



# AWS Lambda

## (specifically container loading and snapstart)

Marc Brooker

VP/Distinguished Engineer

[mbrooker@amazon.com](mailto:mbrooker@amazon.com)

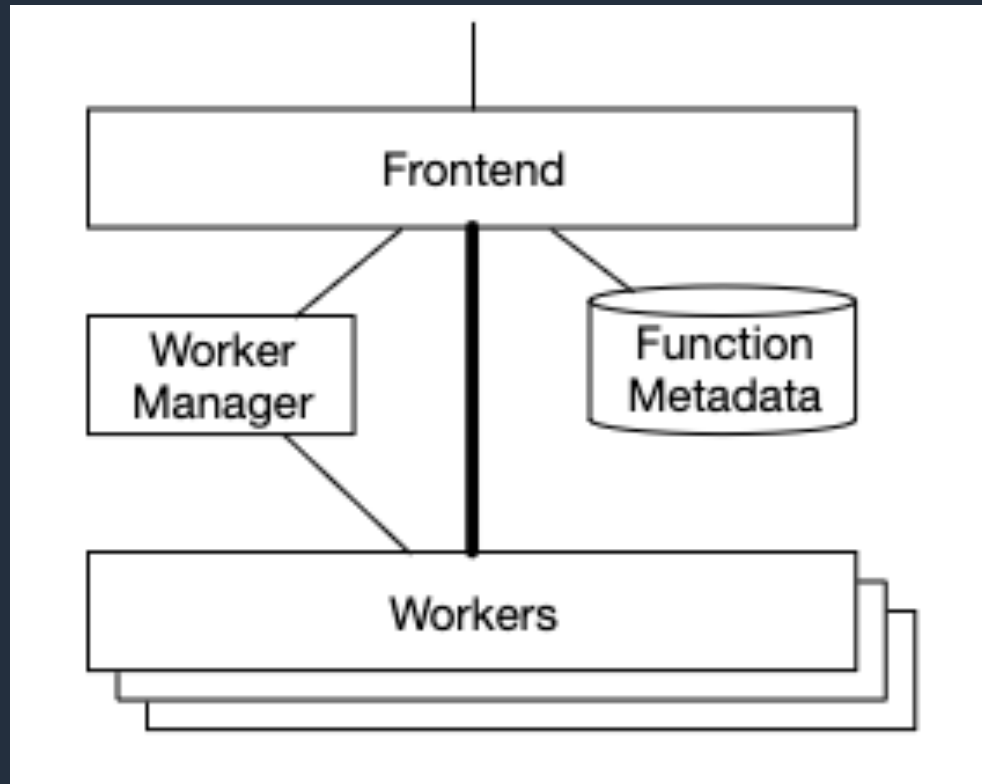
“AWS Lambda is a compute service that runs your code in response to events and automatically manages the compute resources.”

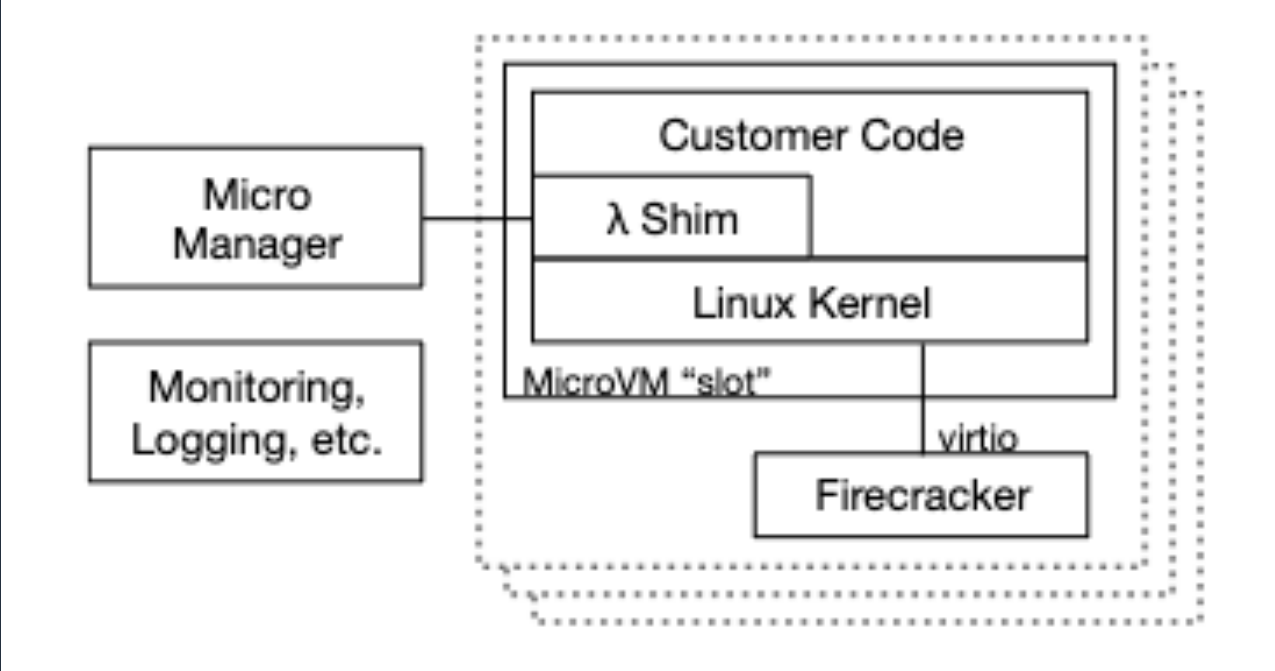
```
import json

def lambda_handler(event, context):
    return {
        'statusCode': 200,
        'body': json.dumps('Hello 6.5840')
    }
```

# Why?

- Building scalable, fault-tolerant systems is hard.
- Driving high hardware utilization is hard.
  - Multi-tenancy makes it easier.
- The cloud needs *glue*.
- *Code close to data is more efficient?*





```
import json

def lambda_handler(event, context):
    return {
        'statusCode': 200,
        'body': json.dumps('Hello 6.5840')
    }
```

“... simple ...”

“... easy ...”

“... fast ...”

“... complex ...”

“... complicated ...”

“... difficult ...”

