# ViPR Distributed Storage System

Shashwat Srivastav, Sr Director Engg.

Kamal Srinivasan, Principal Prod Manager

EMC²

# Agenda

- ViPR Overview

- ViPR Architecture for scalability

- Geo distributed storage

- Demo

- Q&A

EMC²

# ViPR Overview

**EMC²**

# Storage Systems Today

## Storage silos

- Impede development of applications
- Requires movement of data from one to another (e.g. File to HDFS)

## Enterprise scale

- Can't economically scale for cloud
- Lack of elasticity

## Not ready for modern apps

- Choice of API and HW
- Consistency semantics

4

**EMC²**

# Deliver Storage On Commodity Platforms with ViPR
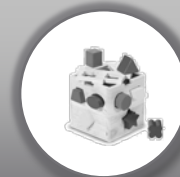


**ViPR Services**

- FILE STORAGE
- HDFS STORAGE
- OBJECT STORAGE

**Commodity Platforms**

5

EMC²

# Scalable Architecture

EMC²

# ViPR Architecture

**BLOCK STORAGE**   **HDFS STORAGE**   **OBJECT STORAGE**

## ViPR Storage Engine
Active-Active read/write support with strong consistency
No single point of failure
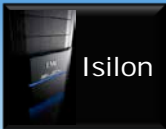Performance and efficiency for small and large objects
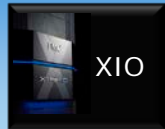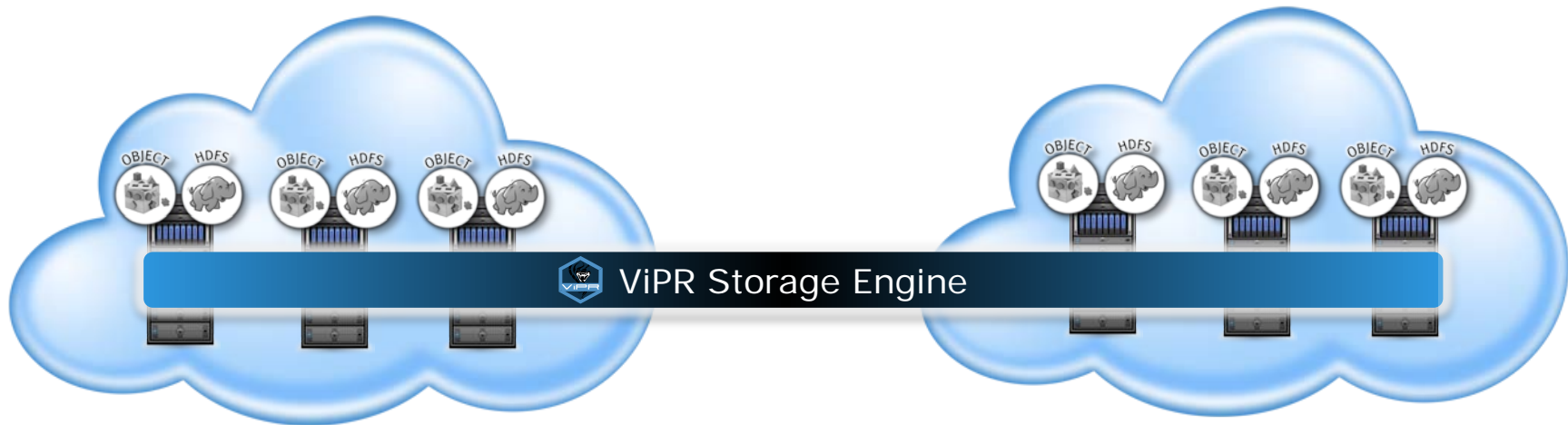
### Storage Arrays
VMAX   VNX   Isilon   XIO   3rd Party

