

InfiniCache: Exploiting Ephemeral Serverless Functions to Build a Cost-Effective Memory Cache

Ao Wang* and Jingyuan Zhang*, Xiaolong Ma†, Ali Anwar‡, Lukas Rupprecht‡,

Dimitrios Skourtis‡, Vasily Tarasov‡, Feng Yan†, Yue Cheng*

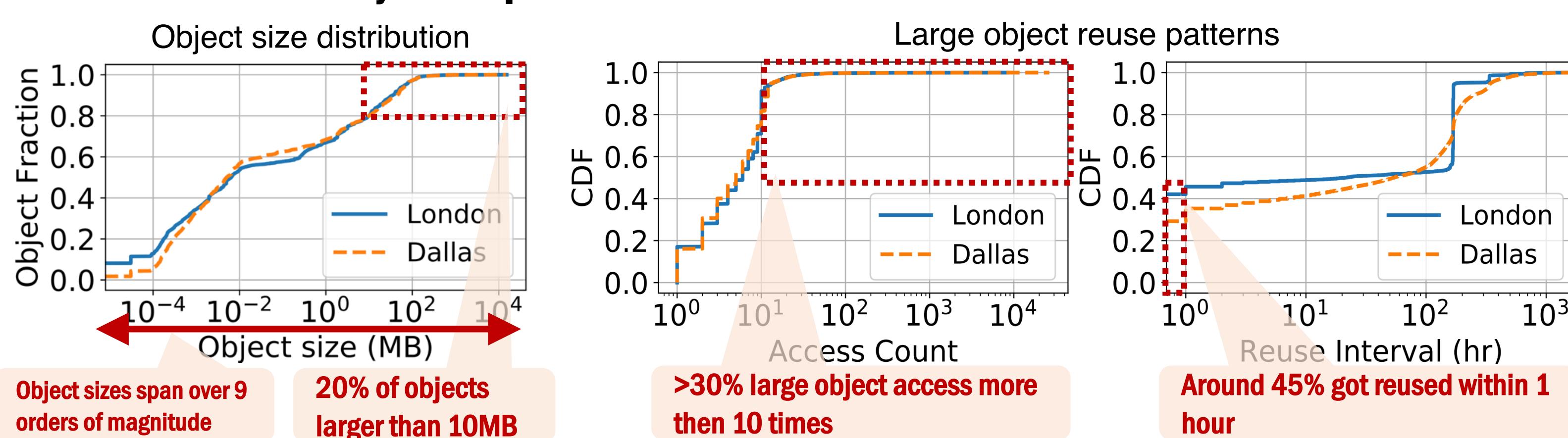
*George Mason University †University of Nevada, Reno ‡IBM Research-Almaden

