

# **Distributed Storage Systems**

Object Storage

# Different Data Storage Systems

-

# Different Data Storage Systems

- Relational databases (SQL)
-

# Different Data Storage Systems

- Relational databases (SQL)
- NoSQL
-

# Different Data Storage Systems

- Relational databases (SQL)
- NoSQL
- Block storage
-

# Different Data Storage Systems

- Relational databases (SQL)
- NoSQL
- Block storage
  - Oldest and simplest
  - Data stored in fixed-sized chunks called “blocks”
  - Block typically houses a portion of the data
  - Address: only identifying part of a block — no block metadata
  - Good performance when app & storage are local, but can lead to more latency the further apart they are
  - Addressing requirements limit scalability

# Different Data Storage Systems

- Relational databases (SQL)
- NoSQL
- Block storage
- File Storage
-

# Different Data Storage Systems

- Relational databases (SQL)
- NoSQL
- Block storage
- File Storage
  - User names files/data, places them in folders, & can nest them under more folders (hierarchy)
  - File has a limited set of metadata associated with it
  - Works well with smaller files, some issues when retrieving large amounts of data
  - Scaling is a problem: it becomes harder to find Information
  - Unique addresses→ finite number of files you can store
  - Sharing via NAS: great locally, issues over wide area networks



# Different Data Storage Systems

- Relational databases (SQL)
- NoSQL
- Block storage
- File Storage
- Object storage

# What is an Object?

**File:**

15x6x3\_10x00x41741x64561  
\_21x58x511\_n.jpg



# What is an Object?

## File:

15x6x3\_10x00x41741x64561  
\_21x58x511\_n.jpg

## Metadata:

Image size: 1516x2048

Date taken: 2013-12-27 13:19

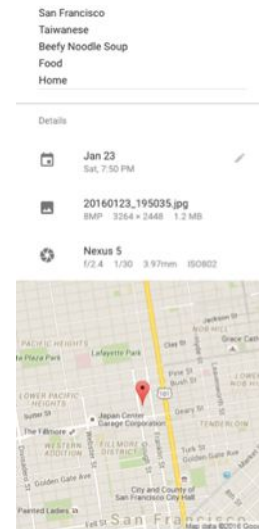
Tags: Asia,  
Japan,  
Bamboo forest,  
Arashiyama

GPS: 35.016520, 135.670436

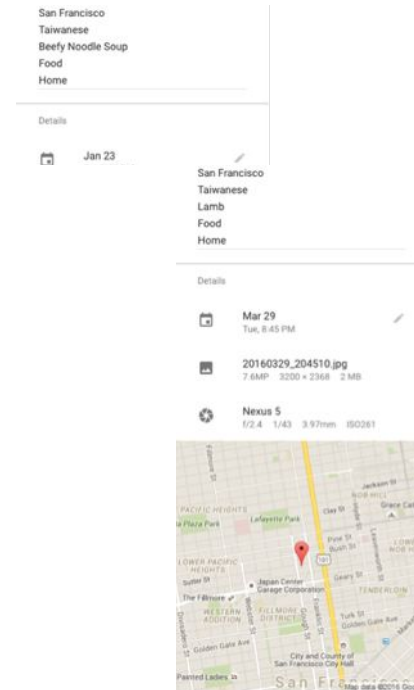


● **Object = File + Metadata**

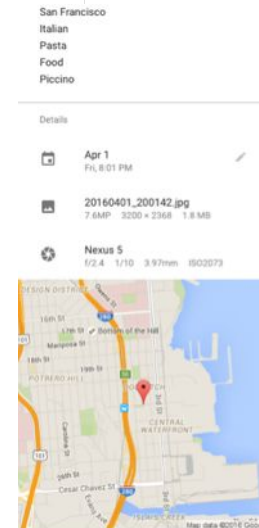
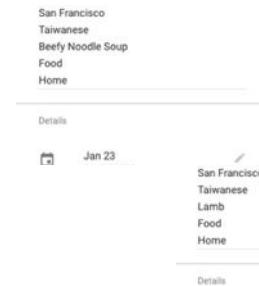
# We have a lot of objects



# We have a lot of objects

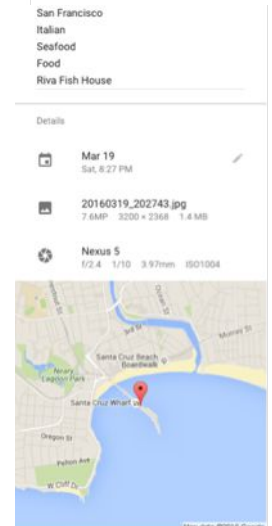
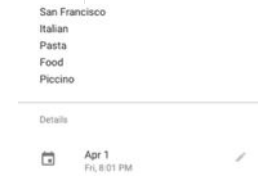
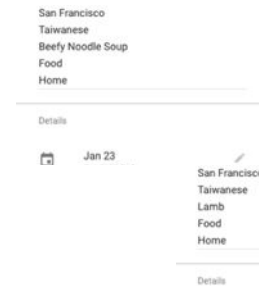


# We have a lot of objects





# We have a lot of objects



# What is for Dinner?

Metadata Search:

Taiwanese Beefy Noodle

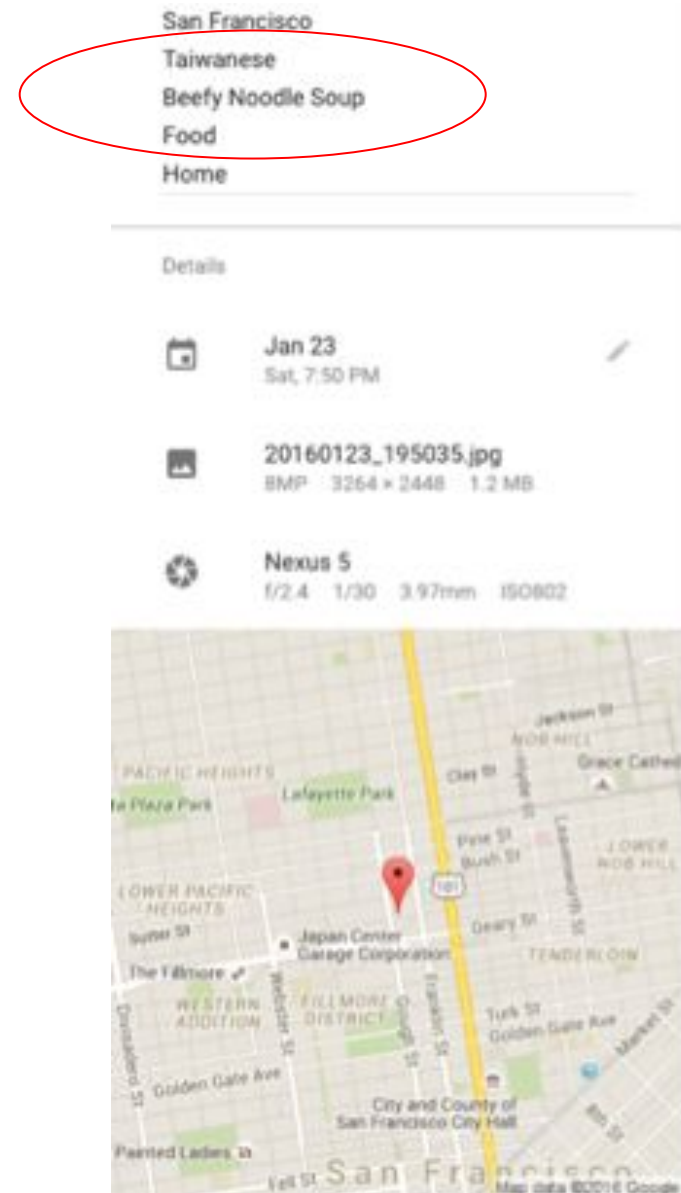
Metadata info is stored with:

