

# Too Big to Eat: Boosting Analytics Data Ingestion from Object Stores with Scoop

Yosef Moatti<sup>1</sup>, Eran Rom<sup>2</sup>, Raúl Gracia-tinedo<sup>3</sup>, Dalit Naor<sup>1</sup>, Doron Chen<sup>1</sup>, Josep Sampé<sup>3</sup>, Marc Sánchez-Artigas<sup>3</sup>, Pedro García-López<sup>3</sup>, Filip Gluszek<sup>4</sup>, Eric Deschodt<sup>4</sup>, Francesco Pace<sup>5</sup>, Daniele Venzano<sup>5</sup>, Pietro Michiardi<sup>5</sup>  
 IBM Research<sup>1</sup>, OpenStack Storlets Project<sup>2</sup>, Universitat Rovira i Virgili<sup>3</sup>, GridPocket<sup>4</sup>, Eurecom<sup>5</sup>

## Problem

- Large Volumes of Unstructured Data Objects
- Disaggregated Compute and Storage Clusters
- Costly Data Ingestion When Running SQL Queries

- Analytics Jobs Delegate SQL Tasks to The Storage
- The Storage Cluster Runs These Tasks To Cut Ingestion
- Framework To Orchestrate Cooperation between Analytics Jobs and Storage

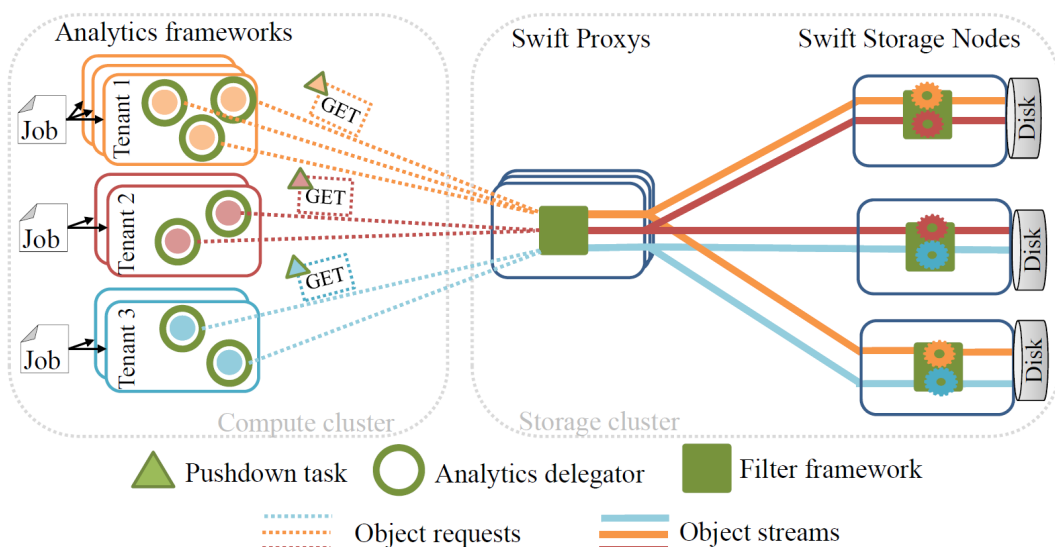
## Results

- Open-Source Prototype Implementation
- GridPocket Use-Case: Energy Data Management
- Evaluation on a 63-Machine Cluster with Zoe – Analytics on demand