Introduction

Web applications feature heterogeneous I/O patterns

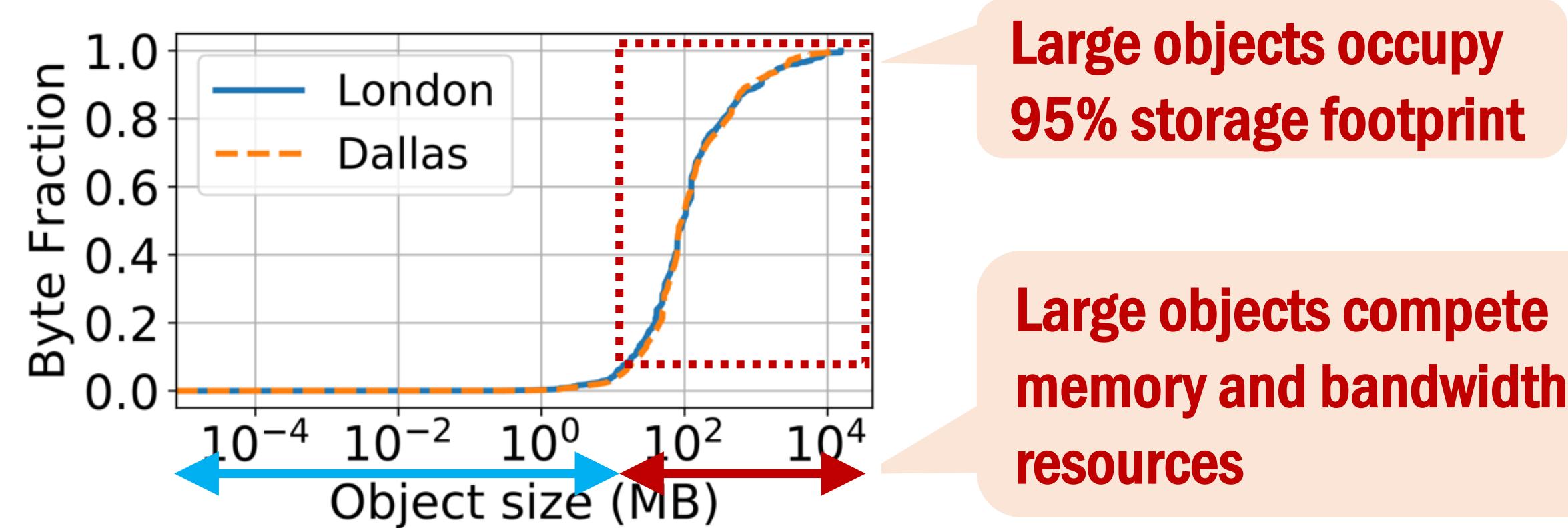


To quantify the I/O behavior in those storage-intensive applications, we conduct a case study:



Caching large objects is beneficial

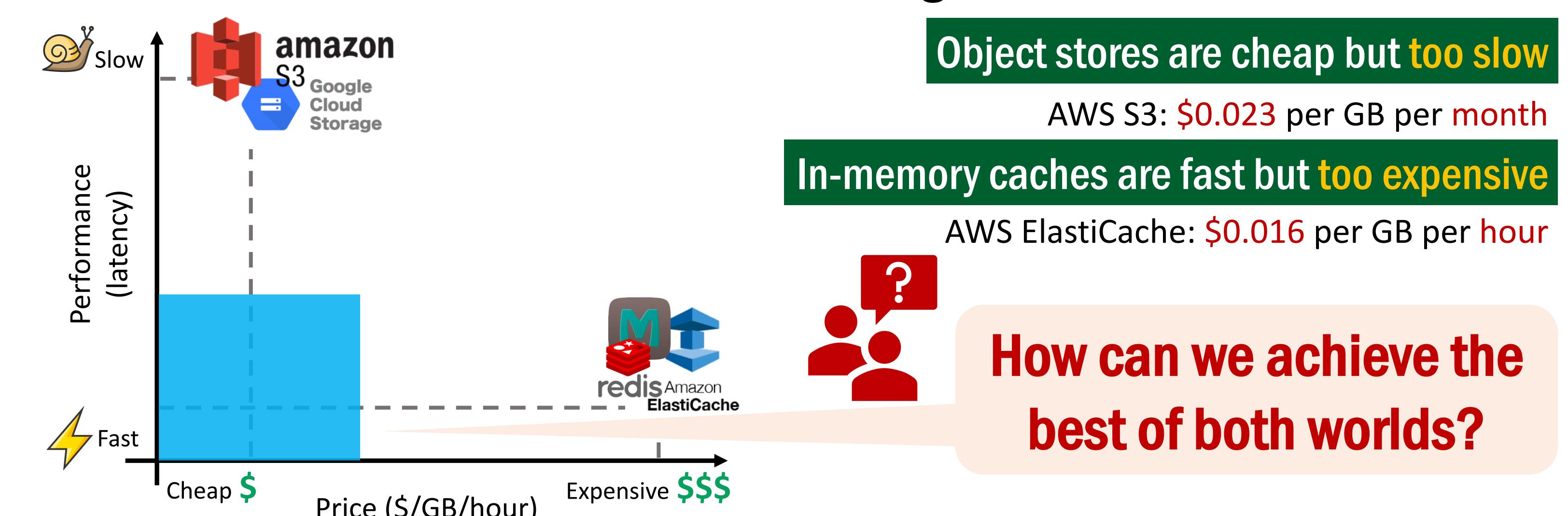
To cache both small and large objects, we analyze the storage footprint of objects:



