



# The Pros and Cons of Erasure Coding & Replication vs. RAID in Next-Gen Storage Platforms

---

Abhijith Shenoy – Engineer, Hedvig Inc.  
@hedviginc

# The need for new architectures



Business executives



Developers



IT infrastructure / DevOps

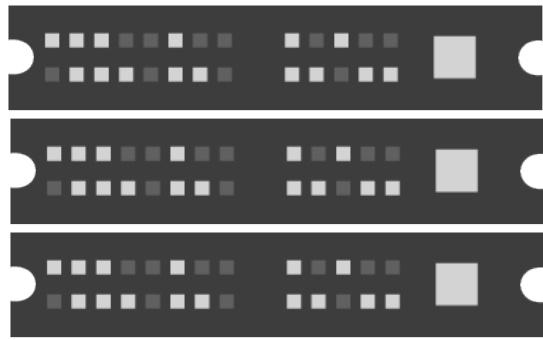
# Modern apps need. . .

- Scale
- Flexibility
- Self-service
- Automation

To achieve this, the world is moving to a software-defined, distributed systems approach

# Software-defined Storage

# Software-defined storage



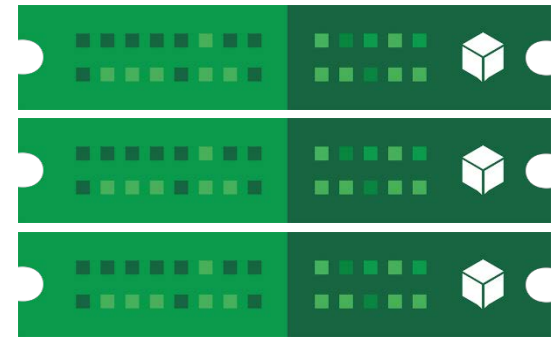
commodity servers

+



software

=



software-defined  
storage

# Common software-defined storage elements



Storage  
software

Forms elastic storage cluster with commodity servers or cloud infrastructure

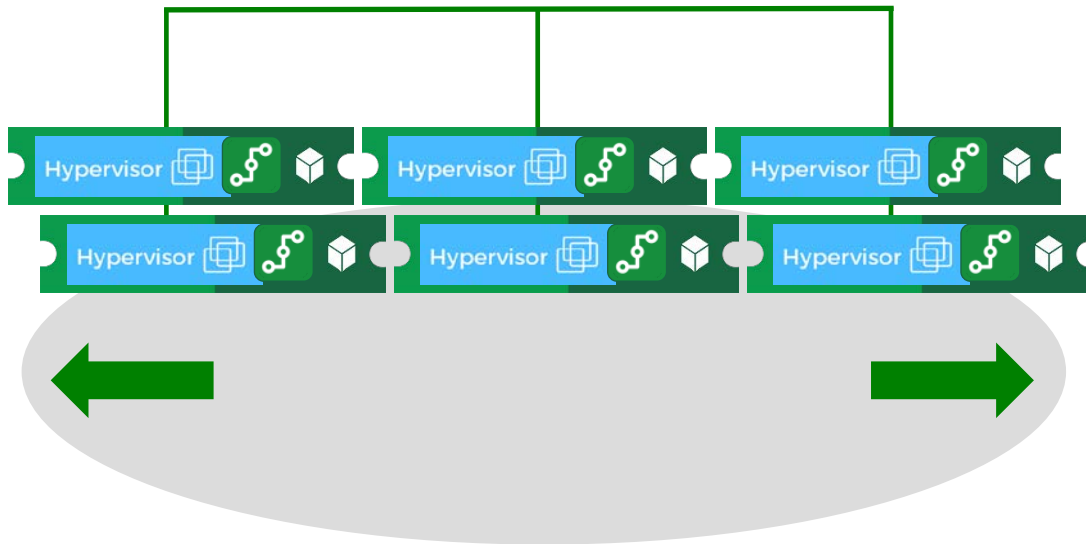


Proxy/client

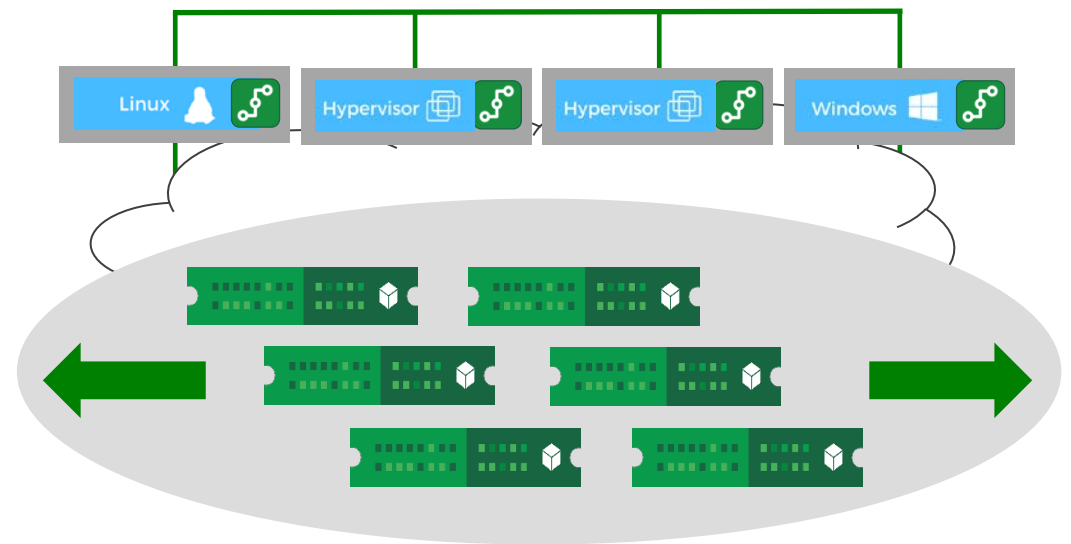
Provides storage access to application compute environment via common protocols

