



Data Management with Skyhook

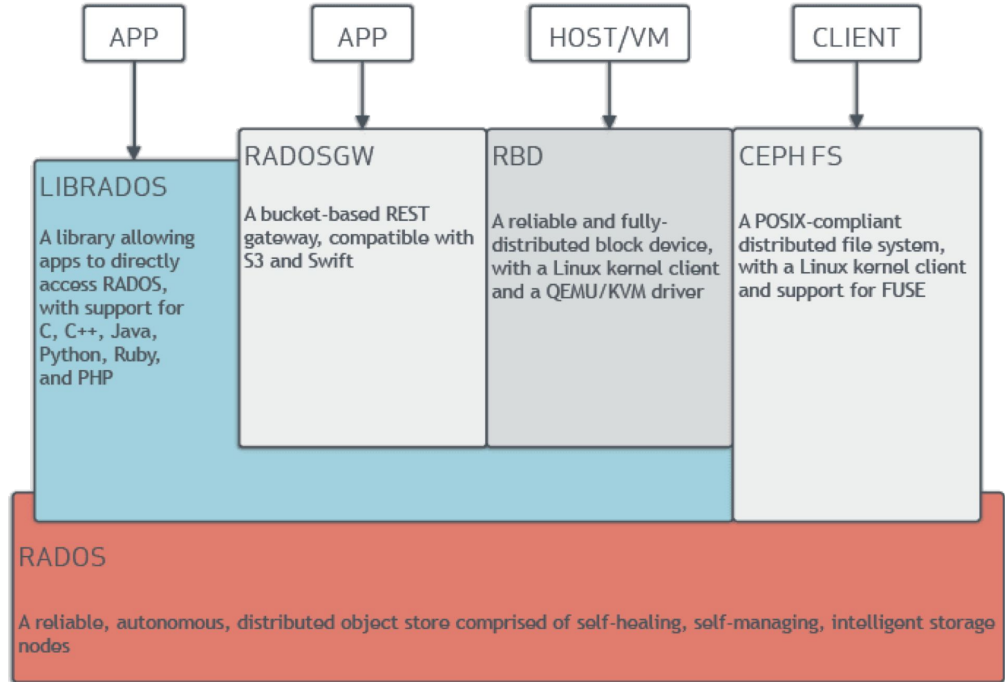
Jayjeet Chakraborty, Carlos Maltzahn, Jeff LeFevre, Ivo Jimenez,
Oksana Shadura, Alex Held, Fengping Hu, Brian Bockelman

Ceph

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

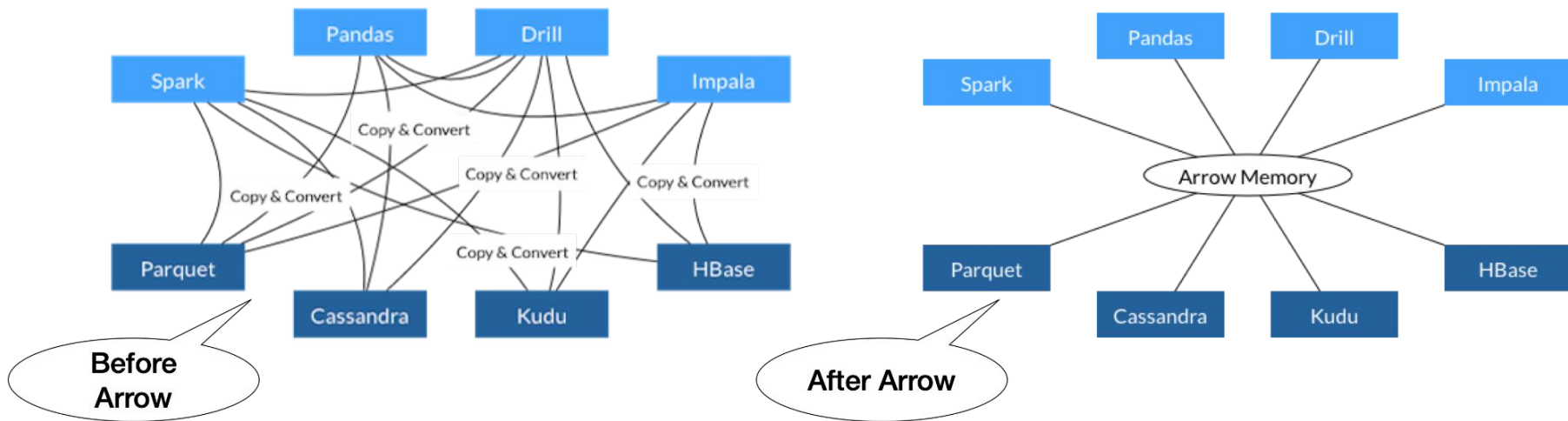
No central point of failure. Uses
CRUSH maps that contains object -
OSD mapping. A CRUSH map in each
client. Client talks directly to the OSDs.