### Commodity Platforms

EMC²

# Common Geo Functionality

GET https://account/bucket/object
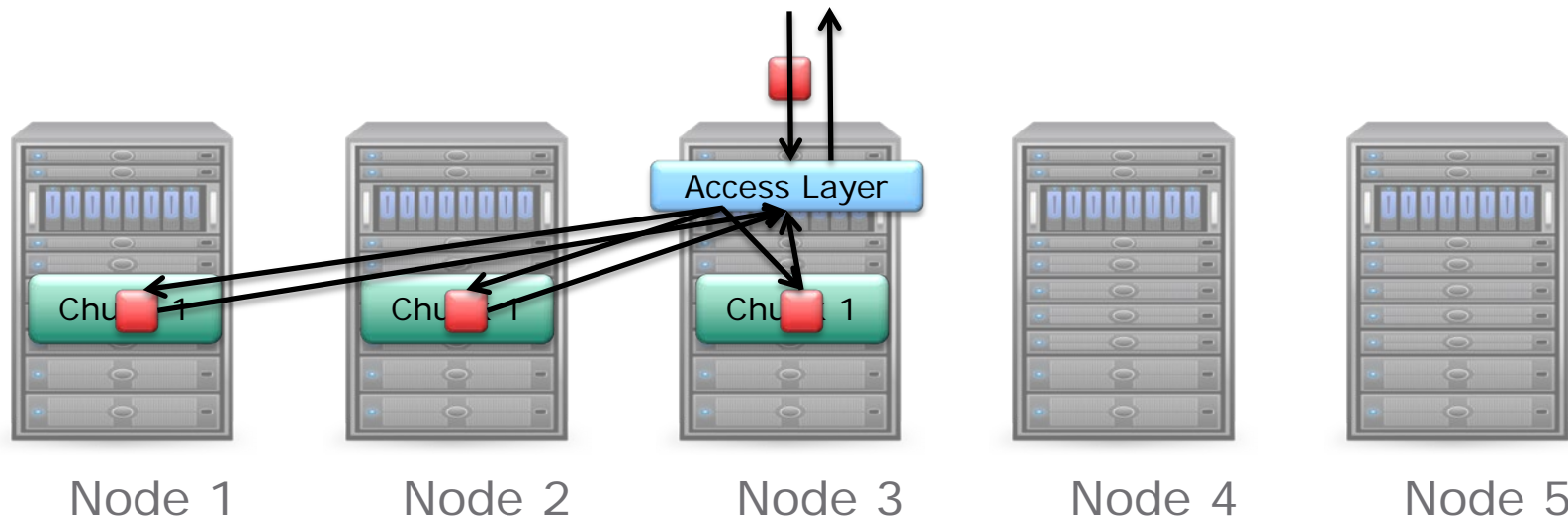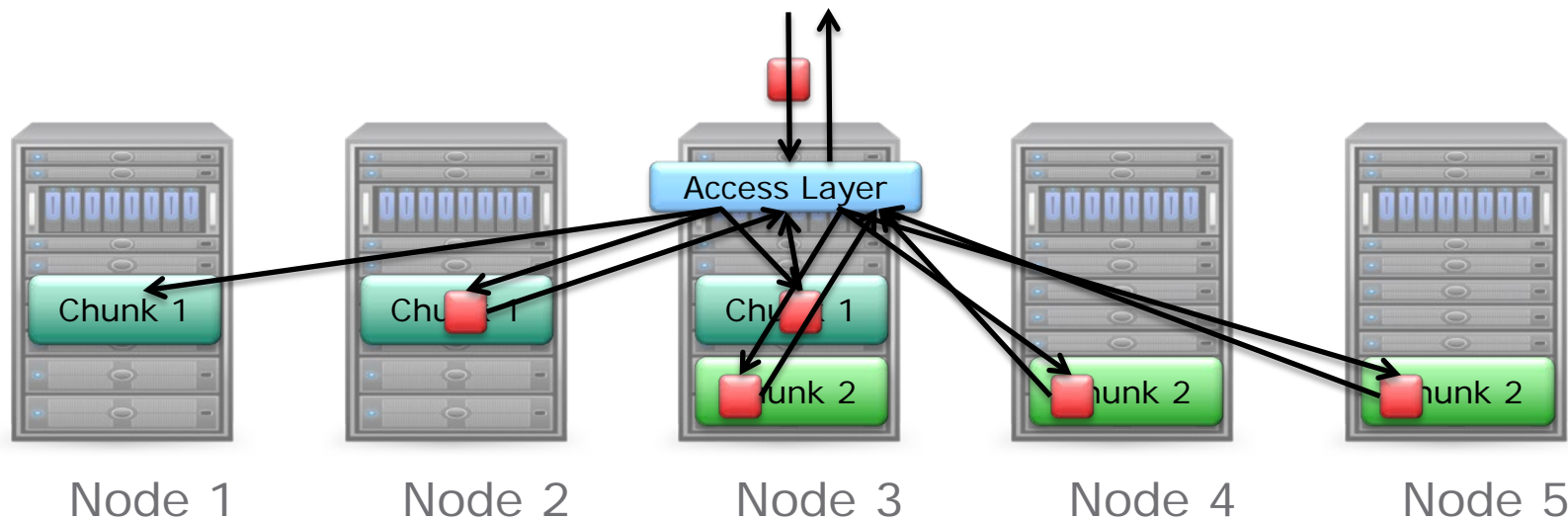
ViPR Storage Engine