# Modern software-defined storage architectures

## Hyperconverged

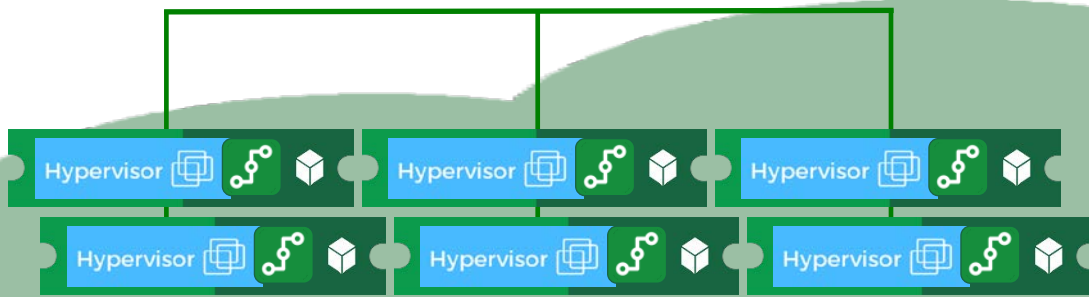


## Hyperscale

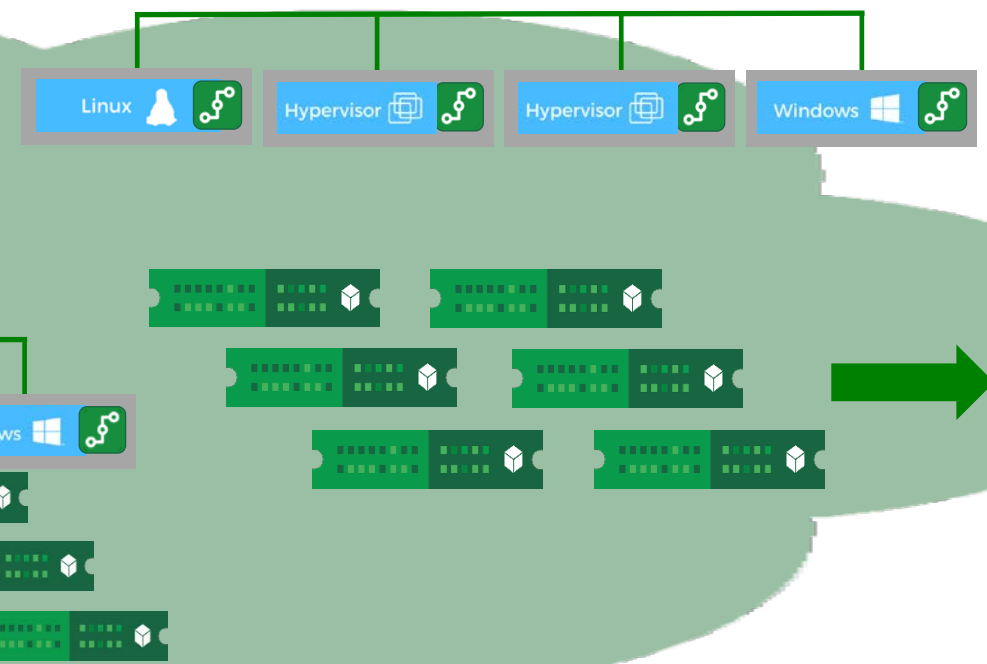


# Spanning multiple DCs and clouds

## Hyperconverged



## Hyperscale

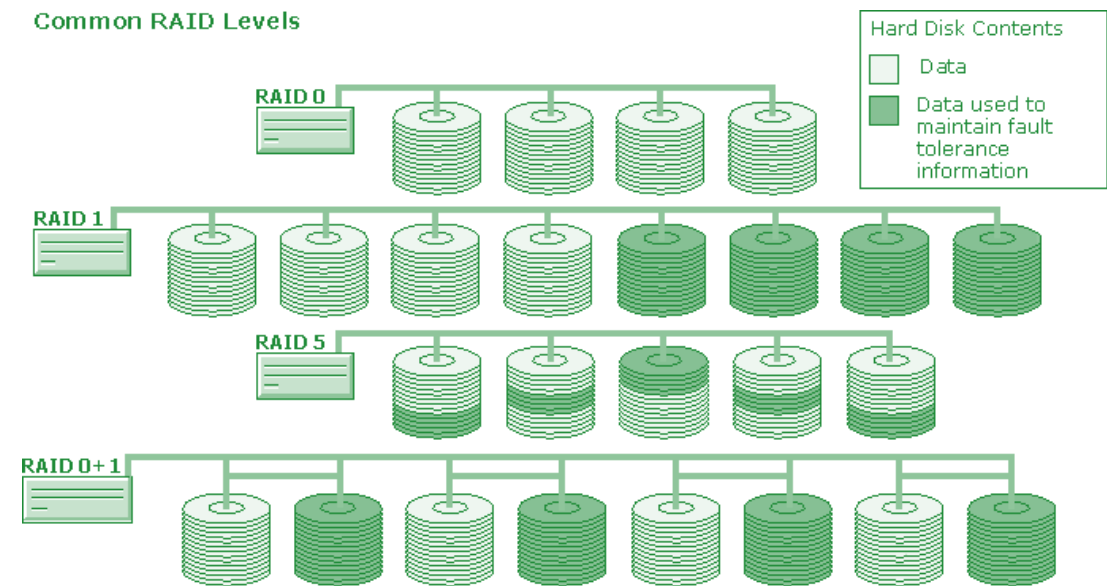




Data protection with SDS:  
RAID, Erasure Coding  
& Replication

# Protecting stored data: RAID

- Redundant Array of Independent Disks
  - Divides or replicates data across multiple drives to deliver performance and fault tolerance
  - Commonly used: RAID 0, RAID 1, RAID 5, RAID 10
- Pros
  - Trusted protection solution in the traditional array world
  - Known performance delivery
- Cons
  - High-capacity drive (8TB+) rebuilds can take days or even weeks
  - RAID controllers add complexity for requisite performance