Highly extensible Object storage layer
via the Ceph Object Classes SDK.



Apache Arrow

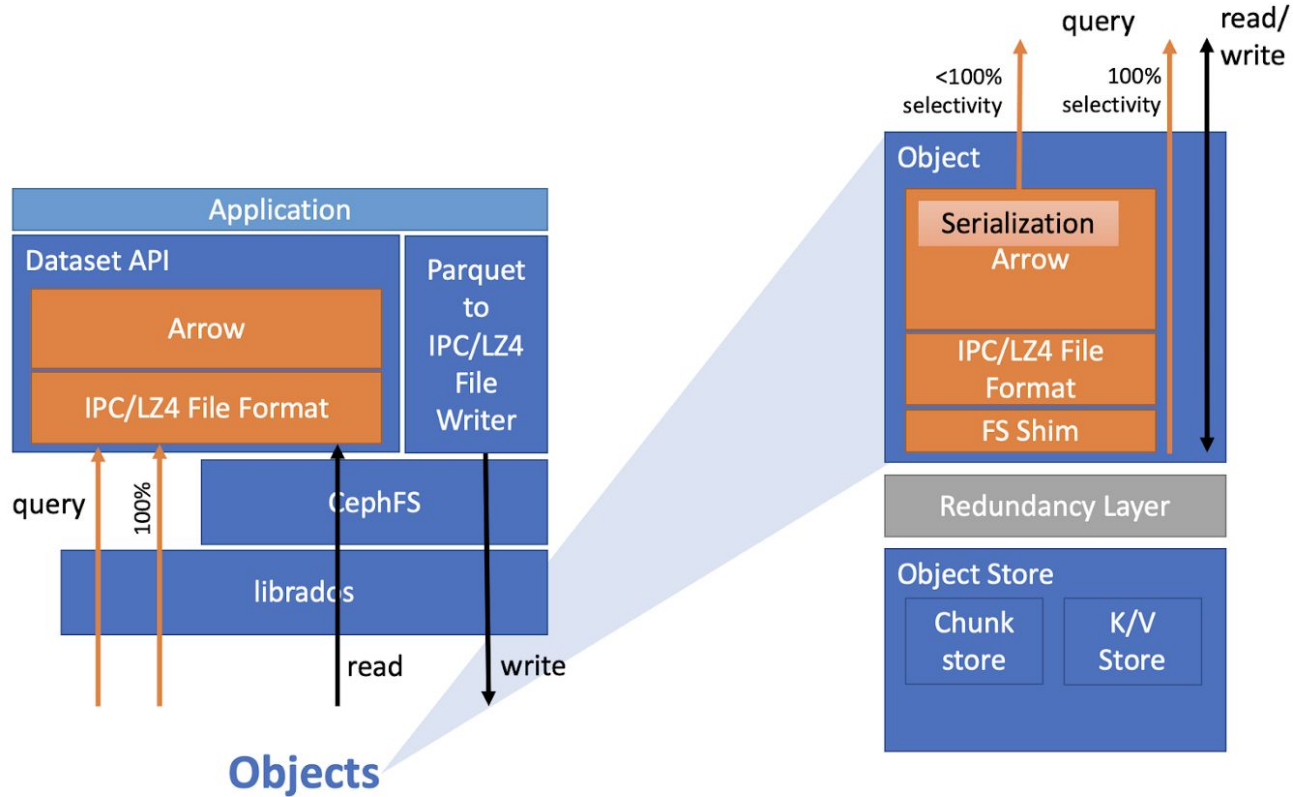
- Language-independent columnar memory format for flat and hierarchical data, optimized for efficient analytic operations on modern hardware.
- Share data between processes without serialization overhead.