EMC²

# Chunks

- ViPR stores all types of data and index in "chunks"

- Chunks are:
  - Logical containers of contiguous space (128MB)
  - Written in an append-only pattern

- All data protection operations are done on chunks

#EMCVIPR
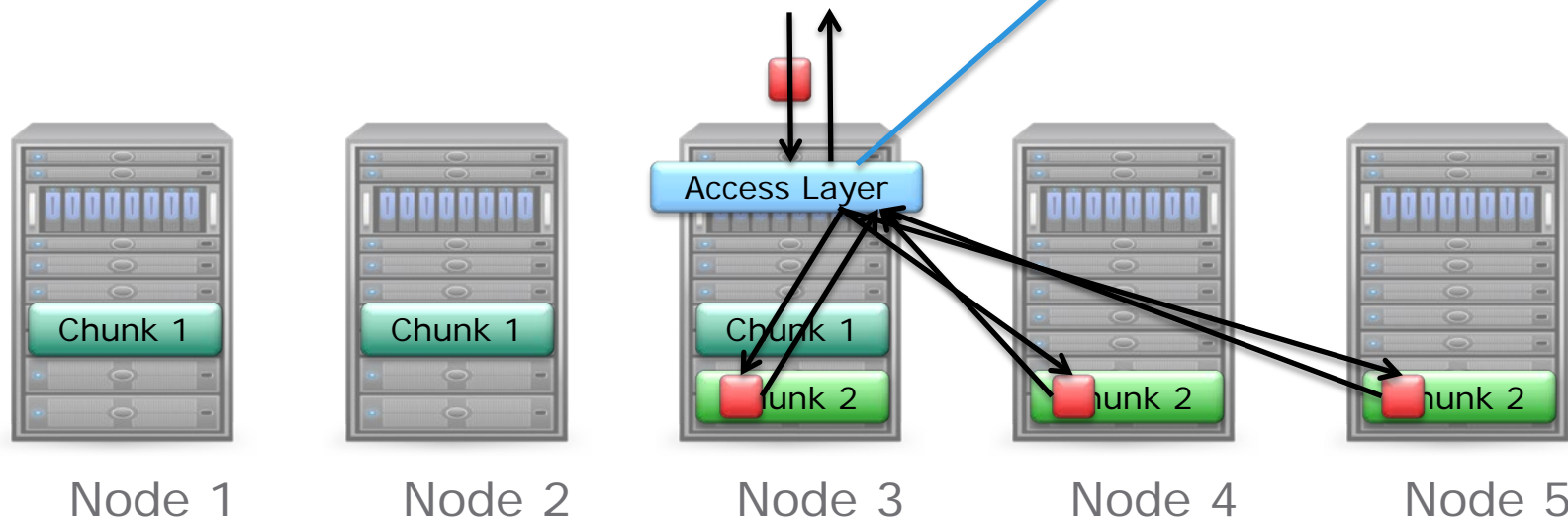
**EMC²**

# Chunk Write



Node 1    Node 2    Node 3    Node 4    Node 5

EMC²

# Chunk Write

# Index

| Content Name | Chunk Location |
|---|---|
| Image1.jpg | Chunk 2 (offset:100) |

