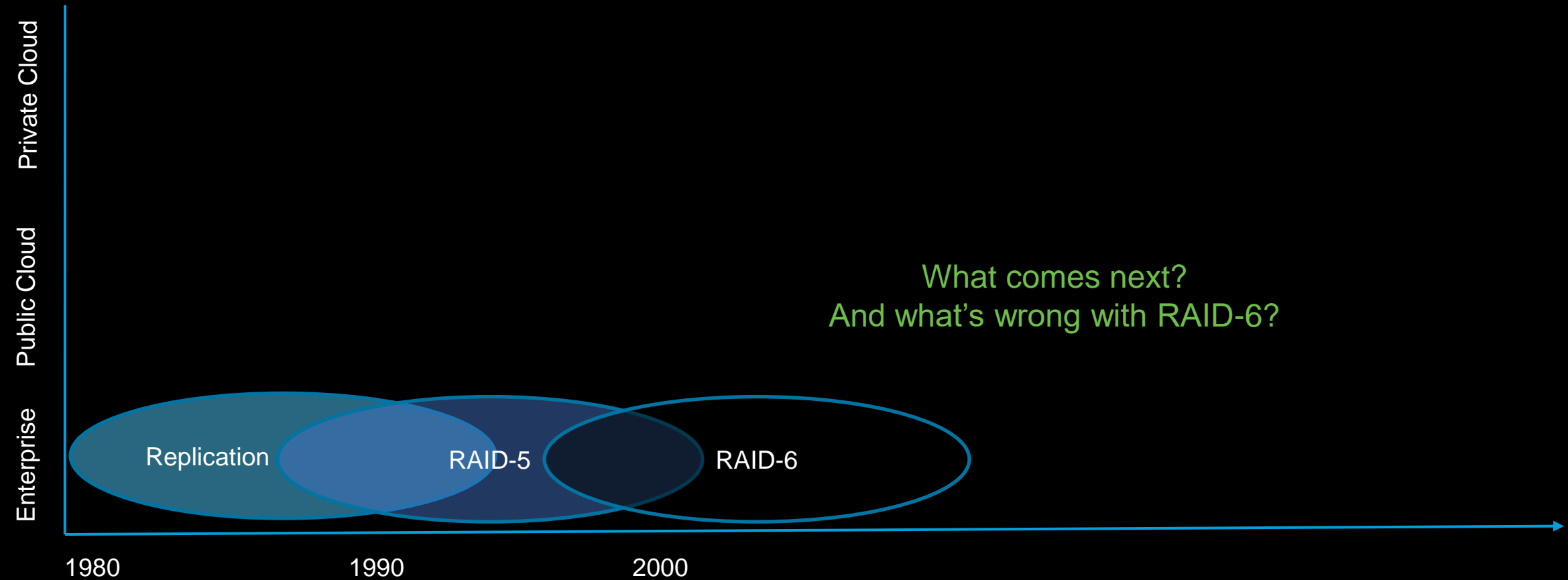
RAID6 Logic? High School Algebra! Two Equations, Two Unknowns.

* System understanding != Math understanding



$$\begin{array}{l} 62 = 17 + x + 15 + y \\ 44880 = 17 * x * 15 * y \end{array} \xrightarrow{\text{(math)}} x = 22 ; y = 8$$

The Evolution of Data Durability in Data Centers



$$t(\textit{write entire drive}) = \textit{capacity} / \textit{bandwidth}$$

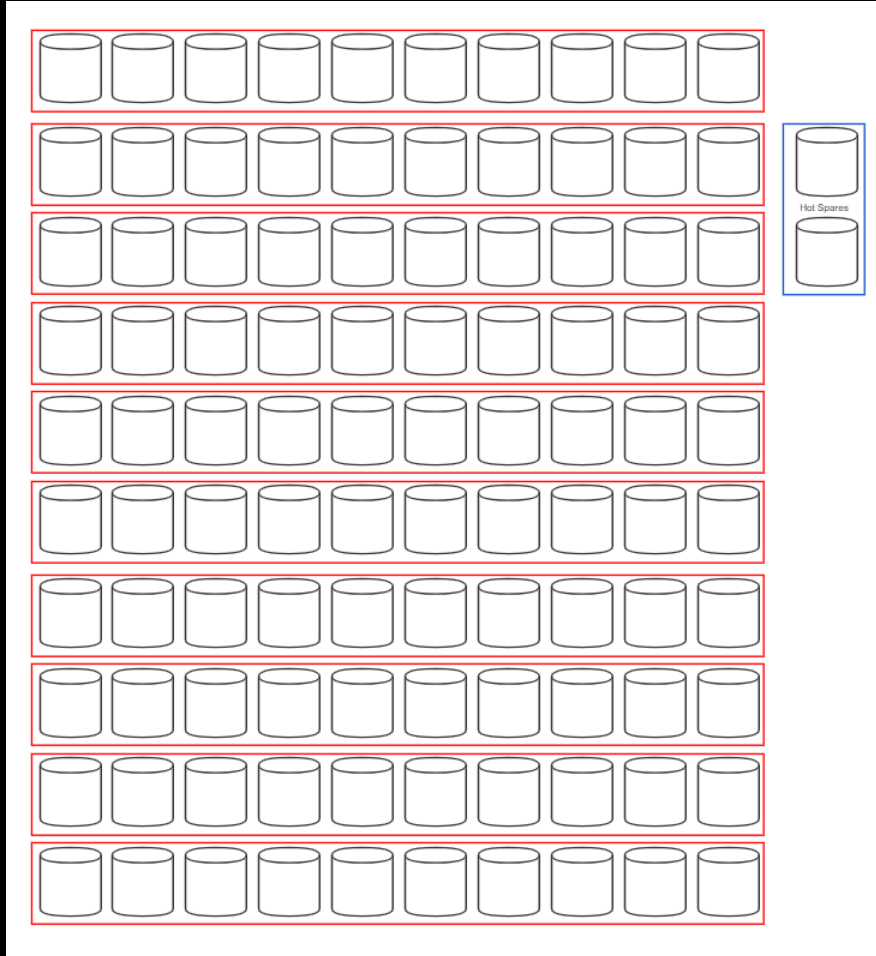
PLUS

$$\textit{growth}(\textit{bandwidth}) \approx \sqrt{\textit{growth}(\textit{capacity})}$$