- Object
- Searchable Index

This allows for meta data search



# What is for Dinner?



## Metadata Search: Taiwanese Beefy Noodle

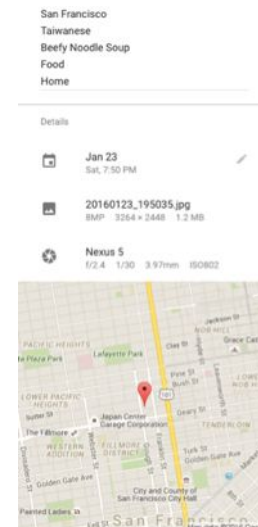
Metadata info is stored with:

- Object
- Searchable Index

This allows for meta data search

# What is for Dinner?

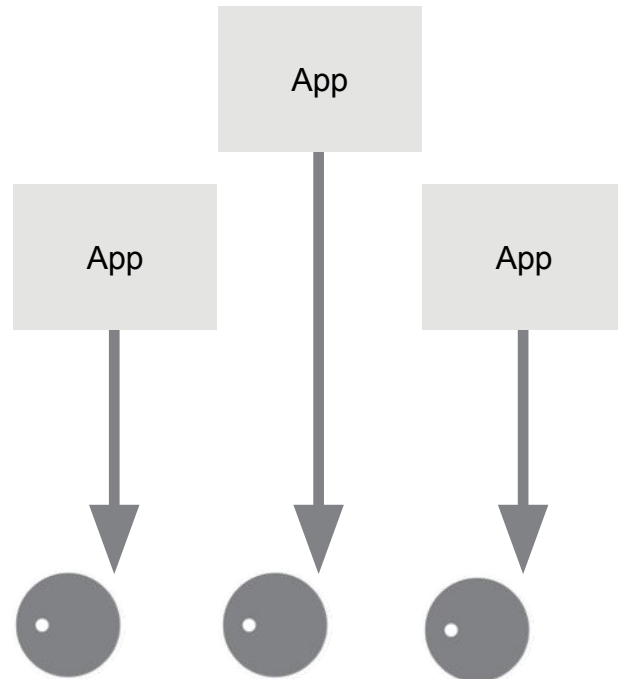
## Metadata Search: Taiwanese Beefy Noodle



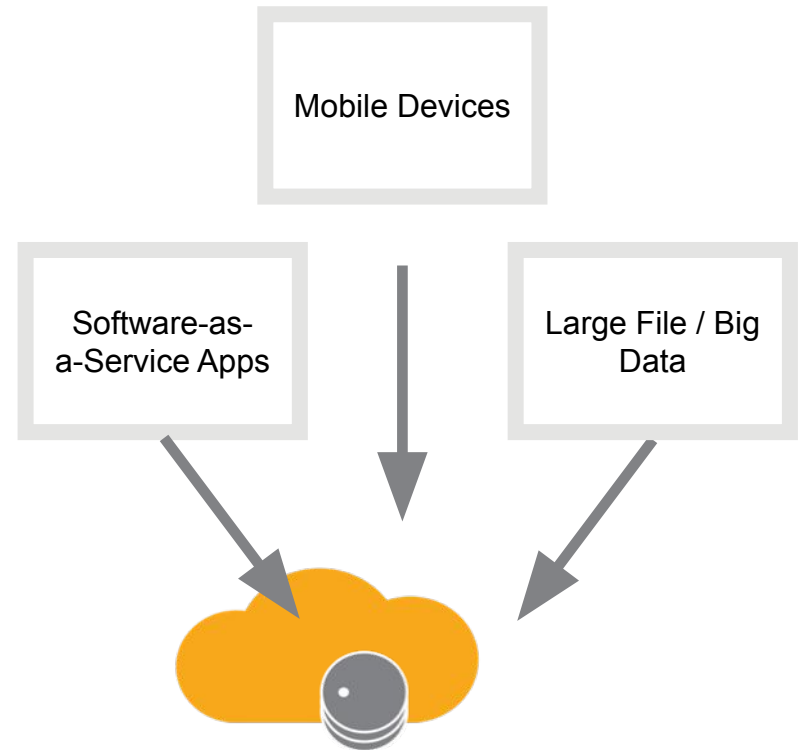
# **Why Object Storage?**

# Why would we need Objects?

Traditional - File-based



Object Storage - HTTP Namespace

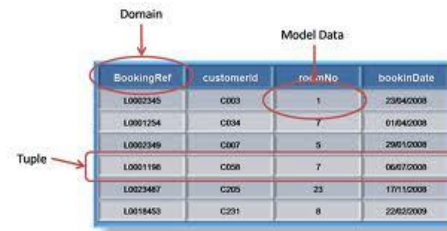


# Right Tool for the Job in the Cloud

	Object	File	Block
Price	\$	\$\$	\$\$\$
Data-type	Semi + Unstructured	Structured, Semi, Unstructured	Structured
Scales	Billions of Objects and Multi-Exabytes	Millions of Files and 1 Petabyte (ish)	Thousands of files and many TB
Protocols	HTTP REST APIs (e.g. S3, OpenStack Swift)	File Protocols (e.g. CIFS/NFS)	Block protocols (e.g. iSCSI, FC)
Optimized for	<ul style="list-style-type: none"> <li>Capacity</li> <li>Scalability</li> <li>Eventual consistency</li> <li>Web accessibility</li> <li>Broad geo distribution</li> </ul>	<ul style="list-style-type: none"> <li>Fast file-sharing</li> <li>Fast file-serving</li> <li>File-locking</li> <li>Canonical true files</li> </ul>	<ul style="list-style-type: none"> <li>Fast random lookups</li> <li>Reads and writes on small records</li> <li>Storage for hypervisors</li> </ul>
Typical apps	Everything under “File” + gene sequences, video, log files, photos, web content, large data sets	Office automation, design automation, collaborative engineering, word processing docs, presentation graphics	Transactional apps such as ERP, CRM, databases
Approach	An object = a File + all associated metadata + a globally unique identifier	File is stored in directory structure. Limited metadata stored in the file system itself (separate from file)	File is written in “blocks” on spinning media

# What Object Storage is Not

- Distributed File System
  - Does not provide POSIX file system API support
- Relational Database
  - Does not support ACID semantics
- NoSQL Data Store
  - Not built on the Key-Value/Document/Column-Family model
- Block Storage System
  - Does not provide block-level storage service



The diagram shows a table with four columns: BookingRef, customerId, roomNo, and bookinDate. Annotations include: 'Domain' pointing to the BookingRef column, 'Model Data' pointing to the roomNo column, and 'Tuple' pointing to a row. The table data is as follows:

BookingRef	customerId	roomNo	bookinDate
L0002345	C003	1	23/4/2008
L0001254	C034	7	01/04/2008
L0002349	C007	5	29/1/2008
L0001190	C008	7	06/7/2008
L0023487	C205	23	17/1/2009
L0018453	C231	8	22/03/2009