Access Layer

Chunk 1

Chunk 1

Chunk 1

Chunk 2

Chunk 2

Chunk 2

Node 1

Node 2

Node 3

Node 4

Node 5

EMC²

# Index Design

| Partition | Node | e | Chunk Location |
|-----------|------|---|----------------|
| Partition 1 | Node 2 | | Chunk 2 |
| Partition 2 | Node 3 | | |
| Partition 3 | Node 5 | | |
| … | … | | |

Partition 1

Partition 2

Partition 3

Chunk 1

Chunk 1

Chunk 1

Chunk 2

Chunk 2

Chunk 2

Node 1

Node 2

Node 3

Node 4

Node 5

EMC²

# Index Design

| Partition | Node |
|-----------|--------|
| Partition 1 | Node 2 |
| Partition 2 | Node 4 |
| Partition 3 | Node 5 |
| ... | ... |



Node 1      Node 2      Node 3      Node 4      Node 5

EMC²

# Transaction

| Content Name | Chunk Location |
|---|---|
| Image1.jpg | Chunk 2 (offset:100) |



Node 1     Node 2     Node 3     Node 4     Node 5

EMC²

# Chunk Info



| Content Name | Chunk Location |
|---|---|
| Image1.jpg | Chunk 2 |

| Chunk Id | Location |
|---|---|
| Chunk 1 | Node 1; Node 2; Node 3 |
| Chunk 2 | Node 3; Node 4; Node 5 |

| Chunk Id | Location |
|---|---|
| Chunk 3 | Node 2; Node 3; Node 4 |

Node 1    Node 2    Node 3    Node 4    Node 5

EMC²

# High Scalability Technique

Content Data
(stored in chunks)

Chunk Management
(stored in chunks)

Chunk Management
(stored in special
replicated store)
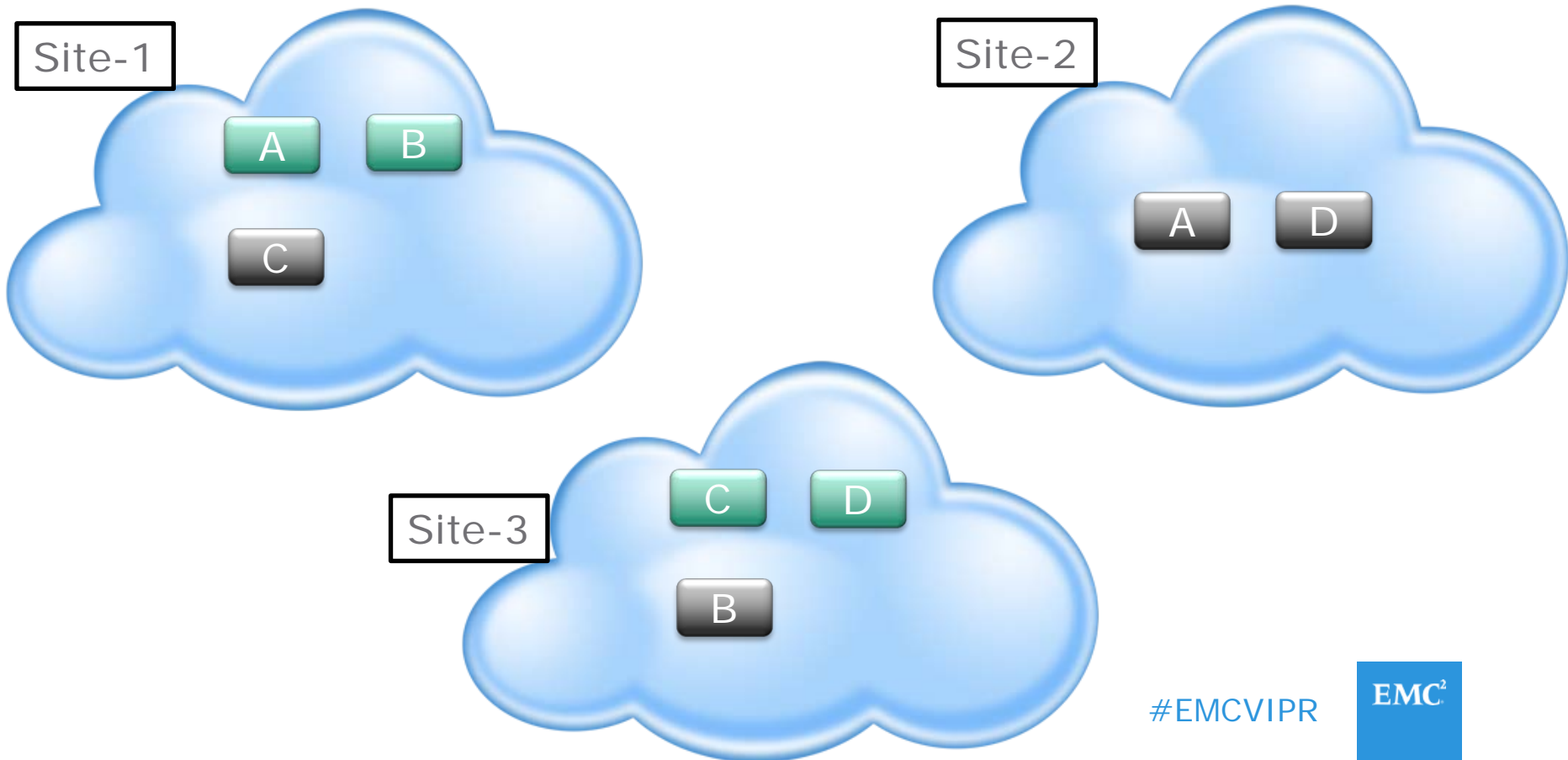
EMC²

