




# Swift架构与实践

@SinaAppEngine



杨雨/2012-08-11

# Content

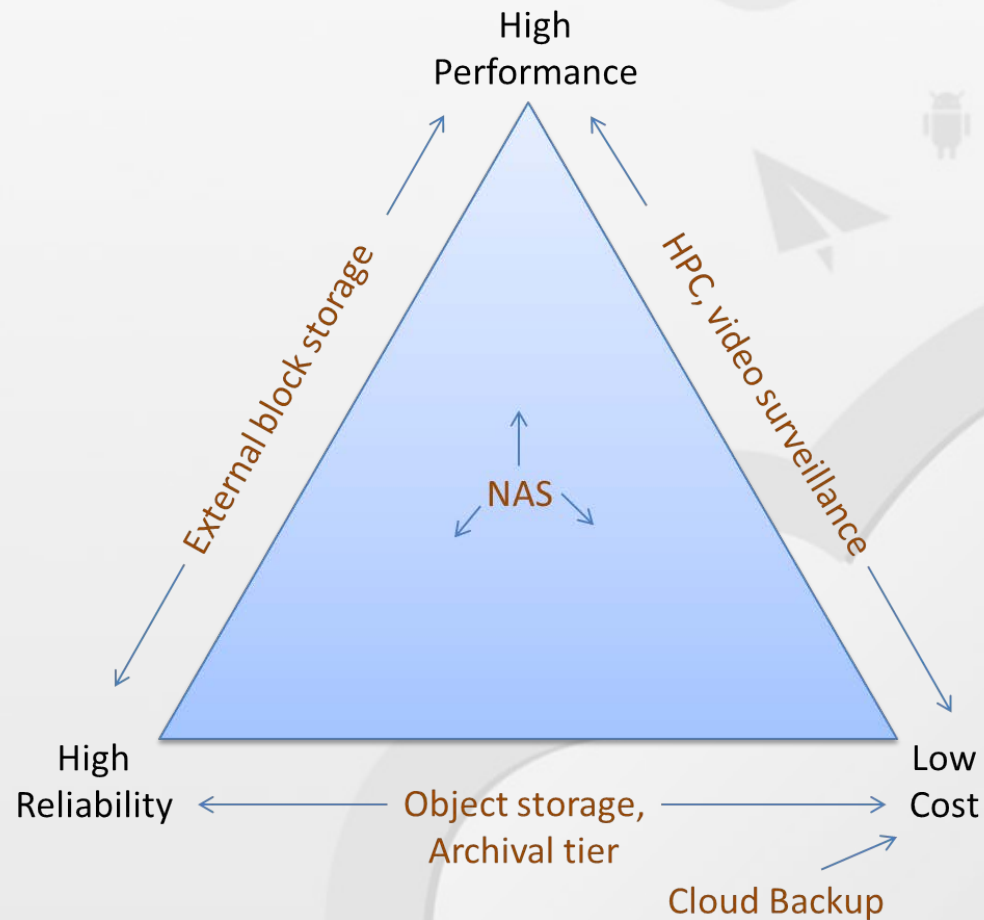
---

- 1 Principles and Architecture
- 2 The Practice of Swift @SinaAppEngine
- 3 Problems and Improvements

# Storage Types

<i>Types</i>	<i>Protocol</i>	<i>Application</i>
<i>Block Storage</i>	<i>SATA, SCSI, iSCSI</i>	<i>SAN, NAS, EBS</i>
<i>File Storage</i>	<i>Ext3/4, XFS, NTFS</i>	<i>PC, Servers, NFS</i>
<i>Object Storage</i>	<i>HTTP, REST</i>	<i>Amazon S3, Google Cloud Storage, Rackspace Cloud Files</i>
<i>Specific Storage</i>	<i>Specific protocol based on tcp</i>	<i>MySQL, MongoDB, HDFS</i>

# Storage Types



From: <http://www.buildcloudstorage.com/2012/08/is-openstack-swift-reliable-enough-for.html>

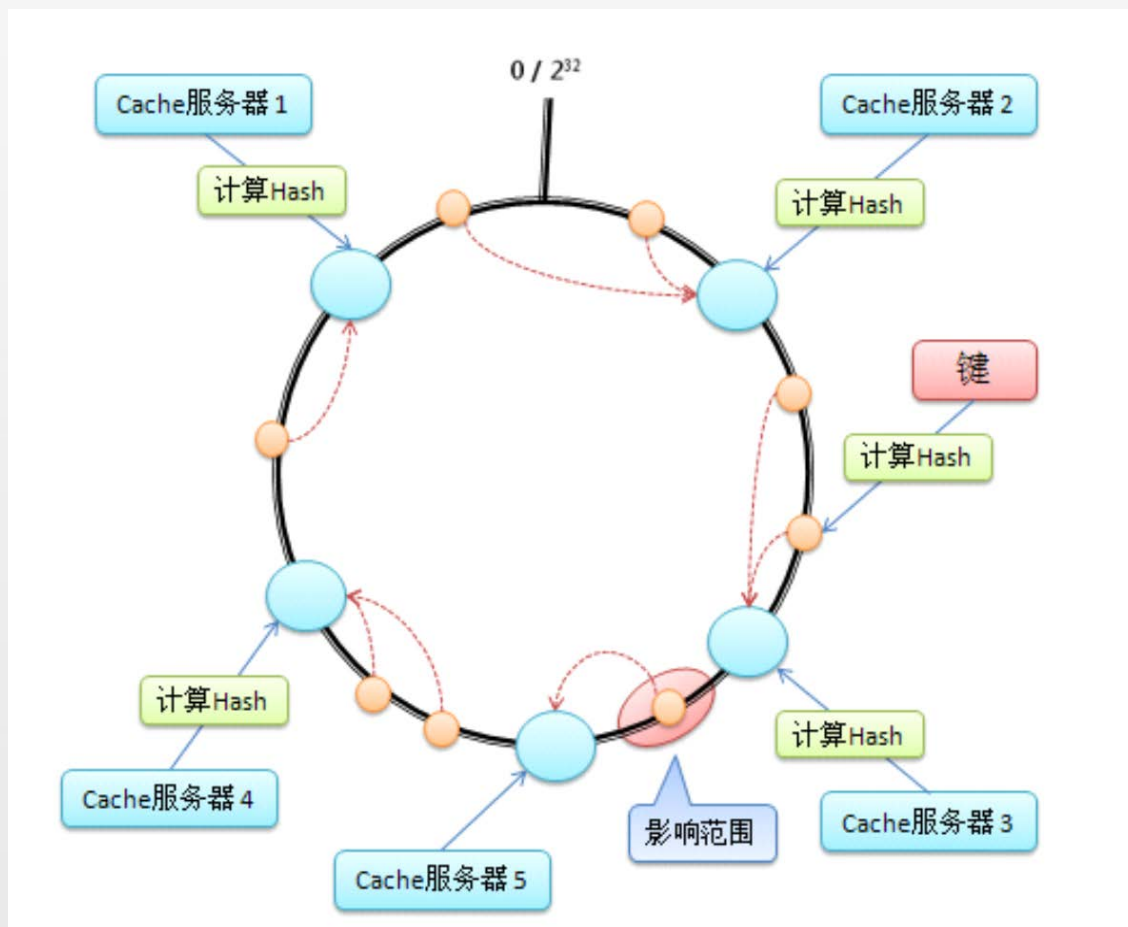
# Targets

---

- High Reliability = High Durability + High Availability
- High Durability - Replicas and Recovery
- High Availability - Replicas and Partition
- Low Cost - Commodity Hardware
- Scale Out - No Single Bottleneck, Share Nothing