***Not a "One Size Fits All" Storage Solution***

# Examples of Object Storage

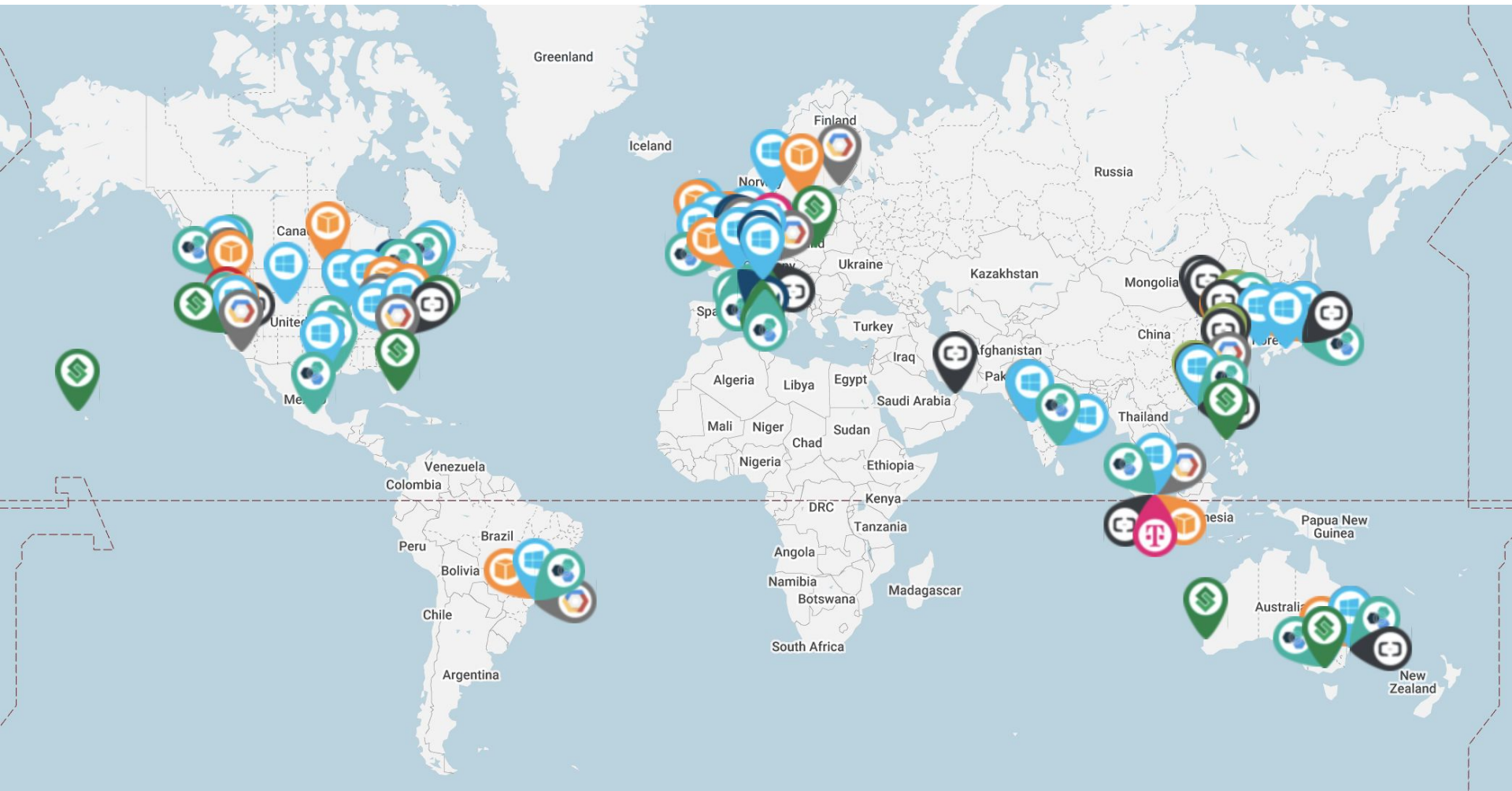
## Object Storage Services:

- Amazon S3
- Google Cloud Storage
- Microsoft Azure
- Rackspace Cloud Files
- HP Cloud Object Storage
- IBM Bluemix Object Storage
- Oracle Cloud Storage
- OVH Cloud Storage
- ...

## Open Source systems:

- OpenStack Object Storage  
("Swift")
- Ceph
- Minio

# (Some) Public Object Storage service regions



Publicly available Object Storage service regions of 10 global cloud providers (early 2020)





# OpenStack Swift



# OpenStack Swift

- Core OpenStack Service

- One of the original 2 projects
- 100% Python
- ~ 40K LOC application, > 80K LOC test code



Top contributing companies include: SwiftStack, Intel, RedHat, IBM, HP, Rackspace, Hitachi, Fujitsu

# OpenStack Swift & Swift API

## **Swift API accessible via HTTP**

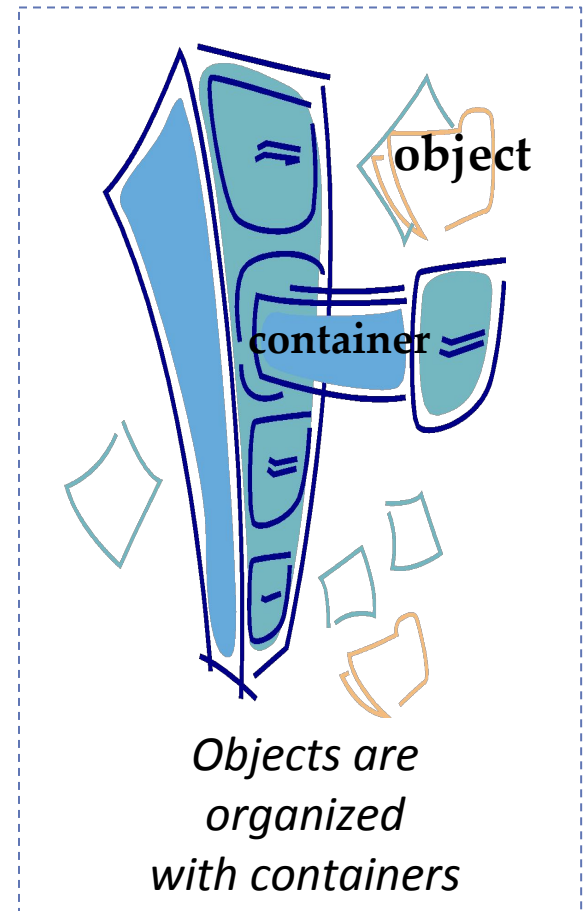
- Applications can consume storage from anywhere
- Larger range of functionality
- Puts the developer in control

## **Swift Object Storage superior to Traditional Storage**

- Designed to scale from TB -> PB -> EB
- Replication is automatic across Nodes, Zones, and Regions
- Supports both Replicas & Erasure Coding
- No Single Points of Failure -> Lose Nodes, Racks or Data Centers
- Ingress / Egress data from all Proxies - No Masters
- Balance Heterogeneous Commodity Hardware

# OpenStack Swift Overview

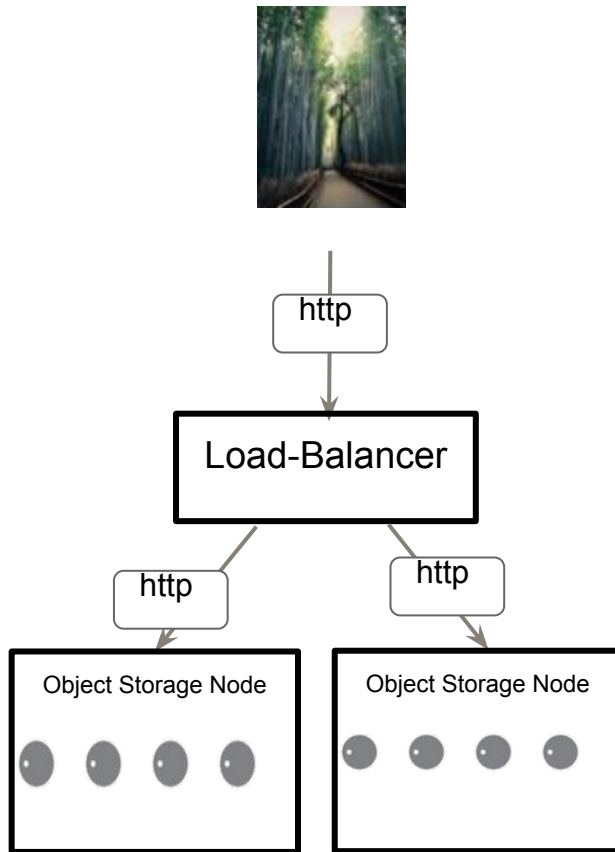
- Uses container model for grouping objects with like characteristics
  - Objects are identified by their paths and have user-defined metadata associated with them
- Accessed via RESTful interface
  - GET, PUT, DELETE
- Built upon standard hardware and highly scalable
  - Cost effective, efficient
- Eventually consistent
  - Designed for availability, partition tolerance