Extreme tension between small and large objects

Motivation

Performance-cost trade-off in existing solutions



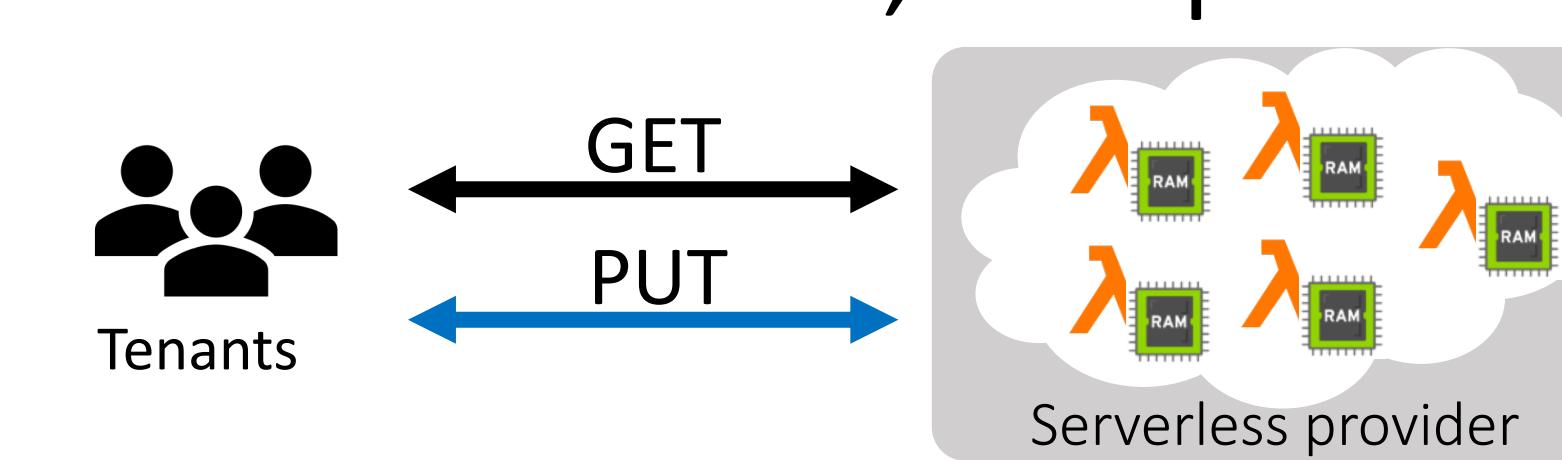
InfiniCache: A first-of-its-kind in-memory caching solution atop Serverless Computing platform

Insights:

Serverless functions' <CPU, Mem> resources are pay-per-use → Cost Effectiveness
Serverless providers offer "free" function caching for tenants → High Performance