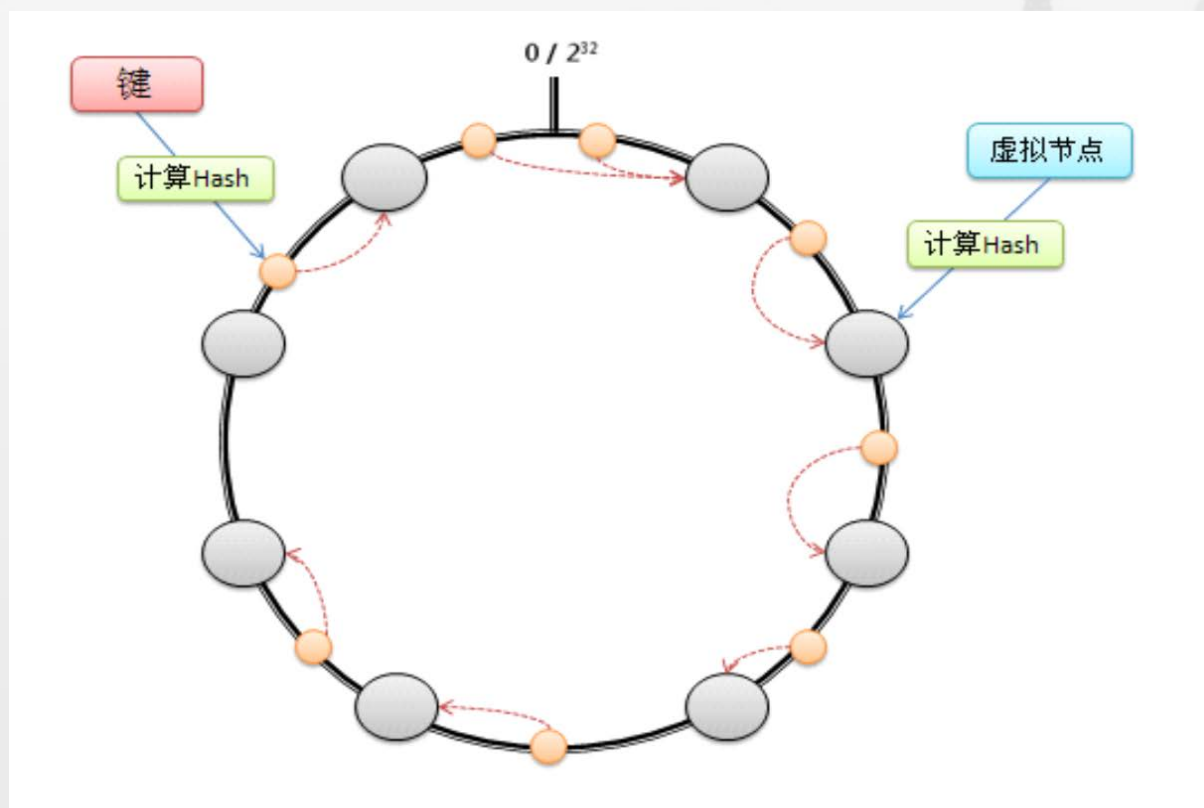
# Reliability

## Consistent Hashing



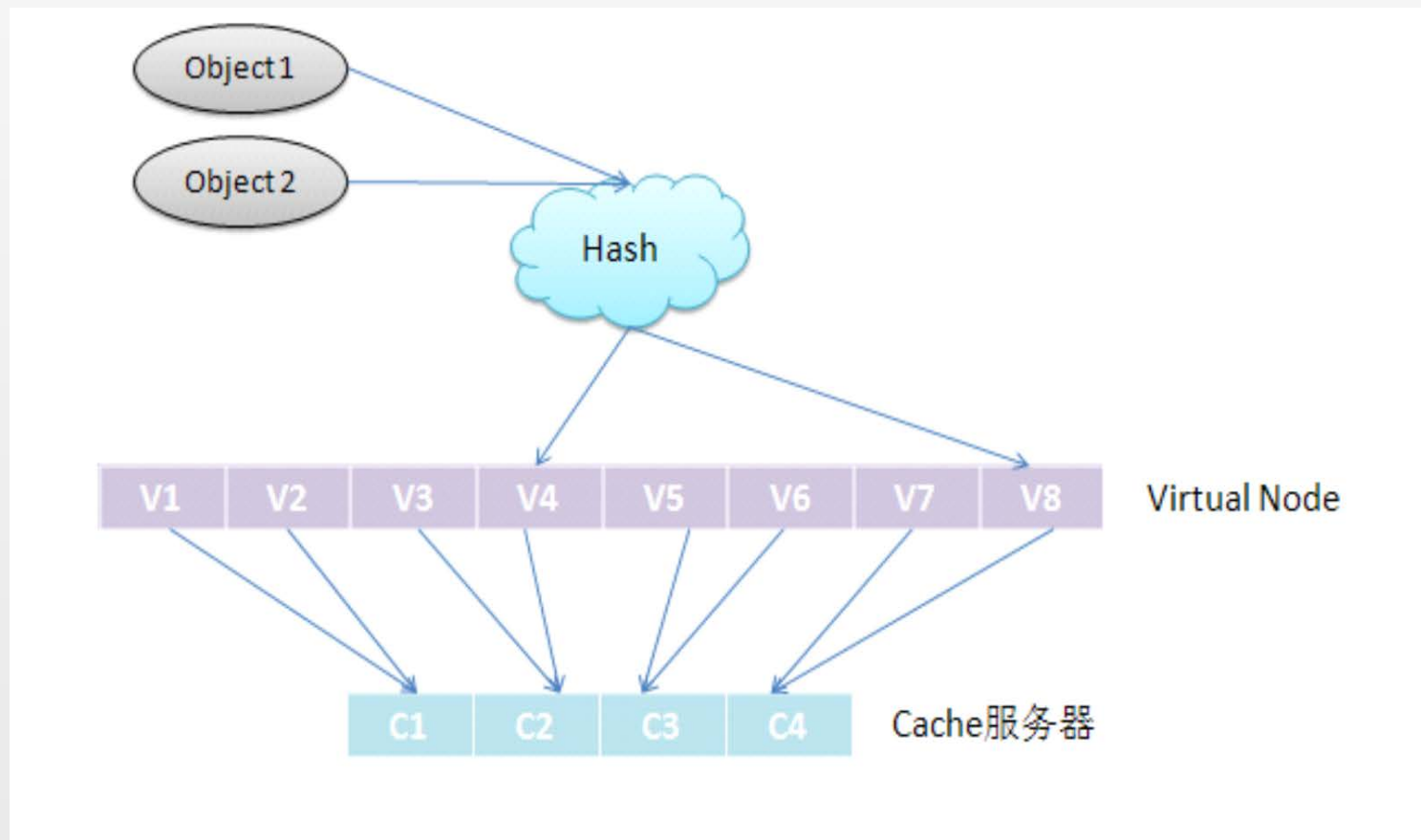
# Reliability

## Consistent Hashing with Virtual Node



# Reliability

## Consistent Hashing with Virtual Node





# Reliability

---

The advantages of consistent hashing:

1. Metadata is small;(AWS S3 has 762 billion objects)
2. Distribution uniformity;
3. Peer to peer communication;
4. Load balance.