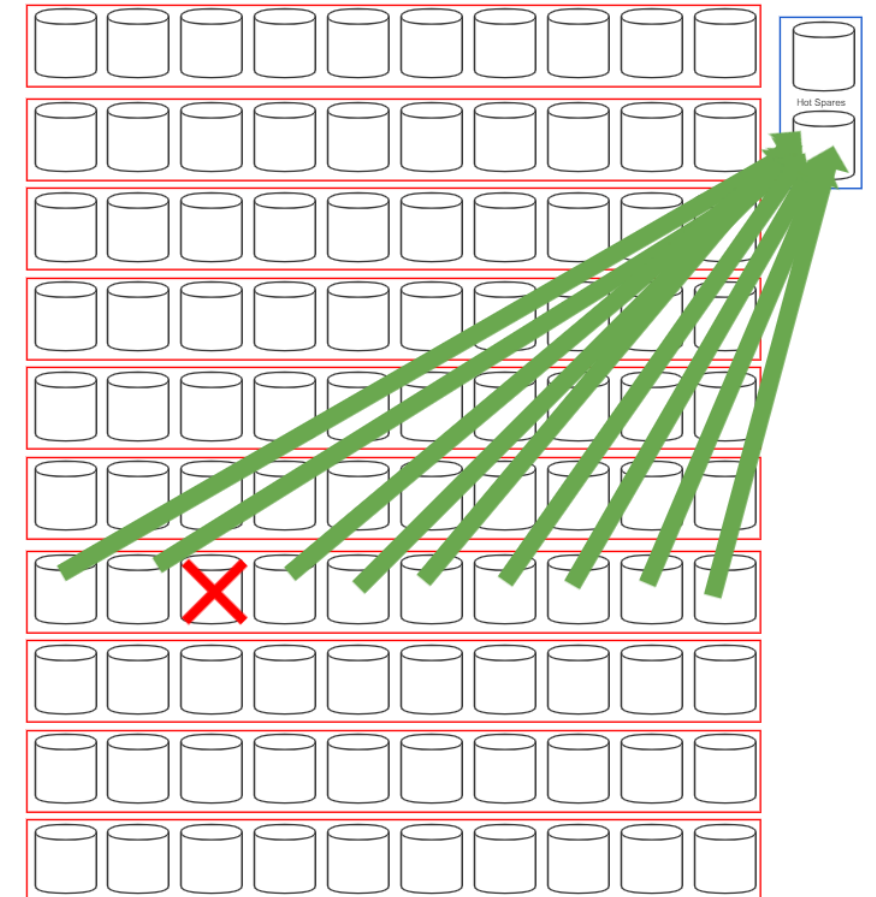
EQUALS

$$\textit{Probability}(\textit{DataLoss}, \textit{RAID6}) \Rightarrow 1$$

Finally, Onto Declustered! (wait, not quite)

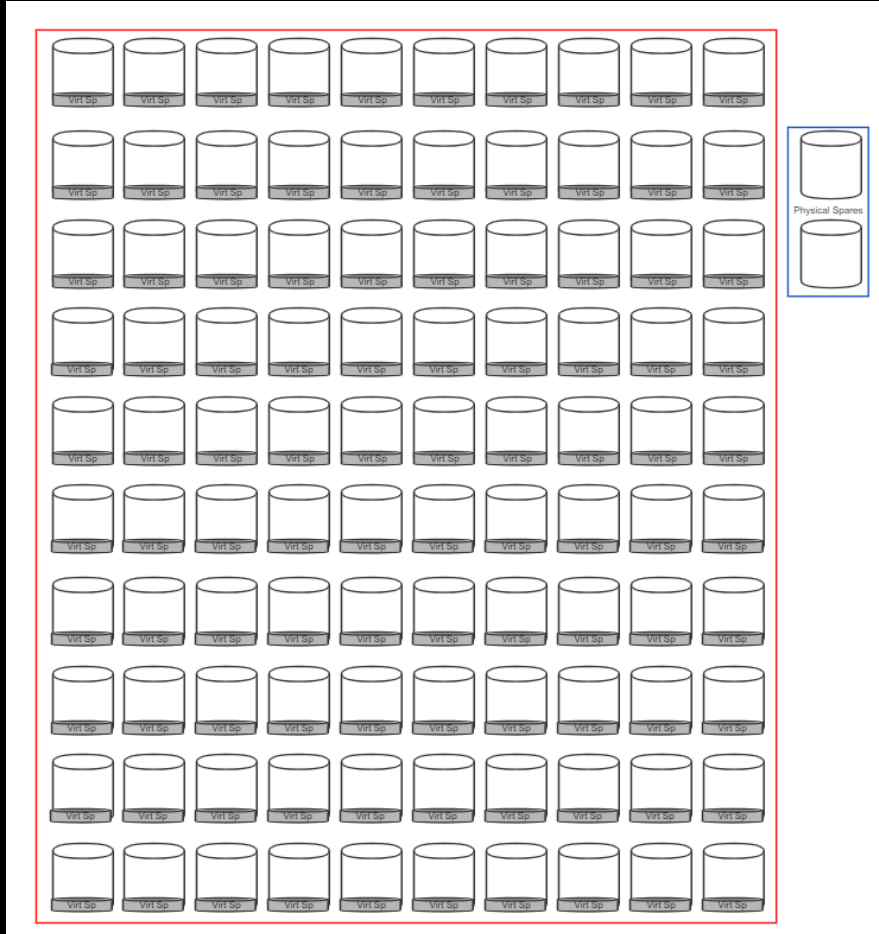


FAILURE!



100 disks, fixed groups, 2 spares

Finally, Onto Declustered!



Declustered Parity Key Concepts

Virtual spares

- $N+K \rightarrow N+K+S$

Replay Phases

- Following the rebuild Phase

Pseudo-Random Layout

- Every stripe goes to a different random $N+K$ drives

Rebuild is all-all

- E.g. 99 to 99

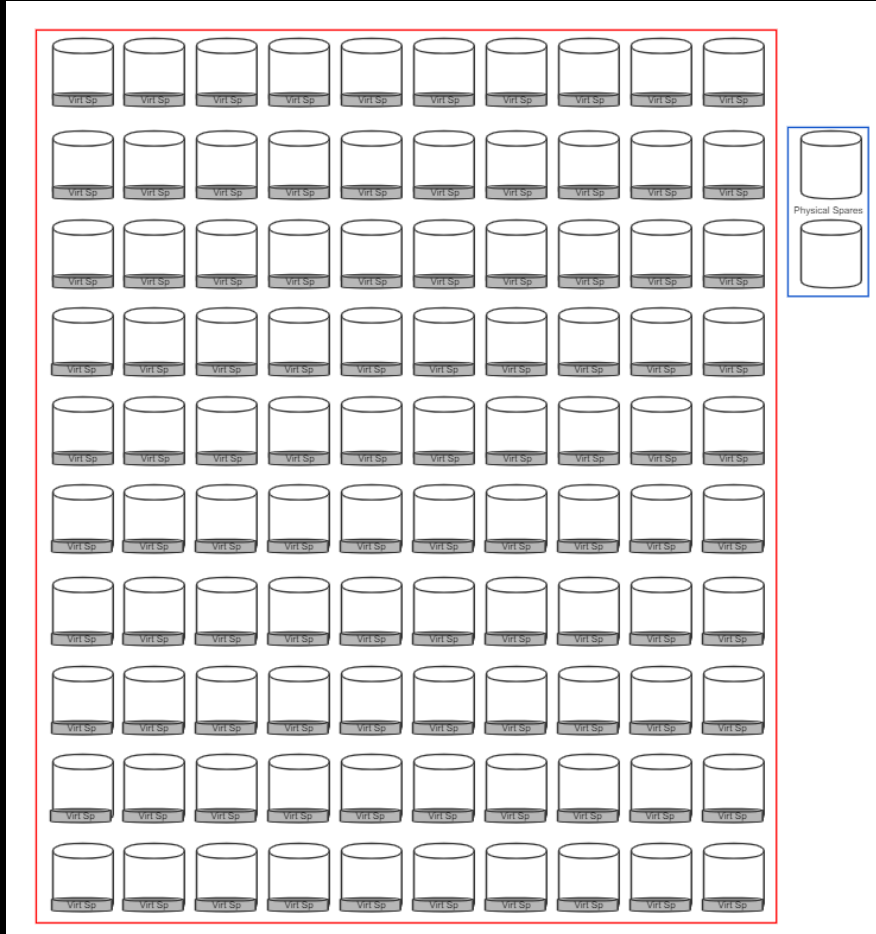
Rebuild time is time to write a disk's worth of data across all