A primer on Serverless Computing

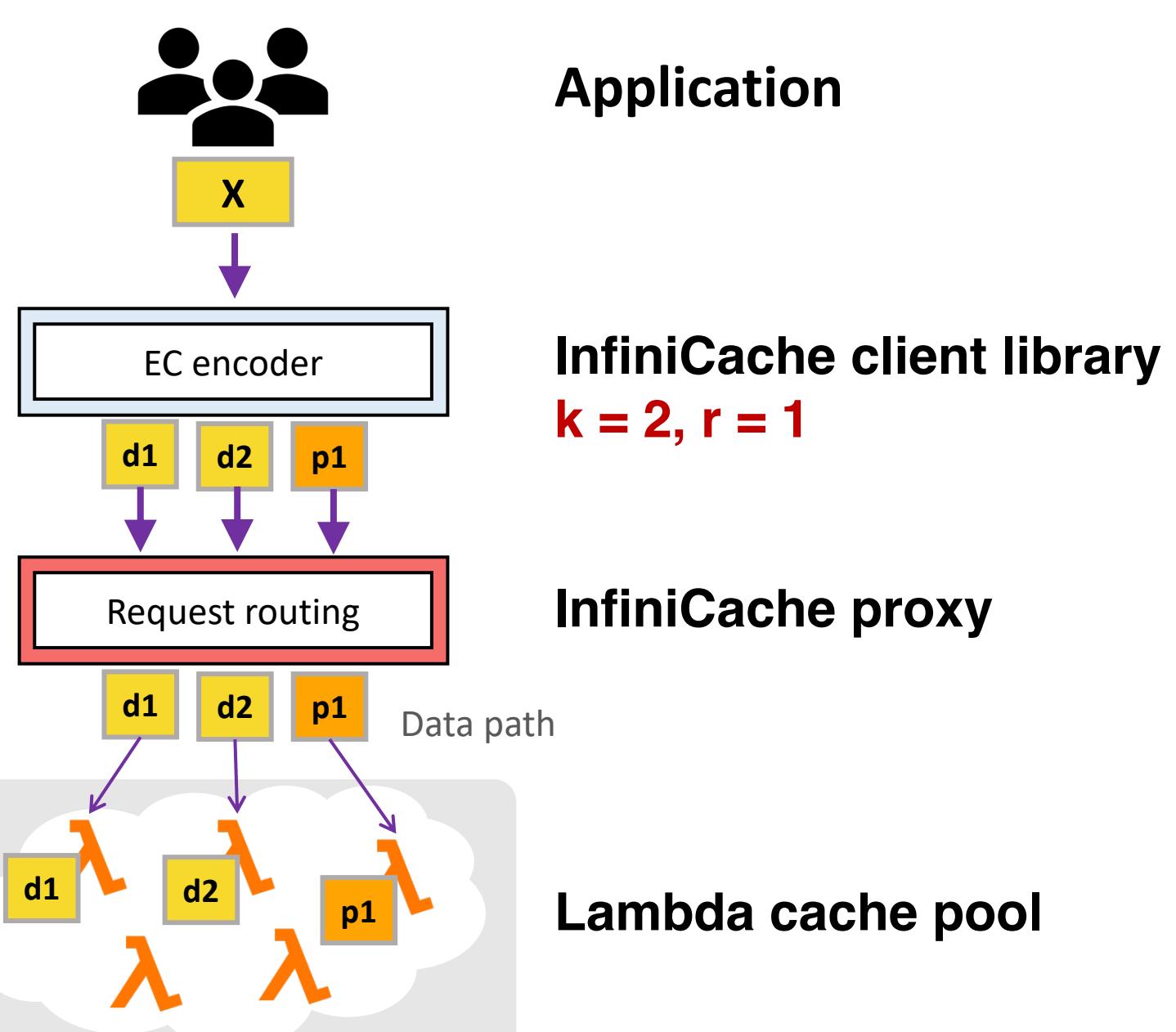
- * Serverless computing enables cloud tenants to launch short-lived tasks (i.e., Lambda functions) with **high elasticity** and **fine-grained resource billing**
- * Function: basic unit of deployment. Application consists of multiple serverless functions
- * Popular use cases: Backend APIs, data processing...



