

## Fast, Elastic Storage for the Cloud

Ana Klimovic



January 2019

The amount of **data** we are **generating** and **analyzing** is growing **exponentially** 

# The amount of **data** we are **generating** and **analyzing** is growing **exponentially**

We are storing 30% to 40% more data every year → data stored is doubling every 2 to 3 years! The amount of **data** we are **generating** and **analyzing** is growing **exponentially** 

We are not just storing, we are analyzing data.  $\rightarrow$  5 ZB of data is predicted to be analyzed in 2025. Users are increasingly storing and processing data in the **cloud** 





#### Cloud computing offers high...

Elasticity

**Cost-efficiency** 





#### Performance

...only if each application receives the storage and compute resources it needs.

Allocating resources to achieve these goals is hard.

#### Why is resource allocation difficult?

Consider the resource utilization of a large-scale service at Facebook, normalized over 6 months: [EuroSys'16]



Compute and storage requirements vary dynamically *over time*.

### Why is resource allocation difficult?

Consider the resource utilization of a large-scale service at Facebook, normalized over 6 months: [EuroSys'16]



At datacenter scale...

#### ~10,000s of servers ~10s of Megawatts

# ...underutilizing resources is extremely wasteful.

Huge opportunity to improve resource utilization! <u>Challenge</u>: maintain high performance How should we build **resource-efficient** and **high performance** storage systems for diverse, elastic applications in the cloud?

## Key insights

#### Decouple storage from compute resources

- Improve resource allocation flexibility by enabling fast, predictable access to remote data
  - Efficiently integrate network & storage software processing
  - $_{\circ}$   $\,$  Introduce token-based I/O scheduling for quality of service



## Automate storage resource management

• Allocate resources dynamically based on application requirements, learned via hints or a machine learning model

#### Decouple storage from compute



Provide **flexibility** to allocate as many CPU cores an application needs, independent of its storage capacity and storage throughput requirements.



#### Decouple storage from compute



Provide **flexibility** to allocate as many CPU cores an application needs, independent of its storage capacity and storage throughput requirements.



#### Decouple storage from compute



Provide **flexibility** to allocate as many CPU cores an application needs, independent of its storage capacity and storage throughput requirements.



<u>Goal</u>: enable any CPU to use any storage device in a cloud facility that has spare capacity & bandwidth.

### Share storage to improve utilization



Share storage devices among applications to increase utilization and improve **cost-efficiency**.



<u>Goal</u>: enable any CPU to use any storage device in a cloud facility that has spare capacity & bandwidth.



To improve resource allocation flexibility while maintaining high performance, we need **fast, predictable** access to remote data.





#### Local access to storage











To improve resource allocation flexibility while maintaining high performance, we need **fast, predictable** access to remote data.



How do we decide which resources to allocate to each application and where to place data?

#### Automate storage resource management



- Automatically decide how many storage devices to allocate, where to place data, and on what type of storage technology
- Dynamically adjust allocations as application load varies



# **Requirement #2:** Learn application resource requirements to automate resource allocation

Learn the application's characteristics (e.g., capacity, throughput, latency requirements) and translate to a performance-cost efficient resource allocation.

- Leverage high-level hints about application characteristics
- $\rightarrow$  use heuristic policies to decide resource allocation

- Leverage performance data collected across application runs
- $\rightarrow$  build a machine learning model that recommends a near-optimal resource allocation

## Two major requirements for cloud storage

1. Fast access to remote data



2. Automatic, elastic allocation of storage resources



## Contributions

1. Fast access to remote data



[OSDI'14] IX: Dataplane OS for Fast Networking
[EuroSys'16] Flash Storage Disaggregation
[HotStorage'17] Rack-Scale Disaggregated Storage
[ASPLOS'17] ReFlex: Remote Flash ≈ Local Flash

2. Automatic, elastic allocation of storage resources

[ATC'18a] Understanding Ephemeral Storage for Serverless Analytics

[OSDI'18] Pocket: Elastic Ephemeral Storage for Serverless Analytics

[ATC'18b] Selecta: ML-based Heterogeneous Cloud Storage Configuration

## Focus of this talk

1. Fast access to remote data



[OSDI'14] IX: Dataplane OS for Fast Networking

[EuroSys'16] Flash Storage Disaggregation

[HotStorage'17] Rack-Scale Disaggregated Storage

[ASPLOS'17] ReFlex: Remote Flash ≈ Local Flash

2. Automatic, elastic allocation of storage resources

[ATC'18a] Understanding Ephemeral Storage for Serverless Analytics

[OSDI'18] Pocket: Elastic Ephemeral Storage for Serverless Analytics

[ATC'18b] Selecta: ML-based Heterogeneous Cloud Storage Configuration

#### Fast access to remote flash

## Flash storage



NVMe Flash:

- 70 µs read latency → 100x faster than disk
- $\circ$  1,000,000 IOPS  $\rightarrow$  1000x higher than disk

# What if we use existing remote storage approaches to access remote flash?

#### Issue #1: Existing software is too slow

Existing remote storage solutions have significant overheads.



