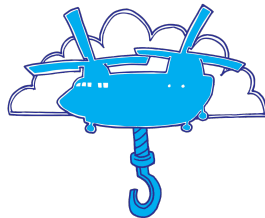# Skyhook Data Management
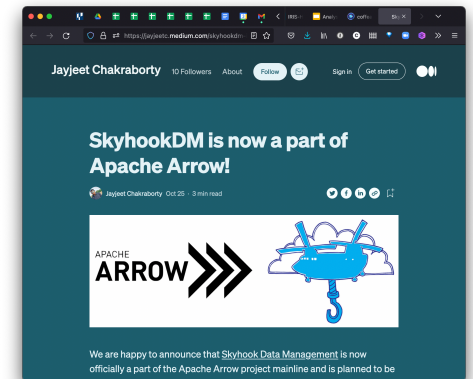
Carlos Maltzahn, 11/4/21

Analysis Grand Challenge Tools 2021 Workshop

# What is Skyhook Data Management?

- Also known as SkyhookDM or just "Skyhook"

- Since 10/22/21 part of Apache Arrow (will be part of 7.0.0)
  - Columnar memory format for flat and hierarchical data
  - Large ecosystem of mapping Arrow data to storage, GPUs, FPGAs

- Offloads Apache Arrow *scans* into a storage system
  - Embeds the Apache Arrow library with minimal changes

- Reduces client-side resource utilization (CPU, memory, network)
  - Faster networks → more CPU and memory BW for data movement
  - Particularly good for data-intensive selection operations

- Storage systems can optimize dataset operations based on *local* info
  - Fewer "magic numbers" applications have to worry about

