# SkyhookDM: An Arrow-Native Storage System

Jayjeet Chakraborty, Carlos Maltzahn

Centre for Research in Open Source Software

UC Santa Cruz

# The Broader Problem

- CPU is the new bottleneck with modern high speed storage and network devices like NVMe and Infiniband networks

- Client-side computation of data and reading from efficient storage formats like Parquet, ORC exhausts the clients CPUs

- Scalability and Latency is severely hampered.

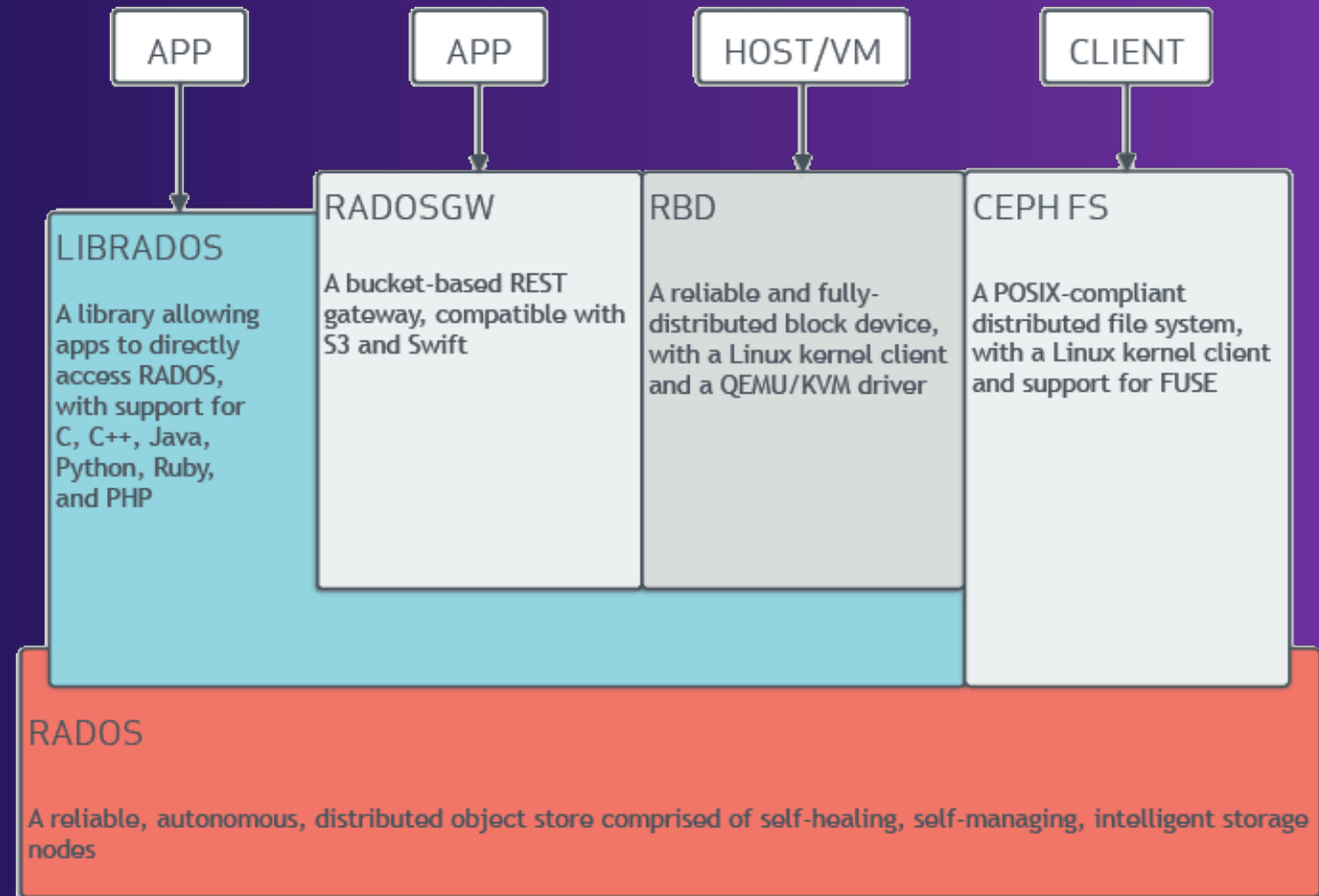STORAGE DEVELOPER CONFERENCE
SDC 21

# Computational Storage as a Solution

- Offload computation from the client to the storage layer as much as possible

- Utilize the idle CPUs of storage systems for increased processing rates and faster queries