## Before

250MB max

Code or .zip

Custom tools

## After

10GB max

Container image

Standard tools

## Before

250MB max

Code or .zip

Custom tools

## After

10GB max

Container image

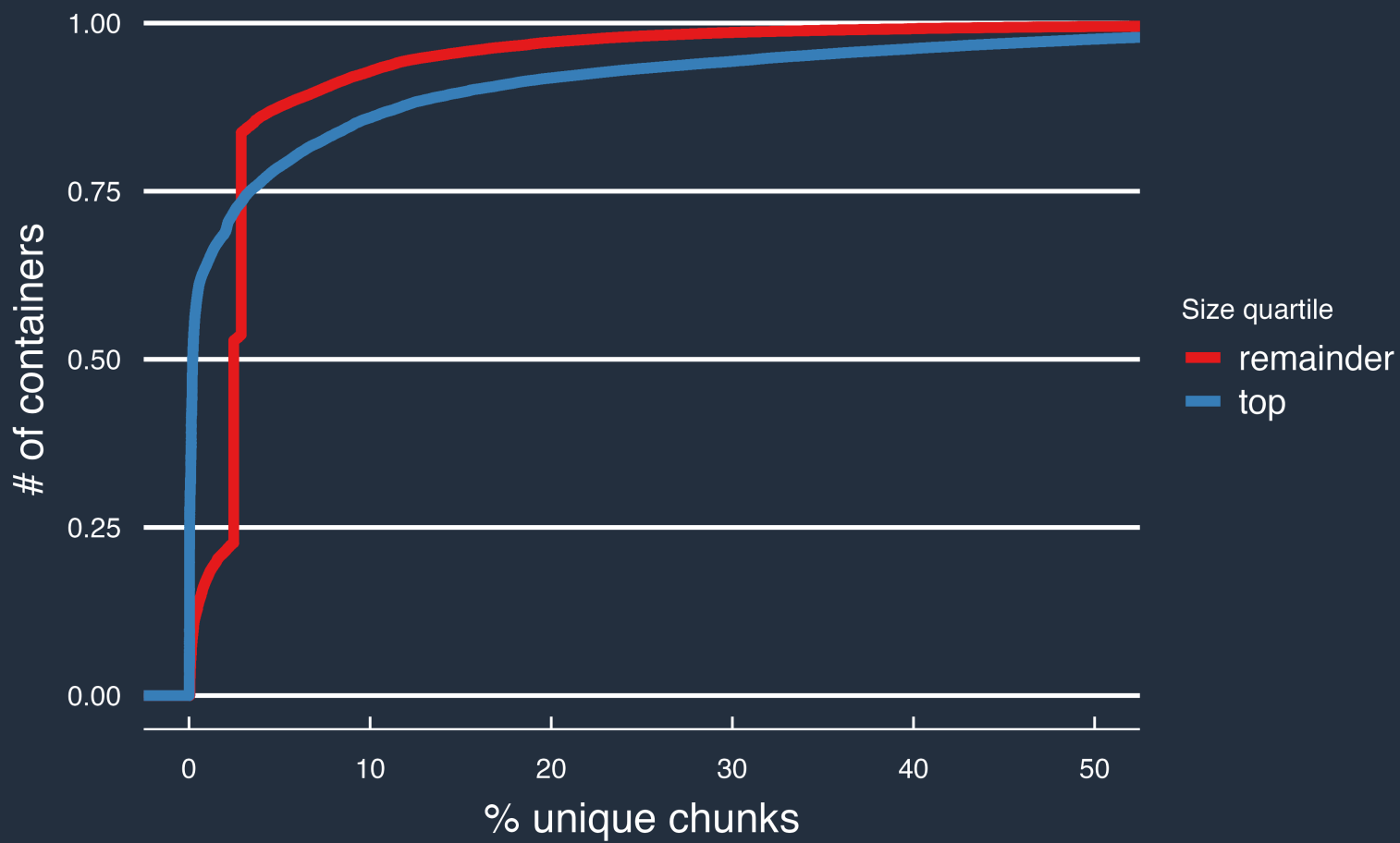
Standard tools

**without slower cold starts!**

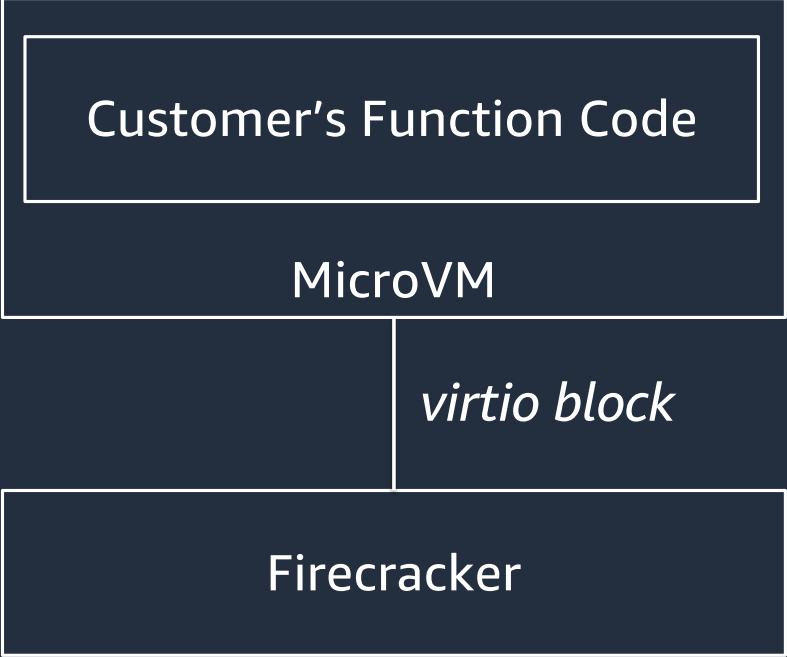
150 Pb/s

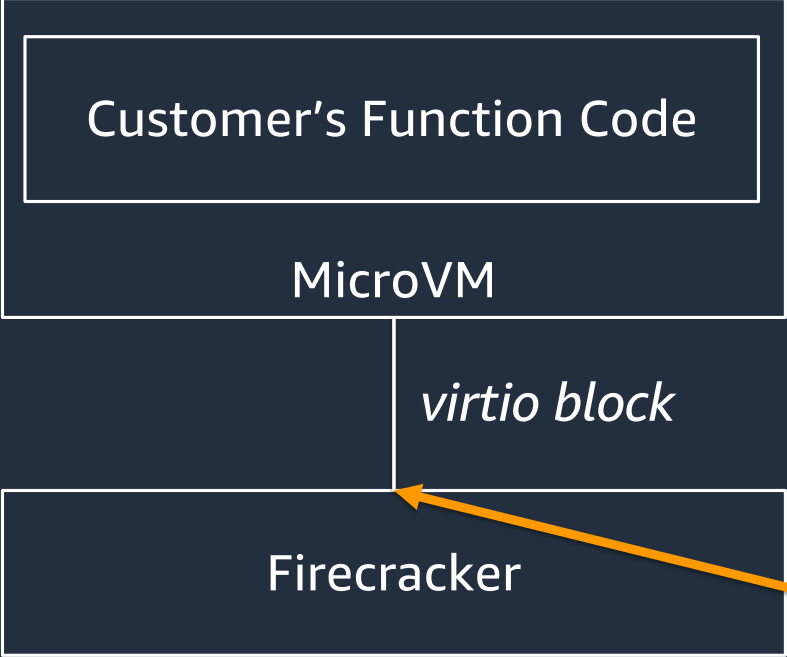
# Harter et al: Only 6.4% of container data is needed at startup!

Harter et al, *Slacker: Fast Distribution with Lazy Docker Containers*, FAST'16



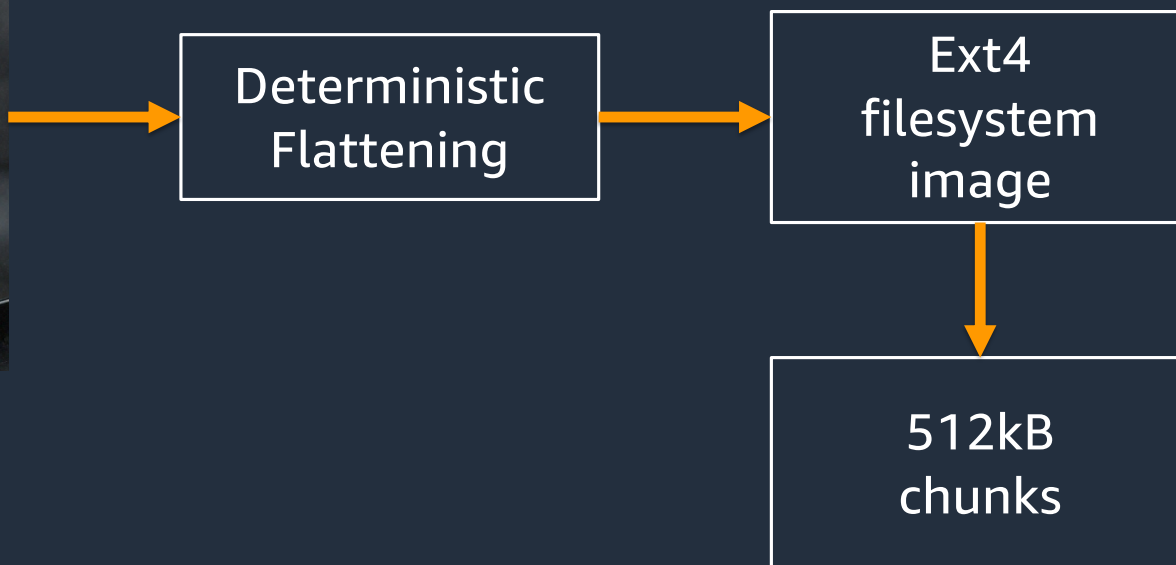
Find a way, invisible to applications, to load data on demand and deduplicate common data.





Insert lies here

But...  
Container images aren't filesystems.



But...

Real storage is really fast ( $<100\mu\text{s}$ ), so loading on demand must be really fast too.

Lambda worker

<1ms

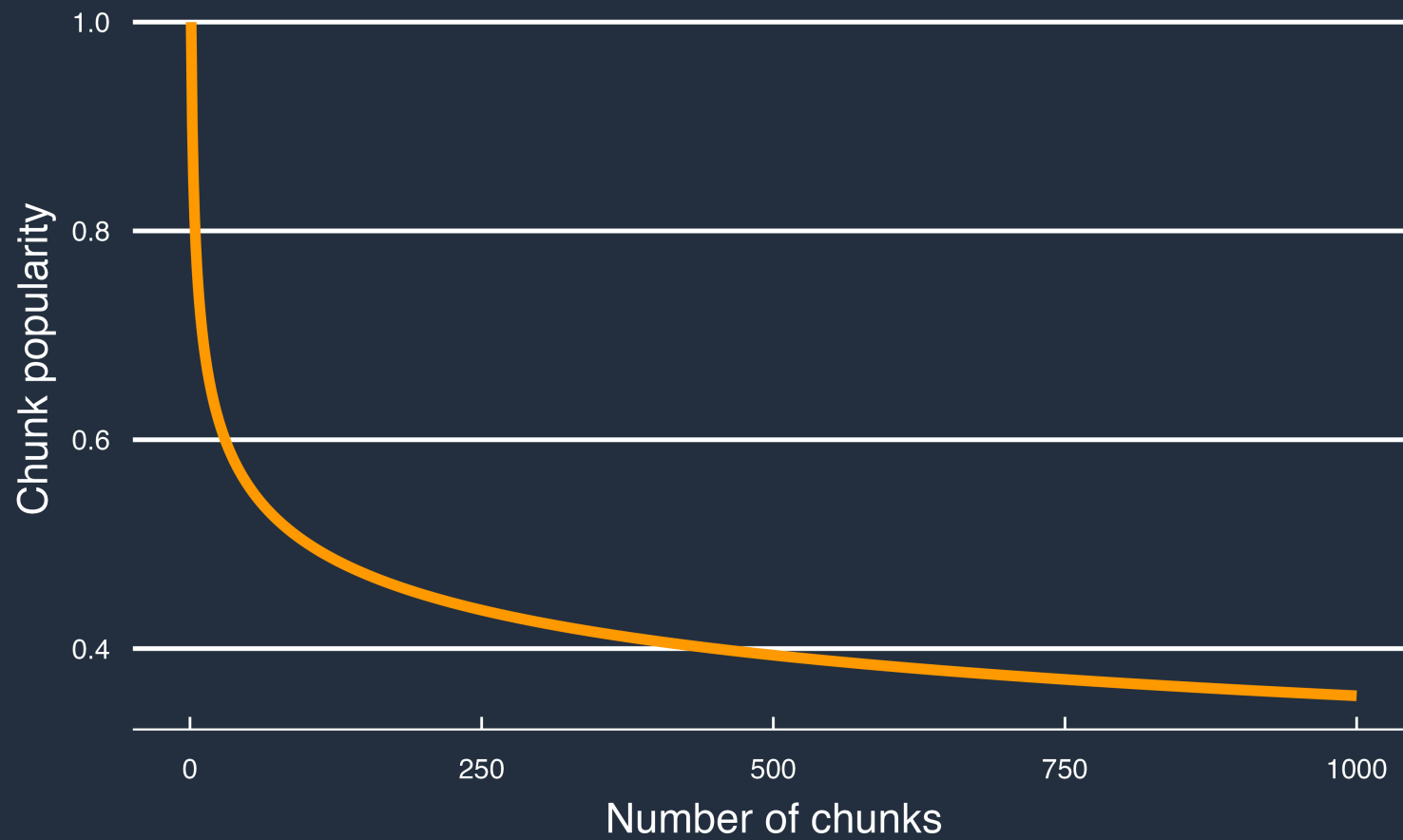
Really fast storage with the  
whole dataset.