# Implementation

```
├── cls
│   ├── cephfs
│   ├── hello
│   ├── journal
│   ├── lock
│   ├── log
│   ├── lua
│   ├── numops
│   ├── rbd
│   ├── refcount
│   ├── replica_log
│   ├── rgw
│   ├── sdk
│   ├── statelog
│   ├── tabular
│   ├── timeindex
│   ├── user
│   └── version
```
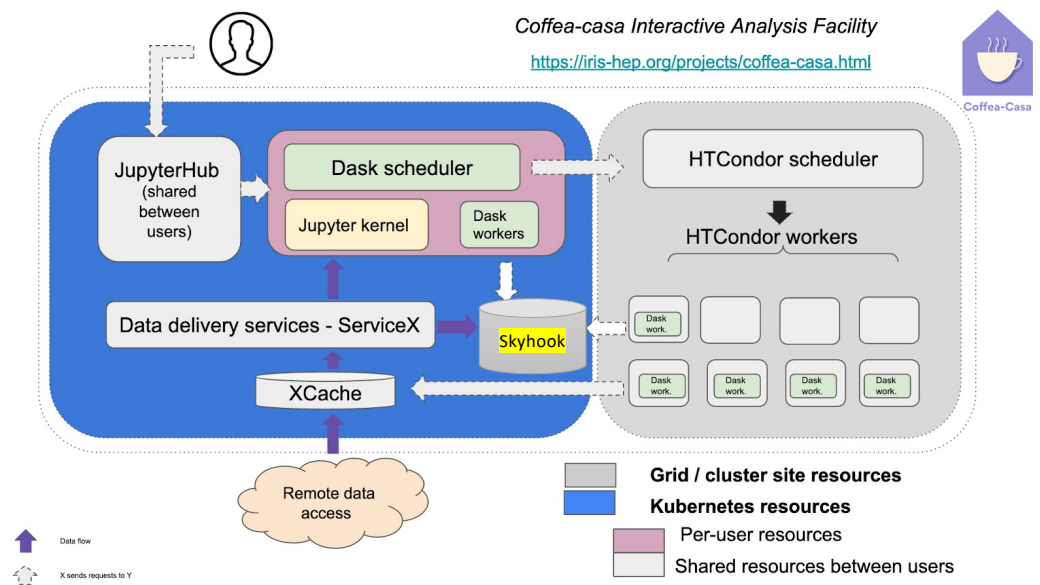
tree -d
ceph/src
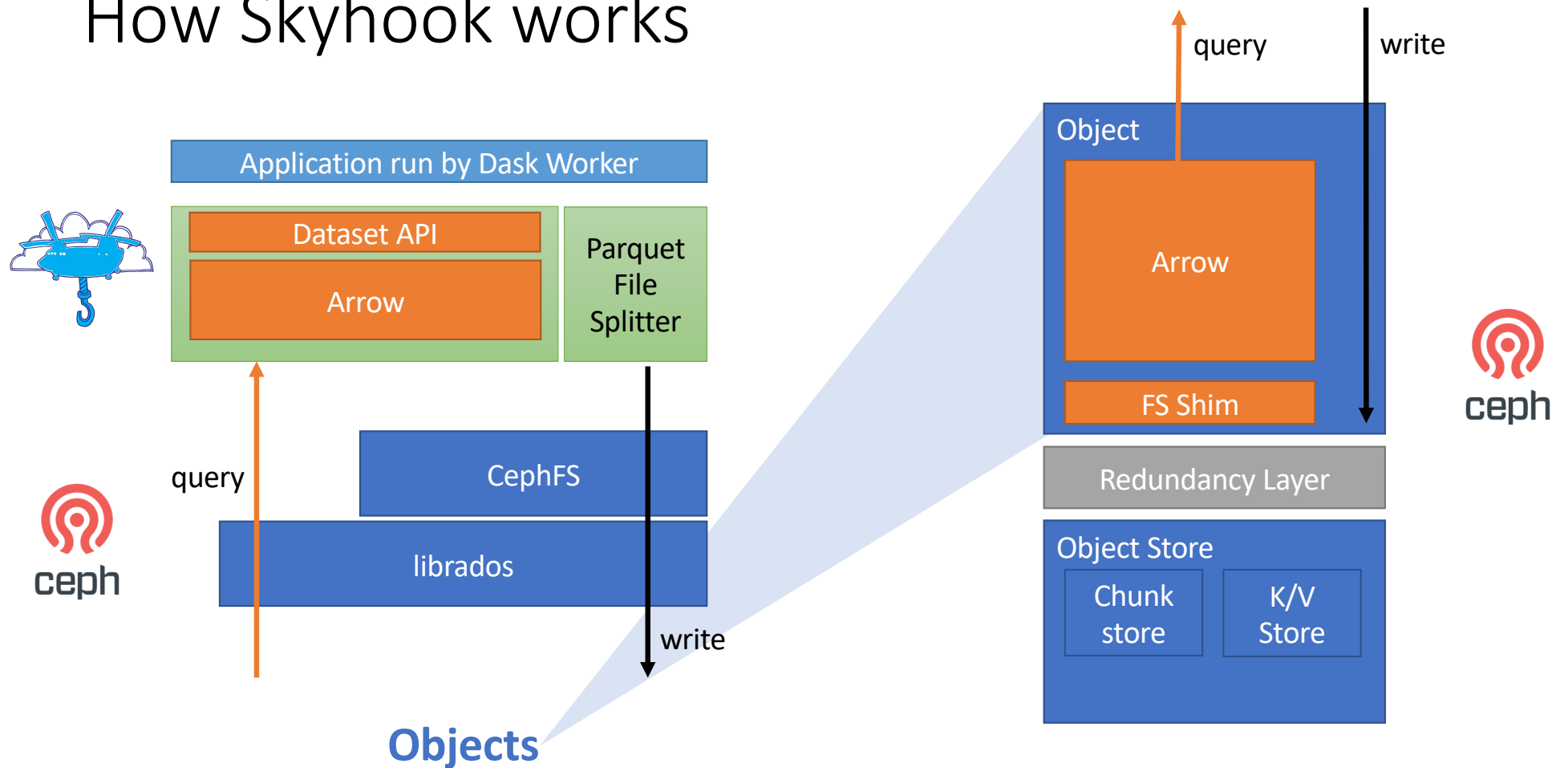
SkyhookDM

An object "class" for Ceph
- No upstream modifications required
- Inherits Ceph's properties now and in the future
- Can use all other object extensions
- **Not a database**

# Role in Analysis Facility

- Data caching
- Database caching
- Persistent dataset views
- Friend Tree management
- Selection by trigger bits



*Coffea-casa Interactive Analysis Facility*
https://iris-hep.org/projects/coffea-casa.html

# How Skyhook works

| Application run by Dask Worker | |
|---|---|
| Dataset API<br><br>Arrow | Parquet File Splitter |

CephFS

librados

**query**

**write**

**Objects**

**ceph**

---

**query** **write**

Object

Arrow

FS Shim

Redundancy Layer

Object Store

| Chunk store | K/V Store |

**ceph**

# (Old) example notebook

https://github.com/CoffeaTeam/coffea/blob/master/binder/nanoevents_pq.ipynb