

# MOS: Taming the Cloud Object Storage

Ali Anwar\*, Yue Cheng\*, Aayush Gupta<sup>†</sup>, Ali R. Butt\*

\**Virginia Tech &* <sup>†</sup>*IBM Research – Almaden* 



#### Cloud object stores enable cost-efficient data storage



Object storage







Windows<sup>®</sup>Azure<sup>®</sup>



#### Cloud object store supports various workloads



Online video sharing

Enterprise backup



#### One size does not fit all



#### Replace monolithic object store with specialized fine-grained object stores each launched on a sub-cluster



## Reason 1: Classification of workloads





## Small objects





## Large objects





## Reason 2: Heterogeneous resources

- Dcenters hosting object stores are becoming increasingly heterogeneous
- Hardware to application workload mismatch
- Meeting SLA requirement is challenging









#### Outline

## Introduction Motivation Contribution Design Evaluation



## Background: Swift object store





## Swift: Proxy and Storage servers





## Swift: Ring architecture





## Benchmark used: CosBench

- COSBench is Intel developed Benchmark to measure Cloud Object Storage Service performance
  - For S3, OpenStack Swift like object store
  - Not for file system or block device system
- Used to compare different hardware software stacks
- Identify bottlenecks and make optimizations



## Workload used

Workload	<b>Workload Characteristics</b> Object Size Distribution		Application scenario
Workload A	1 – 128 KB	<b>G: 90%</b> , P: 5%, D:5%	Web hosting
Workload B	1 – 128 KB	G: 5%, <b>P: 90%</b> , D:5%	Online game hosting
Workload C	1 – 128 MB	<b>G: 90%</b> , P: 5%, D:5%	Online video sharing
Workload D	1 – 128 MB	G: 5%, <b>P: 90%</b> , D:5%	Enterprise backup



#### Experimental setup for motivational study



## Configuration 1 – Default monolithic



![](_page_15_Picture_2.jpeg)

## Configuration 2 – Favors small objects

![](_page_16_Figure_1.jpeg)

![](_page_16_Picture_2.jpeg)

## Configuration 3 – Favors large objects

![](_page_17_Figure_1.jpeg)

![](_page_17_Picture_2.jpeg)

#### Performance under multi tenant environment – Workload A & B

![](_page_18_Figure_1.jpeg)

![](_page_18_Picture_2.jpeg)

#### Performance under multi tenant environment- – Workload A & B

![](_page_19_Figure_1.jpeg)

![](_page_19_Picture_2.jpeg)

#### Performance under multi tenant environment - latency

![](_page_20_Figure_1.jpeg)

![](_page_20_Picture_2.jpeg)

## Key Insights

- Cloud object store workloads can be classified based on the size of the objects in their workloads
- When multiple tenants run workloads with drastically different behaviors, they compete for the object store resources with each other

![](_page_21_Picture_3.jpeg)

#### Outline

## Introduction Motivation

## Contribution Design Evaluation

![](_page_22_Picture_3.jpeg)

## Contributions

- Perform a performance and resource efficiency analysis on major hardware and software configuration opportunities
- We design MOS, Micro Object Storage:
   1) dynamically provisions fine-grained microstores
   2) exposes the interfaces of microstores to the tenants
- Evaluate MOS to showcase its advantages

![](_page_23_Picture_4.jpeg)

#### Outline

## Introduction

## Motivation

## Contribution

![](_page_24_Picture_4.jpeg)

Evaluation

![](_page_24_Picture_6.jpeg)

## Design criteria for MOS

- We studied the effect of three knobs on performance of a typical object store to come up with design rules/ rules of thumb
  - Proxy Server settings
  - Storage Server settings
  - Hardware changes

![](_page_25_Picture_5.jpeg)

#### Effect of Proxy server settings

![](_page_26_Figure_1.jpeg)

#### Effect of Proxy server settings

![](_page_27_Figure_1.jpeg)

## Effect of Storage server settings

![](_page_28_Figure_1.jpeg)

![](_page_28_Picture_2.jpeg)

## Effect of Storage server settings

![](_page_29_Figure_1.jpeg)

![](_page_29_Picture_2.jpeg)

## Effect of Storage server settings

![](_page_30_Figure_1.jpeg)

![](_page_30_Picture_2.jpeg)

#### Effect of hardware settings

![](_page_31_Figure_1.jpeg)

![](_page_31_Picture_2.jpeg)

## Rules of thumb

CPU on proxy serves as the first-priority resource for small-object intensive workloads

Network bandwidth is more important than CPU on proxy for large-object intensive workloads

proxyCores = storageNodes \* coresPerStorageNode

BWproxies = storageNodes \* BWstorageNode

Faster network cannot effectively improve QPS for small-object intensive workloads – use weak network (1 Gbps NICs) with good storage devices (SSD)

![](_page_32_Picture_6.jpeg)

## MOS Design

![](_page_33_Figure_1.jpeg)

![](_page_33_Picture_2.jpeg)

## Resource Provisioning Algorithm

- Initially, the algorithm allocates the same amount of resources to each microstore conservatively then use greedy approach for resource allocation
- Keep track of free set of resources (including hardware configuration, current load served, and the resource utilization such as CPU and network bandwidth utilization)
- Periodically collect monitoring data from each microstore to aggressively increase and linearly decrease resources from each microstore

![](_page_34_Picture_4.jpeg)

#### Outline

## Introduction

## Motivation

## Contribution

![](_page_35_Picture_4.jpeg)

Evaluation -

![](_page_35_Picture_6.jpeg)

#### Preliminary evaluation via simulation – Experimental setup

#### Compute nodes:

- **3** 32 core machines
- **4** 16 core
- 31 8 core machines
- **12** 4 core machines
- Network:
  - **18** 10 Gbps
  - 32 1 Gbps NICs

□ HDD to SSD ratio was 70% to 30%.

![](_page_36_Picture_10.jpeg)

#### Aggregated throughput

![](_page_37_Figure_1.jpeg)

![](_page_37_Picture_2.jpeg)

#### Aggregated throughput

![](_page_38_Figure_1.jpeg)

![](_page_38_Picture_2.jpeg)

#### Aggregated throughput

![](_page_39_Figure_1.jpeg)

![](_page_39_Picture_2.jpeg)

#### Timeline under dynamically changing workloads

![](_page_40_Figure_1.jpeg)

![](_page_40_Picture_2.jpeg)

#### Resource utilization timeline

![](_page_41_Figure_1.jpeg)

## Related Work

- MET proposes several system metrics that are critical for a NoSQL database and highly impacts server utilization's estimation
- CAST and its extension perform coarse-grained cloud storage (including object stores) management for data analytics workloads
- IOFlow solves a similar problem by providing a queue and control functionality at two OS stages – the storage drivers in the hypervisor and the storage server

![](_page_42_Picture_5.jpeg)

## Conclusion

- We performed exhausted study of cloud object stores
- We proposed a set of rules to help cloud object store administrator to efficiently utilize resources
- We presented MOS which can outperform extant object store under multi-tenant environment
- Our analysis shows that it is possible to exploit heterogeneity inherited by modern datacenter to the advantage of object store providers

![](_page_43_Picture_5.jpeg)

![](_page_44_Picture_0.jpeg)

![](_page_44_Picture_1.jpeg)

![](_page_45_Picture_0.jpeg)

http://research.cs.vt.edu/dssl/

Ali Anwar Yue Cheng Aayush Gupta Ali R. Butt

![](_page_45_Picture_3.jpeg)