- Results in lesser data movement, memory copying, and network traffic

STORAGE DEVELOPER CONFERENCE
SDC 21

# Ceph

# Introduction

- Provides 3 types of storage interface: File, Object, Block

- No central point of failure. Uses CRUSH maps that contains Object - OSD mapping

- Extensible Object storage layer via the Ceph Object Classes SDK

APP — LIBRADOS: A library allowing apps to directly access RADOS, with support for C, C++, Java, Python, Ruby, and PHP

APP — RADOSGW: A bucket-based REST gateway, compatible with S3 and Swift

HOST/VM — RBD: A reliable and fully-distributed block device, with a Linux kernel client and a QEMU/KVM driver

CLIENT — CEPH FS: A POSIX-compliant distributed file system, with a Linux kernel client and support for FUSE

RADOS: A reliable, autonomous, distributed object store comprised of self-healing, self-managing, intelligent storage nodes

STORAGE DEVELOPER CONFERENCE
SDC 21

# Object Class Mechanism

- Utilizing Ceph's object class mechanism ("cls")
  - Object storage extension mechanism
  - Present in [ceph/src/cls](ceph/src/cls)
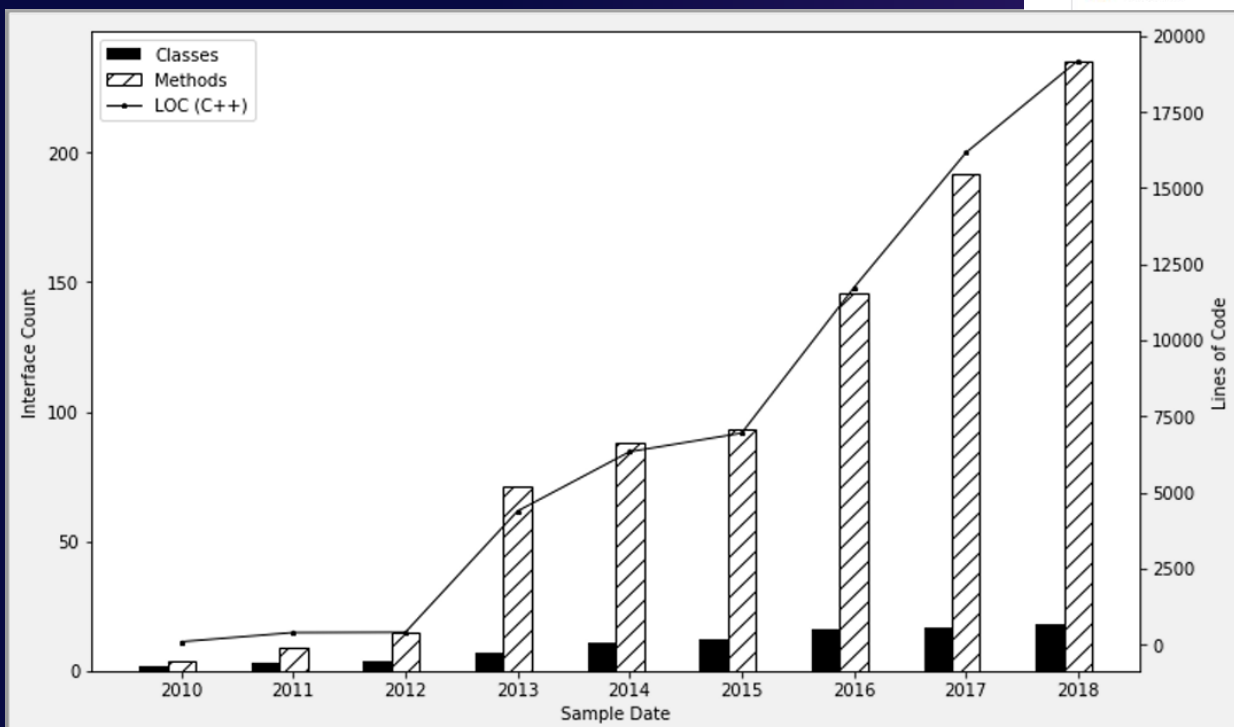- Used by several Ceph internals
  - CephFS, RGW, RBD