Lambda worker



**Amazon S3 Express One Zone**

Delivers consistent  
single-digit millisecond  
request latency on  
hundreds of thousands of  
transactions per second for  
faster data processing

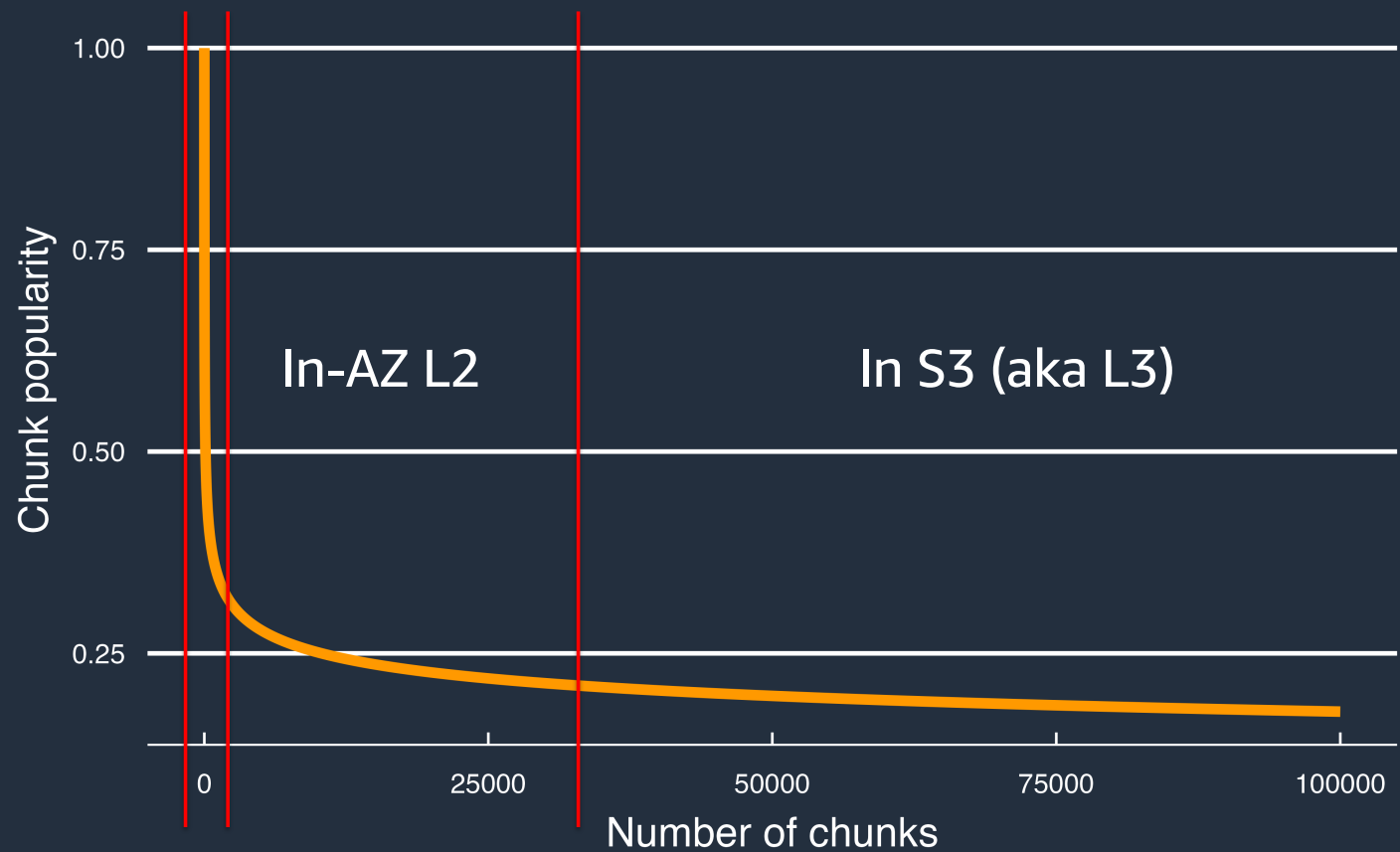


This is fake data, aimed to give the general idea, because I can't share the real numbers.

© 2020, Amazon Web Services, Inc. or its Affiliates. All rights reserved.