#### Issue #2: Unpredictable performance on shared flash

Write requests from one tenant can **interfere** with read requests from another tenant, leading to unpredictable performance on shared flash



### **ReFlex: Remote Flash ≈ Local Flash**

A software system for fast, predictable access to remote flash storage

## **ReFlex: Remote Flash ≈ Local Flash**

A software system for fast, predictable access to remote flash storage

**Issue #1:** Existing software is too slow

→ Custom network-storage OS, designed for low latency and high throughput

- $\checkmark\,$  Run to completion
- ✓ Adaptive batching
- $\checkmark\,$  Reduce data copies
- $\checkmark\,$  Direct access to hardware

## How does ReFlex achieve high performance?



## How does ReFlex achieve high performance?



## How does ReFlex achieve high performance?


## How does ReFlex achieve high performance?



## How does ReFlex achieve high performance?



#### **ReFlex performance: throughput per core**



#### **ReFlex performance: throughput per core**



#### **ReFlex performance: latency**



#### **ReFlex: Remote Flash ≈ Local Flash**

A software system for fast, predictable access to remote flash storage

Issue #1: Existing software is too slow
→ Custom network-storage OS, designed for low-latency and high throughput

Issue #2: Unpredictable performance on shared flash storage
→ Novel I/O scheduler provides tail latency and throughput guarantees

ReFlex schedules requests to satisfy the performance objectives specified by each tenant.



**Goal**: Provide **tail latency** and **throughput** guarantees to tenants sharing flash.

**Challenge**: Account for the performance properties of different request types (e.g., read vs. write, 4 KB vs. 1 MB request)

**Solution**: profile device to derive a request cost model that captures how each type of request impacts tail latency and throughput

Step 1: Build a request cost model

 $\rightarrow$  account for different request types

For this device, write cost = 10 x read cost



- Step 1: Build a request cost model → account for different request types
- Step 2: Schedule requests

Example scenario

#### Tenant A:

- 1ms tail latency
- 200K IOPS

#### Tenant B:

• Best-effort (use slack)





- Step 1: Build a request cost model → account for different request types
- Step 2: Schedule requests

Example scenario

#### Tenant A:

- 1ms tail latency
- 200K IOPS

#### Tenant B:

• Best-effort (use slack)





#### **ReFlex performance guarantees on shared Flash**



Tenants A & B: latency-critical; Tenant C + D: best effort

#### **ReFlex performance guarantees on shared Flash**



Tenants A & B: latency-critical; Tenant C + D: best effort

#### **ReFlex performance guarantees on shared Flash**



Tenants A & B: latency-critical; Tenant C + D: best effort

## **ReFlex impact**

- ReFlex provides high throughput, low latency with the flexibility to use commodity networks
- Broadcom is porting ReFlex to a SoC platform
- ReFlex integrated into Apache Crail storage system  $\rightarrow$  runs on  $rac{1}{2}$  mazon





# Contributions



[OSDI'14] IX: dataplane OS for fast networking
[EuroSys'16] Flash storage disaggregation
[HotStorage'17] Rack-Scale Disaggregated Storage
[ASPLOS'17] ReFlex: Remote Flash ≈ Local Flash

2. Automatic, elastic allocation of storage resources

[ATC'18a] Understanding Ephemeral Storage for Serverless Analytics

#### [OSDI'18] Pocket: Elastic Ephemeral Storage for Serverless Analytics

[ATC'18b] Selecta: ML-based Heterogeneous Cloud Storage Configuration

# Automatic, elastic allocation of storage resources

# **Serverless computing**

Serverless computing is a new cloud service → users can launch thousands of tiny tasks with **high elasticity** and **fine-grain billing** 



Users focus on writing code for their applications → no resource management



Cloud providers automatically allocate and scale resources.

# **Serverless analytics**

Serverless computing is increasingly used for interactive analytics
→ exploit massive task parallelism to get a result in near real-time



# Data sharing in serverless analytics

• Analytics jobs involve multiple stages of execution

**Challenge:** exchange intermediate data efficiently between tasks

ephemeral data

# Data sharing in serverless analytics

- Direct communication between serverless tasks is difficult:
  - Tasks are short-lived and stateless



# Data sharing in serverless analytics

The natural approach for sharing ephemeral data is through a shared remote data store



# **Requirements for ephemeral storage**

- 1. High performance for a wide range of object sizes
- 2. Cost efficiency

## **Requirements for ephemeral storage**

- 1. High performance for a wide range of object sizes
- 2. Cost efficiency
  - Example of performance-cost tradeoff for a serverless video analytics job with different ephemeral data store configurations



# **Requirements for ephemeral storage**

- 1. High performance for a wide range of object sizes
- 2. Cost efficiency
- 3. Fault-tolerance