**Object Classes in Ceph**
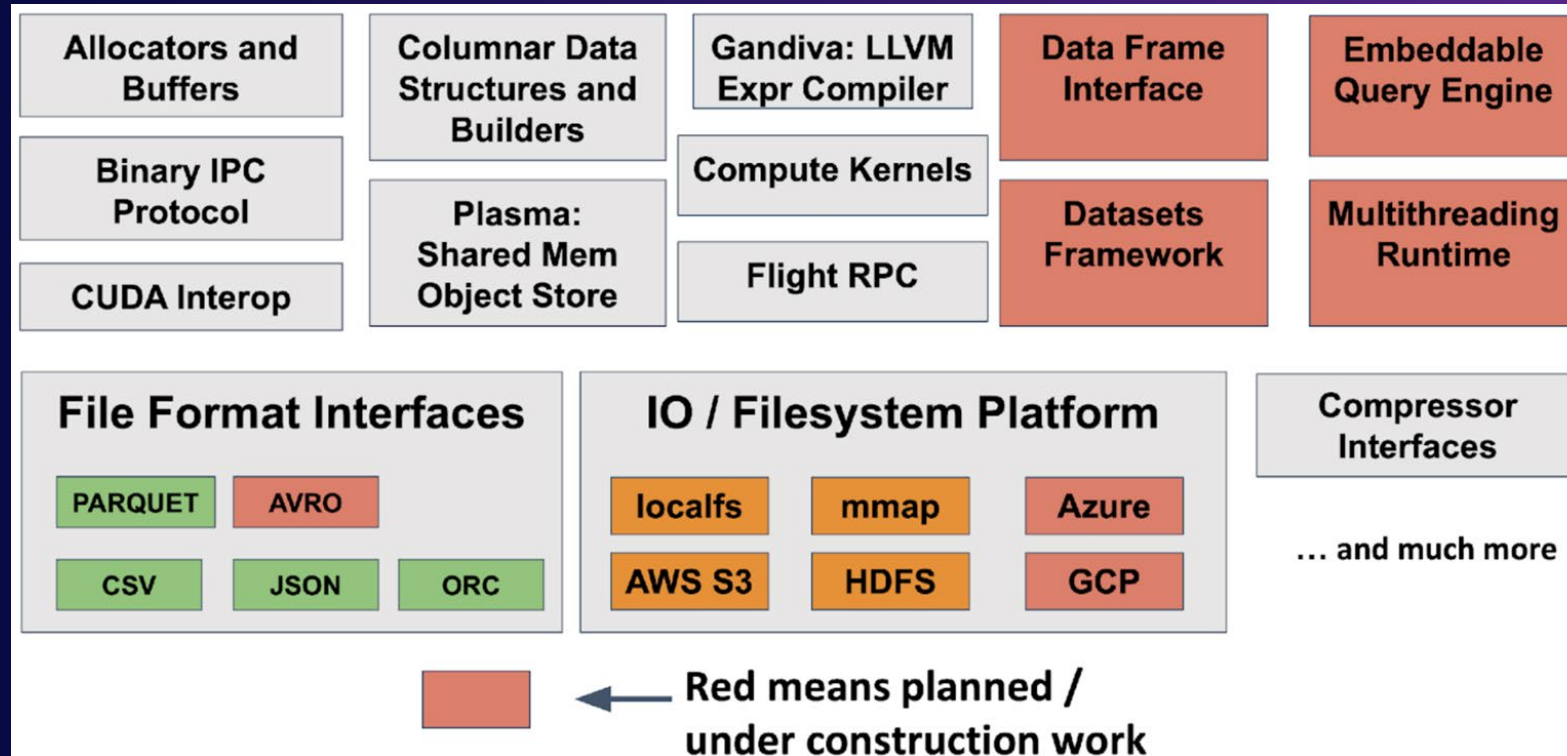
**Growth of Object Classes in Ceph**

# Apache Arrow

- Language-independent columnar memory