# Access Method: RESTful HTTP API

## Object Storage API

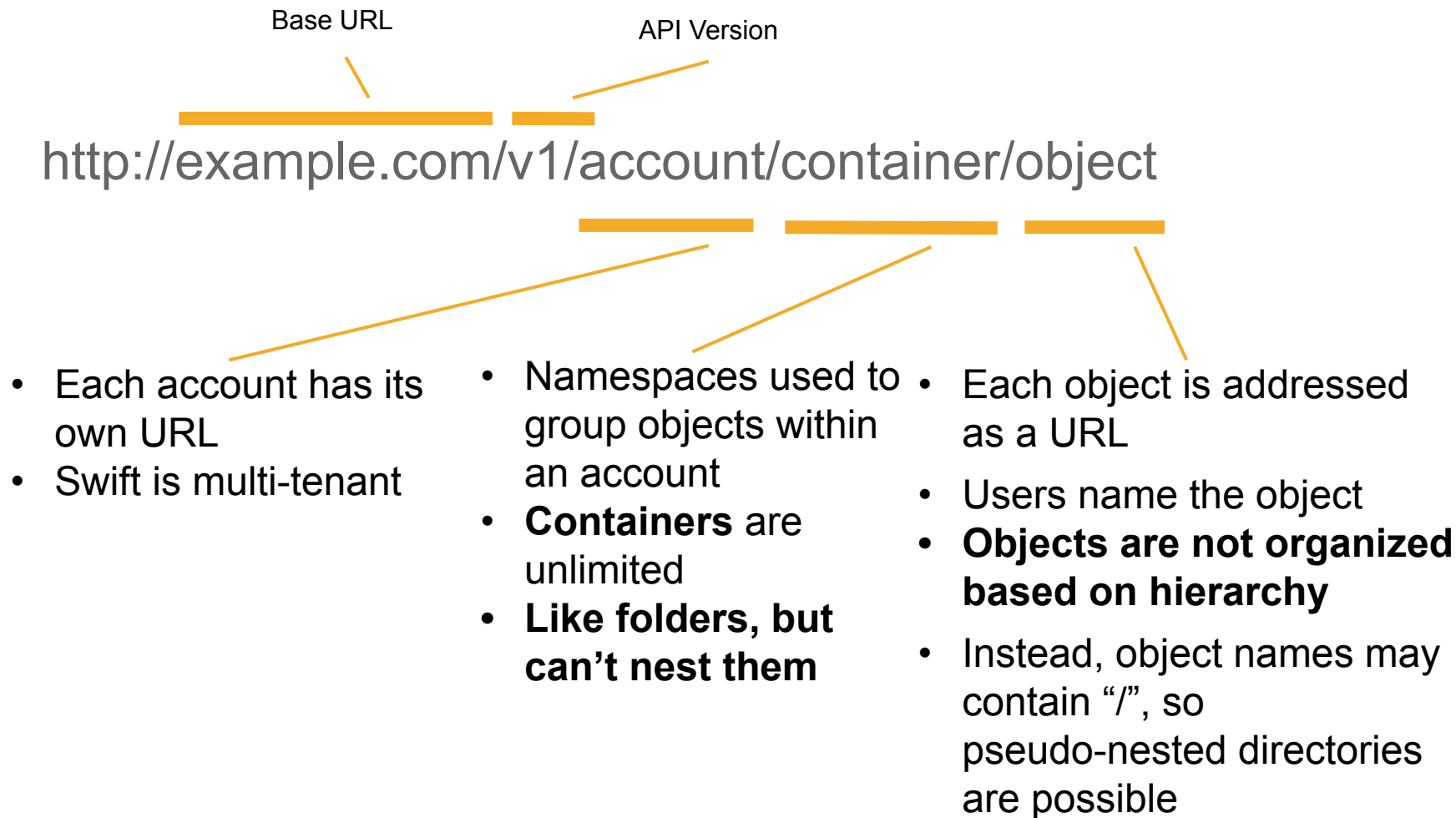
### Operations for Objects:



Method	Description
GET	Downloads an object with its metadata
PUT	Creates new object with specified data content and metadata
COPY	Copies an object to another object
DELETE	Deletes an object
HEAD	Shows object metadata
POST	Creates or updates object metadata

There are also API operations for Accounts and Containers

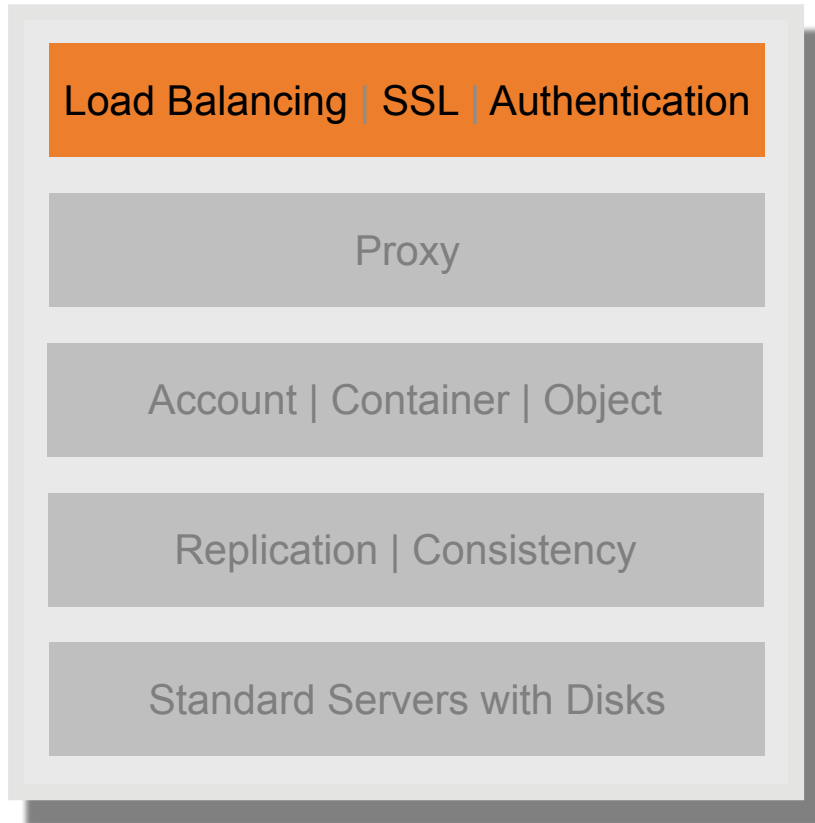
# Every Object Has a URL



# Swift High-Level Architecture

- Load Balancing/Authentication
- Proxy
- Account / Container / Object
- Replication and Consistency
- Standard Hardware