Today's cloud storage systems are not optimized for the requirements of serverless analytics jobs.

e.g., - S3 has low cost, but is optimized for large objects

- Redis offers high performance, but uses DRAM  $\rightarrow$  expensive

61

- A fast and elastic storage service for ephemeral data sharing in serverless analytics
- Pocket achieves high performance and cost efficiency by:
  - 1. Leveraging multiple storage technologies
  - 2. Rightsizing resource allocations for applications
  - 3. Autoscaling storage resources in the cluster based on usage
- Similar performance to Redis, an in-memory data store, while saving ~60% in cost for various serverless analytics jobs



# Pocket design

Job A Job C Job B λλλλλλ λλλλλλλλλ λλλλ λλλλλλ λλλλλλλλλλ λλλλ Controller Metadata server(s) app-driven resource request routing allocation & scaling

Storage server

Storage server

Storage server

Storage server

#### 1. Leverage multiple storage technologies



#### 2. Allocate resources based on job requirements



#### 3. Autoscale resources based on utilization



Job C λλλλλλλλλλ λλλλλλλλλλ













# **Evaluating Pocket**

- Pocket is intended to be managed by cloud providers
- Evaluate Pocket running on Amazon EC2 virtual machines
- Run serverless applications on AWS Lambda:
  - Video analytics (object recognition)
  - MapReduce sort
  - Distributed software compilation



# **Application performance with Pocket**

• Compare Pocket to S3 and Redis, which are commonly used today



# **Application performance with Pocket**

• Compare Pocket to S3 and Redis, which are commonly used today



# **Application storage cost with Pocket**

 Pocket leverages job attribute hints for cost-effective resource allocation and amortizes VM costs across multiple jobs, offering a pay-what-you-use model



## **Autoscaling the Pocket cluster**



Job hints	Job1: Sort
Latency sensitive	False
Ephemeral data capacity	10 GB
Aggregate throughput	3 GB/s

## **Autoscaling the Pocket cluster**



Job hints	Job1: Sort	Job2: Video analytics	Job3: Sort
Latency sensitive	False	False	False
Ephemeral data capacity	10 GB	6 GB	10 GB
Aggregate throughput	3 GB/s	2.5 GB/s	3 GB/s

# Summary

To build resource-efficient, high performance storage systems, we need to:



#### Decouple storage from compute resources

• Improve resource allocation flexibility with fast, predictable access to remote data

[OSDI'14] IX: Dataplane OS for Fast Networking
[EuroSys'16] Flash Storage Disaggregation
[HotStorage'17] Rack-Scale Disaggregated Storage
[ASPLOS'17] ReFlex: Remote Flash ≈ Local Flash



#### Automate storage resource management

• Allocate resources dynamically based on app requirements, learned via hints or ML

[ATC'18a] Understanding Ephemeral Storage for Serverless Analytics

[OSDI'18] Pocket: Elastic Ephemeral Storage for Serverless Analytics

[ATC'18b] Selecta: ML-based Heterogeneous Cloud Storage Configuration
### What's next?

# Machine learning for cloud systems



- Use ML to enhance or replace **heuristics** and user **hints** 
  - E.g., Automate cluster resource allocation

### **ML-based resource allocation**

 Selecta [ATC'18] predicts near-optimal configuration for a job using sparse training data across jobs





### **ML-based resource allocation**



- Selecta [ATC'18] predicts near-optimal configuration for a job using sparse training data across jobs
  - 94% probability of recommending near-optimal performing configuration
  - 80% probability of recommending near-optimal cost configuration



# Machine learning for cloud systems



- Use ML to enhance or replace heuristics and user hints
  - Automate cluster resource allocation
  - Predict resource utilization for opportunistic computing
  - Debug performance and security by detecting anomalies

## Systems for machine learning

#### ML-based applications













#### Hardware









# Systems for machine learning











- Near-storage computing for massive datasets
- Distribute computing between the cloud and edge
- Privacy in the cloud with secure enclaves









## Acknowledgements

Christos Kozyrakis Patrick Stuedi Heiner Litz Yawen Wang Adam Belay Animesh Trivedi Jonas Pfefferle Binu John Eno Thereska George Prekas **Edouard Bugnion** 



## Conclusion

As we continue exponentially generating and analyzing data in the cloud, we need:

1. Fast access to remote data



[OSDI'14] IX: Dataplane OS for Fast Networking
[EuroSys'16] Flash Storage Disaggregation
[HotStorage'17] Rack-Scale Disaggregated Storage
[ASPLOS'17] ReFlex: Remote Flash ≈ Local Flash

2. Automatic, elastic allocation of storage resources



[ATC'18a] Understanding Ephemeral Storage for Serverless Analytics

[OSDI'18] Pocket: Elastic Ephemeral Storage for Serverless Analytics

[ATC'18b] Selecta: ML-based Heterogeneous Cloud Storage Configuration