format for flat and hierarchical data, organised for efficient analytic operations on modern hardware
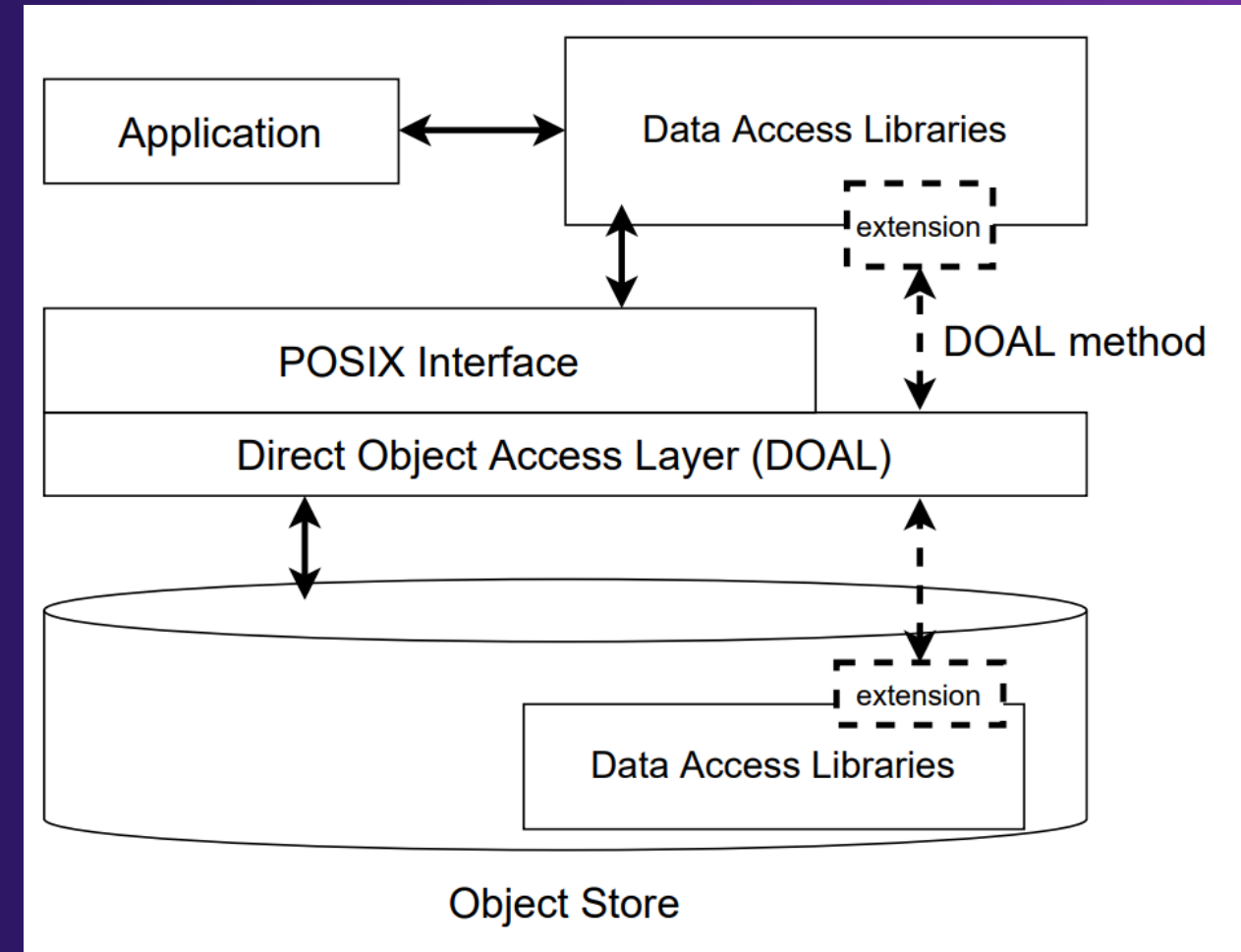
# What else ?