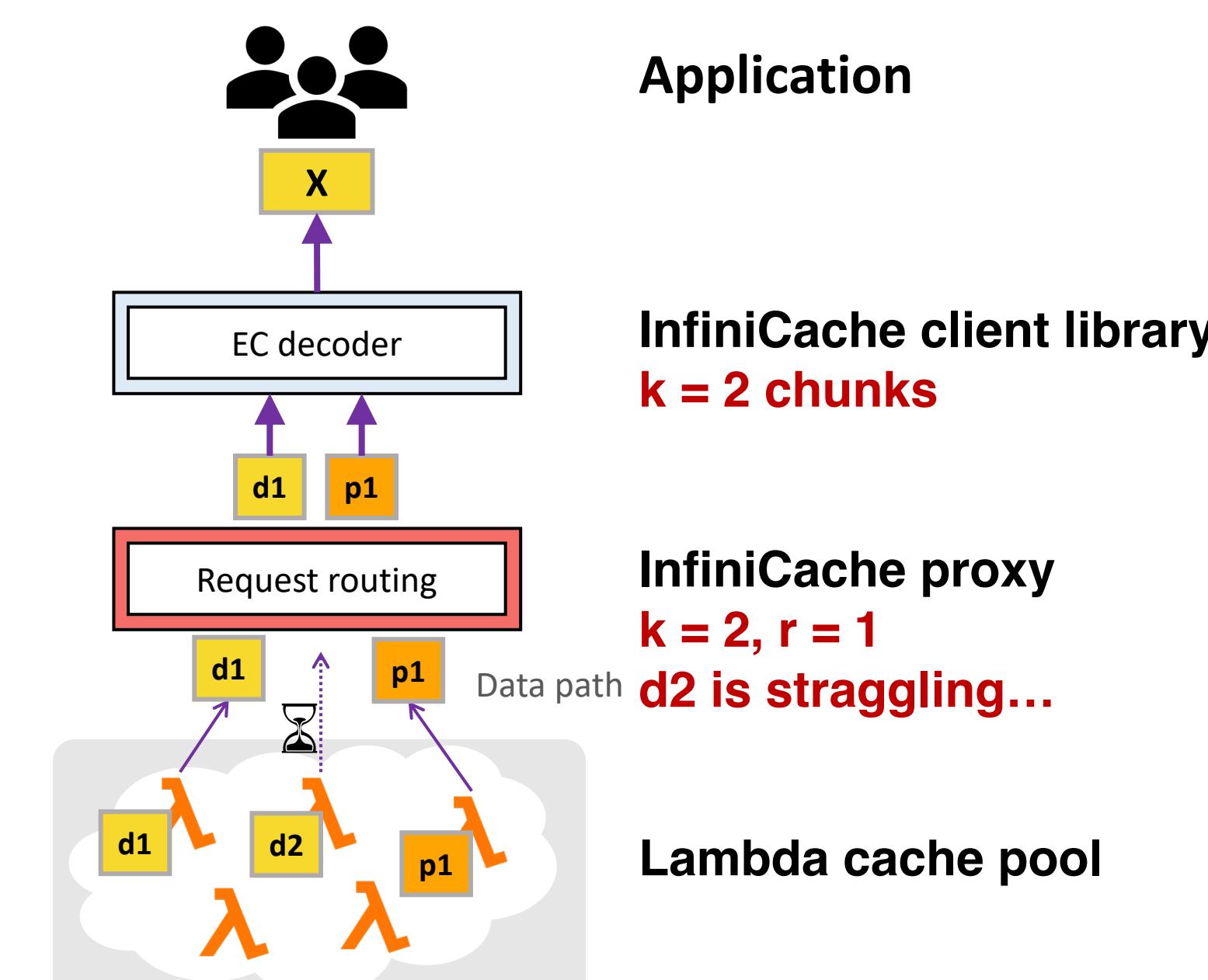
Design

InfiniCache achieves **high performance** by utilizing the aggregated network bandwidth of multiple serverless functions in parallel