# Optimized Geo Protection

# Key Points

- The scheme can tolerate **one site disaster** along with up to **2 node failures** in all the rest of the sites.

- The node failures are repaired using fragments from **local site** without WAN traffic.

- Achieves **~1.8 copies across 4 sites** without having to reconstruct data across the WAN
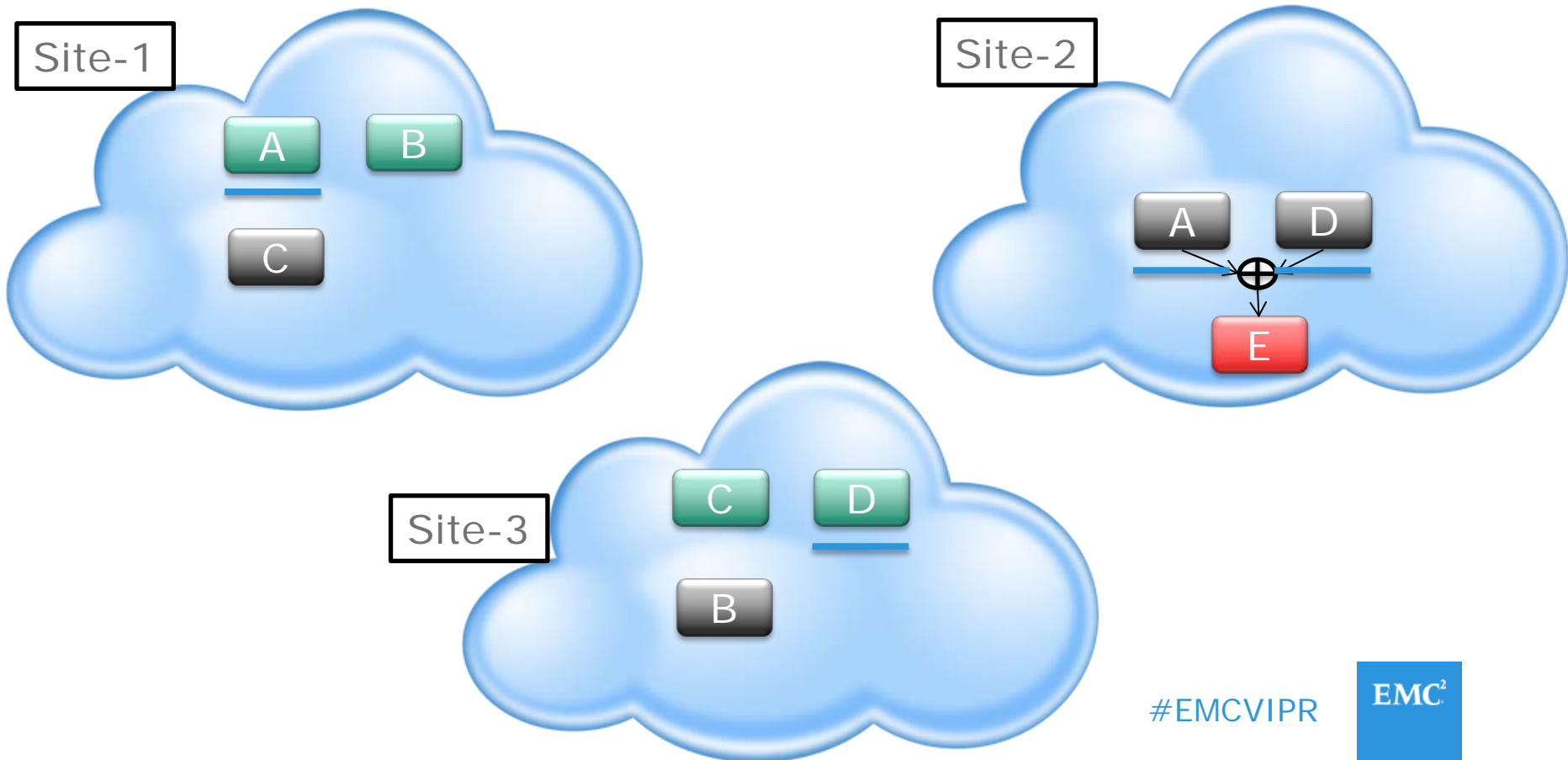
#EMCVIPR