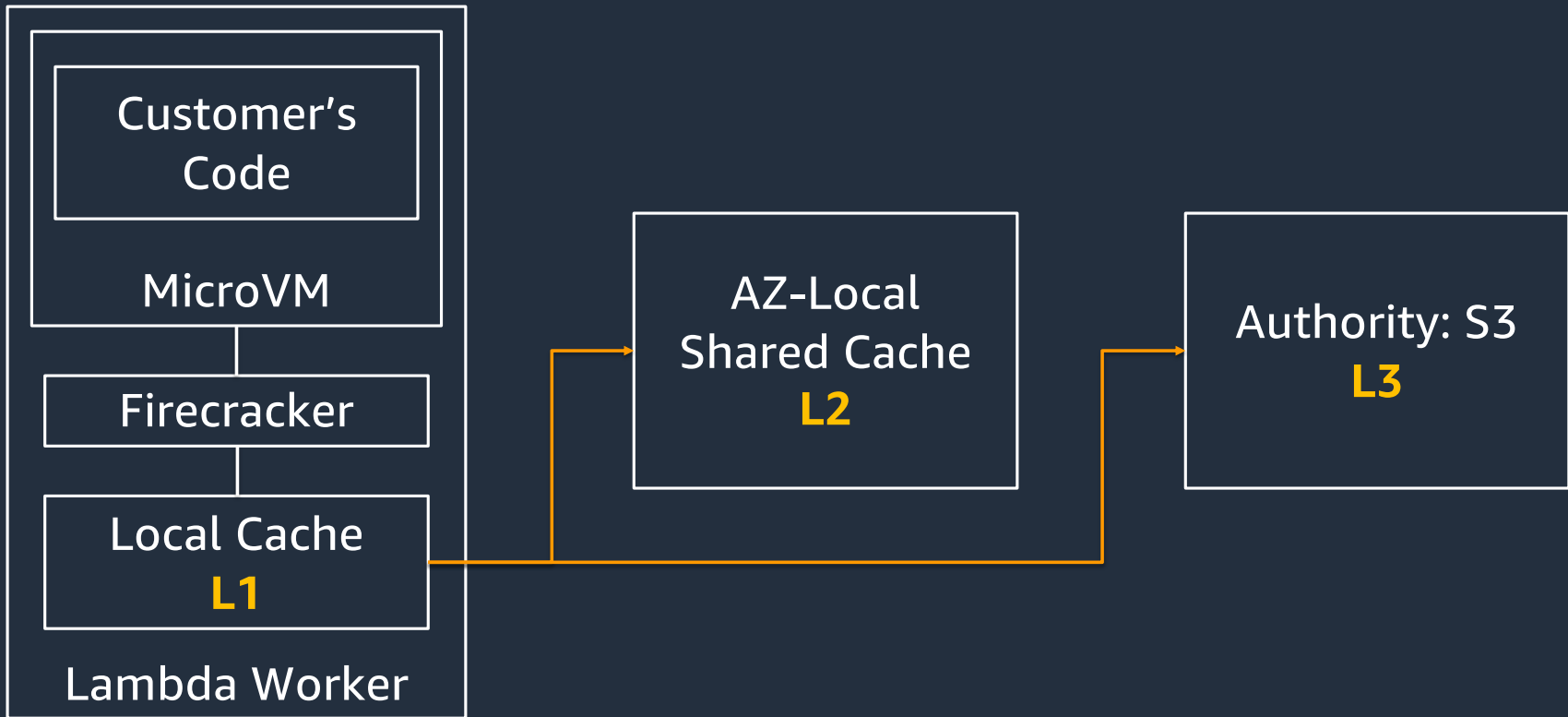


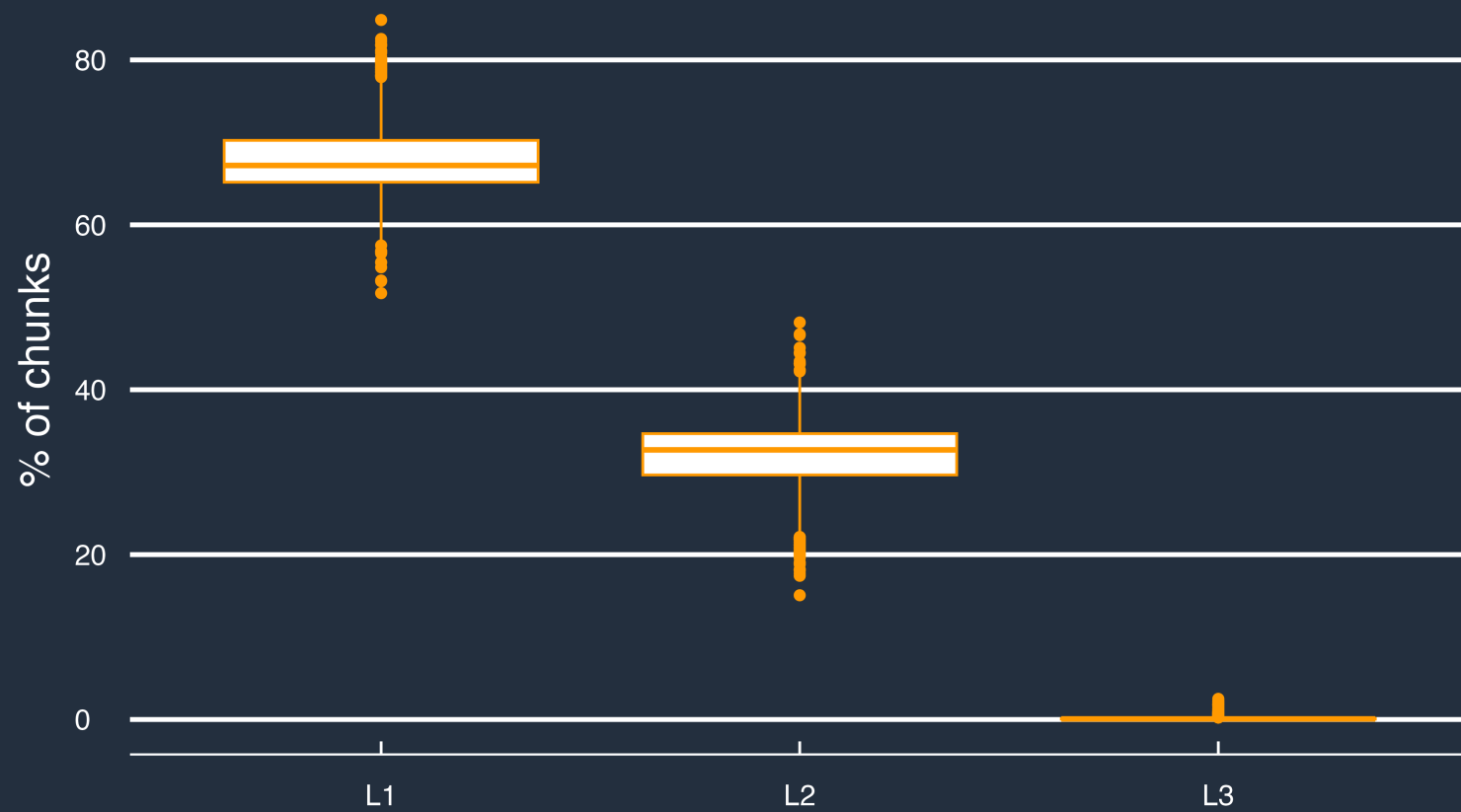
# On-worker L1

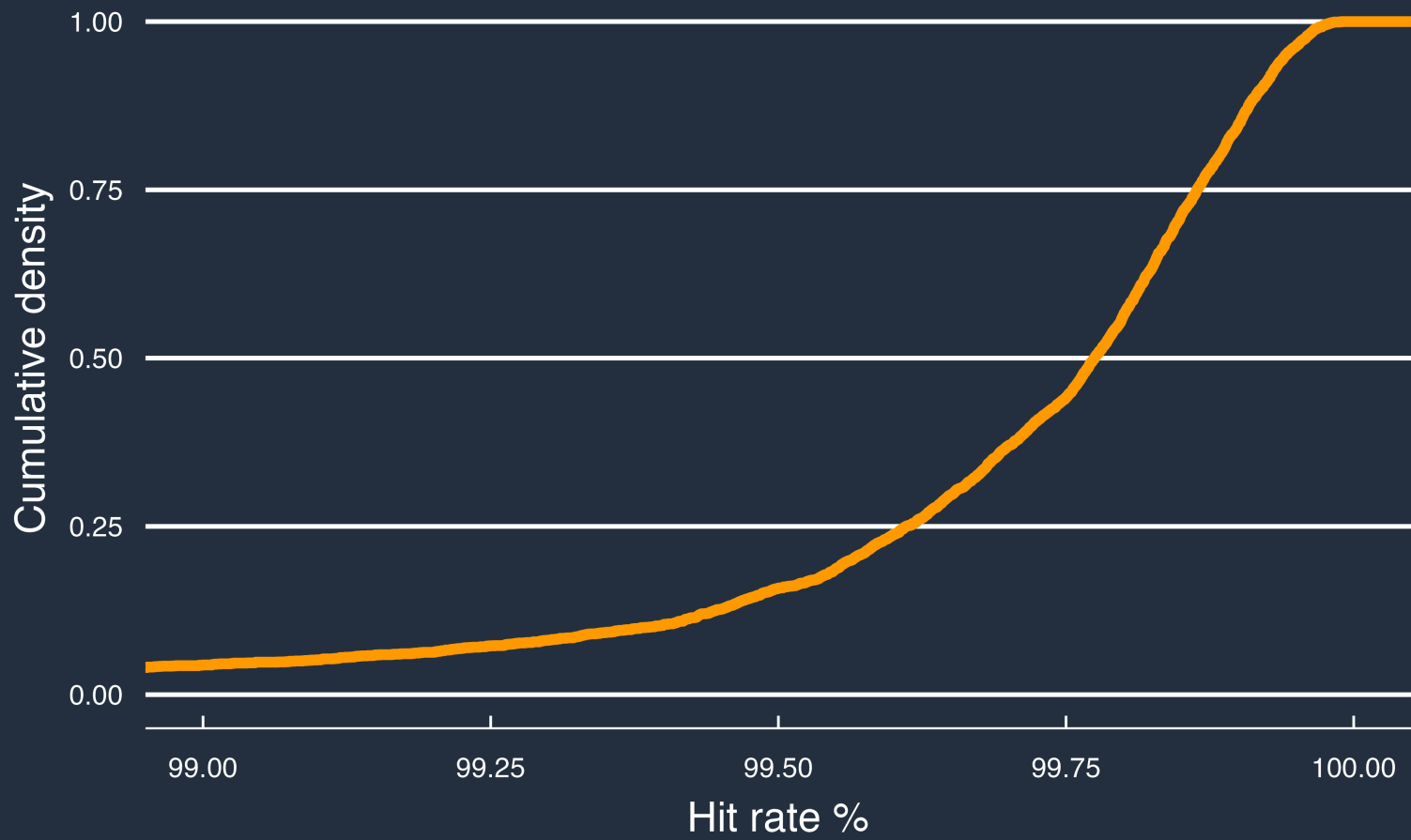


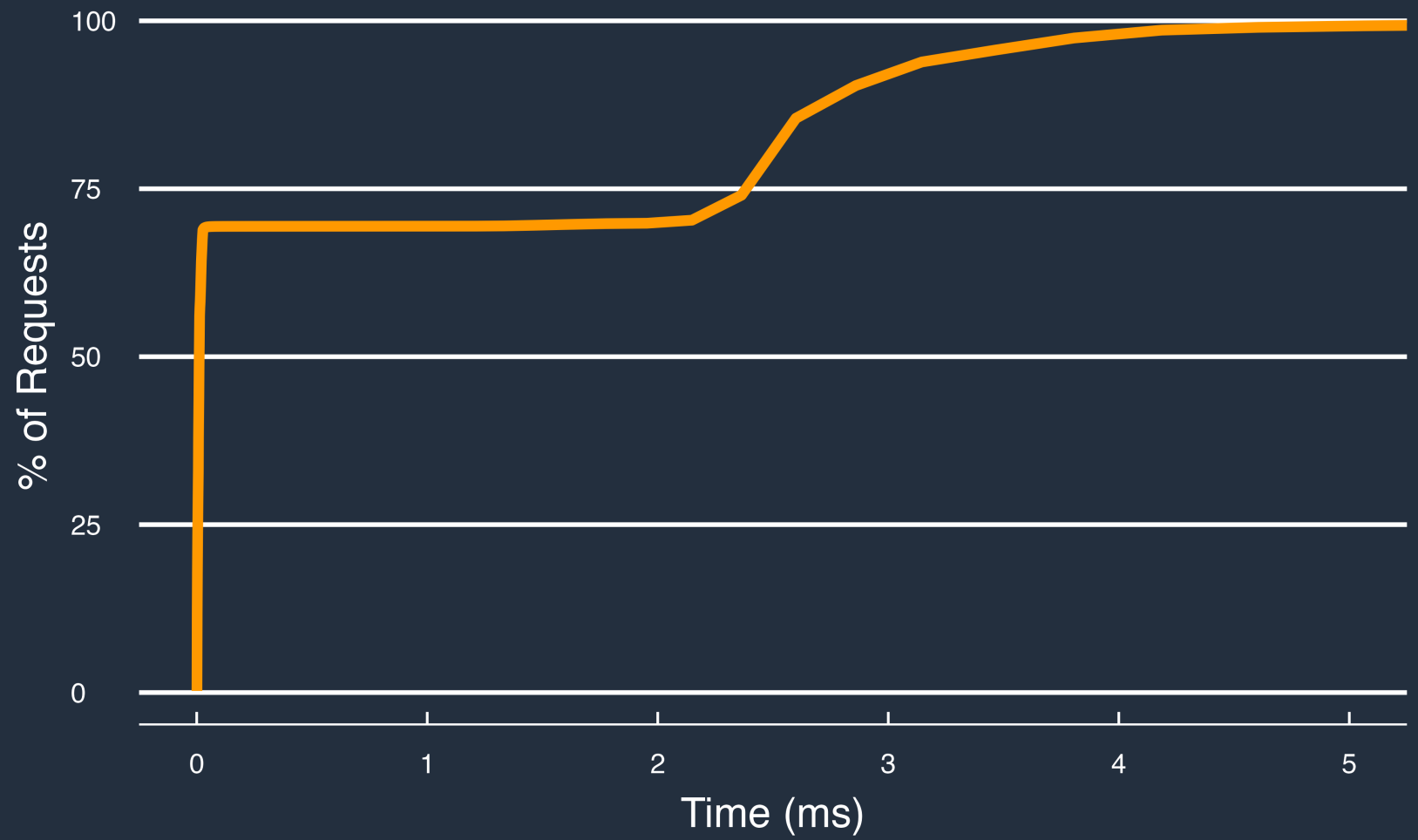
This is fake data, aimed to give the general idea, because I can't share the real numbers.  
© 2020, Amazon Web Services, Inc. or its Affiliates. All rights reserved.