# Reliability

---

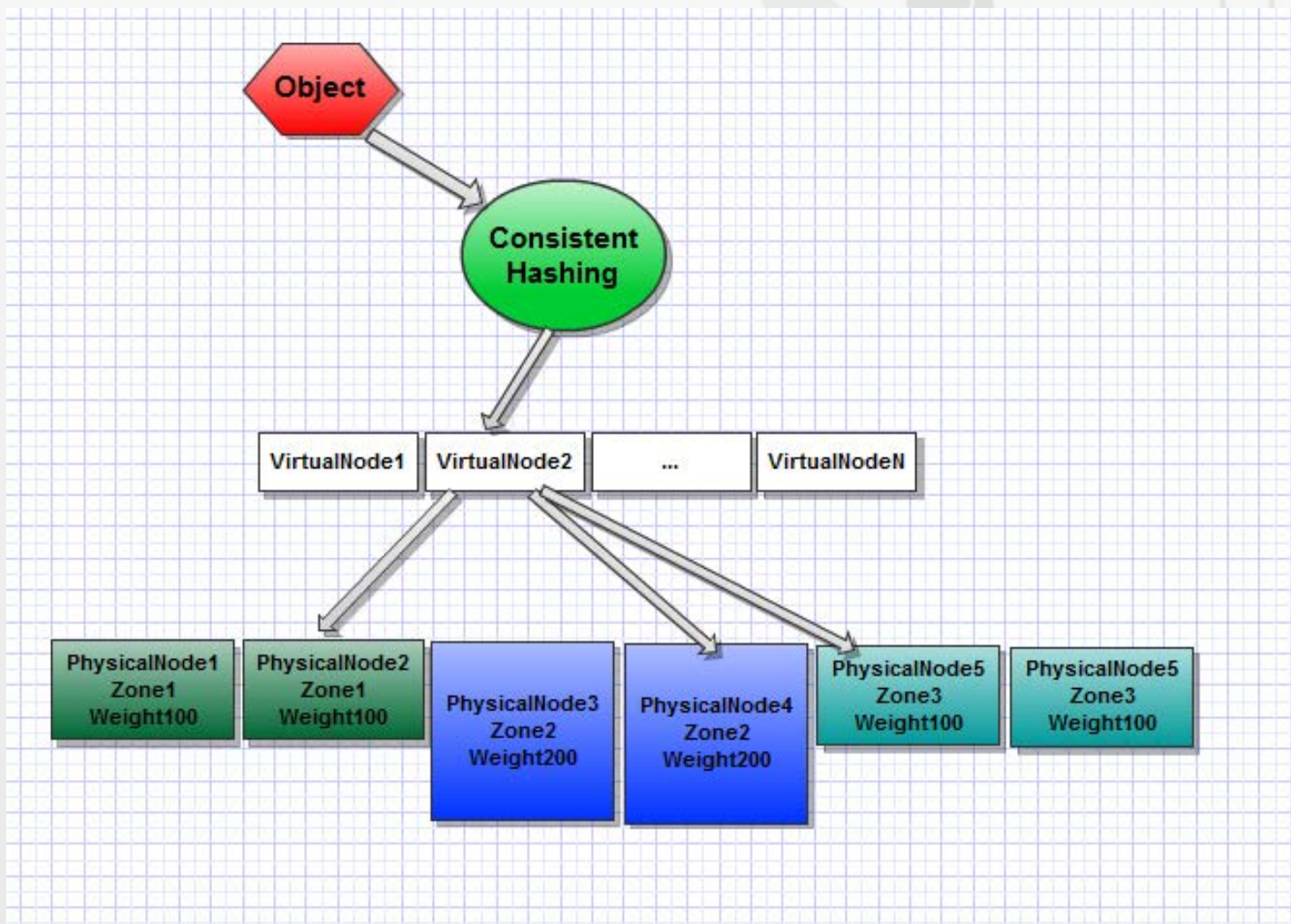
- Virtual node(partition) - Object distribution uniformity
- Weight - Allocate partitions dynamically
- Zone - Partition Tolerance

Zones can be used to group devices based on physical locations, power separations, network separations, or any other attribute that would lessen multiple replicas being unavailable at the same time.

Swift :

Ring - The index file for locating object in the cluster.

# Reliability



# Reliability

---

**Durability: 99.999..9%?** AWS S3 with 99.999999999%

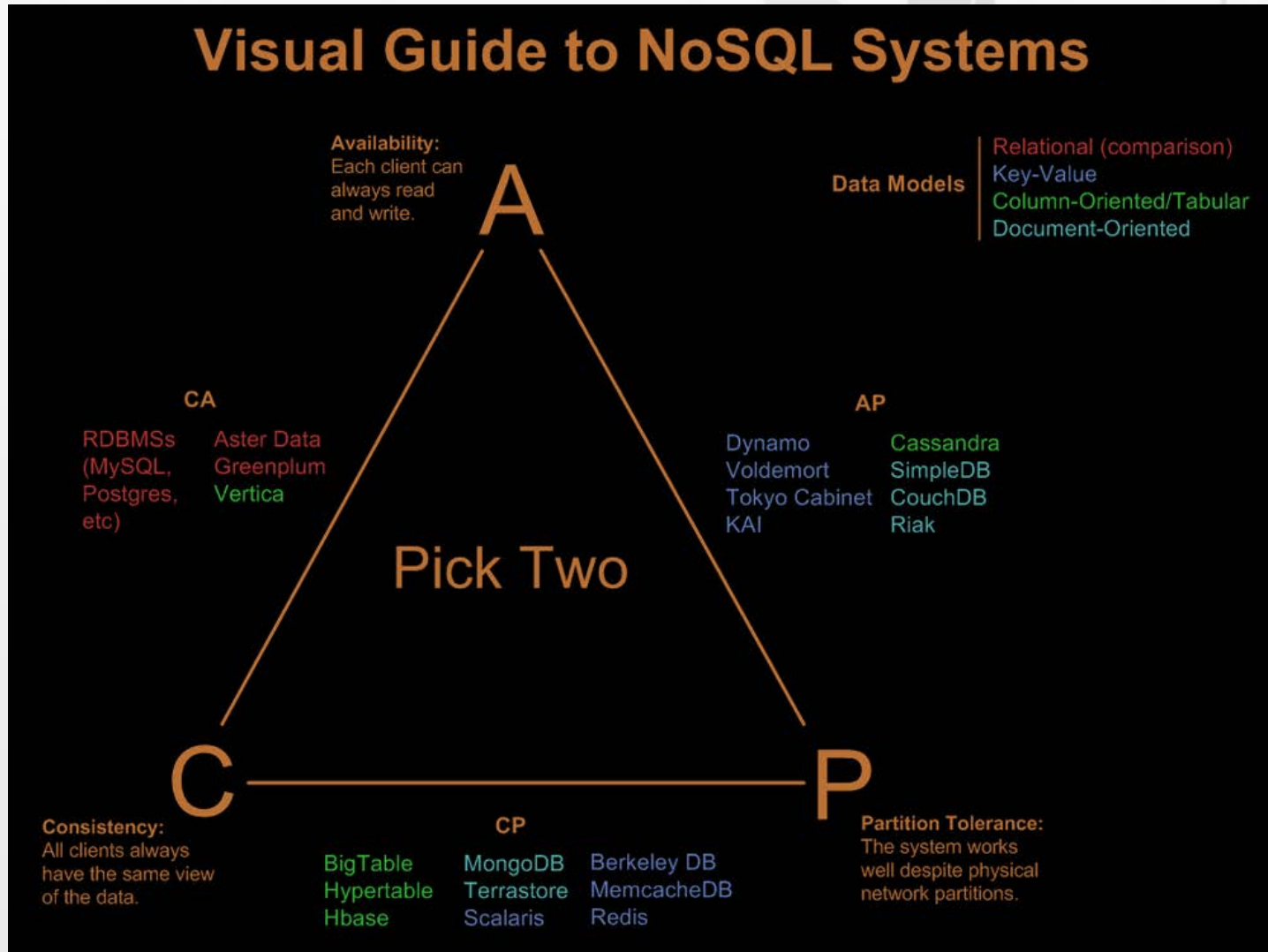
- Annualized failure rate of devices.(disk, network...)
- Bit rot rate and the time to detect bit rot.
- Mean time to recovery.
- More about durability :  
<http://www.buildcloudstorage.com/2012/08/is-openstack-swift-reliable-enough-for.html>

Swift :

Auditor - To detect bit rot;

Replicator-To keep the consistency of object;

# Consistency Model



# Consistency Model

Quorum + Object Version + Async Replication

Quorum Protocol:

$N$ , the number of nodes that store replicas of the data;

$W$ , the number of replicas that need to acknowledge the receipt of the update before the update completes;

$R$ , the number of replicas that are contacted when a data object is accessed through a read operation;

If  $W+R > N$ , then the write set and the read set always overlap and one can guarantee **strong consistency**.