EMC²

# Chunk Backup

Site-1
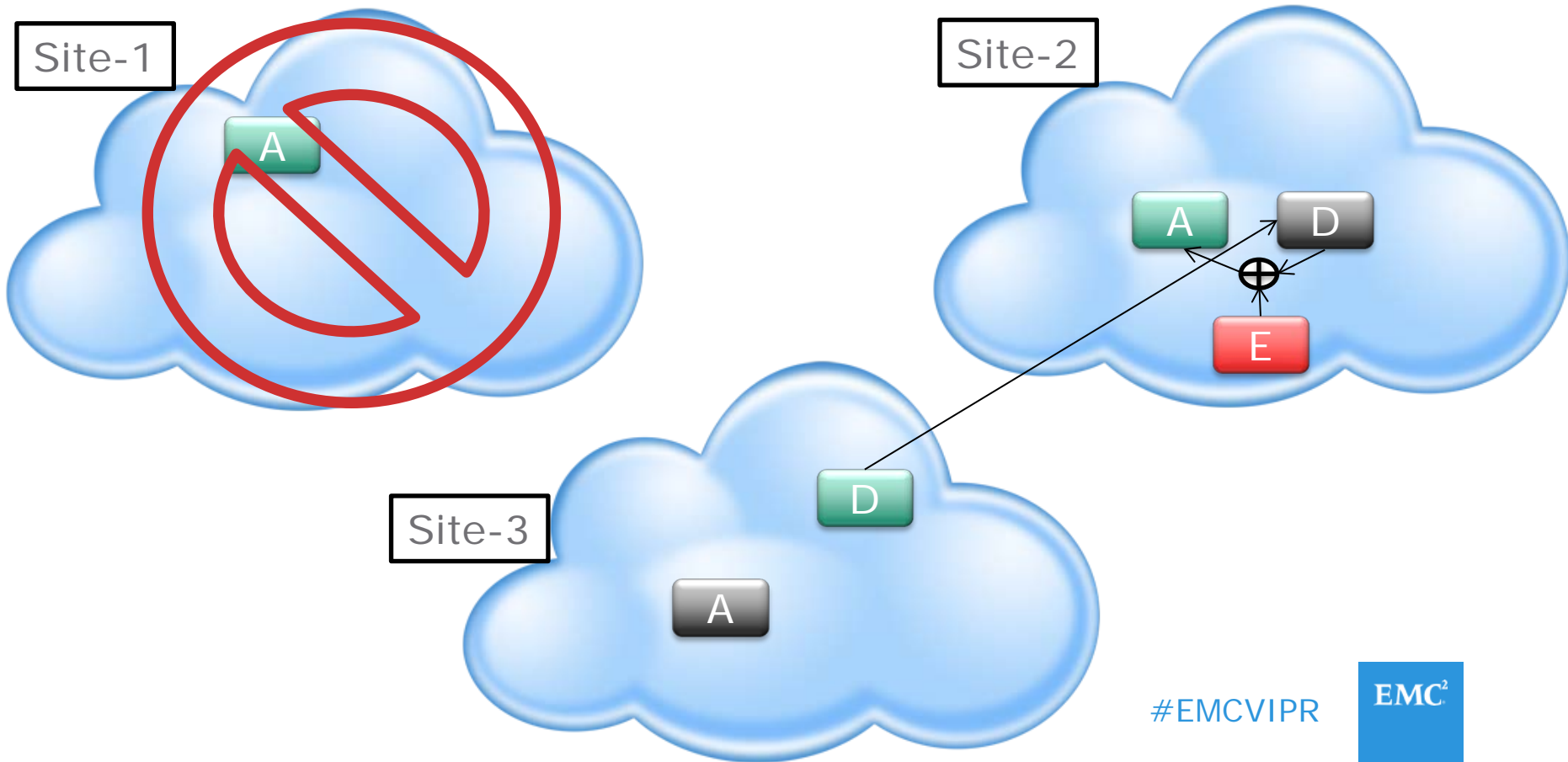
A    B

C

Site-2

A    D

Site-3

C    D

B

#EMCVIPR

EMC²

# Backup Compaction



Site-1

A  B
C

Site-2

A  D
⊕
E

Site-3

C  D

B

#EMCVIPR

EMC²

21

# Chunk Recovery



Site-1

Site-2

Site-3

A

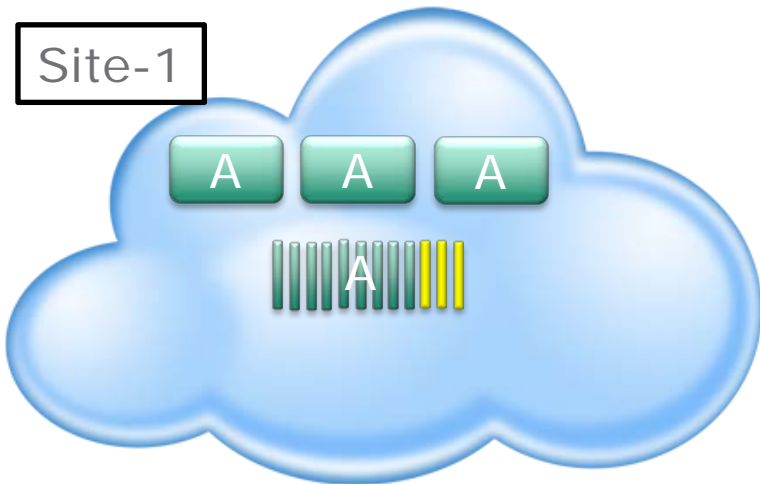A    D

E

D

A

EMC²

# Local Protection



Site-1