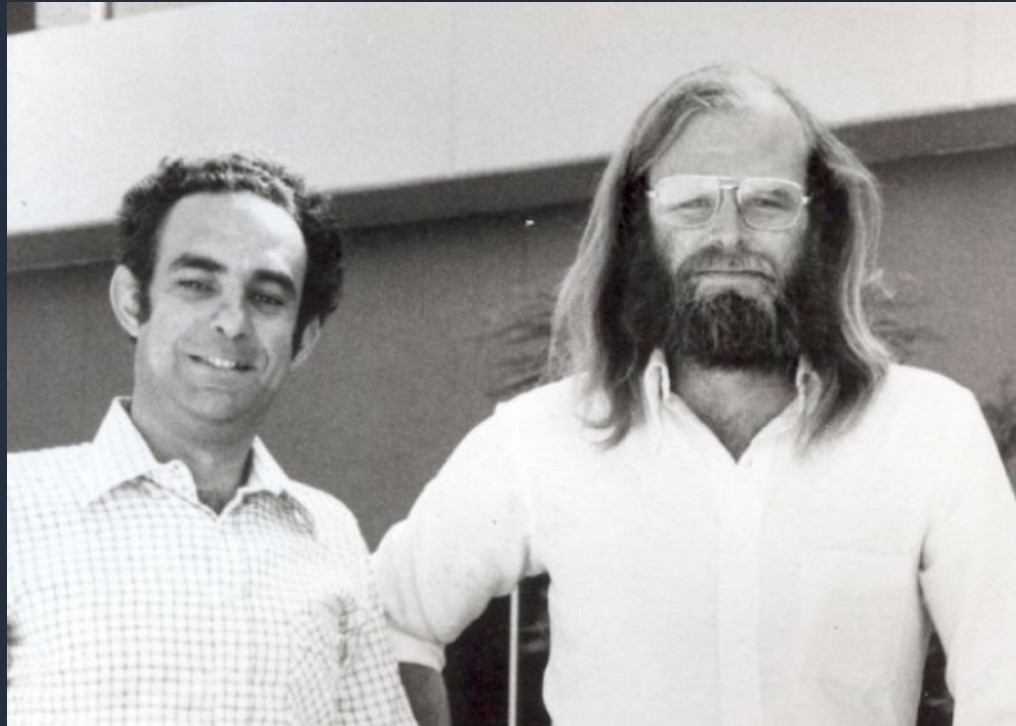








# How Big?



cost of caching = cost of reloading

Jim Gray and Franco Putzolu, IBM Research, 1977 (picture thanks to Pat Helland)  
© 2020, Amazon Web Services, Inc. or its Affiliates. All rights reserved.



Max(size needed for cost,  
size needed for latency)