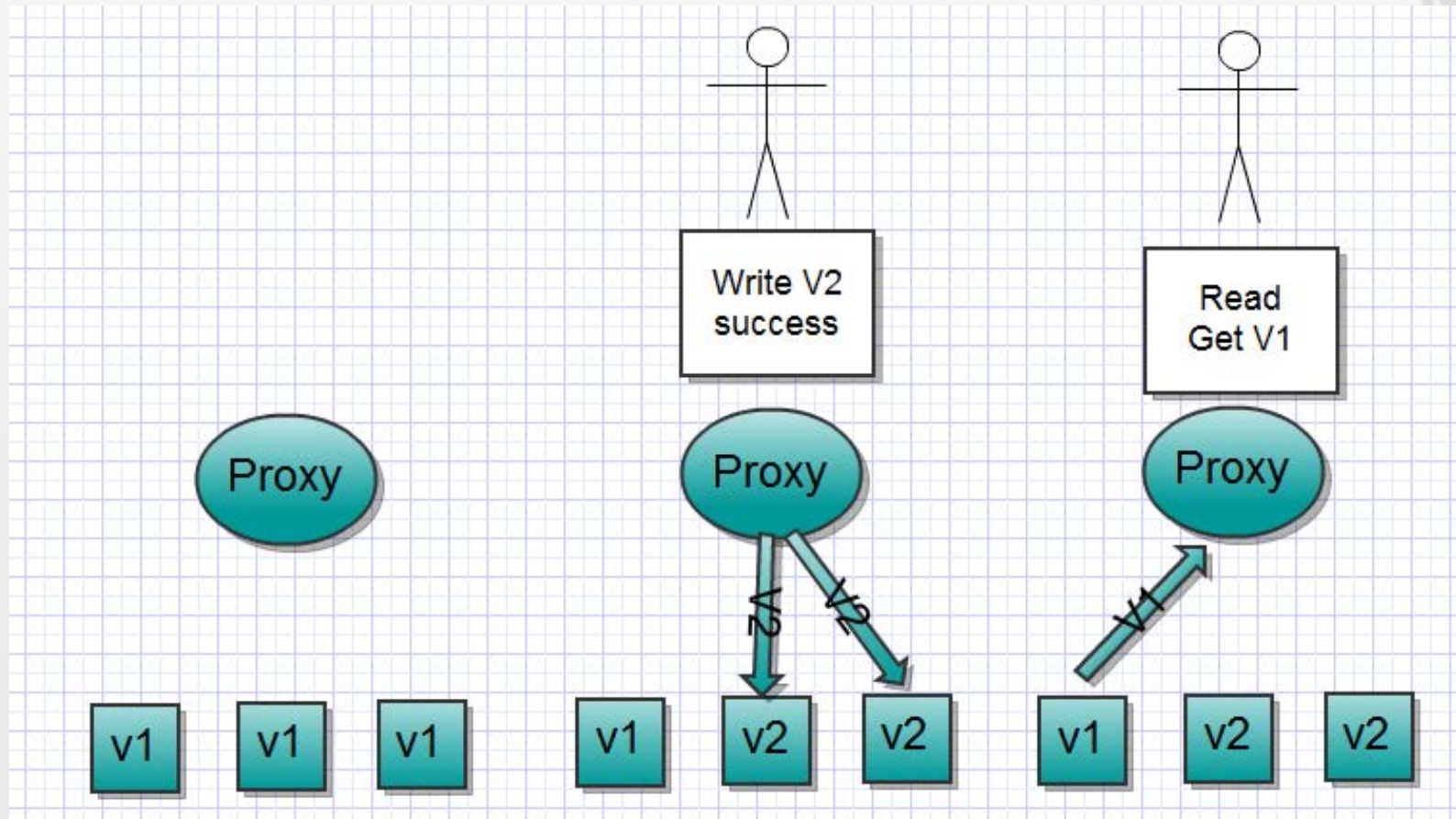
In Swift, NWR is configurable.

General configuration:  $N=3$ ,  $W=2$ ,  $R=1$  or  $2$ , So the swift can provide two models of consistency, strong and eventual.



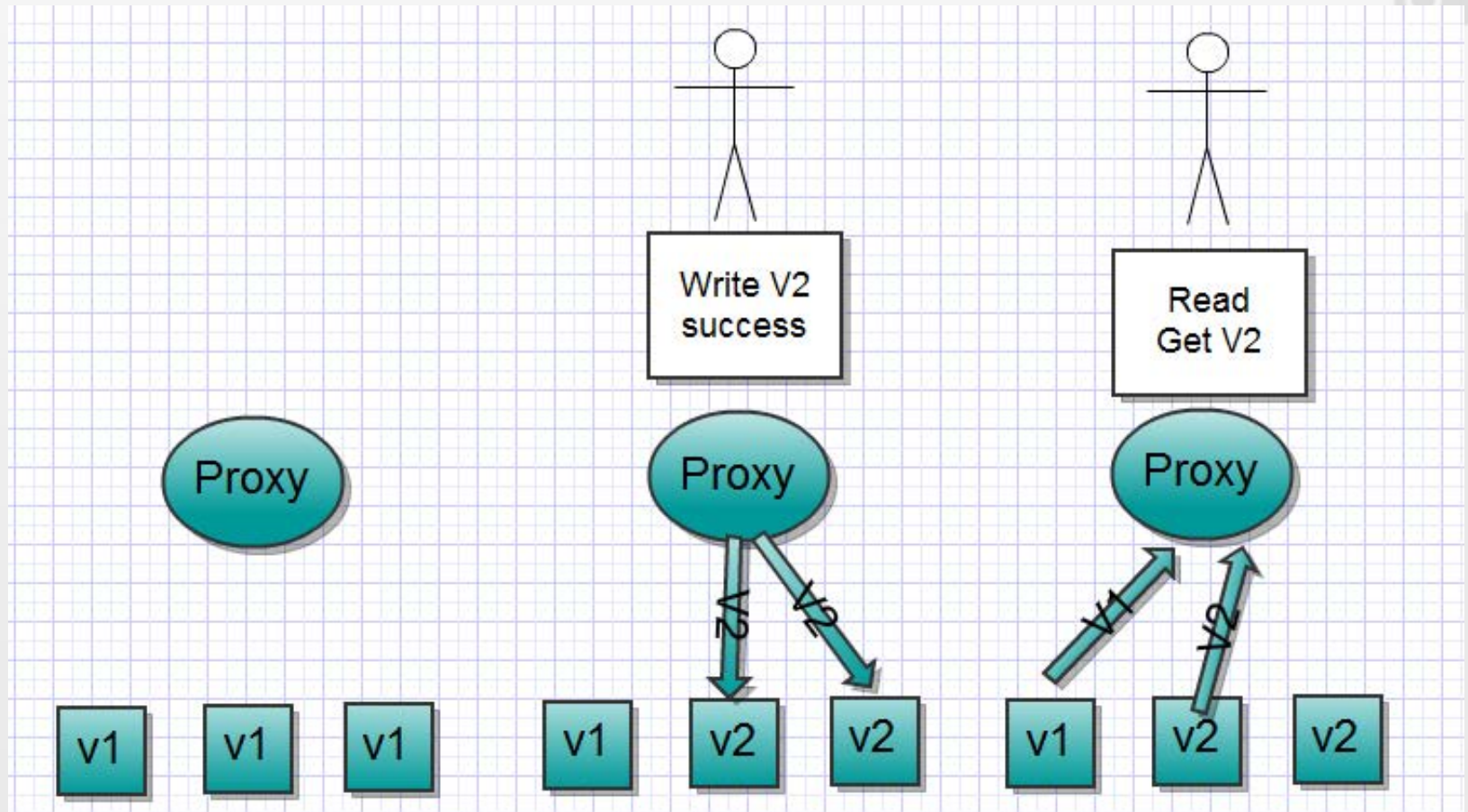
# Consistency Model

Weak Consistency (  $N=3, W=2, R=1$  )