A A A

A

Site-2
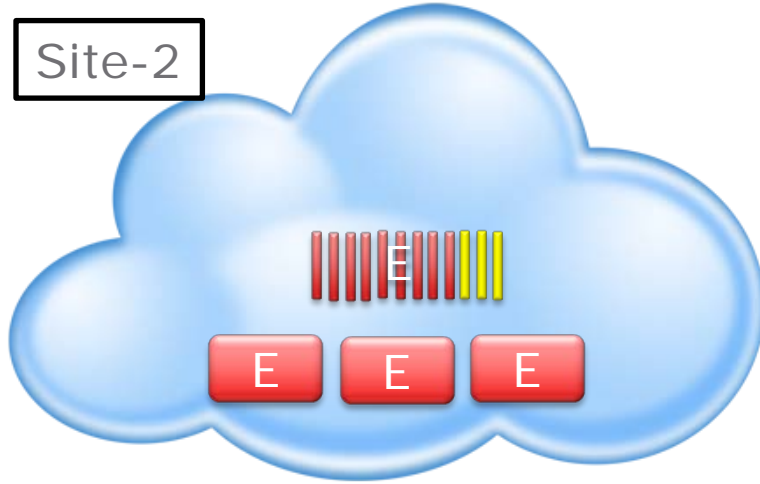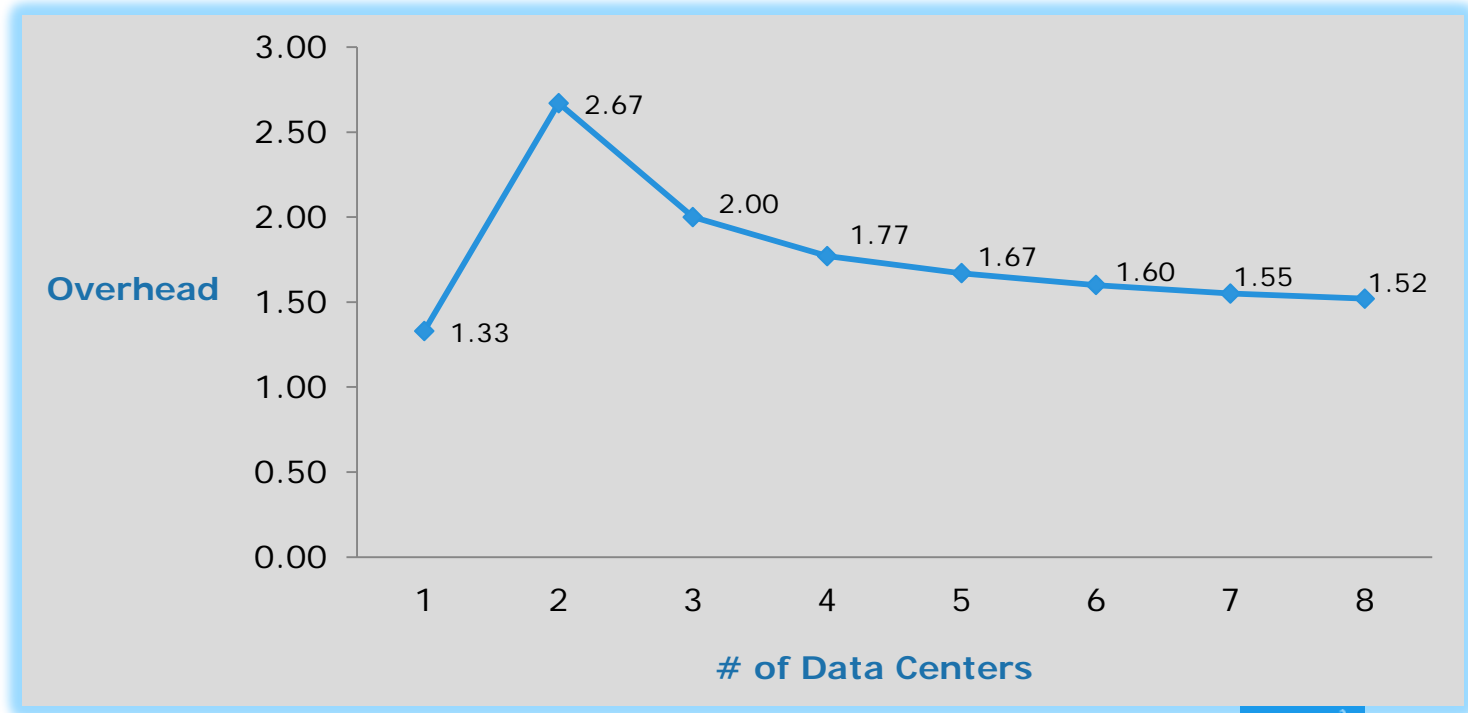
E

E E E

Site-3

#EMCVIPR

EMC²

# Storage Overhead

Optimized data access, protection and efficiency

# Location Agnostic Access With Strong Consistency

**EMC²**

# Industry Solutions

Geo location partitioned namespace

Eventual consistency across geo locations

Sync write all transaction across geo locations

**EMC²**

# ViPR Solution

## Scalable Geo protection

- Each bucket, object, directory, and file is represented as an entity in the index

## Traffic heuristics

- Sense traffic pattern individually for each entity.

## Strong consistency

- Different techniques to avoid WAN round trip

**EMC²**

# EMC²®