PUT request example



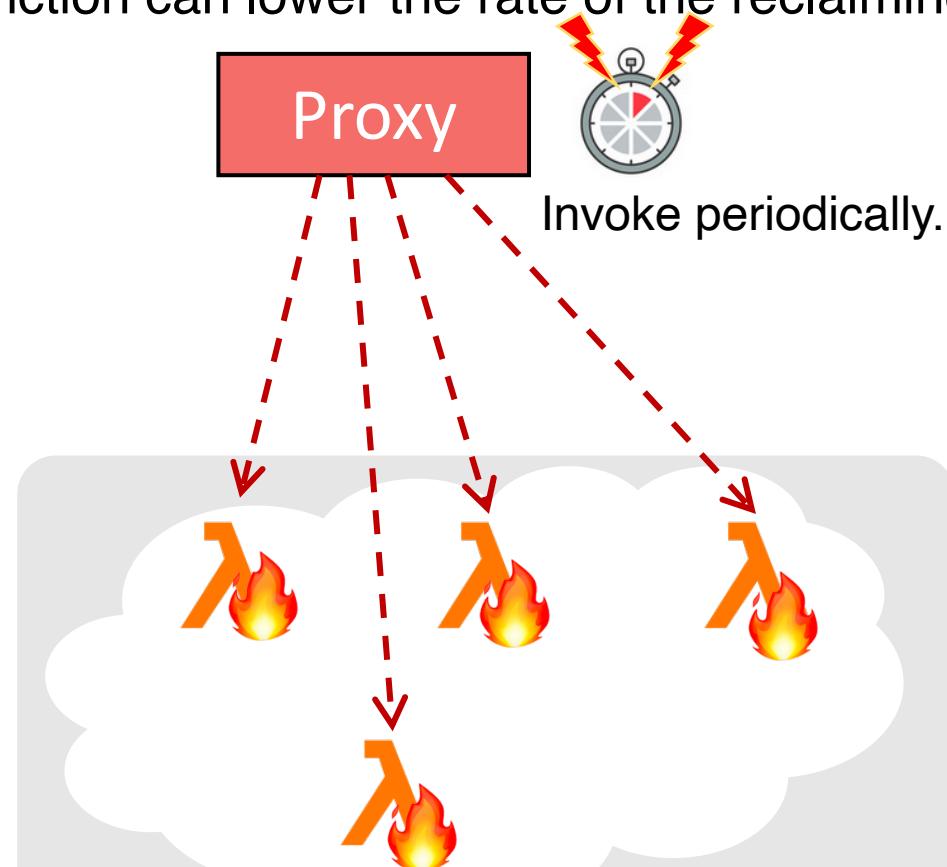
Get request example



InfiniCache achieves **high data availability** by leveraging 1. Erasure coding

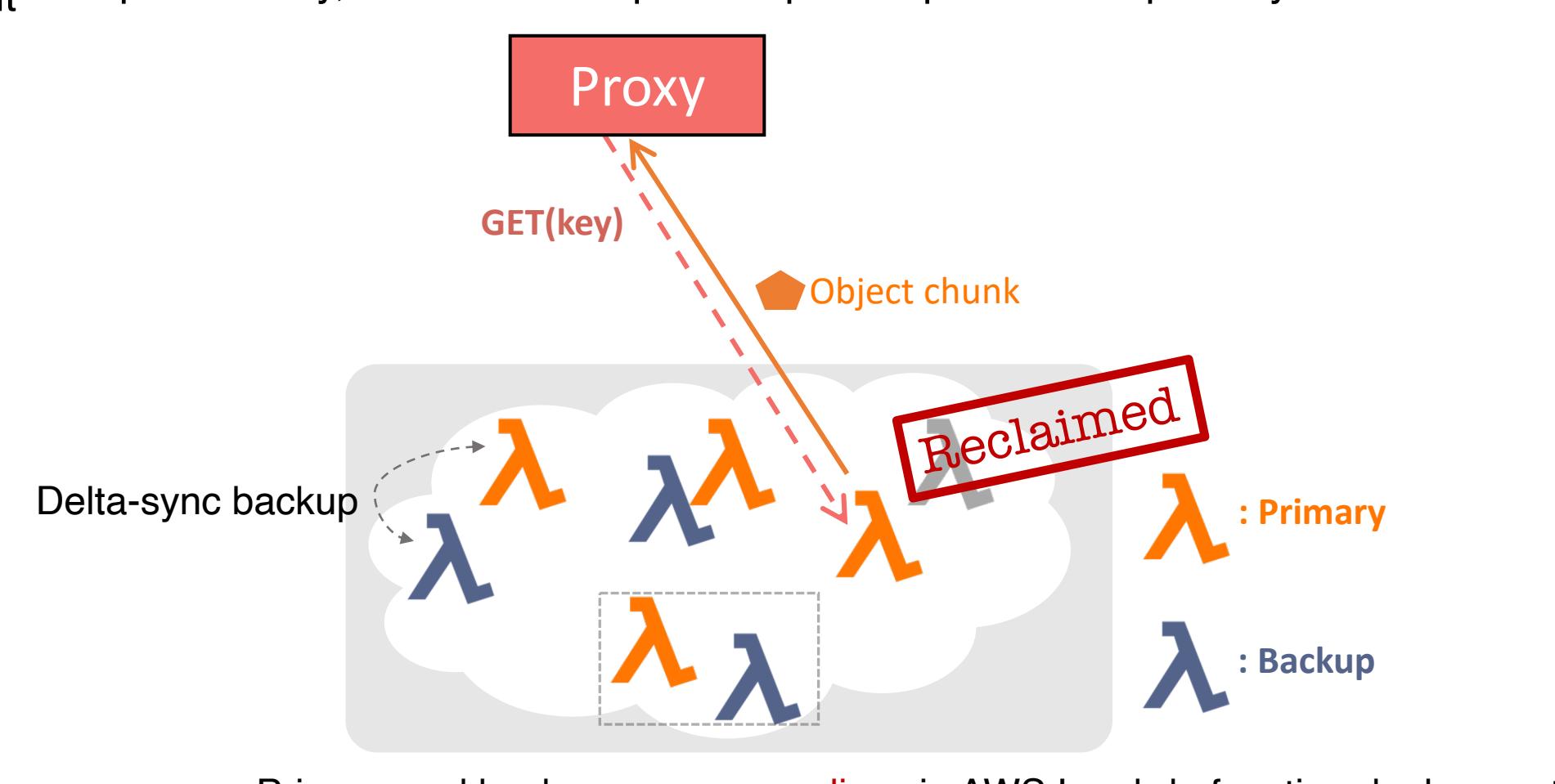
2. Periodic warm-up

AWS may reclaim Lambda functions after idling for a period. Our long-term empirical study shows periodically warming up lambda function can lower the rate of the reclaiming event