- Scales with total population of drives

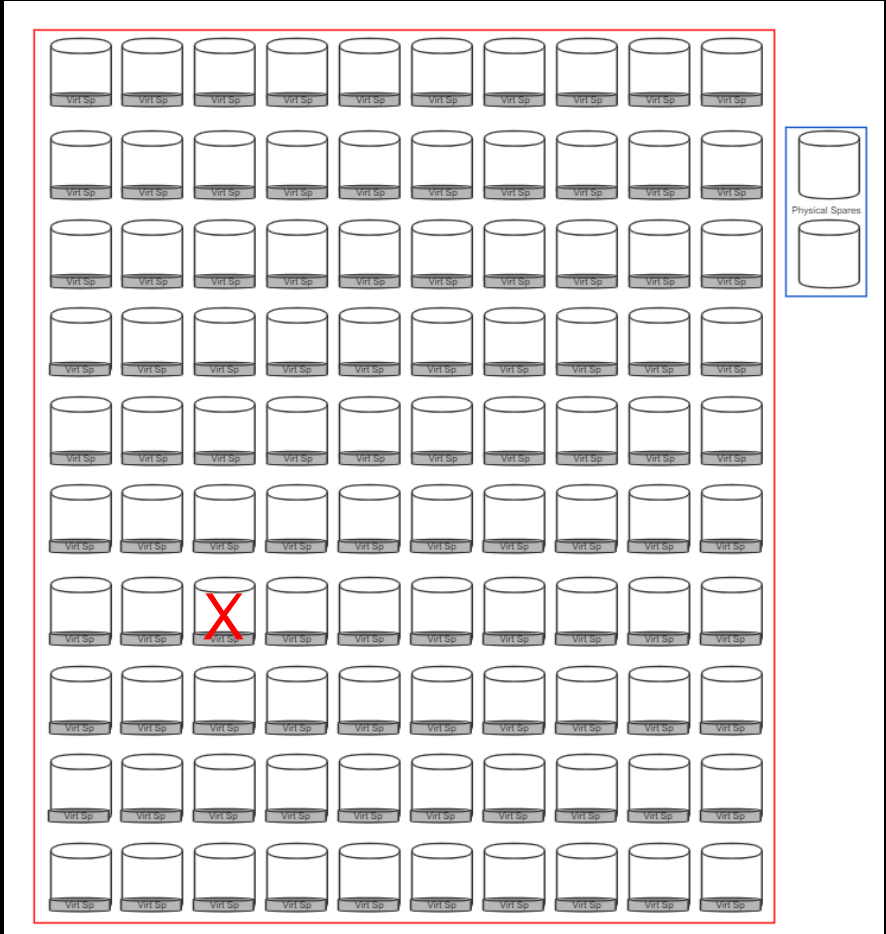
100 disks, one declustered groups, 2 spares

Finally, Onto Declustered!

Not drawing the 9702 arrows! Sorry!

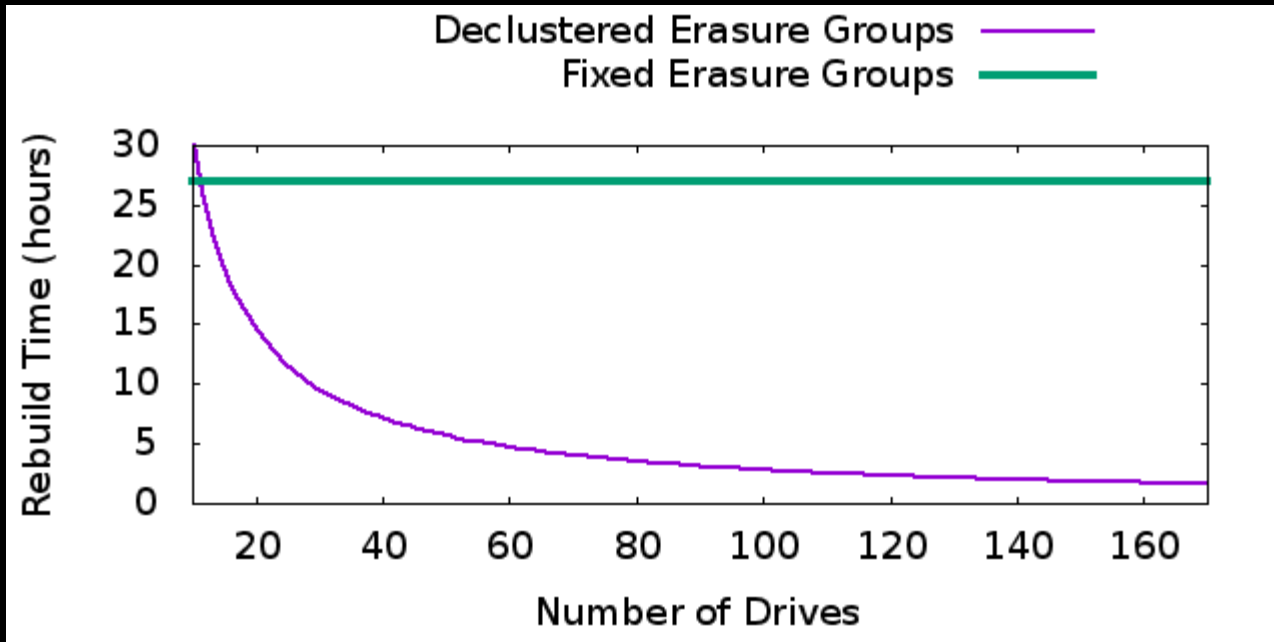


FAILURE!



100 disks, one declustered groups, 2 spares

Rebuild Time Scales with Population



(<https://cortex.link/declustered>)

```
set terminal png font "Calibri,14" size 640,320
```

```
capacity_tb=20  
capacity_mb=capacity_tb*1000*1000  
hdd_mbps=200  
seconds_in_hour=3600
```