# Swift Architecture



## Load Balancing

- Requests are load balanced across all nodes running proxy server processes

## SSL

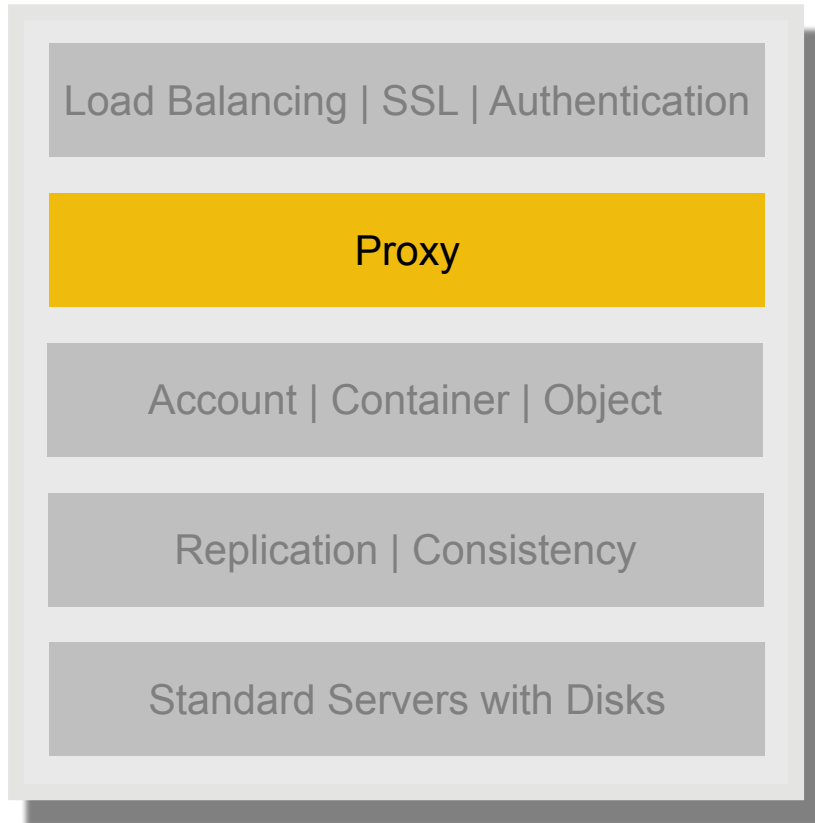
- Optionally, SSL termination can be enabled to encrypt data in flight

## Authentication

- OpenStack Swift has a pluggable authentication system
- Modules include API/UI-driven, LDAP, AD, Keystone



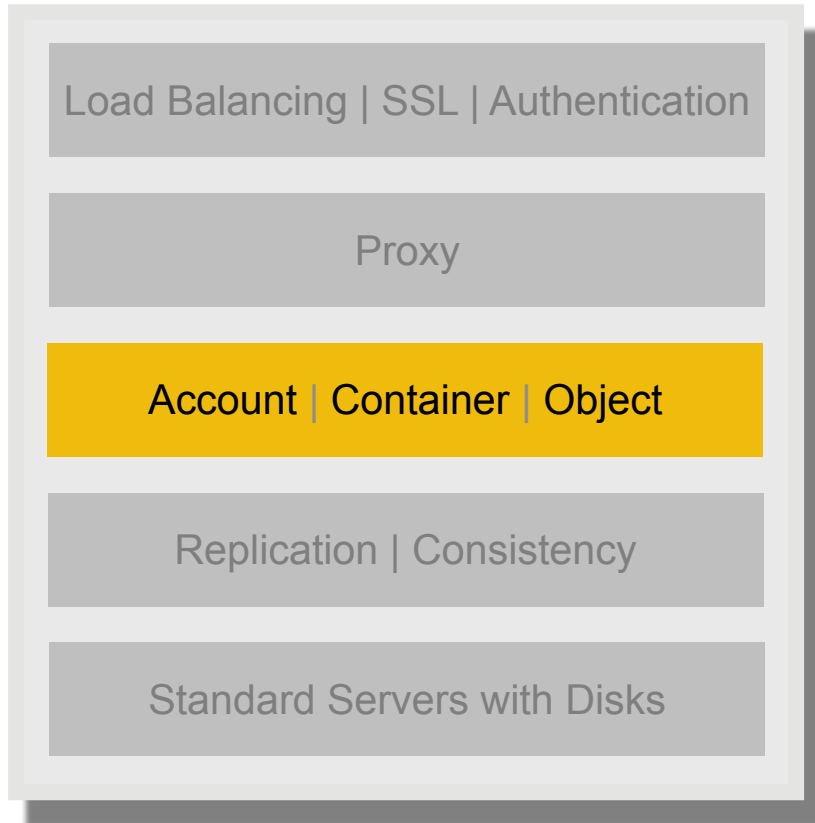
# Swift Architecture



## Proxy

- Only part of the cluster that “talks” to external clients
- Primarily with HTTP RESTful Swift API
- Routes requests from clients to disk
- Three replicas are simultaneously written
- Quorum required
- Uses single replica for reads
- Routes around failures
- Enforces ACLs set by user

# Swift Architecture



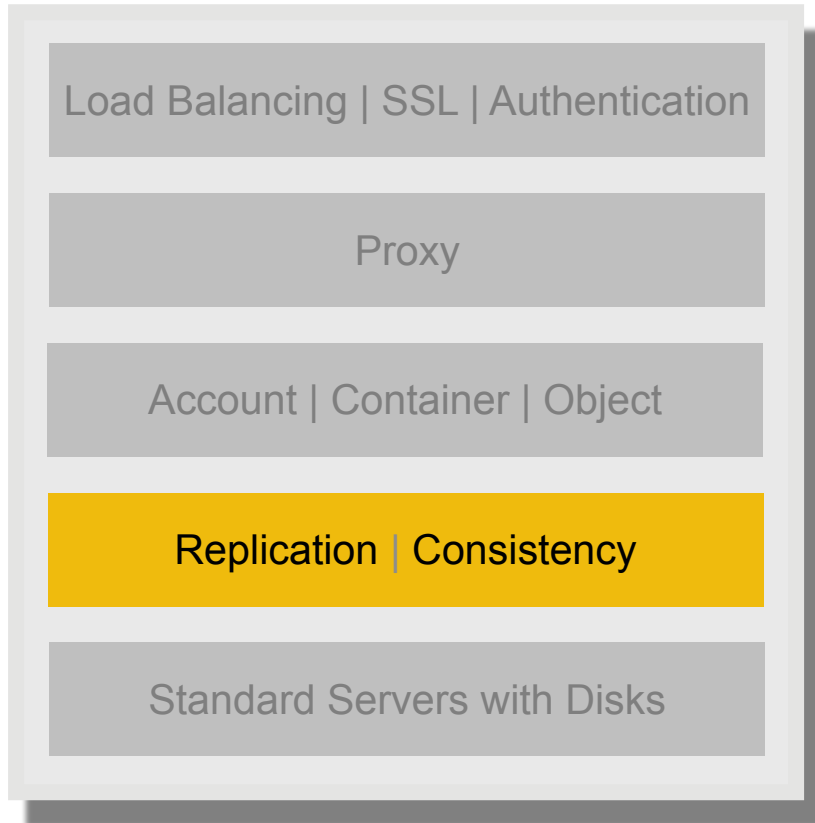
## Account / Container

- Accounts keep records of containers
- Containers keep records of objects

## Object

- The object servers store the data on disk
- Metadata is stored with the data
- Uses standard filesystem (XFS)

# Swift Architecture



## Replication

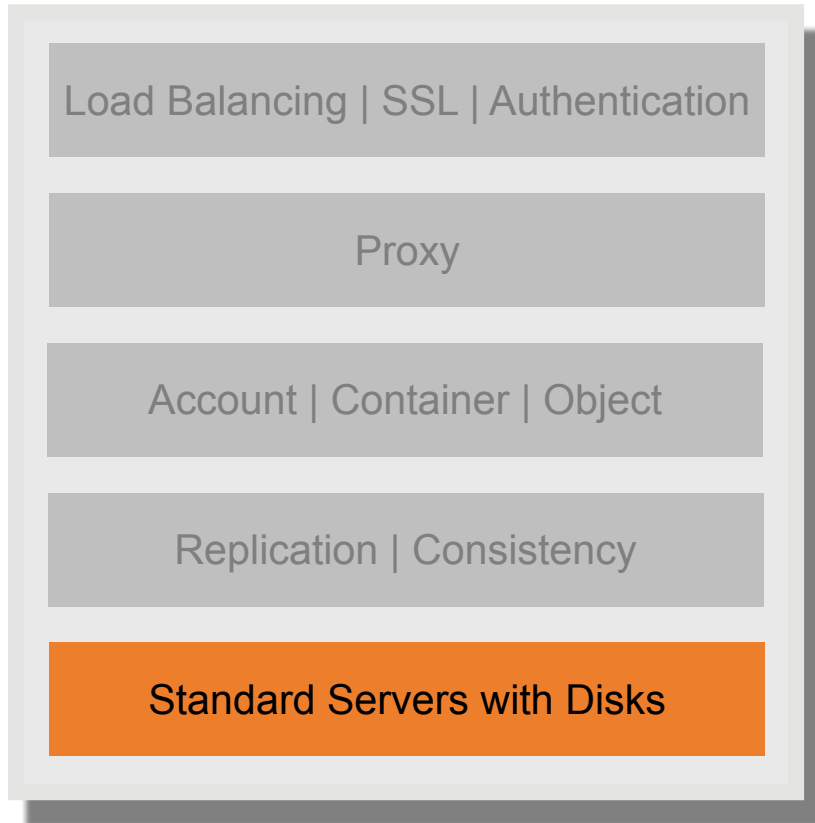
- Constantly checking for replicas status
- Only updates other replica sites, does not pull in newer versions of objects

## Consistency

- Constantly 'scrubbing' data to check for bad data

Note: These processes run in the background on Nodes where account, container, or object server processes are running.