3. Periodic backup & Seamless failover

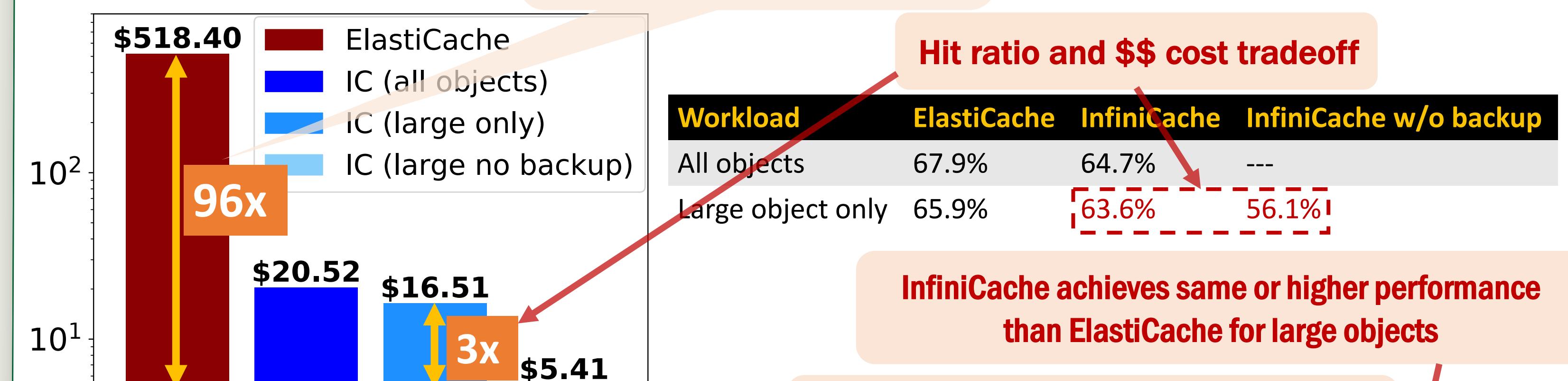
The failover is done seamlessly by exploiting the auto-scaling (elasticity) feature of AWS Lambda. Note each cache node (primary) get backed up periodically, and the backup is the peer replica of the primary



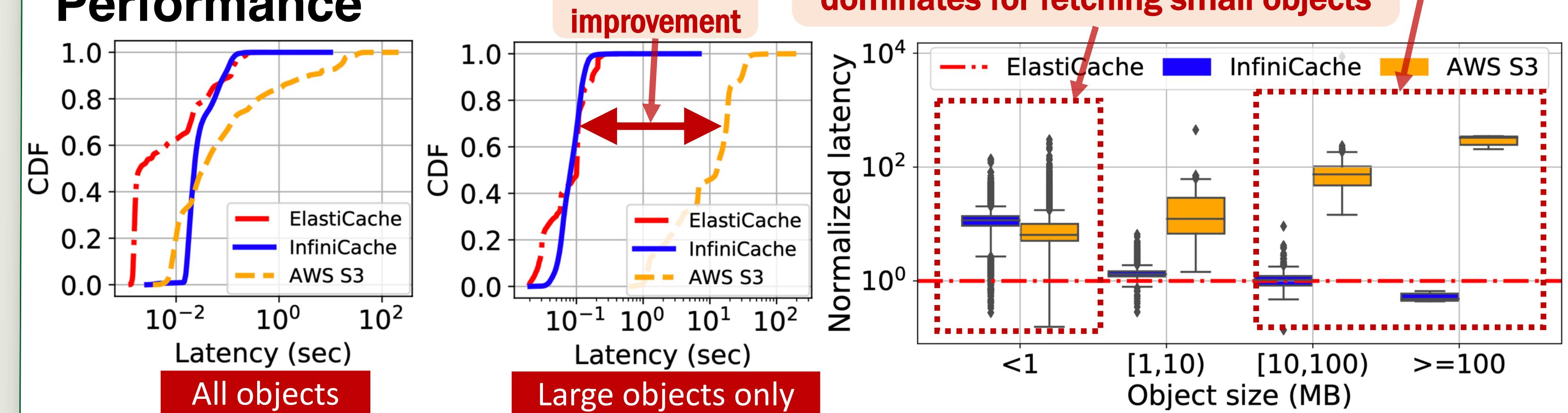
Primary and backup are peer replicas in AWS Lambda function deployment