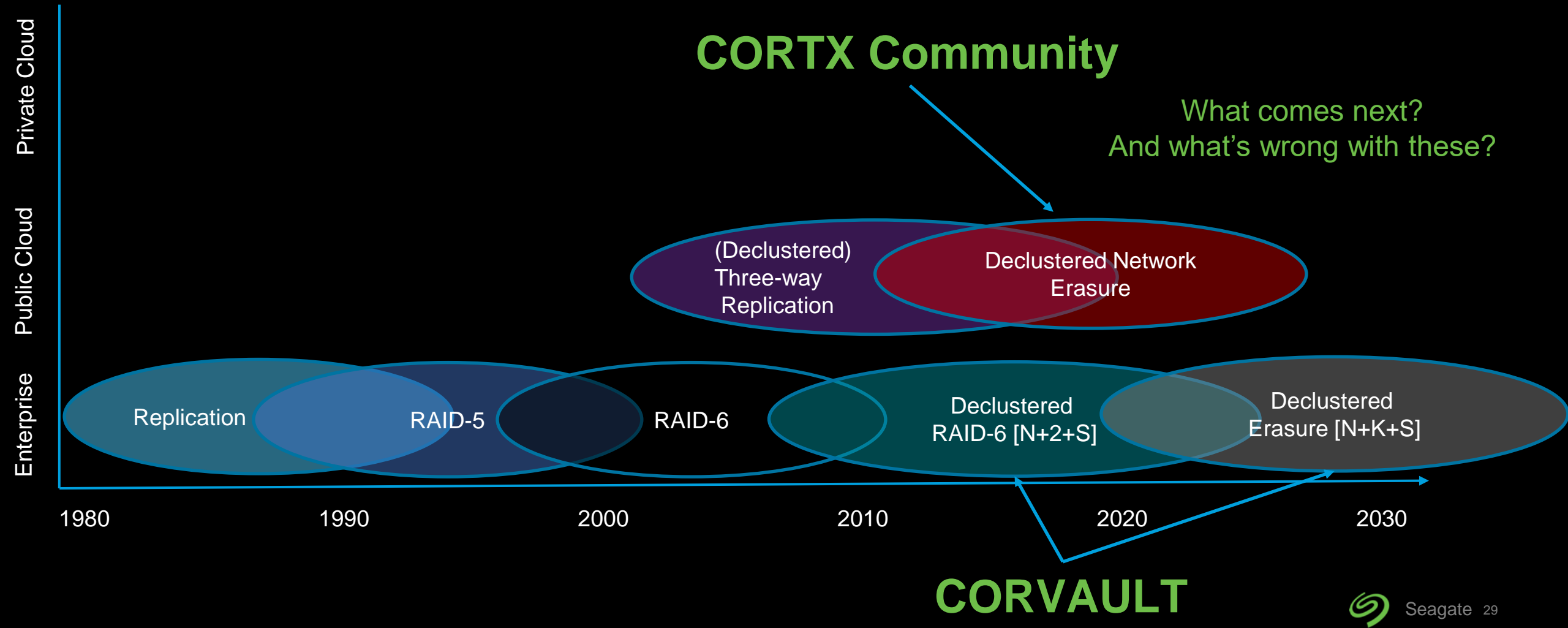
```
amplification=10 # need to read 8 pieces and then write one  
# for legacy, ignore amplification since you are just bound to one drive
```

```
dcr_rebuild_rate(x)=hdd_mbps*(x-1) # subtract out the failed drive  
dcr_rebuild_hours(x)=((capacity_mb*amplification)/dcr_rebuild_rate(x))  
legacy_rebuild_hours(x)=(capacity_mb/hdd_mbps)/3600
```

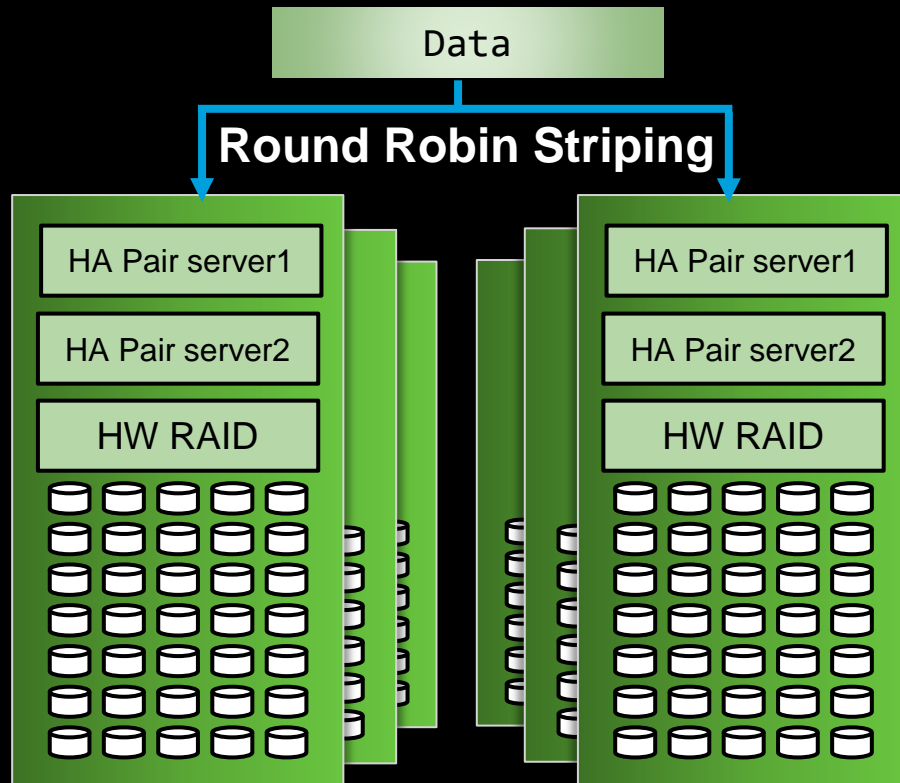
```
set yrange [0:30]  
set xrange [10:170]
```

```
set key above  
unset title  
set ylabel 'Rebuild Time (hours)'  
set xlabel 'Number of Drives'  
set output 'declustered_rebuild.png'  
plot dcr_rebuild_hours(x) lw 2 t "Declustered Erasure Groups", legacy
```

The Evolution of Data Durability in Data Centers

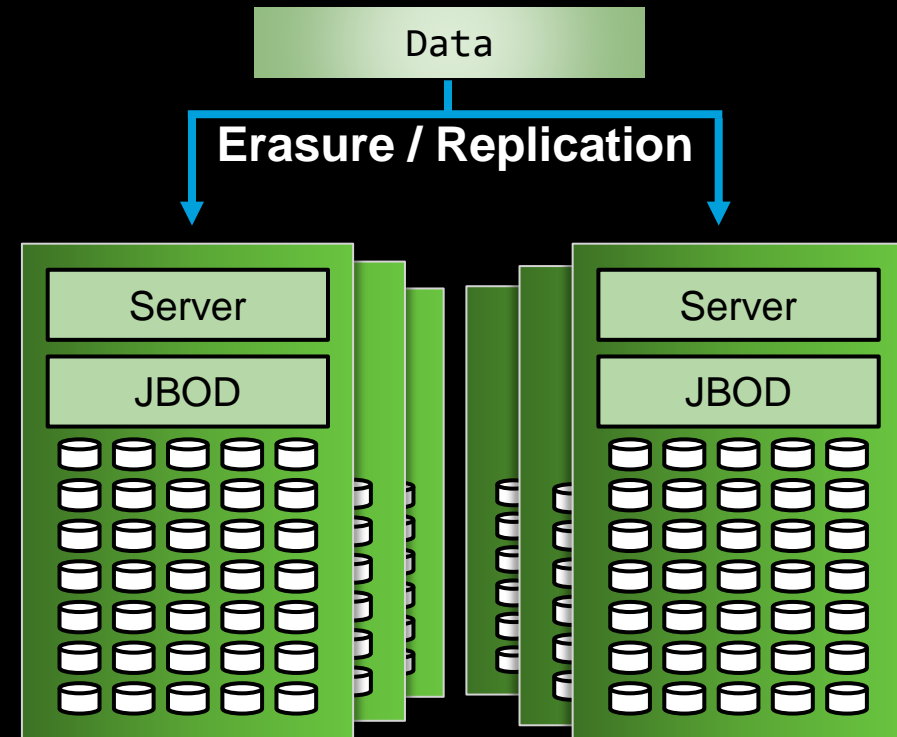


Two Existing Approaches for Data Durability and Availability



HW Local Declustered Erasure

Lustre, PVFS, BeeGFS, etc.