# Swift Architecture



## Standard Servers

- Runs on standard server hardware
- SATA or SAS disks
- No RAID
  - Visibility to hardware beneath
  - Swift is already providing data redundancy
- Reduce cost

# Swift Architecture Summary

Load Balancing | SSL | Authentication

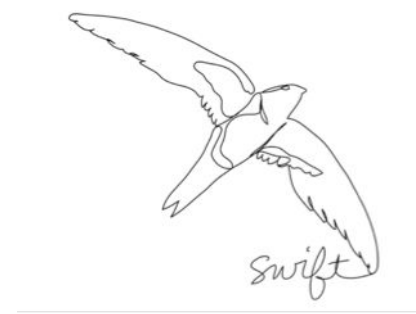
Proxy

Account | Container | Object

Replication | Consistency

Standard Servers with Disks

- Native HTTP API
- Scales Linearly
- No Single-point of Failure
- Standard Servers and Linux
- Extremely Durable
- Resilient to hardware failures
- Routes around network failures
- Consistency Model Enables Multi-Data Center



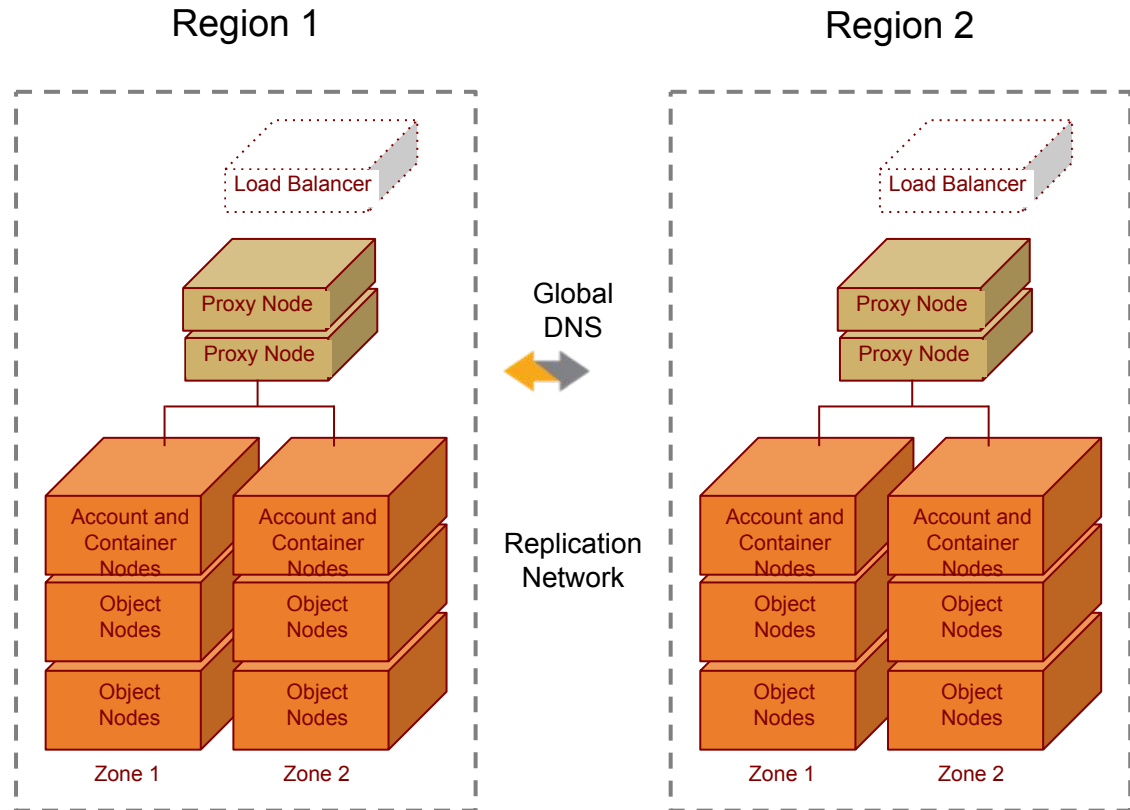
# OpenStack Swift

Multi-Region and  
Global Clusters

# Regions and Zones

## Regions

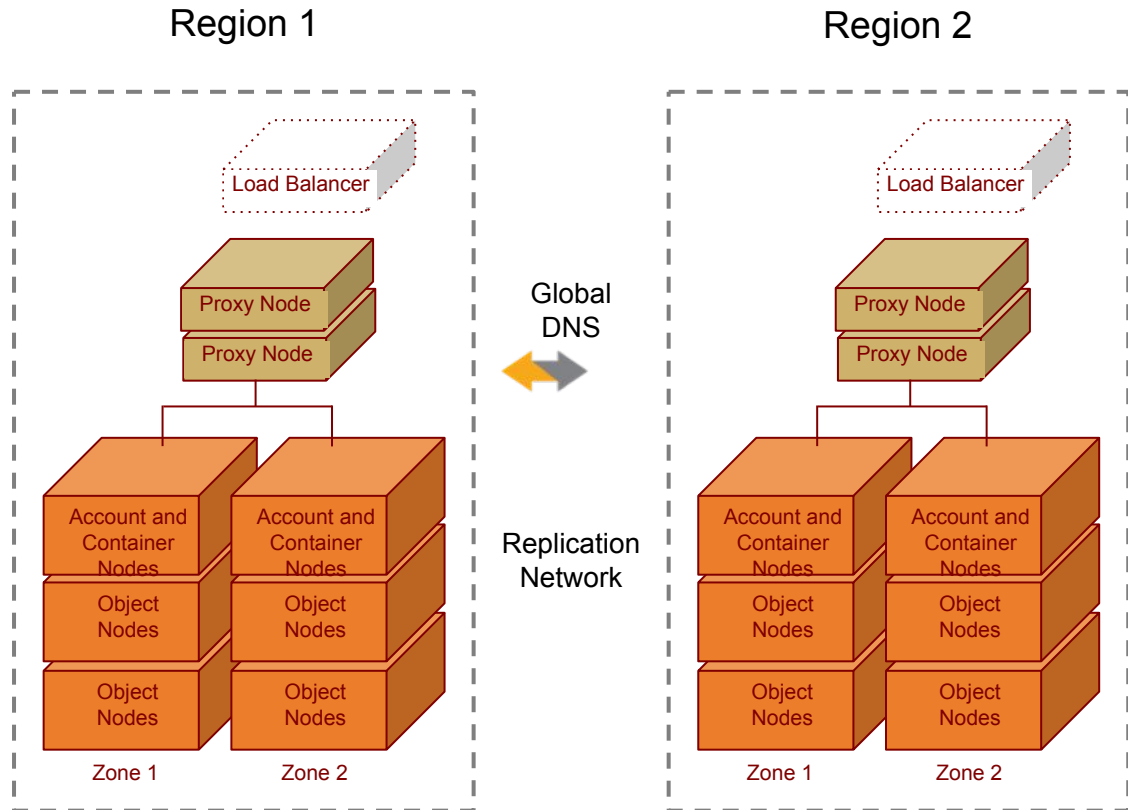
- Physically separate, often defined by geographical boundaries
- Minimum: 1 Region