Evaluation

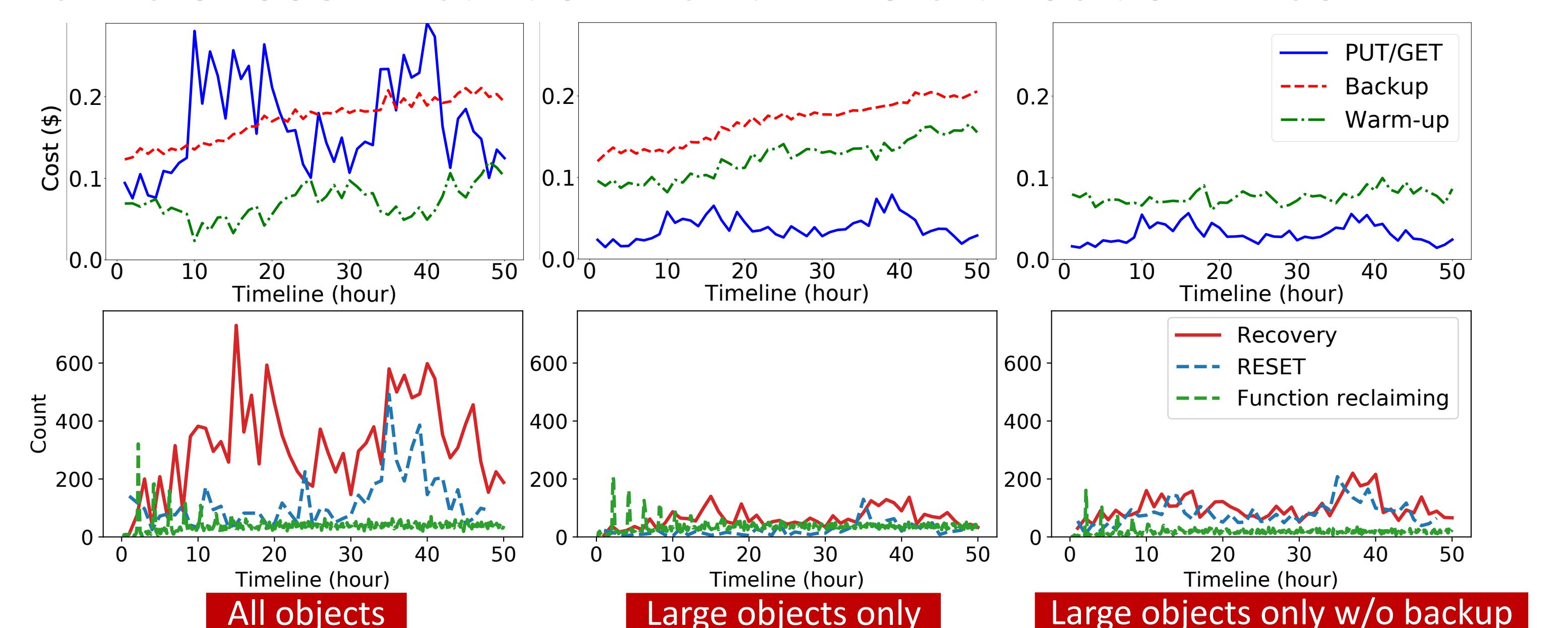
Cost Effectiveness



Performance



Timeline of cost breakdown & fault tolerance activities



Fork InfiniCache on GitHub! <https://github.com/mason-leap-lab/infinicache>