Skyhook ([CCGrid '22](#))

- A programmable storage system to offload selections and projections to the storage system
- Reduce data movement across the system
 - Especially, when servers are in different racks with wimpy network
- Mitigates CPU scarcity on the client by giving compute scalable access to storage layer CPUs

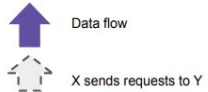
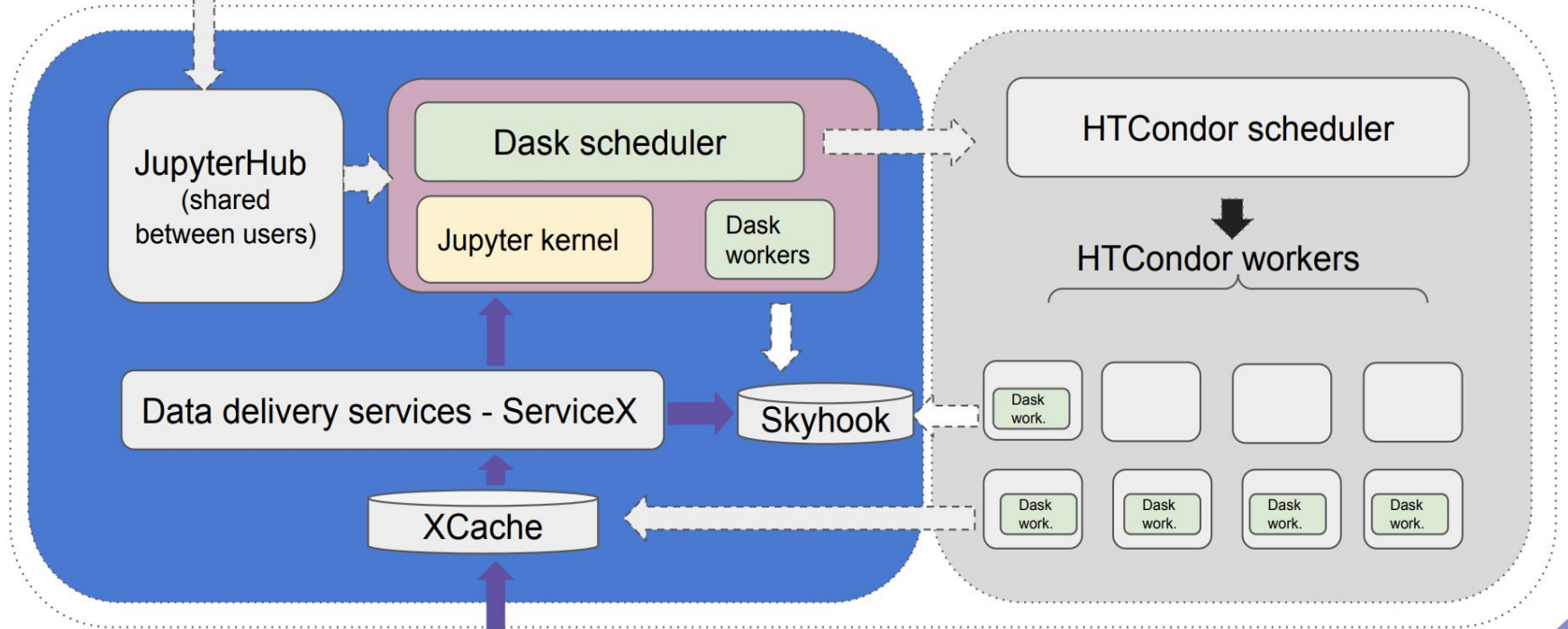
Architecture



Skyhook integrates with..

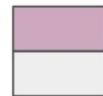
- **Dask:** parallel computing framework similar to Spark
 - Enable offloading dataframe operations in cluster computing systems
- **Coffea:** high energy physics analysis framework
 - Offload operations such as cut and zip on Nano events in HEP data analytics

Where does Skyhook fit in the Analysis
Facility ?



Grid / cluster site resources

Kubernetes resources



Per-user resources

Shared resources between users



Demo