SW Network Declustered Erasure

Cloudian, ActiveScale, SwiftStack, Ceph, HDFS, etc.

Extreme Scale Requires Extreme Protection

Availability in Globally Distributed Storage Systems

Daniel Ford, François Labelle, Florentina I. Popovici, Murray Stokely, Van-Anh Truong*,
Luiz Barroso, Carrie Grimes, and Sean Quinlan
{ford, flab, florentina, mstokely}@google.com, vatrung@ieor.columbia.edu
{luiz, cgrimes, sean}@google.com
Google, Inc.

SPATIAL FAILURE BURST

Multiple simultaneous drive failures within a single rack.
Protect against these with erasure across enclosures.
Parity within enclosure is insufficient.

ASPATIAL FAILURE BURST

Multiple simultaneous drive failures across multiple racks.
Protect against these with erasure within enclosures.
Erasure across enclosures is insufficient..

*No single tier of parity can protect against all these failures.
Google knows about this and presumably had a team of PhD's
implement a solution. Private cloud needs our help to solve this;
we can do so with CORVAULT & tiered erasure.*

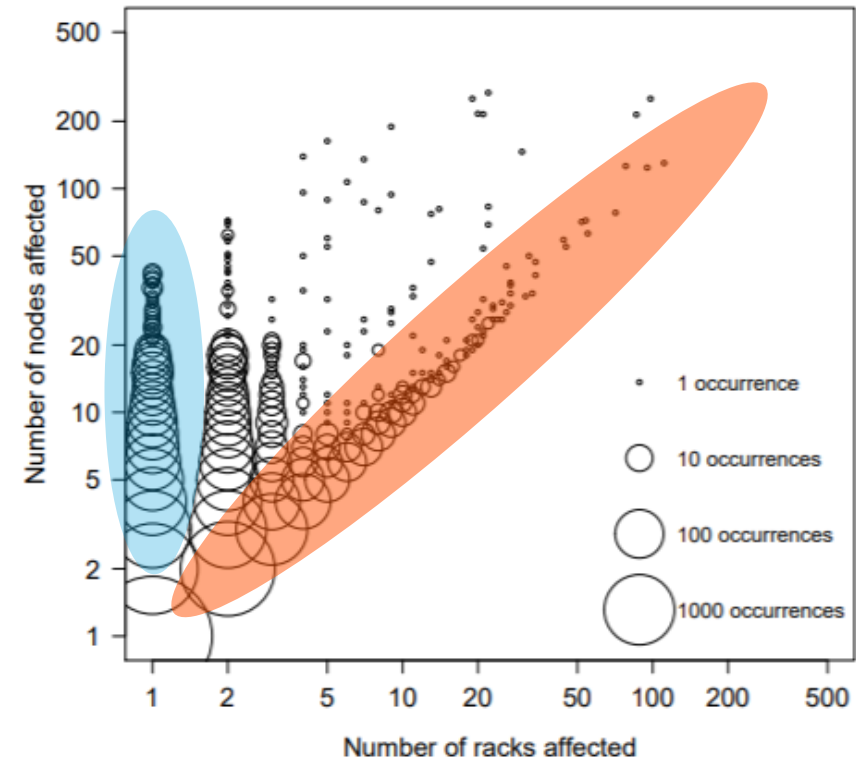
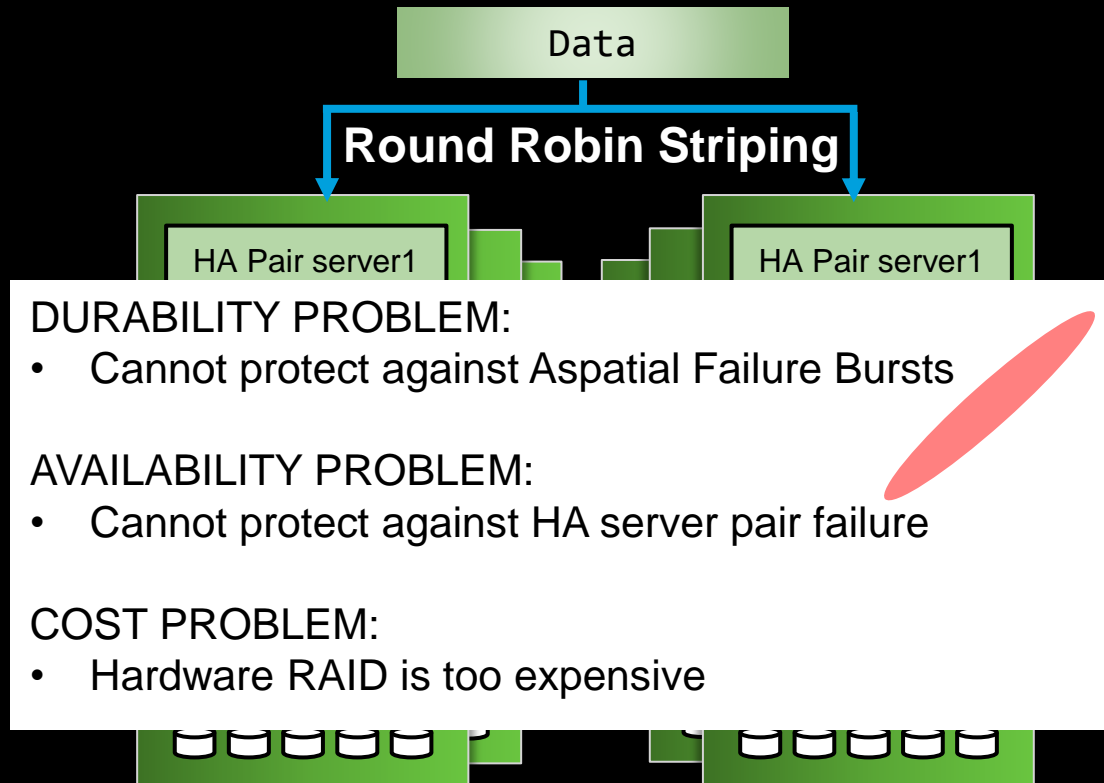