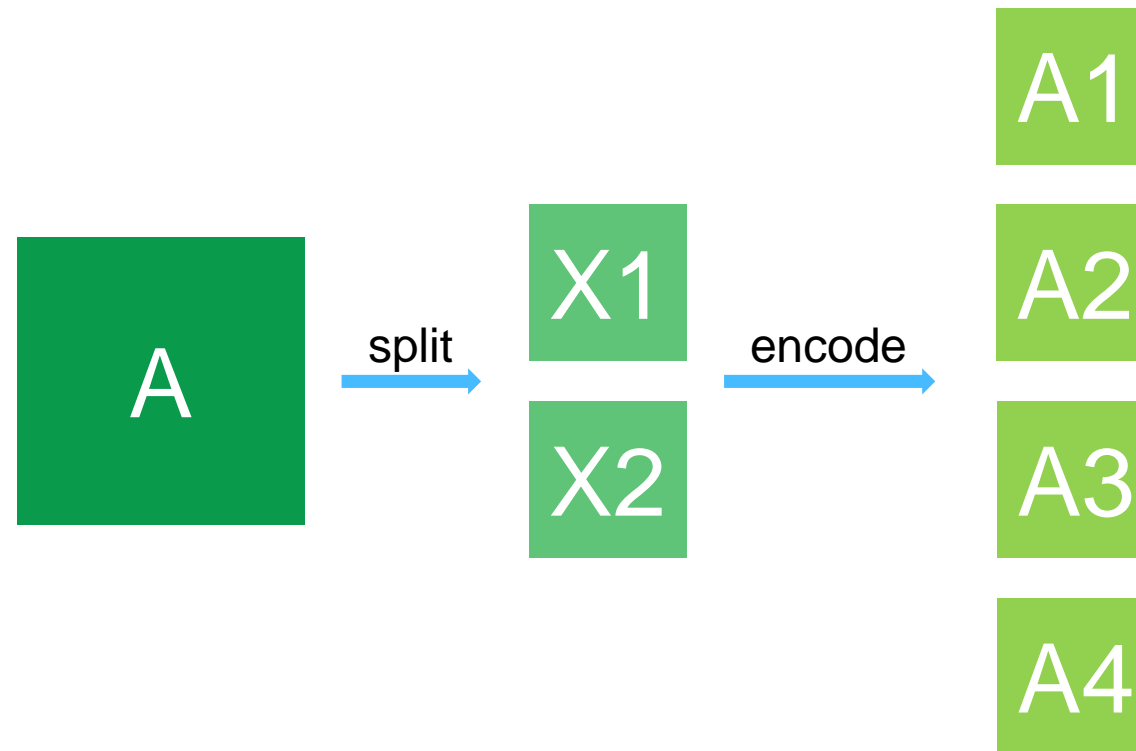


# Protecting stored data: Erasure coding

- A parity based protection technique
  - Data broken into fragments and encoded
  - Stored across different locations with a configurable number of redundant pieces
- Pros
  - Consumes less storage than replication – good for cheap/deep
  - Allows for the failure of two or more elements of a storage system
- Cons
  - Parity calculation is CPU-intensive
  - Increased latency can slow production writes and rebuilds

# How erasure coding works

- Split a file into  $n$  chunks and code into  $m$  parity blocks



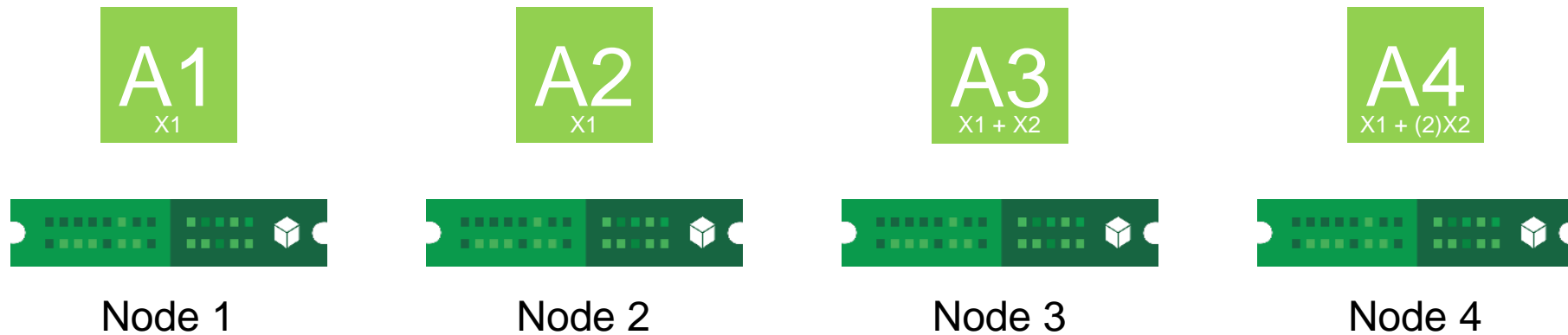
# How erasure coding works

- Tolerate  $m$  erasures (failures)

$$\begin{array}{l} \boxed{A1} \\ \boxed{A2} \\ \boxed{A3} \\ \boxed{A4} \end{array} = \begin{array}{l} \boxed{X1} \\ \boxed{X2} \\ \boxed{X1} \\ \boxed{X1} \end{array} + \begin{array}{l} \\ \\ \boxed{X2} \\ 2 \boxed{X2} \end{array}$$

# How erasure coding works

- In a distributed system, chunks are spread across nodes
- In this example, 2 nodes can fail and data can still be rebuilt