# What about failures?

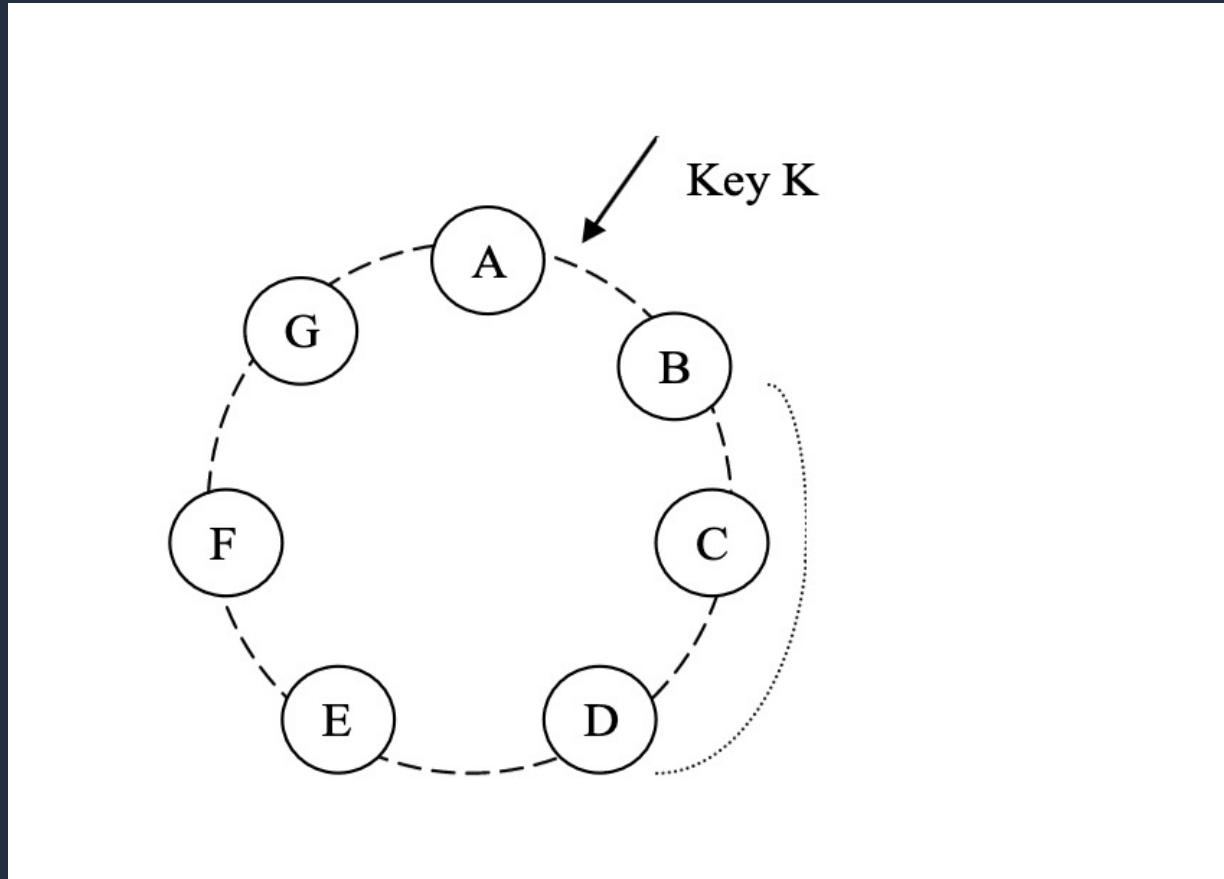
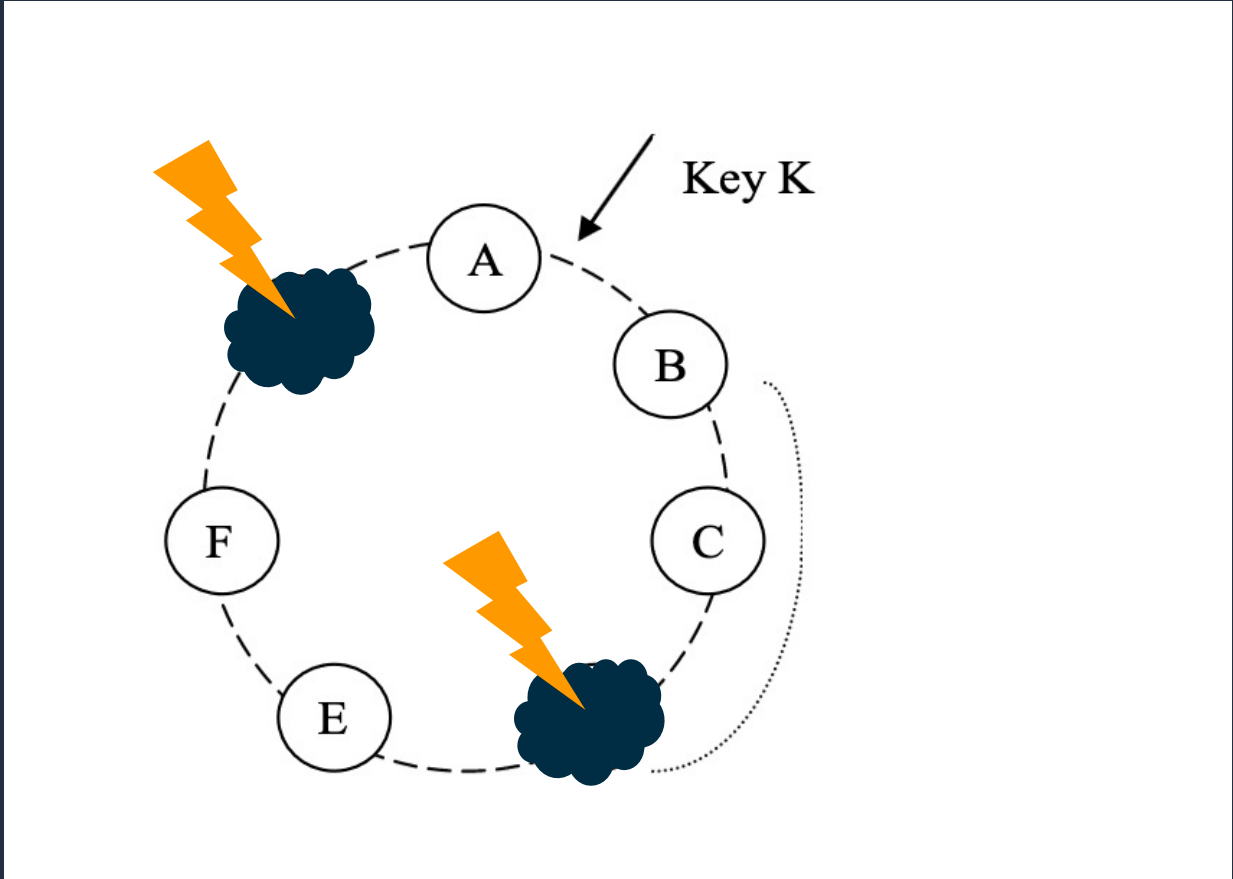
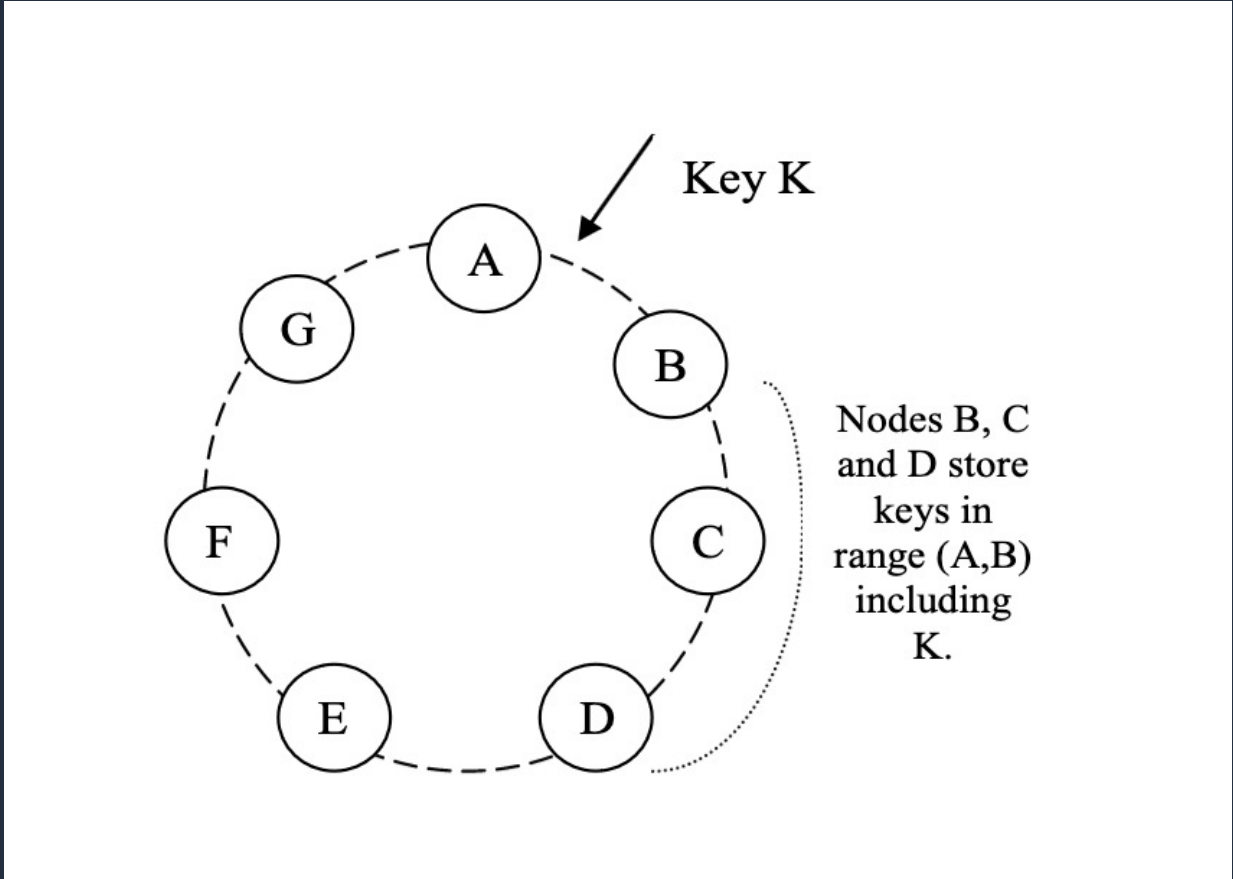
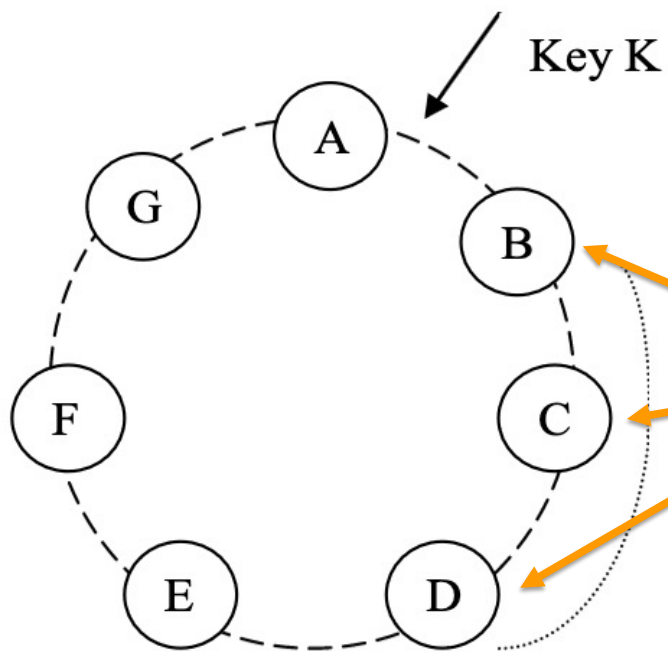


Figure shamelessly borrowed from Decandia et al, *Dynamo: Amazon's Highly Available Key-value Store*, SOSP'07