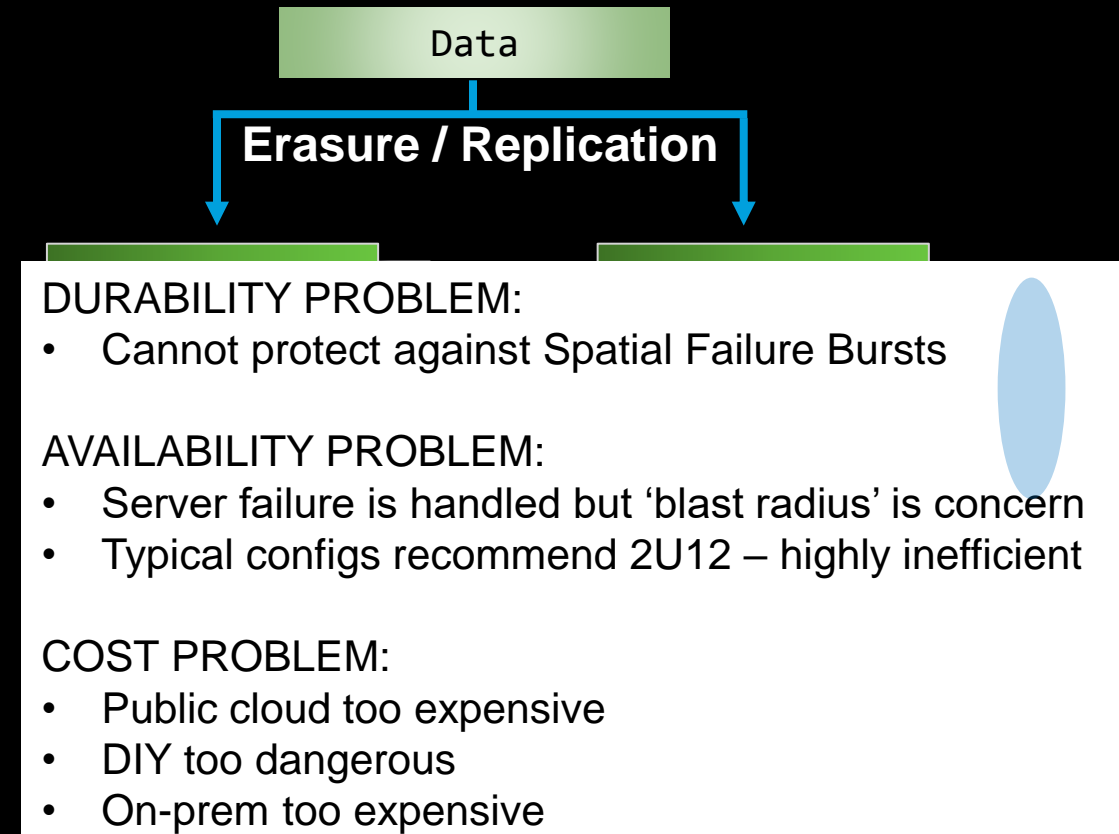
Figure 8: Frequency of failure bursts sorted by racks and nodes affected.

Two Existing Approaches for Data Durability and Availability



HW Local Declustered Erasure

Lustre, PVFS, BeeGFS, etc.

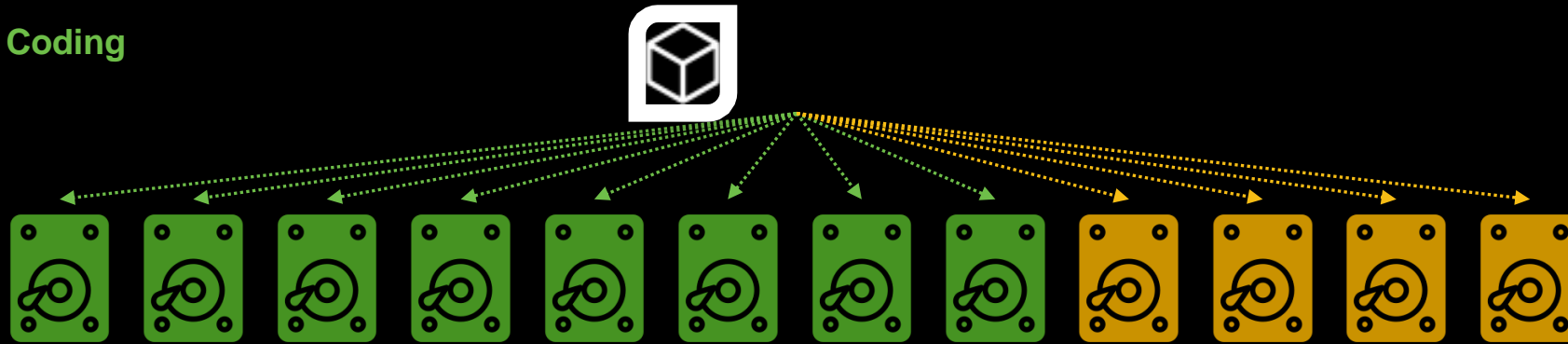


SW Network Declustered Erasure

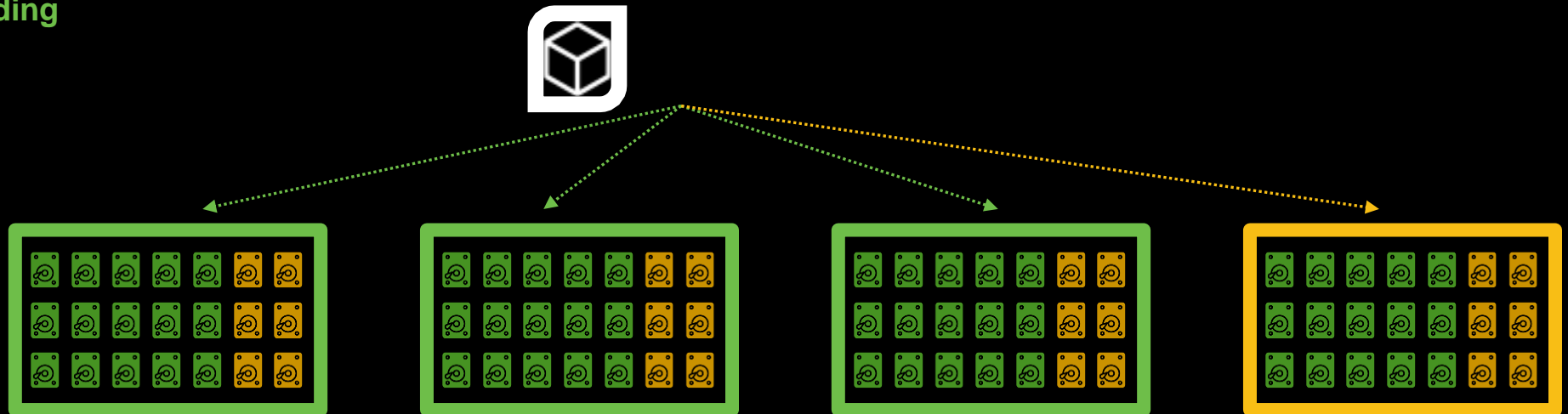
Cloudian, ActiveScale, SwiftStack, Ceph, HDFS, etc.