- Rich collection of pluggable components for building data processing systems
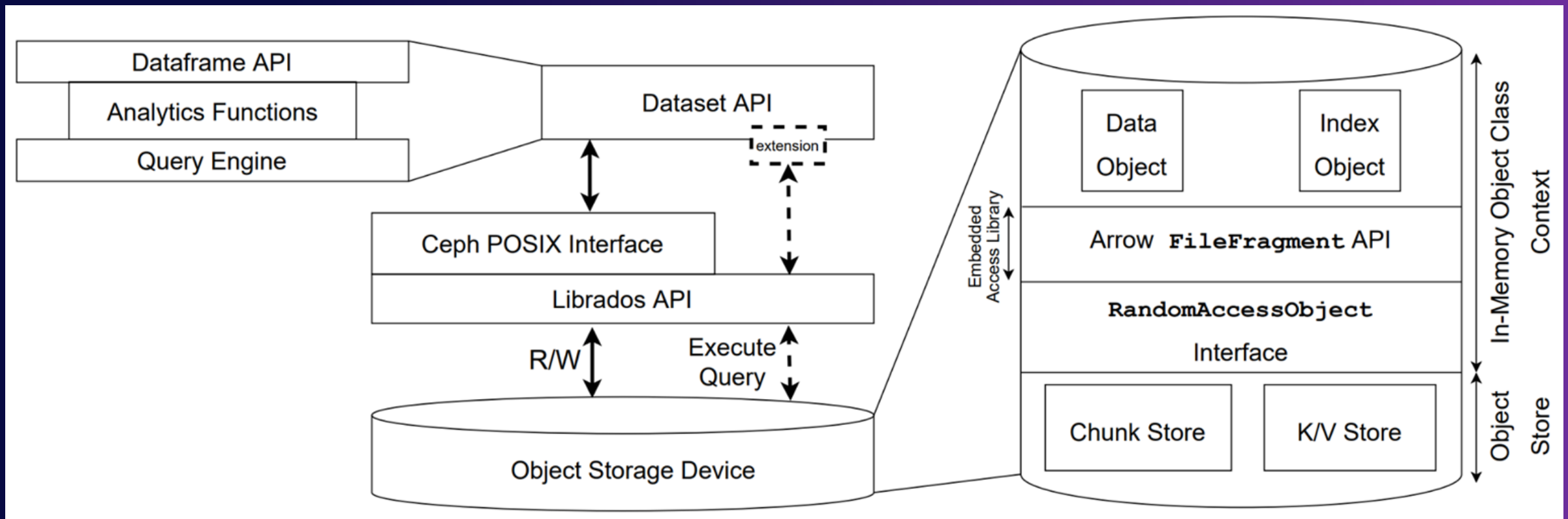
# Design Paradigm

- Extend client and storage layers of programmable storage systems with data access libraries

- Embed a FS shim inside storage nodes to have file-like view over objects

- Allow direct interaction with objects in an object store while bypassing the filesystem layer utilising FS metadata
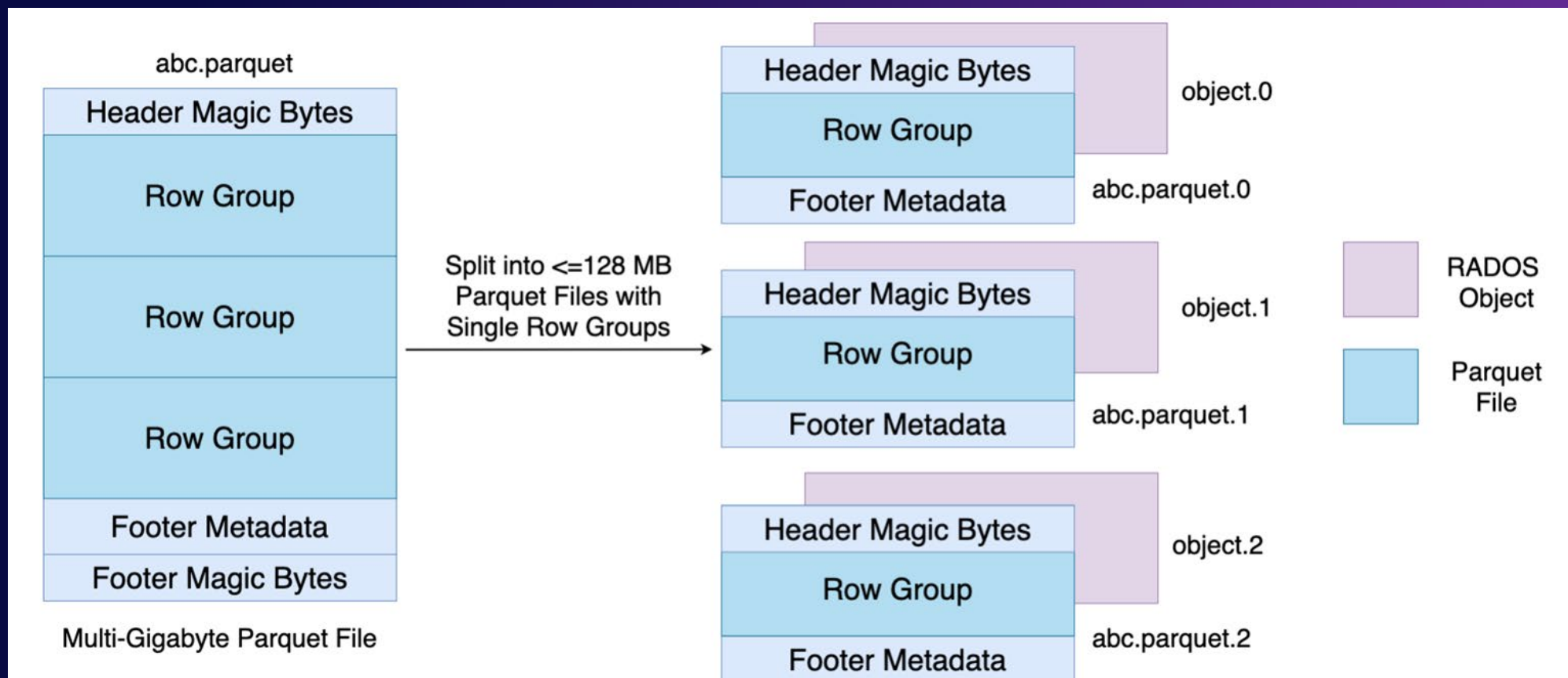
# Architecture