# Regions and Zones

## Zones

- Regions contain one or more zones
- Designate a group of nodes sharing a set of physical hardware
- Also referred to as failure domains

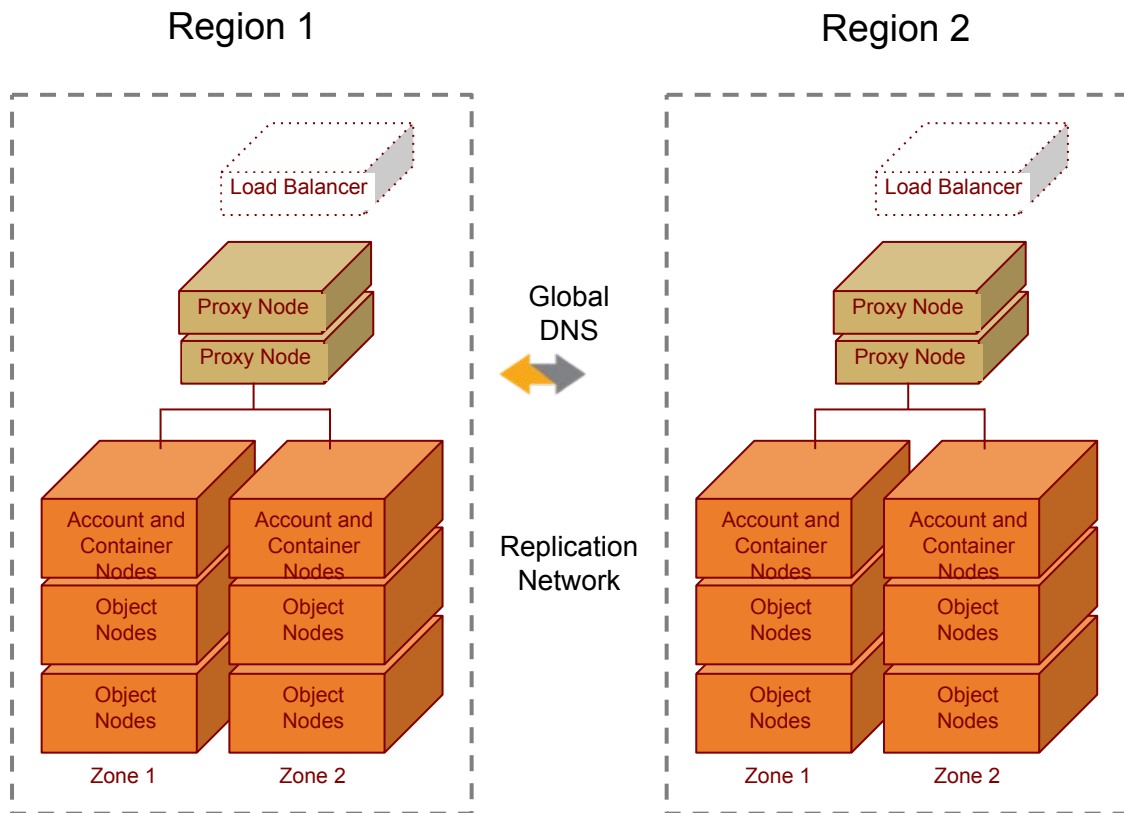




# Regions and Zones

## Multi-Region Clusters

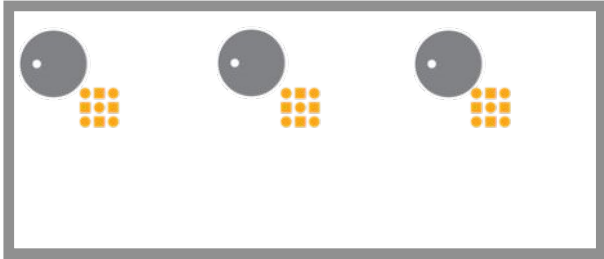
- Goal: Tolerate failures across regions and zones
- Policies control placement across regions
- Send requests to “closest” region



# Data Placement: “Unique as Possible”

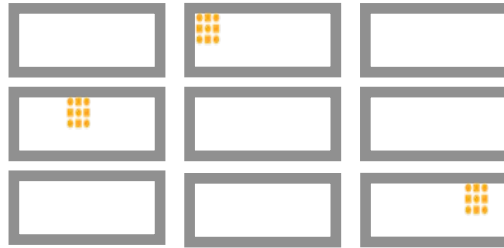
## Single Node Cluster

Disks are “as-unique-as-possible”



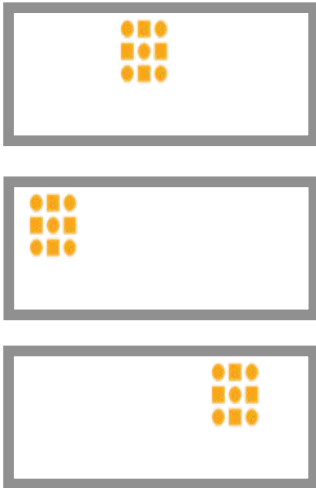
## Large Cluster

Storage Racks are “as-unique-as-possible”



## Small Cluster

Storage Nodes are “as-unique-as-possible”



## Multi-Region

Distributed data centers are “as-unique-as-possible”



# Storage Policies

## Benefits

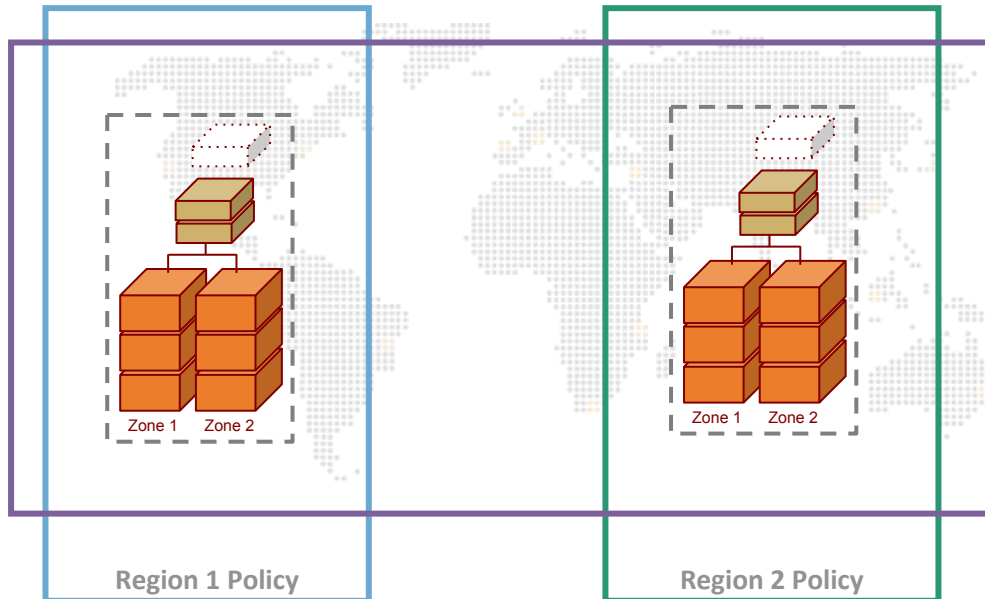
- Optimizes storage for applications and users
- Consolidates storage tiers under one system
- Simplifies management
- Lowers storage infrastructure TCO