Multilevel Erasure Coding

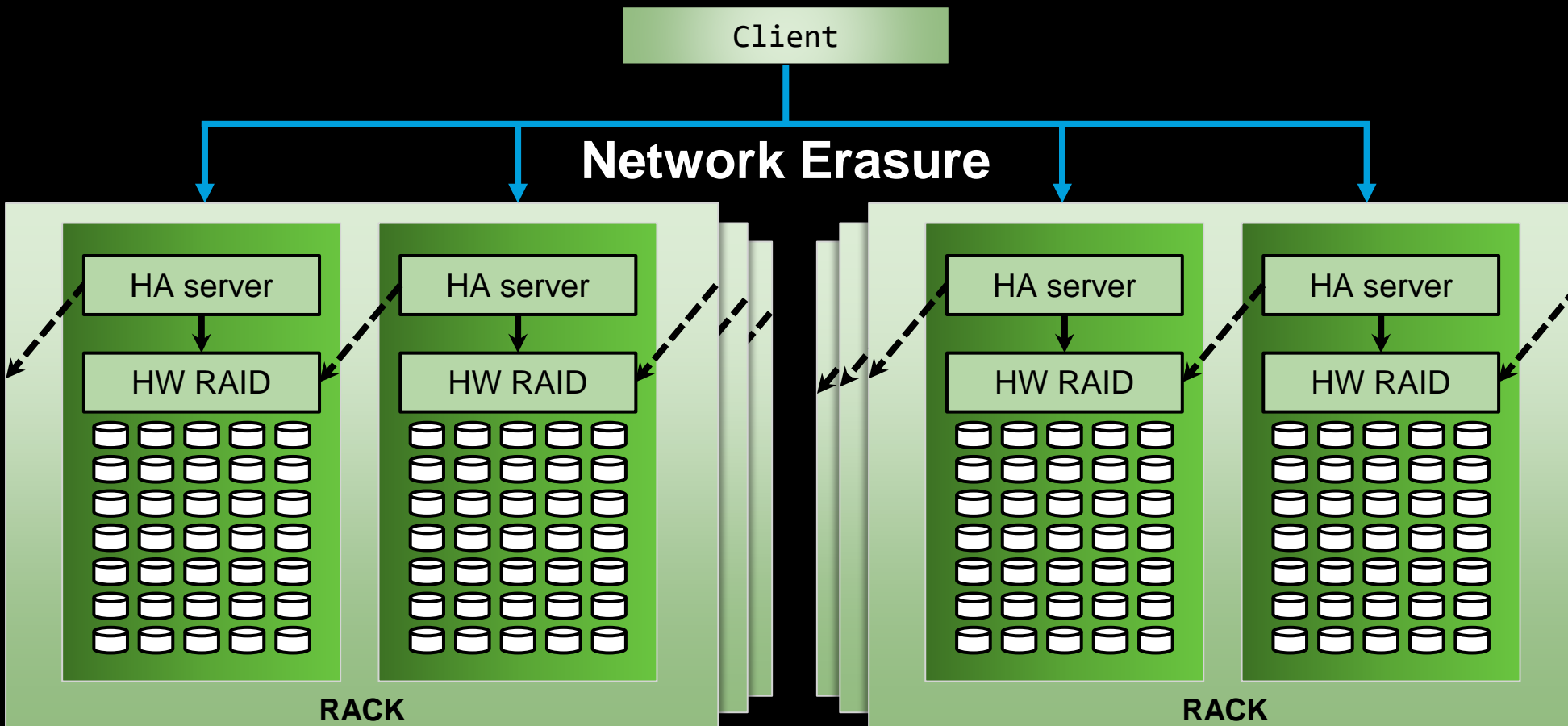
Single Level Erasure Coding



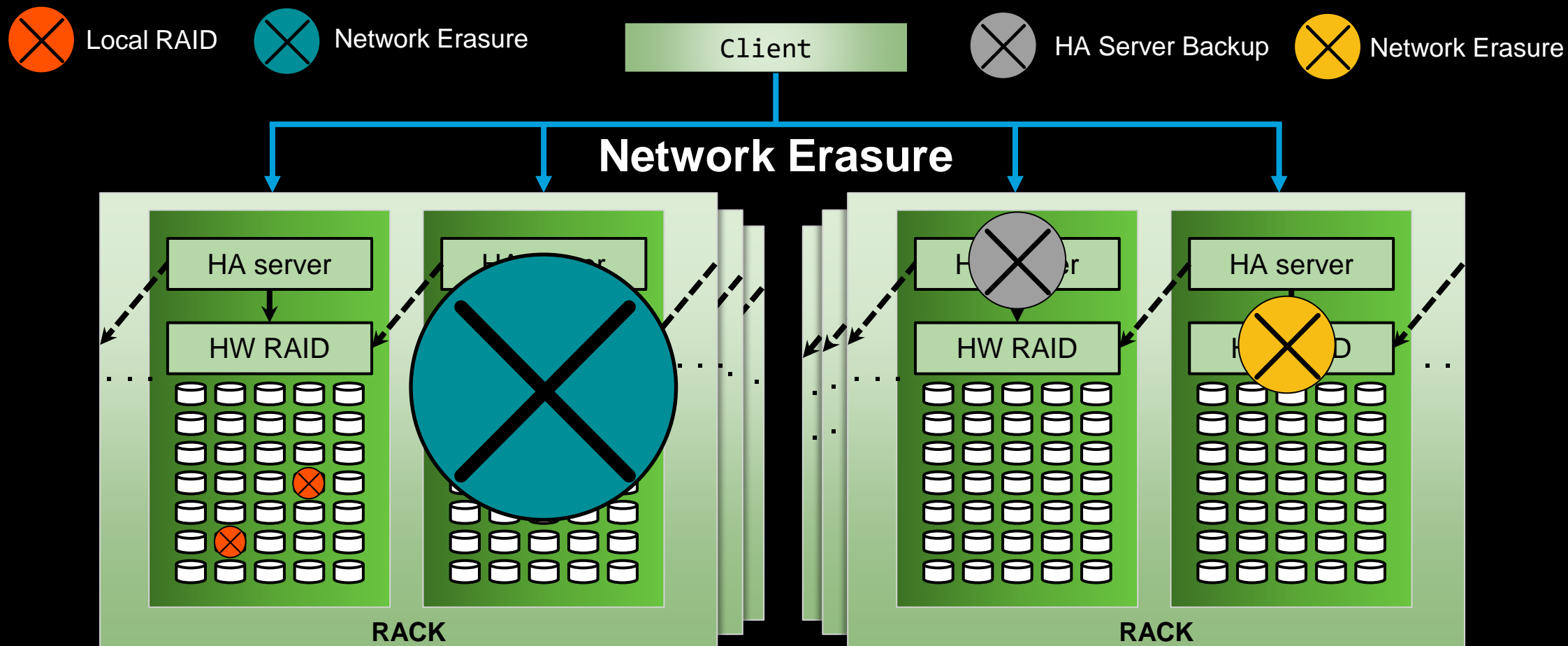
Hierarchical Erasure Coding



Hybrid Approach: Two Layers of Erasure

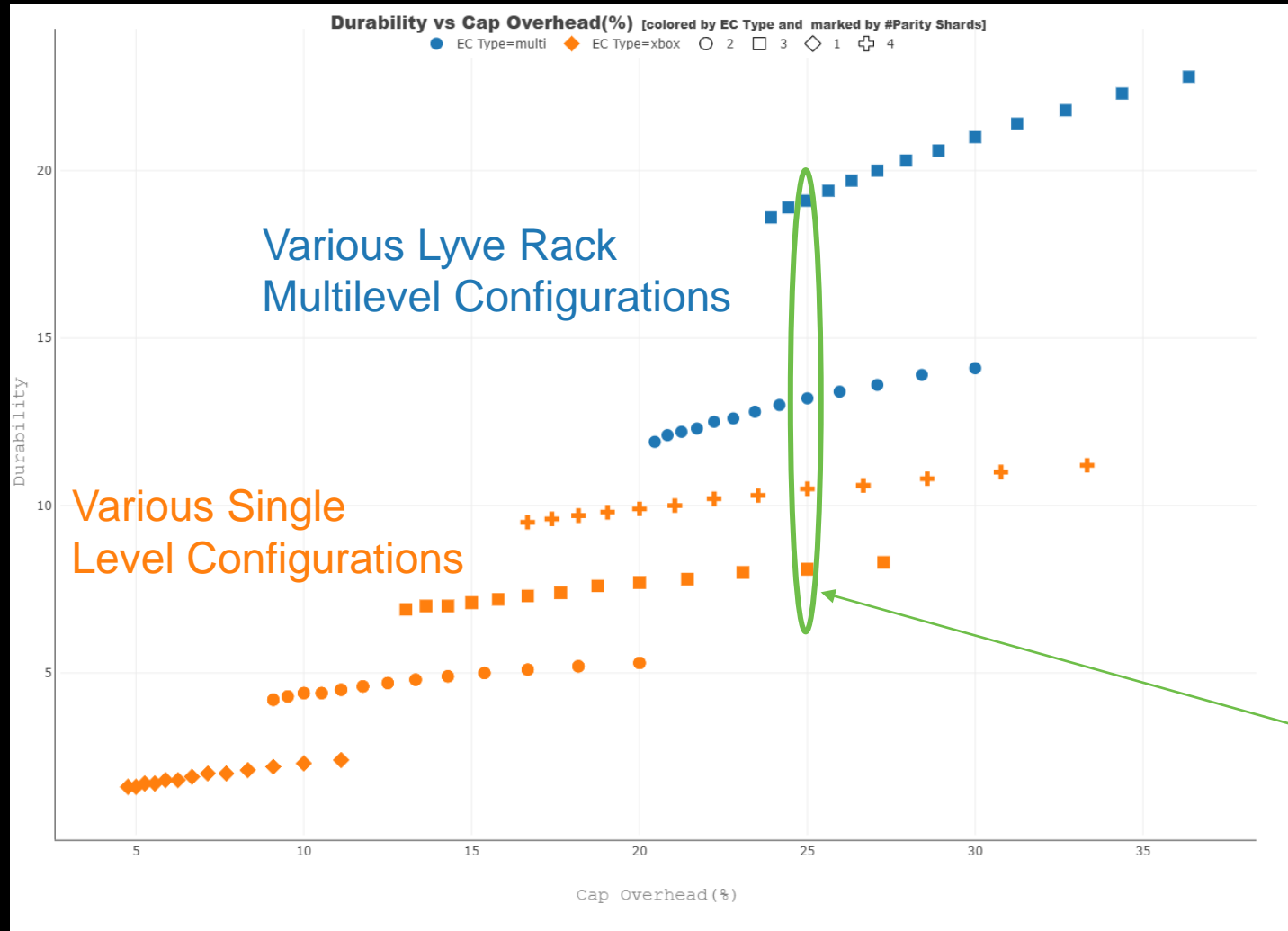


Hybrid Approach: Two Layers of Erasure



CORTX Hybrid Approach Maximizes Durability, Availability, and Density

The Value of Multilevel Erasure Coding