# Consistency Model

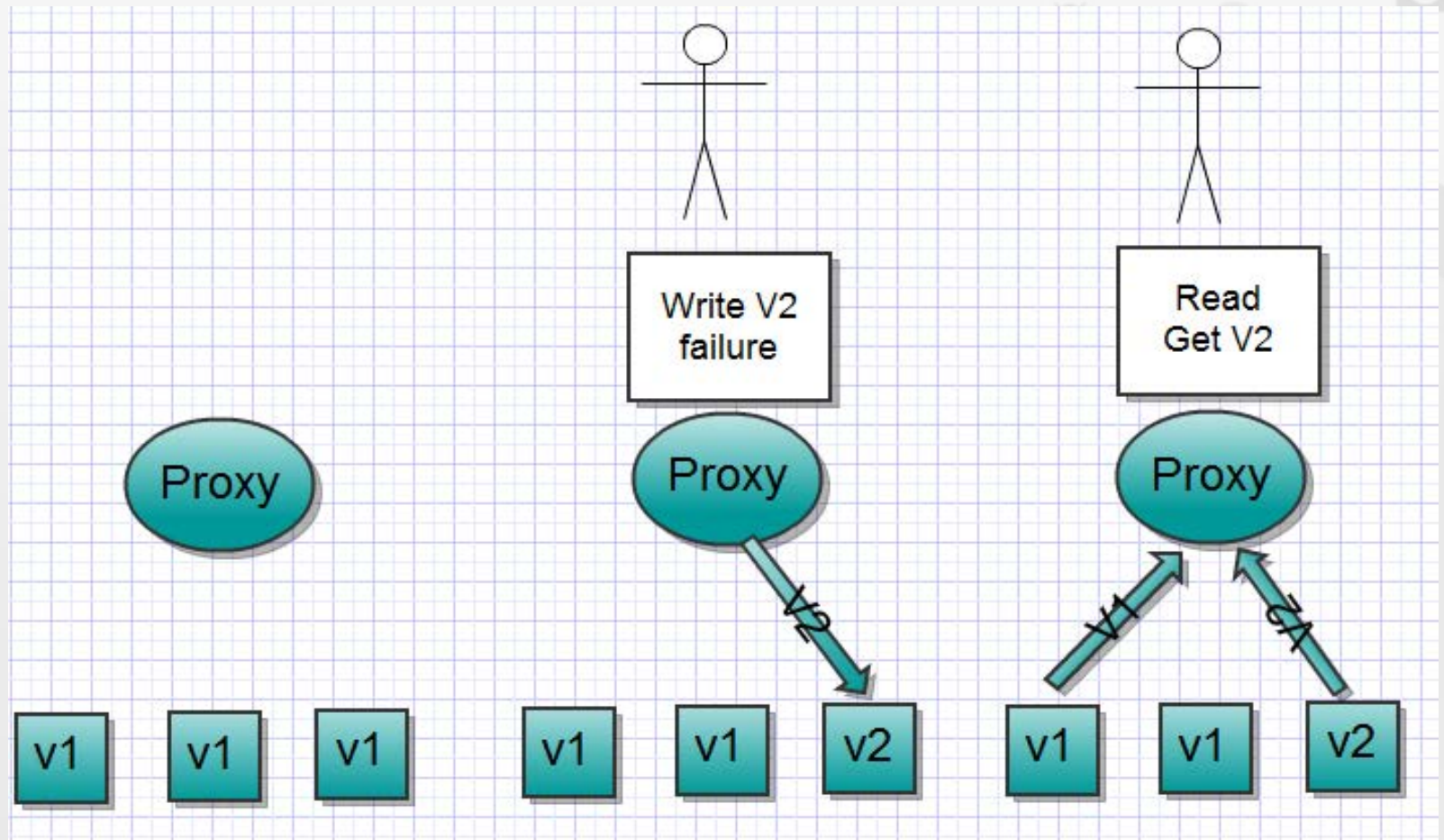
Strong Consistency (N=3,W=2,R=2)



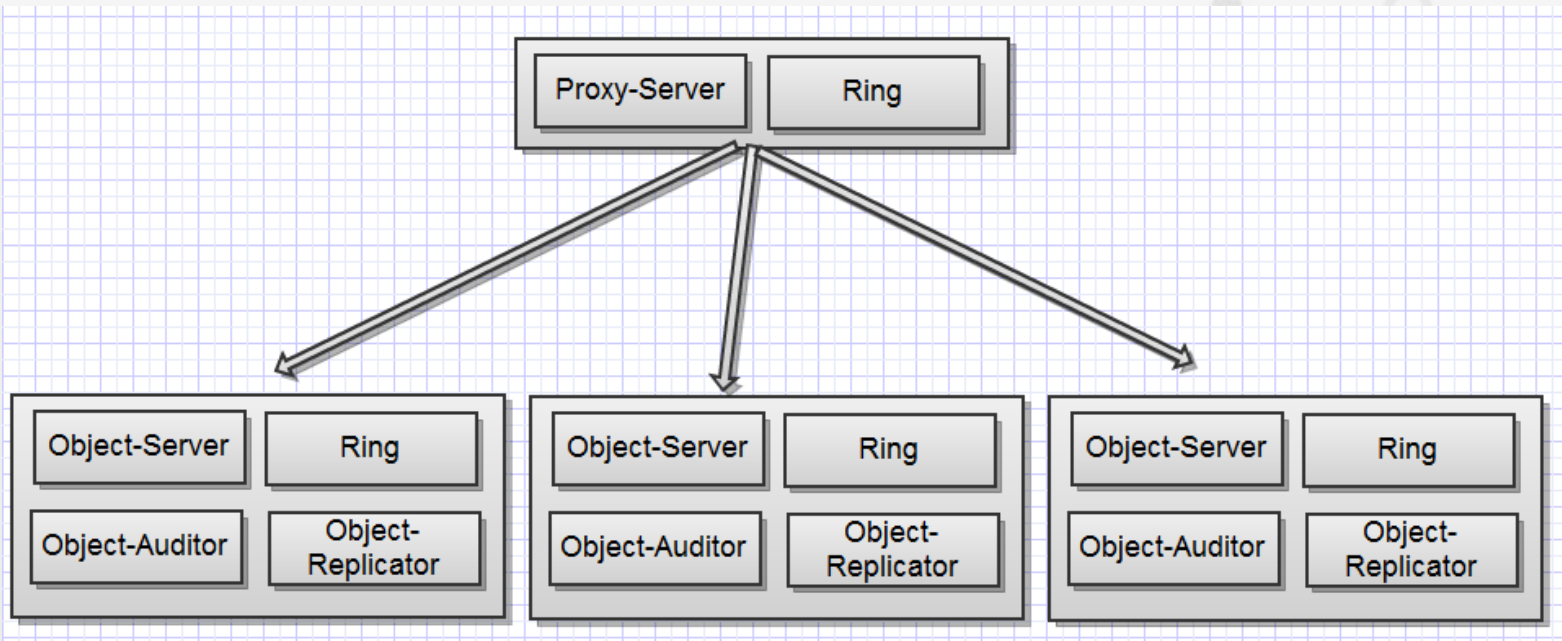


# Consistency Model

Special Scene : dirty read



# Architecture Prototype



# Metadata

---

```
account
|-- container1
|-----obj1
|-----obj2
|-- container2
|-----objN
```

How to store the relationship of account, container and object?

- Relation Database
- NoSQL, Cassandra, MongoDB
- Relation Database with Sharding

# Metadata

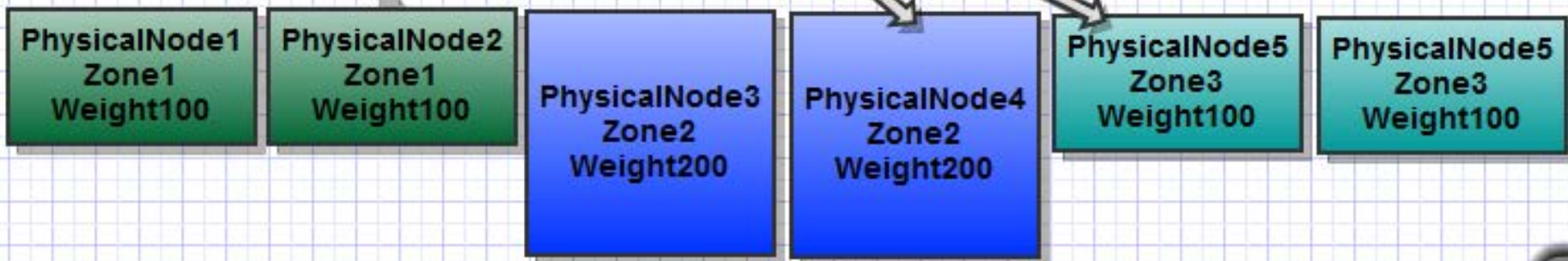
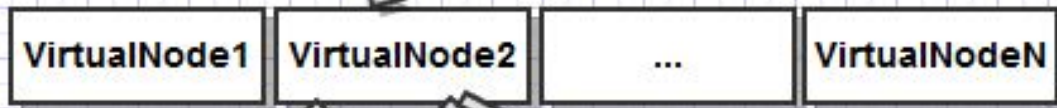
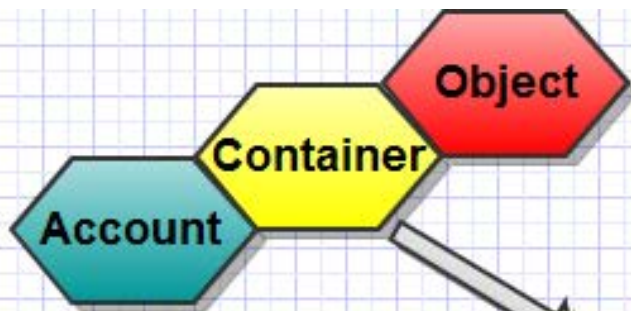
---