STORAGE DEVELOPER CONFERENCE

≋SD C 21

- Arrow data access libraries embedded inside Ceph OSDs to allow scanning data fragments in the Ceph storage layer

- Extend Arrow Dataset API with `SkyhookFileFormat` to expose the offload capability

# File-Layout Design

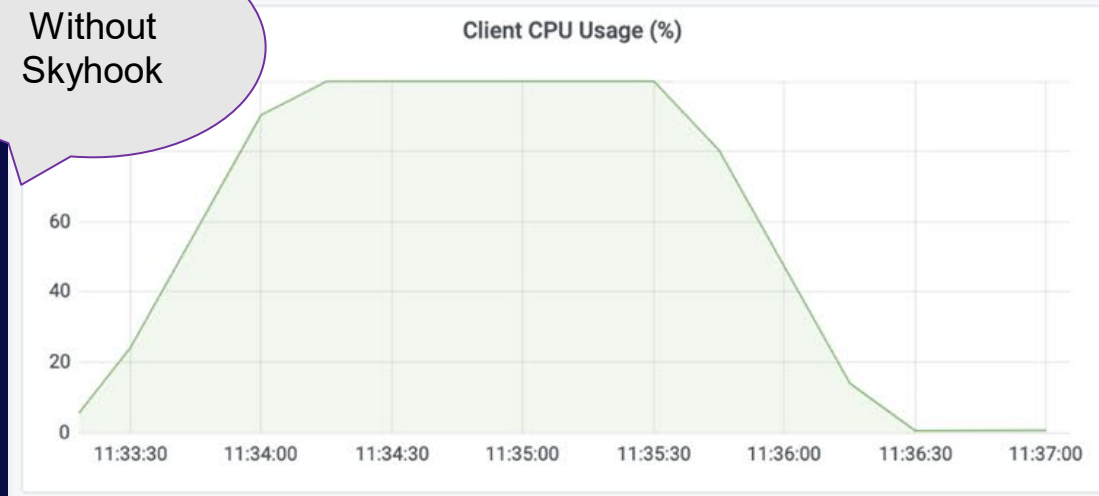STORAGE DEVELOPER CONFERENCE

SDC 21

- 16MB is the preferred file size in SkyhookDM as found out from several experiments with different file sizes.

- Files larger than 16MB are splitted into smaller files of ~16MB and each file is stored in a single RADOS object.

- Due to Arrow Dataset API being the data access library, a wide range of file formats like IPC, Parquet, CSV are supported out of the box.
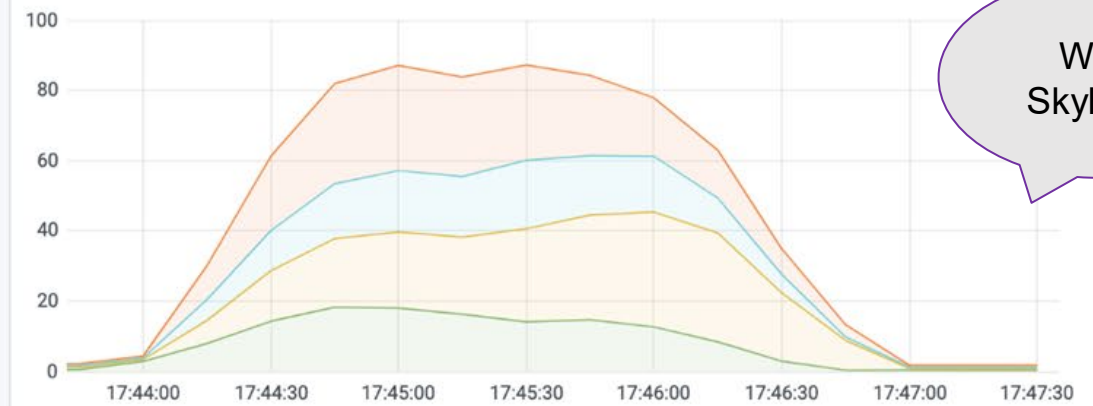
STORAGE DEVELOPER CONFERENCE
SDC 21