## Policies Encompass:

- Storage media
- Number of replicas
- Erasure Codes / Replicas

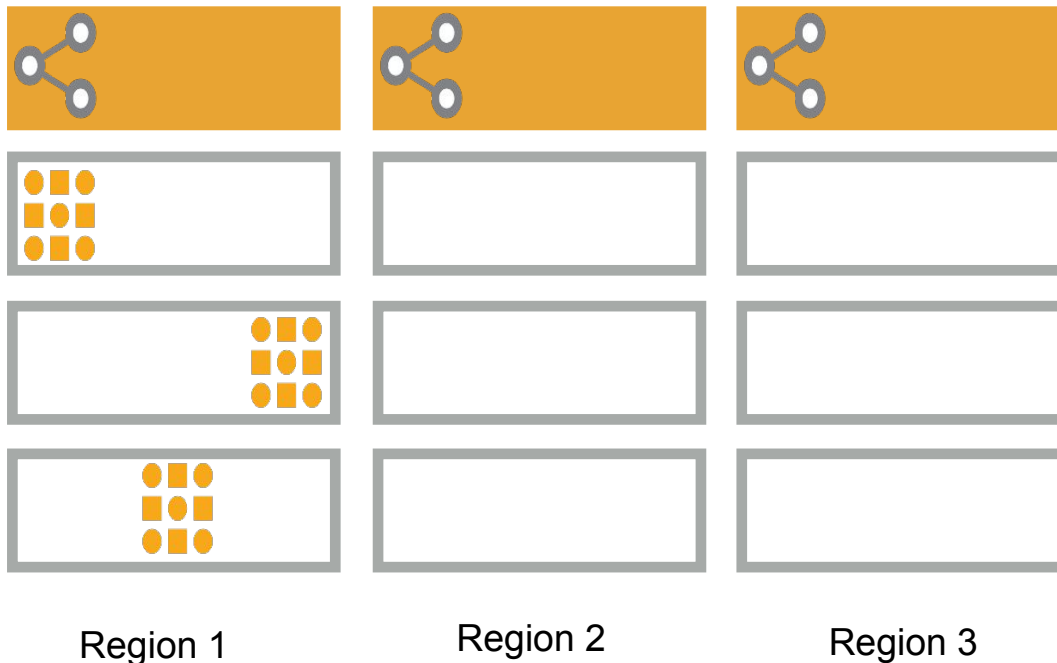
## Common Policy Groupings:

- According to performance
- According to geography or political boundary
- According to data protection scheme, e.g., number of replicas or erasure coding



# MRC: Write Affinity

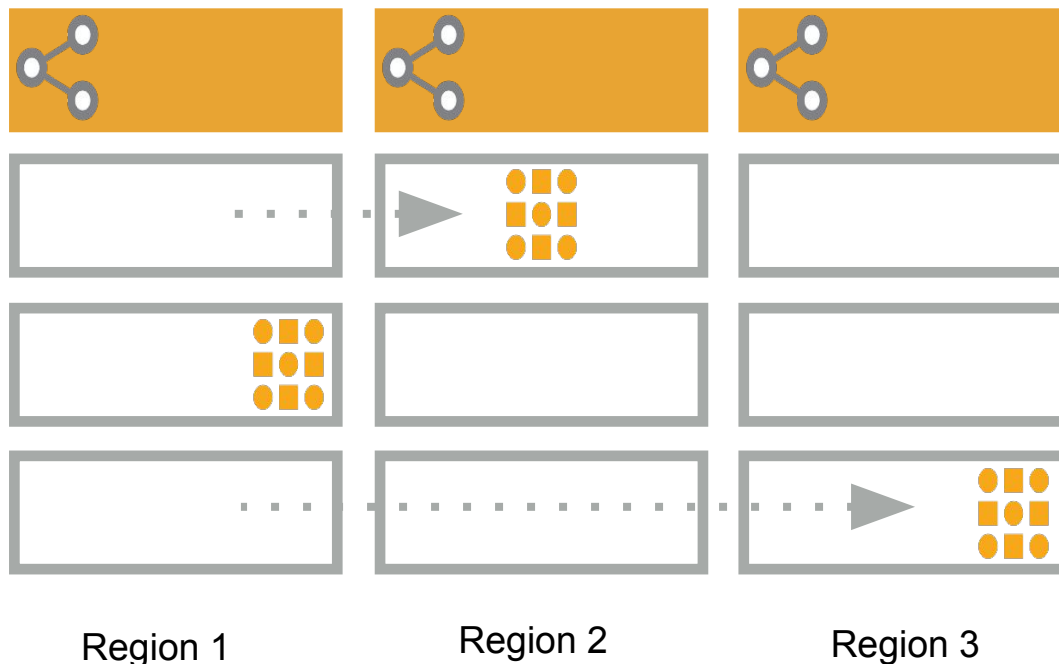
## Proxy Write Affinity



- By default, **objects are written to all the locations simultaneously**
- **Write Affinity:** writes all copies locally then transfer asynchronously to other regions

# MRC: Write Affinity

## Proxy Write Affinity



- By default, **objects are written to all the locations simultaneously**
- **Write Affinity:** writes all copies locally then transfer asynchronously to other regions

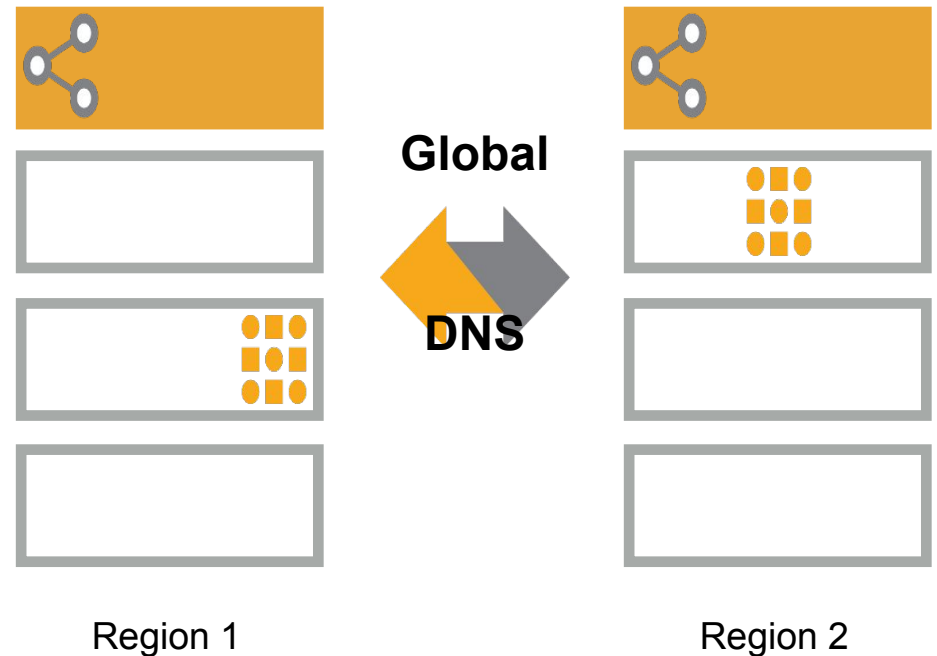
# MRC: Read Affinity

## Prioritizes “Nearby”

- Zone/Region, where am I?
- Latency to Storage Node

## DNS Routes User

- Each proxy pool has its own hostname
- Routes user to closest region



# Use Case – Backup

## Backup

- Enterprise: Workstations
- MSP: Cloud backup service

## Behavior

- Write optimized
- High throughput
- Low concurrency

# Use Case – Big Data / MapReduce

Hadoop / Spark

Philosophy: Let HDFS do what it's best at:  
Serving data where you want it when you need it

Swift for warm and cold data storage. Why?

- Durability and reliability guarantees
- Managed capacity: Grow to and beyond petabytes
- Easier integration with various data input sources
- Share results using a common storage platform

Example: Run MapReduce job

- Read input data from SwiftStack
- Write transient results to HDFS
- Write result to Swift