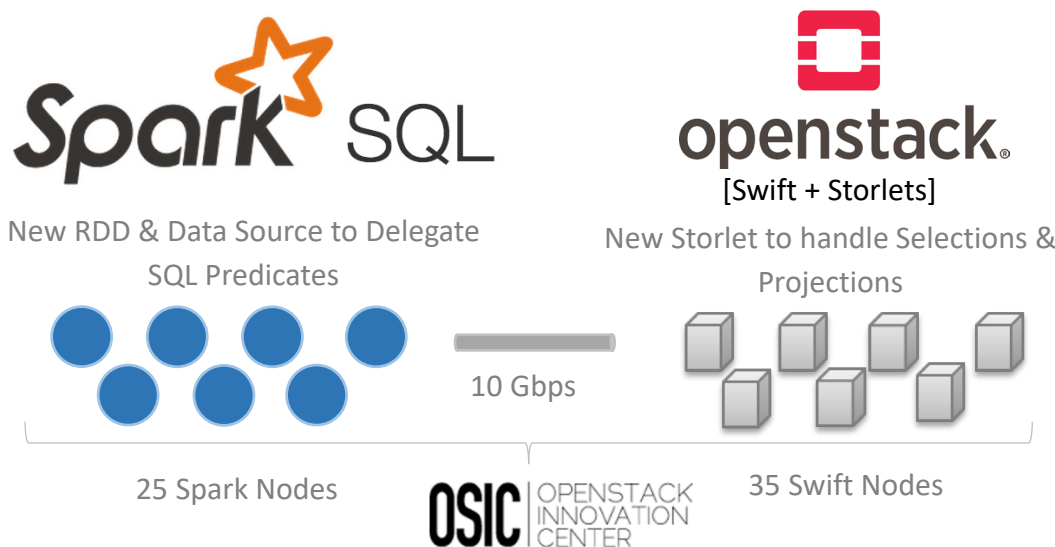
## Goals

## Platform & Experiments

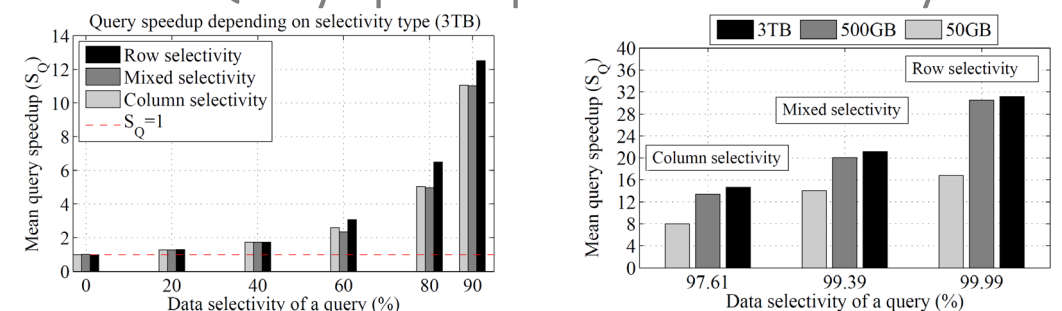
### Conceptual Design of Scoop



### Implementation/Deployment

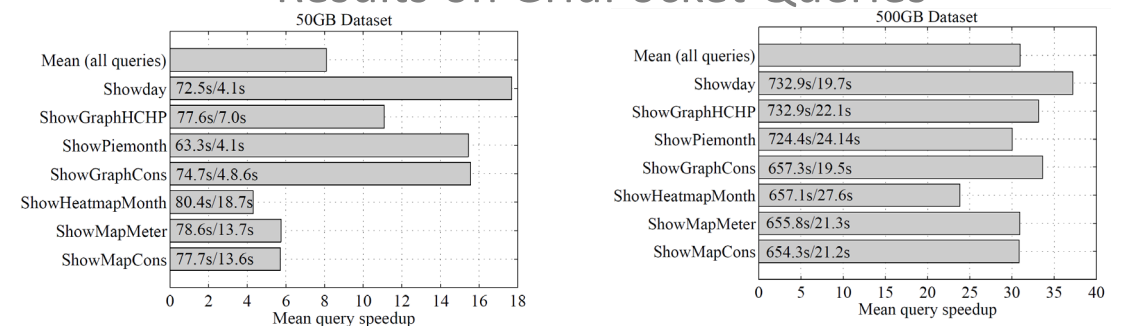


### Query Speedup vs Data Selectivity



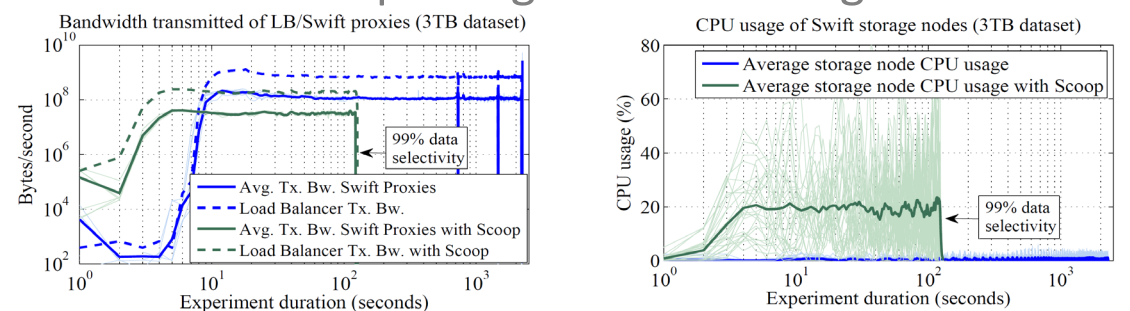
High Query Speedup: As Query Data Selectivity Grows, Job Speedup Exhibits a Non-linear Improvement

### Results on GridPocket Queries



Solving GridPocket Problems: Typical GridPocket Data Intensive Queries Show Speedups Up To x32

### Inspecting Resource Usage



Nice trade-off: We Avoid Consuming 99% Transfer BW by Using 20% of (Mostly Idle) Storage Compute Power



Find US!

<http://iostack.eu>  
[@iostackproject](https://twitter.com/iostackproject)  
<https://github.com/iostackproject>