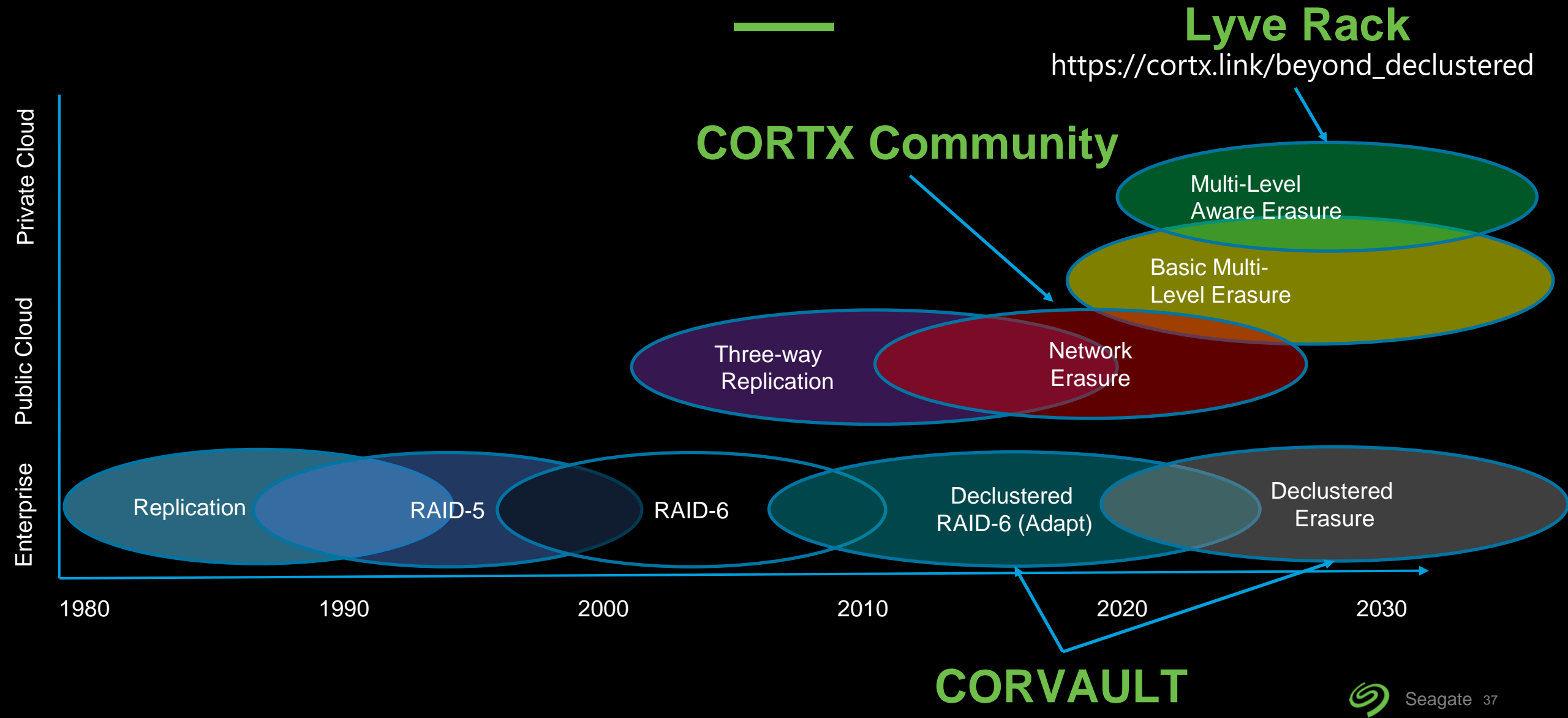


Basic rule: add more overhead to improve durability

Advanced rule: use multilevel to improve durability *without* increasing overhead

At comparable capacity overhead, the multilevel configuration provides many more 9's of durability.

The Evolution of Data Durability in Data Centers





Thanks!

john.bent@seagate.com

<https://cortex.io>