# Erasure coding use case: Archival storage

- Goal
  - Need long-term storage of PBs of files
  - Minimizing storage costs critical to business profitability
- Solution
  - Software-defined storage + erasure coding
- Results
  - Store and protect archival data in 1.5x disk space
  - Performance adequate for workload
  - Rebuilds slower than desired, but capacity savings outweigh latency

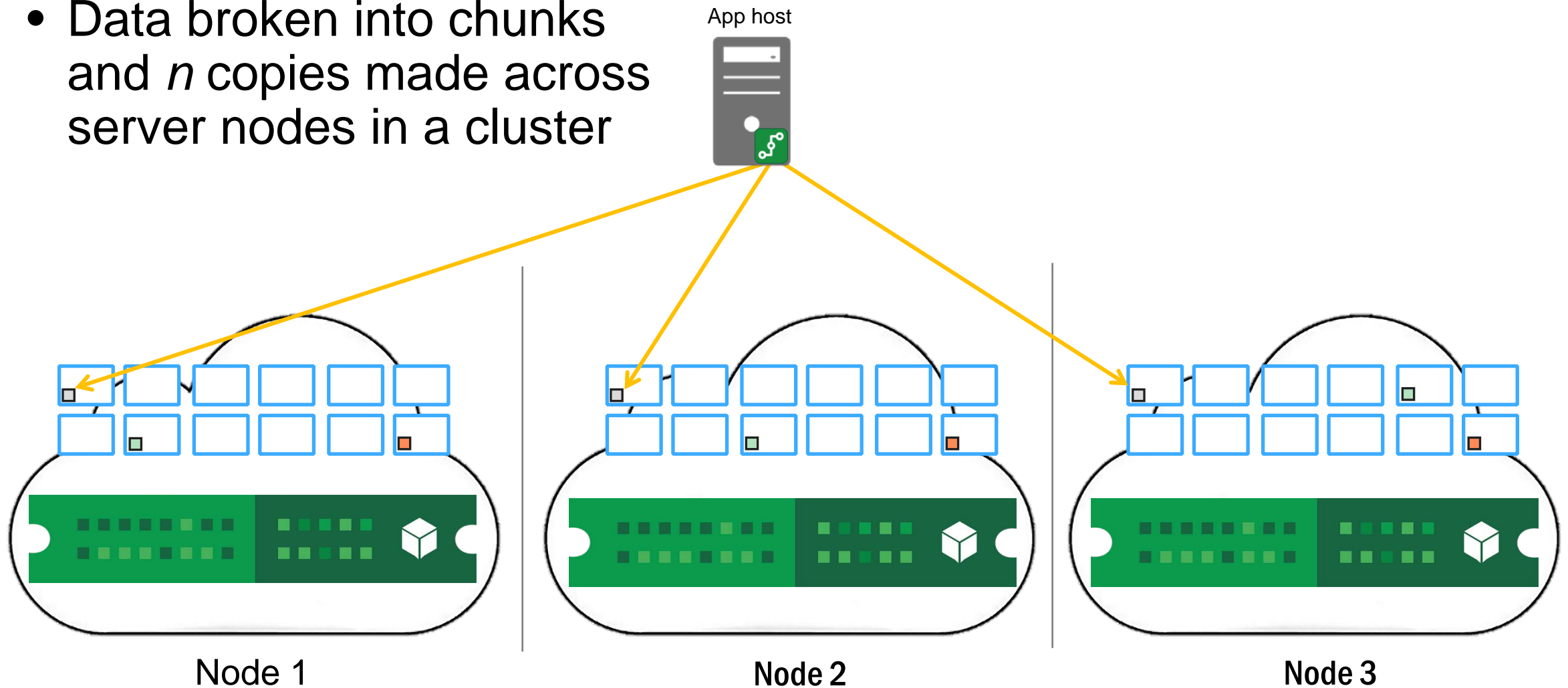
# Protecting stored data: Replication

- The creation of data copies across different locations of the storage system
  - Typically 2 or 3 copies, configurable based on accepted risk level
  - If a drive fails, data is recreated on another drive from replica(s)
- Pros
  - Less CPU intensive = faster write performance
  - Simple restores = faster rebuild performance
- Cons
  - Requires 2x or more the original storage space