The target of swift: no single failure, no bottleneck, scale out

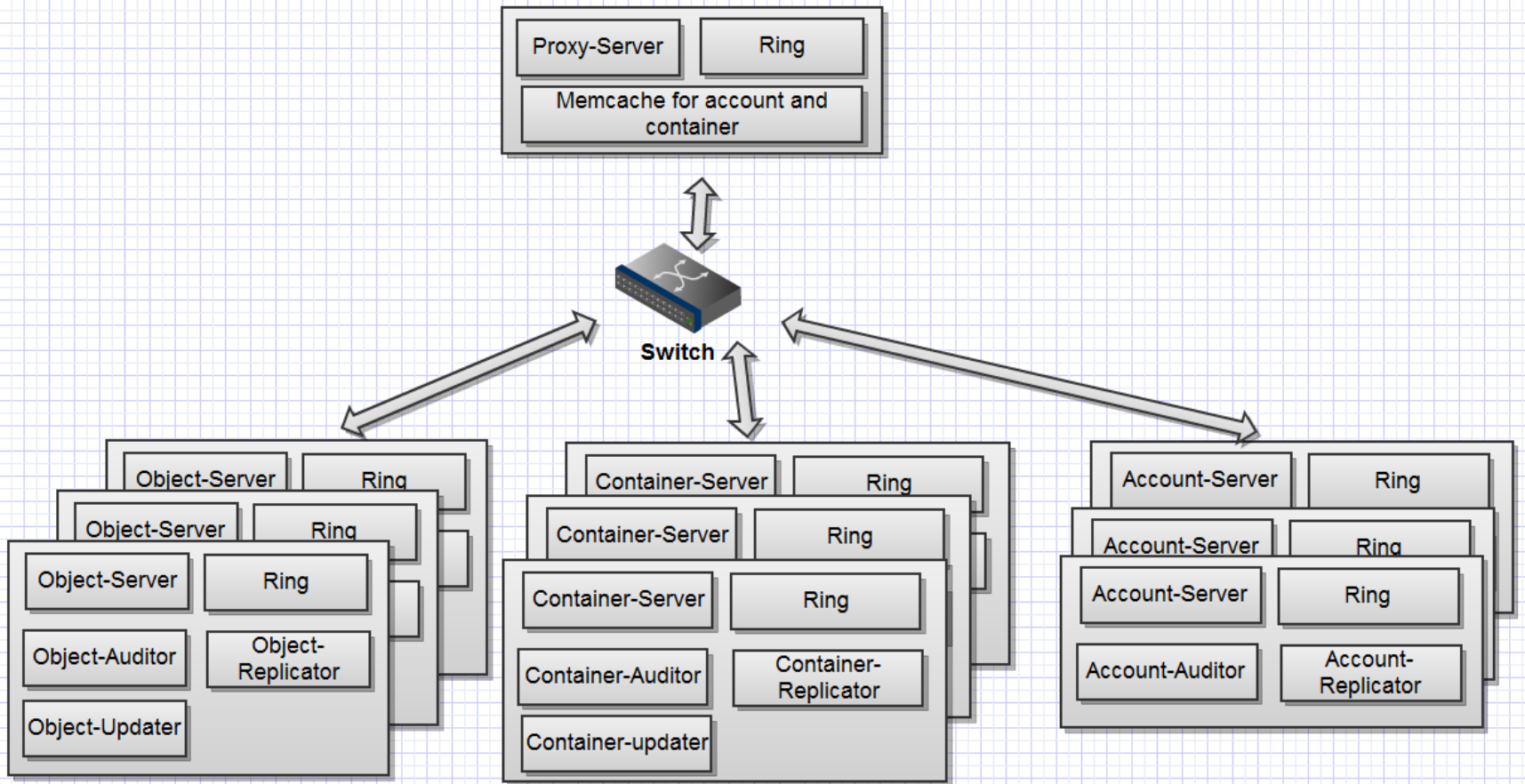
The Swift way: sqlite + consistent hashing + quorum

A sqlite db file is an object.

So the database is HA, durable and with eventual consistency.