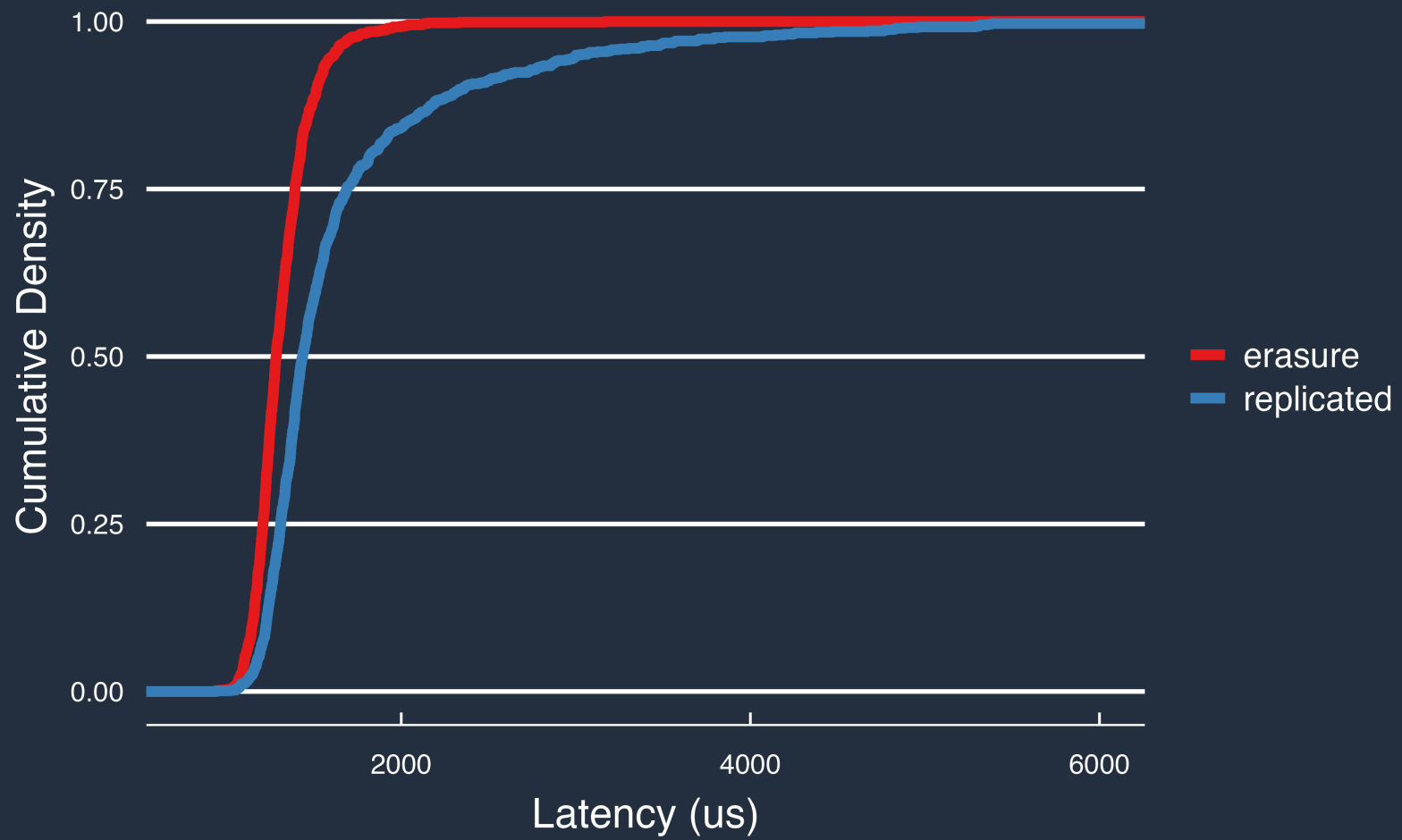


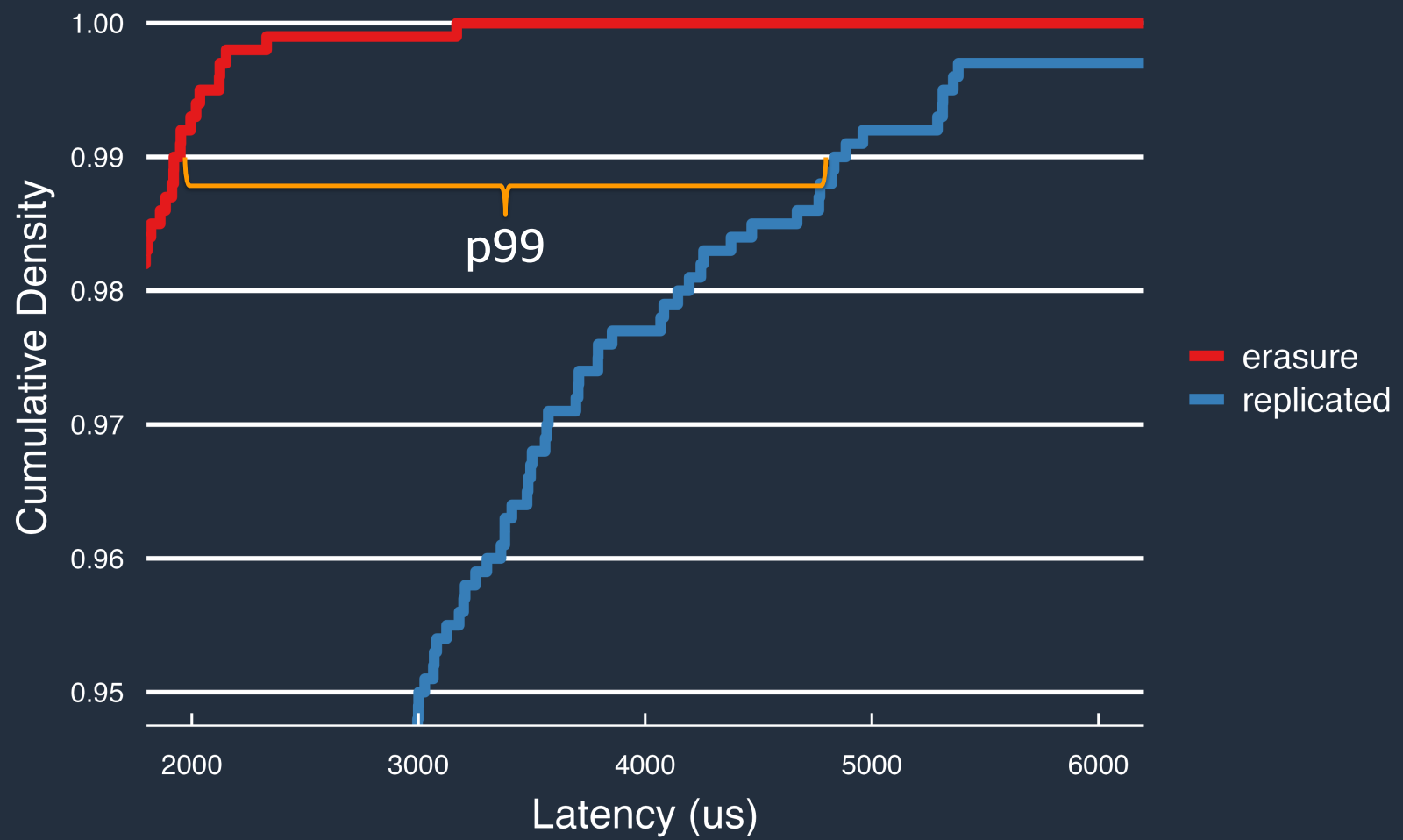


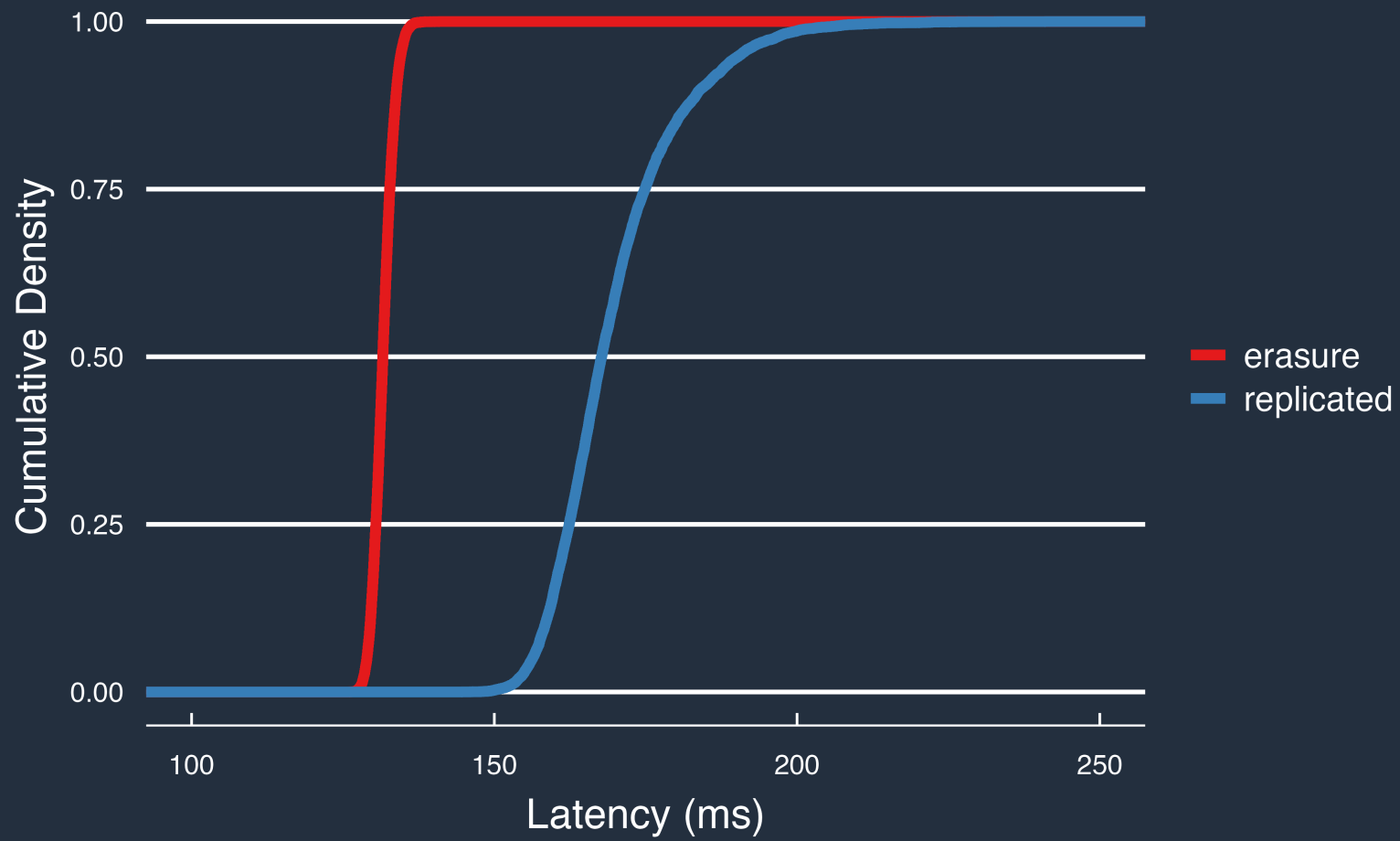
- One Copy
  - Not durable enough.
- Two Copies
  - 2x the cost! (or, half the effective cache size)



Erasure Code  
"Any 4 of 6"