# Results

STORAGE DEVELOPER CONFERENCE

SDC 21

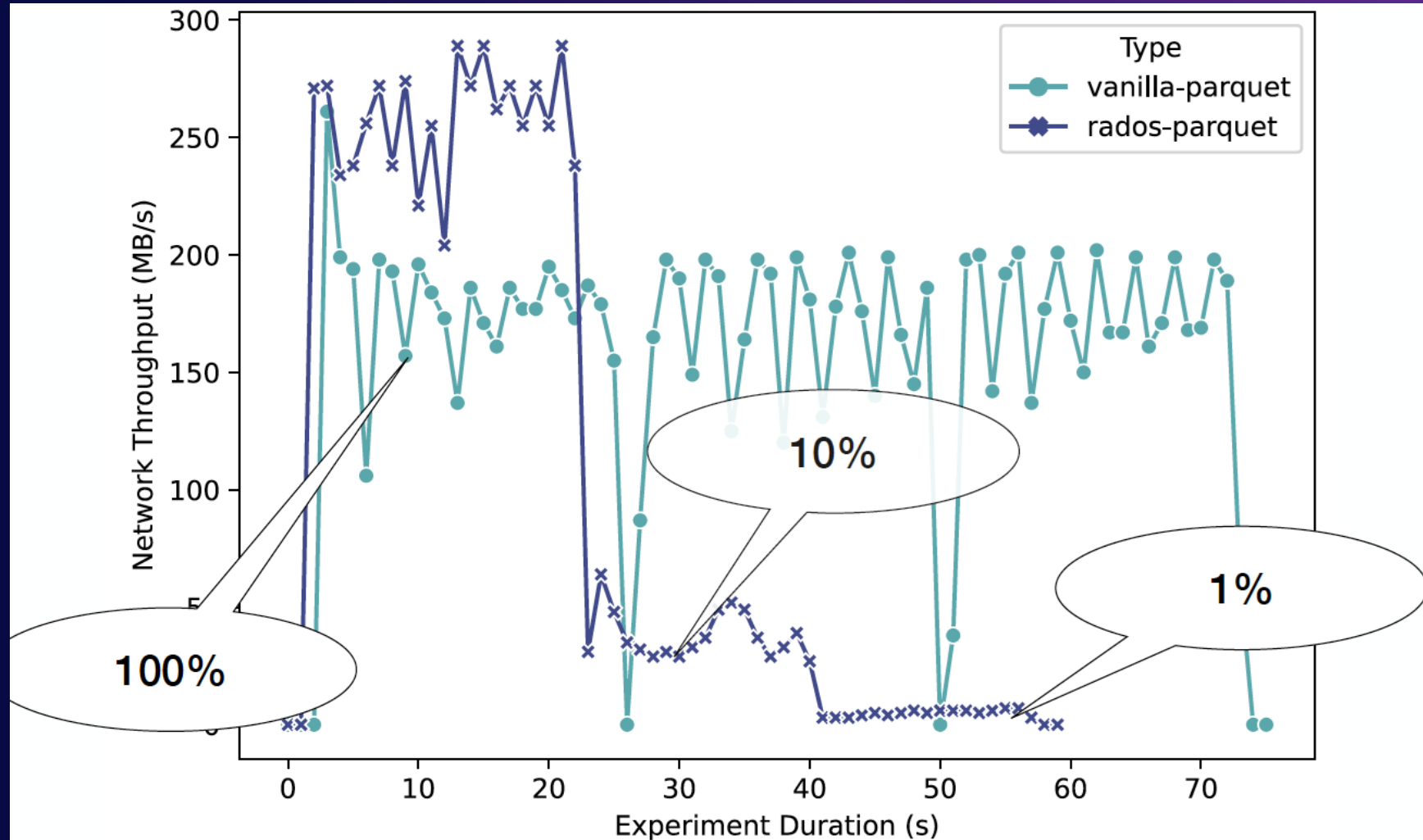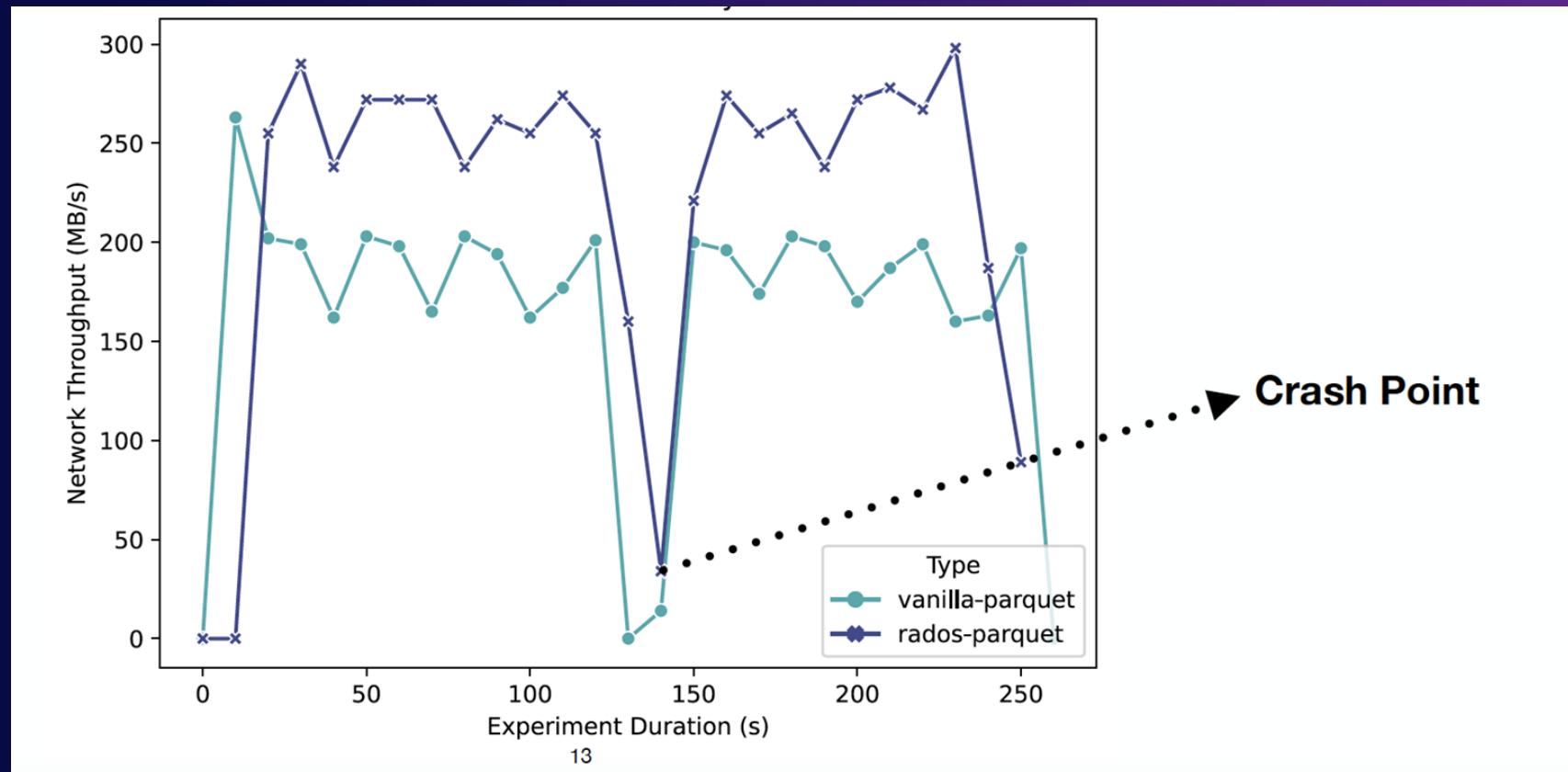# Offloaded CPU usage

# Reduced Wastage of Network Bandwidth

# Automatic Failure Recovery

Since, compute is colocated with storage nodes, the failure recovery and consistency semantics of the storage system apply naturally to the query processing layer

# Please take a moment to rate this session

STORAGE DEVELOPER CONFERENCE

SD C 21

# Thank You !