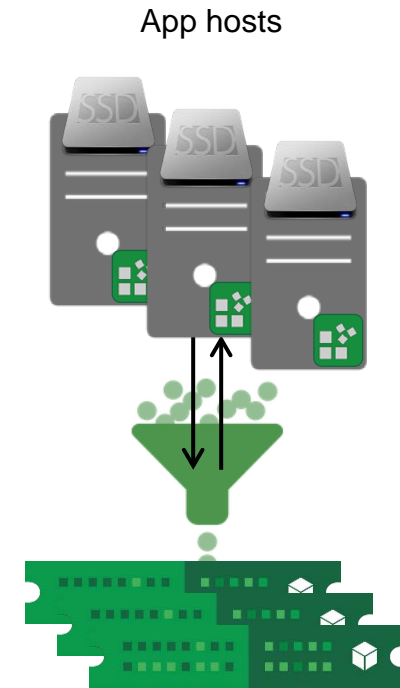
# How replication works with software-defined storage

- Data broken into chunks and  $n$  copies made across server nodes in a cluster



# Offsetting replication overhead

- Compression
  - ~ 2:1 reduction
- Deduplication
  - ~5:1 or higher reduction
- Low disk cost
  - HDD and flash economics declining
  - Overhead of replication more tolerable



# Doesn't have to be one-size-fits-all

- Modern solutions provide per-volume choice
- Choose protection type based on workload

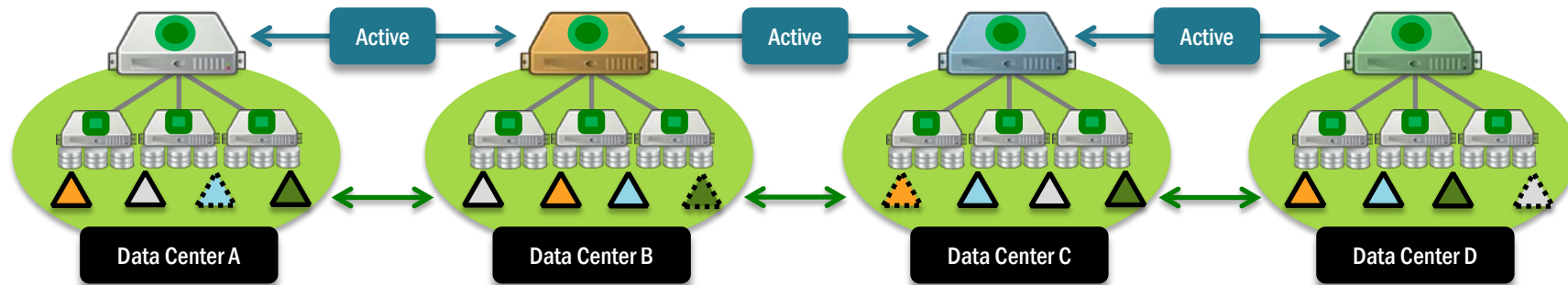
The screenshot displays the 'Virtual Disk Management' interface with a modal window titled 'Add New Virtual Disk'. The interface includes a top navigation bar with 'Cluster Watch' and 'Virtual Disk Management', a user profile 'malfoy | v 0.95', and a search filter. The modal window contains the following configuration options:

- Batch:**
- Name:** [Text input field]
- Size:** 10 GB
- Disk Type:** BLOCK (highlighted with a callout: **Block (iSCSI) | NFS**)
- Residence:**  HDD  Flash (highlighted with a callout: **512 bytes – 64k**)
- Block Size:** 4096 (highlighted with a callout: **512 bytes – 64k**)
- Replication Policy:** Agnostic (highlighted with a callout: **Agnostic | Rack-aware | Datacenter-aware**)
- Replication Factor:** 3 (highlighted with a callout: **2-6 copies**)
- Other options:**  RDM,  Enable Deduplication,  Clustered File System,  Compressed,  Client-side Caching.

A 'Run' button is located at the bottom left of the modal window.

# Replication use case: Primary data storage

- Company: **Large financial organization**
- Situation
  - Hosting 500TB of data across four datacenters in two countries
  - Want maximum availability and recoverability
- Solution
  - Deployed software-defined storage with 4-way replication
- Results
  - Achieve high-performance, high-availability, and quick rebuilds



# Summary

- Protection technologies are evolving along with architectures
- RAID has met its limitation with large capacity drives
- Erasure coding is a good option for latency tolerant, large capacity stores
- Replication provides protection in demanding performance and availability environments
- Software-defined storage offers choice and flexibility to deploy each protection technology where it makes sense

Thank You!