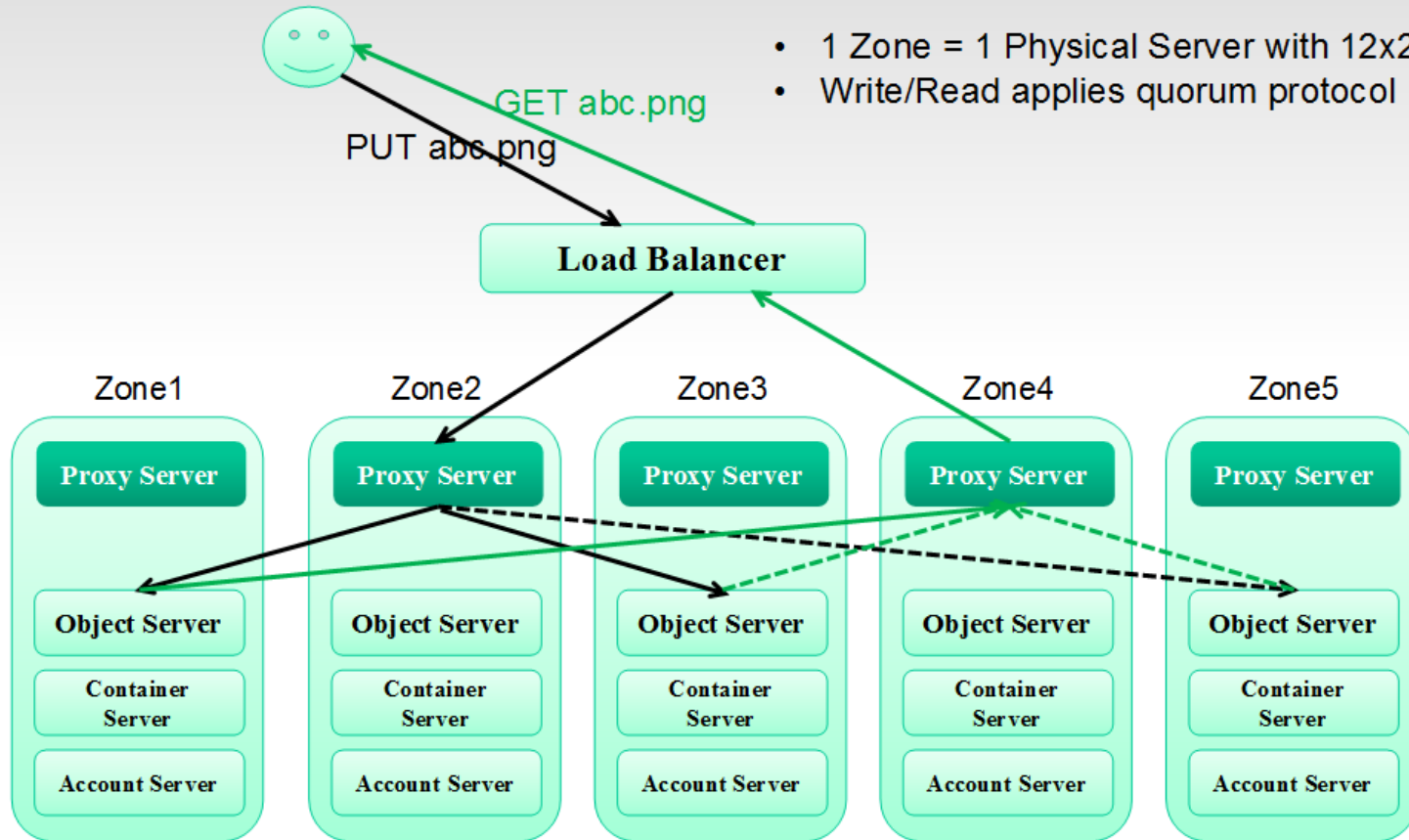


# Architecture

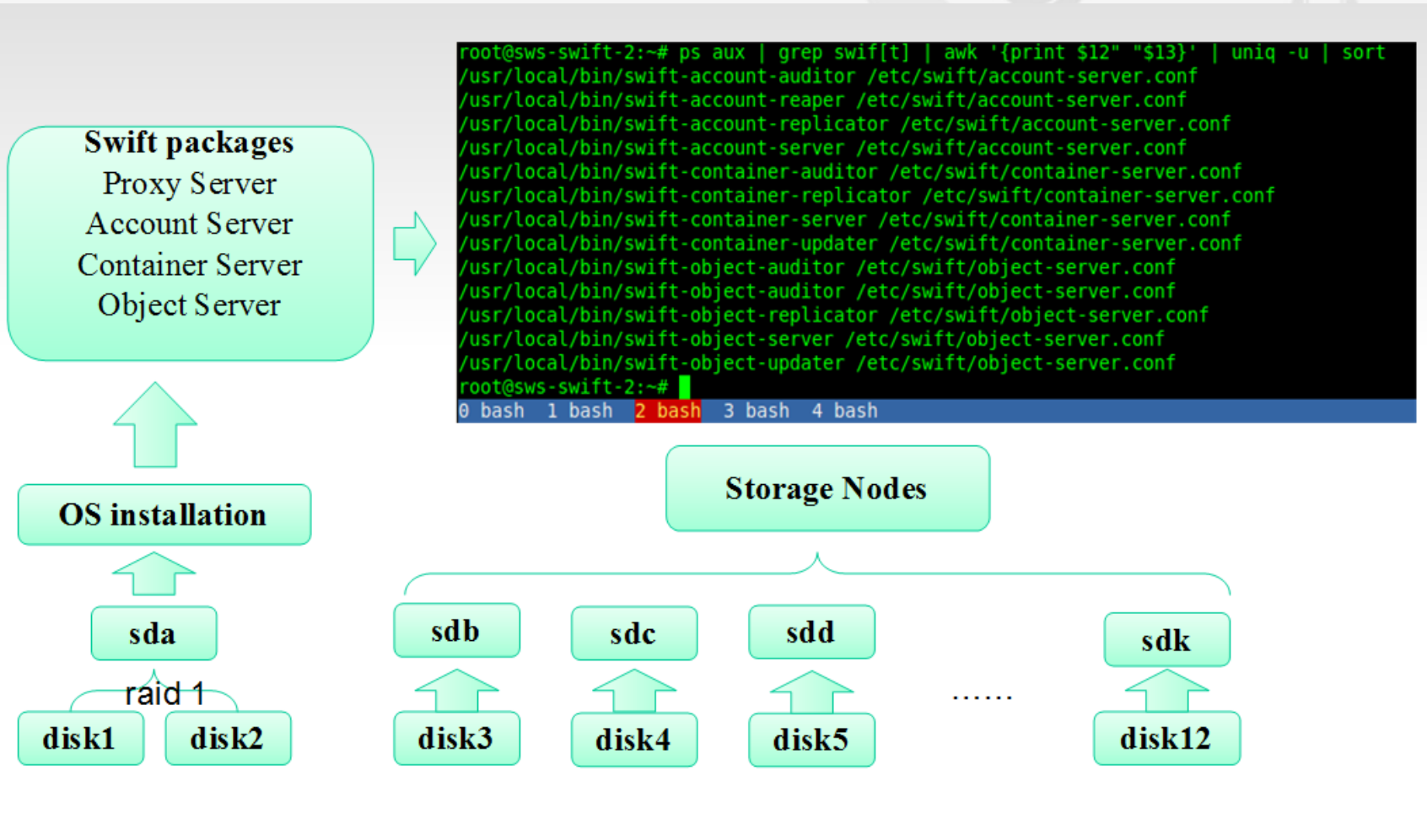




# Swift Practice@SinaAppEngine



# Swift Practice@SinaAppEngine





# Our Works

---

## •Swift as the SAE Storage

- Auth module for SAE(Key-Pair)
- Keystone for SAE(Token)
- HTTP Cache-Control module
- Quota(limit the number of containers and objects, limit the storage usage)
- Domain remap  
app-domain.stor.sinaapp.com/obj To  
sinas3.com/v1/SAE\_app/domain
- Rsync with bwlimit
- Billing
- Storage Firewall

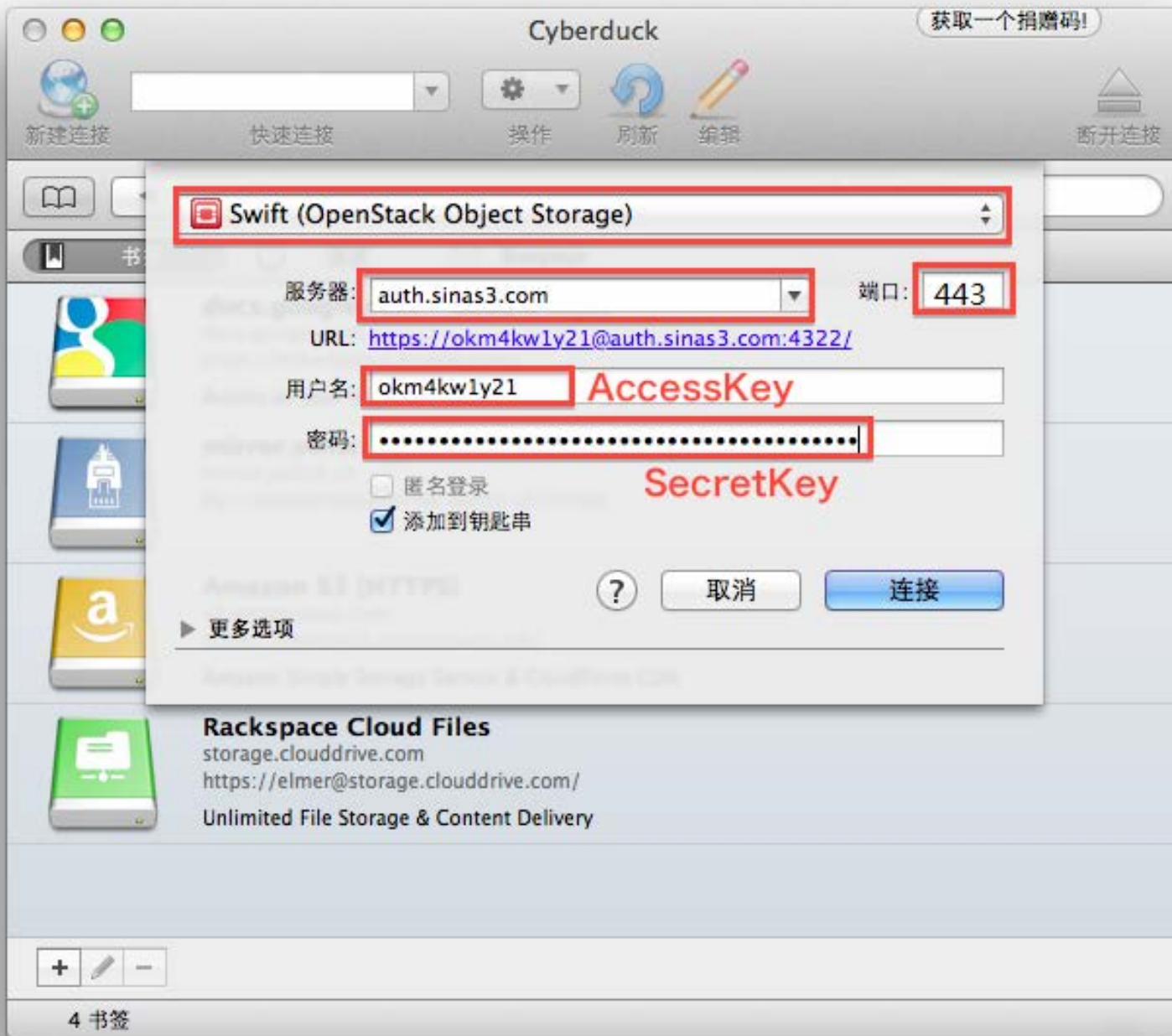
## •Swift as the SWS Simple Storage Service