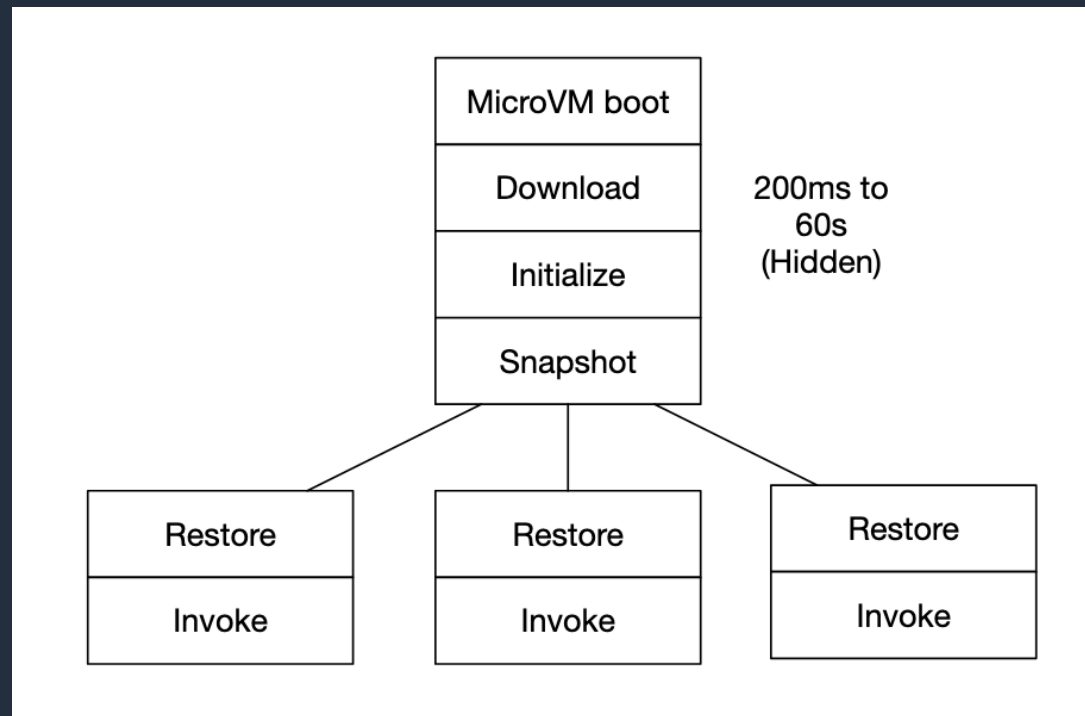
What about chunks of other stuff?  
Like memory?

MicroVM boot	~200ms
Download	20ms to 60s “cold start”
Initialize	
Invoke	<10ms “Warm start”
Invoke	
...	
Death	

Brooker, et al, *Restoring Uniqueness in MicroVM Snapshots*, <https://arxiv.org/pdf/2102.12892>

© 2020, Amazon Web Services, Inc. or its Affiliates. All rights reserved.





Brooker, et al, *Restoring Uniqueness in MicroVM Snapshots*, <https://arxiv.org/pdf/2102.12892>

© 2020, Amazon Web Services, Inc. or its Affiliates. All rights reserved.



**userfaultfd** - create a file descriptor for handling  
page faults in user space

# Not only cold starts!

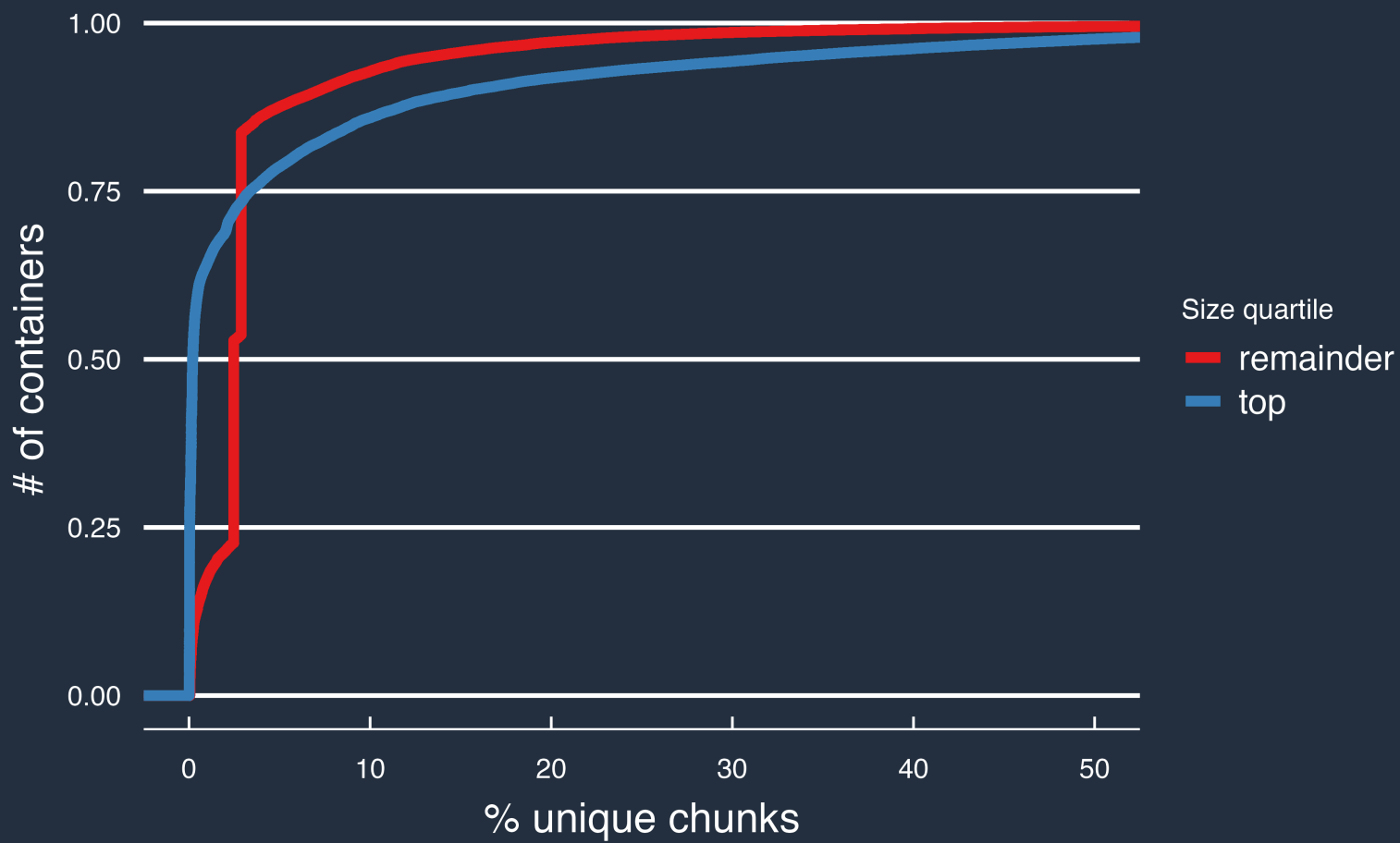
## Tenant Isolation (for multi-tenant services handling untrusted code).

## Session isolation (e.g. for AI agents).



# Questions

# Erasure Coding VS Complexity



- Deterministic Flattener
  - Can only put chunks, not read them.
- Worker
  - Can only read chunks for the functions it is running.
- Everybody else
  - Can do basically nothing.

Chunk name =

SHA2( Chunk data )

Chunk name =

SHA2(AES( Chunk data Key IV )))

# Chunk ciphertext =

$$\text{AES}_{\text{gcm}}(\text{Chunk data}, \text{key} = f(\text{Chunk data}), \text{iv} = [0, \dots])$$

Chunk name =

SHA2( chunk ciphertext )

Chunk name =

SHA2( chunk ciphertext )

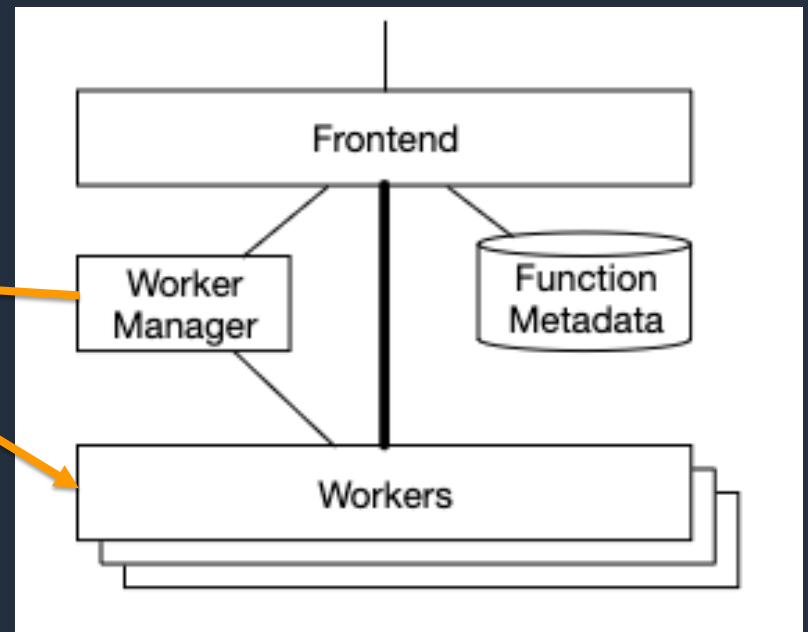
Q: Why not use the GCM *tag*?

chunk 1, name, key  
chunk 2, name, key  
...  
chunk N, name, key

Customer's KMS Key

chunk 1, name, key  
chunk 2, name, key  
...  
chunk N, name, key

Customer's KMS Key



- Deterministic Flattener
  - Can only put chunks, not read them.
- Worker
  - Can only read chunks for the functions it is running.
- Everybody else
  - Can do basically nothing.