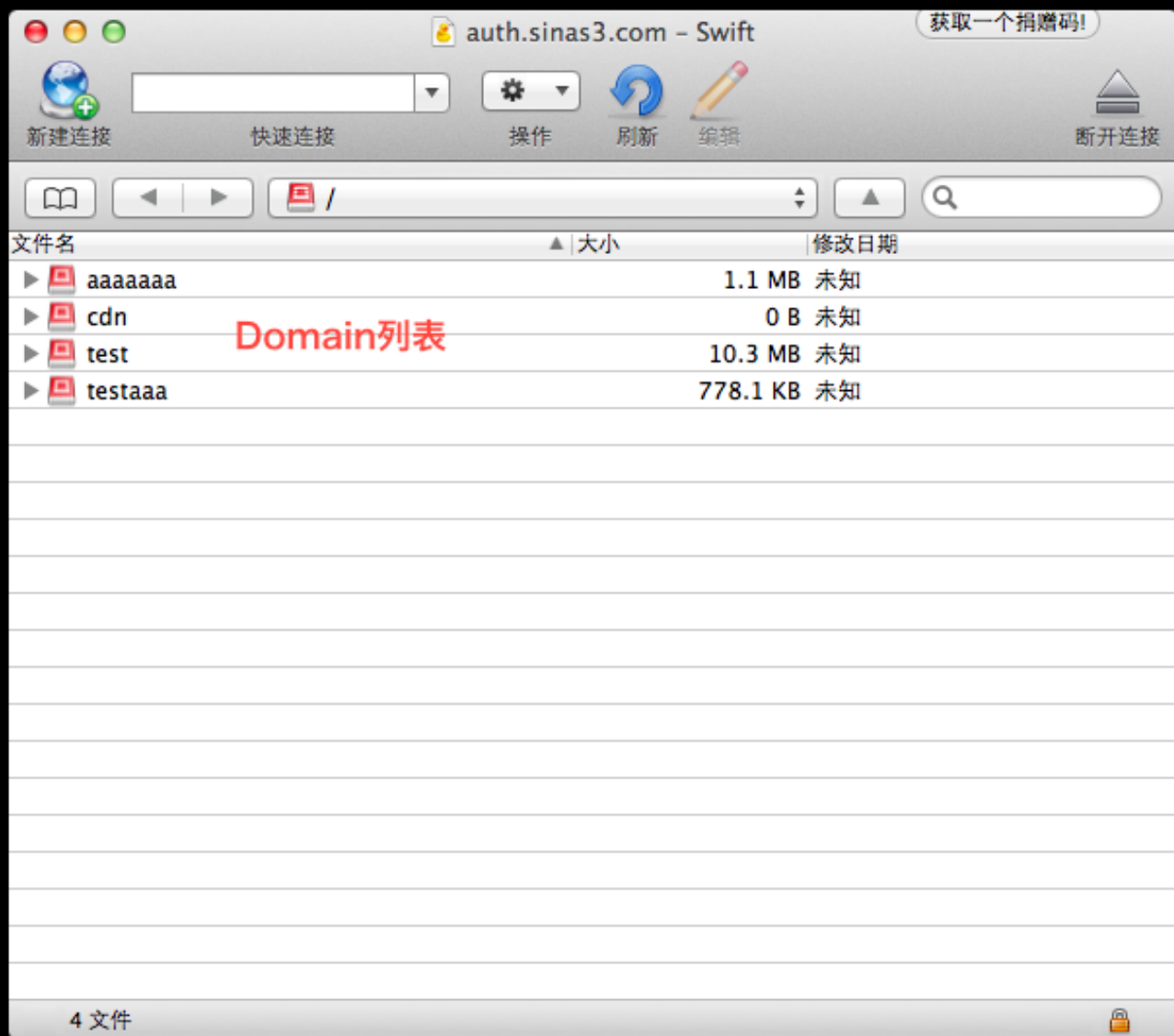
- Container unique module  
container.sinas3.com/object
- Keystone middleware for auth protocol converting

# Our Steps for Switching SAE Storage

---

- Bypass test one month online(旁路读写测试)
- Switch step by step one month(灰度切换)
- Ops
  - Monitoring: I/O, CPU, Memery, Disk
  - LogCenter: syslog-ng, statistics and analytics





# shoplocket

BETA



Sell anything from anywhere, in minutes.

 Connect with Facebook

or

Enter your email address

Choose a password

Try it for free

# Problems&Imporvements

The async processor for keeping eventual consistency is inefficient.

Replicator, auditor, container-updater

**Logic:** loop all objects or dbs on disk, query replica's server to determine whether to sync.

**Results:** High I/O; High CPU usage; Stress on account/container/object-servers, impact the availability; Long time for eventual consistency, impact the durability; The list operation is not consistent.

How to improvements?

1. Runing replicator, auditor and container updater during idle time;
2. An appropriate deployment;
3. A new protocol for keeping replica's consistency;  
(based on log and message queue)
4. Adding new nodes, scale out.

# Problems&Imporvements

---

The performance of sqlite.

Quota for objects and containers

Running sqlite on the high performance I/O devices

The bandwidth of rsync is not under control.

Out-of-band management

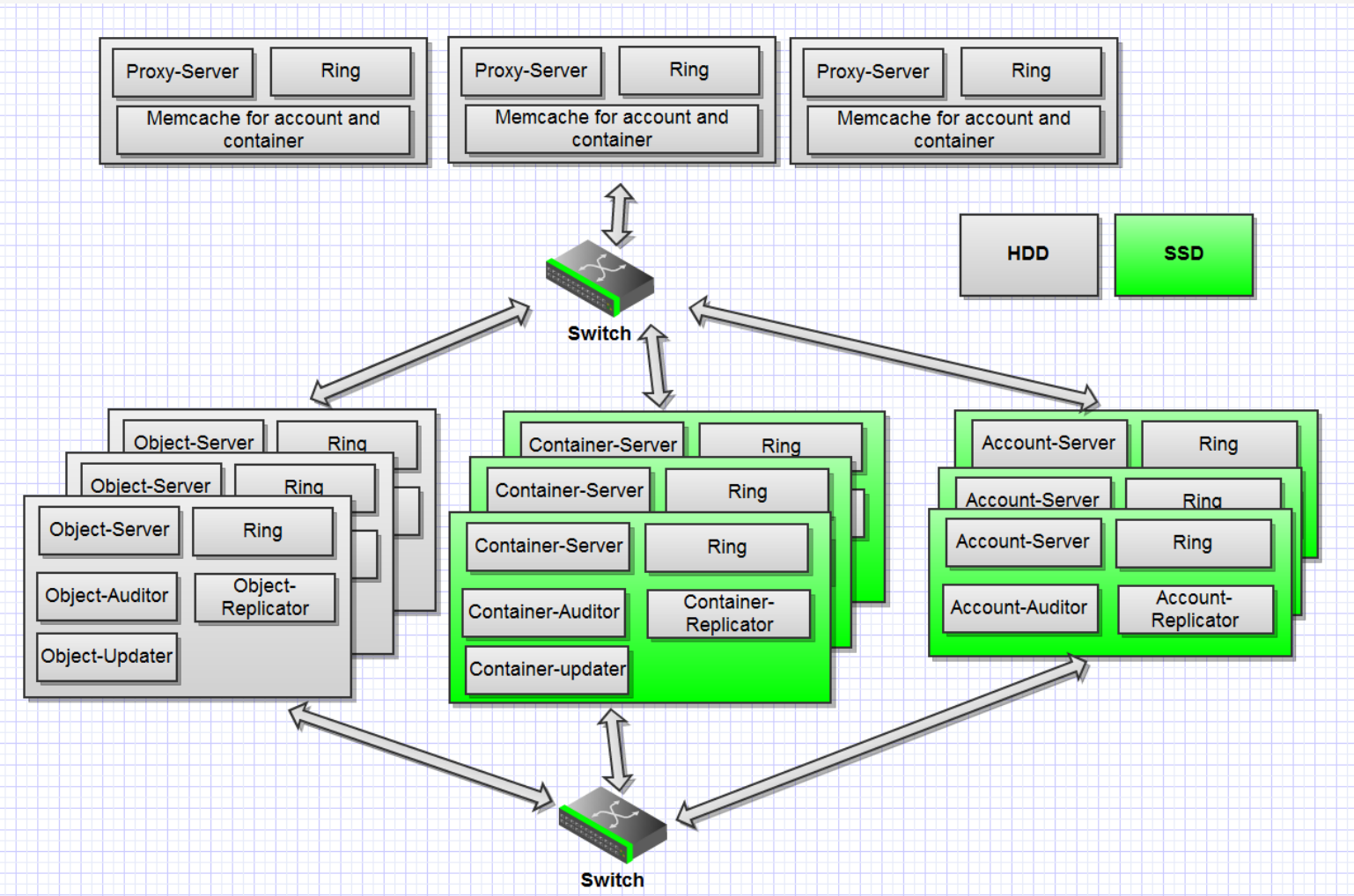
Add bandwidth limitations for rsync

Database centralized or distributed?



# Problems&Imporvements

An appropriate deployment





# Q&A

---



Weibo: @AlexYang\_Coder

Email: yangyu4@staff.sina.com.cn

GTalk: alex890714@gmail.com

Blog: <http://alexyang.sinaapp.com>