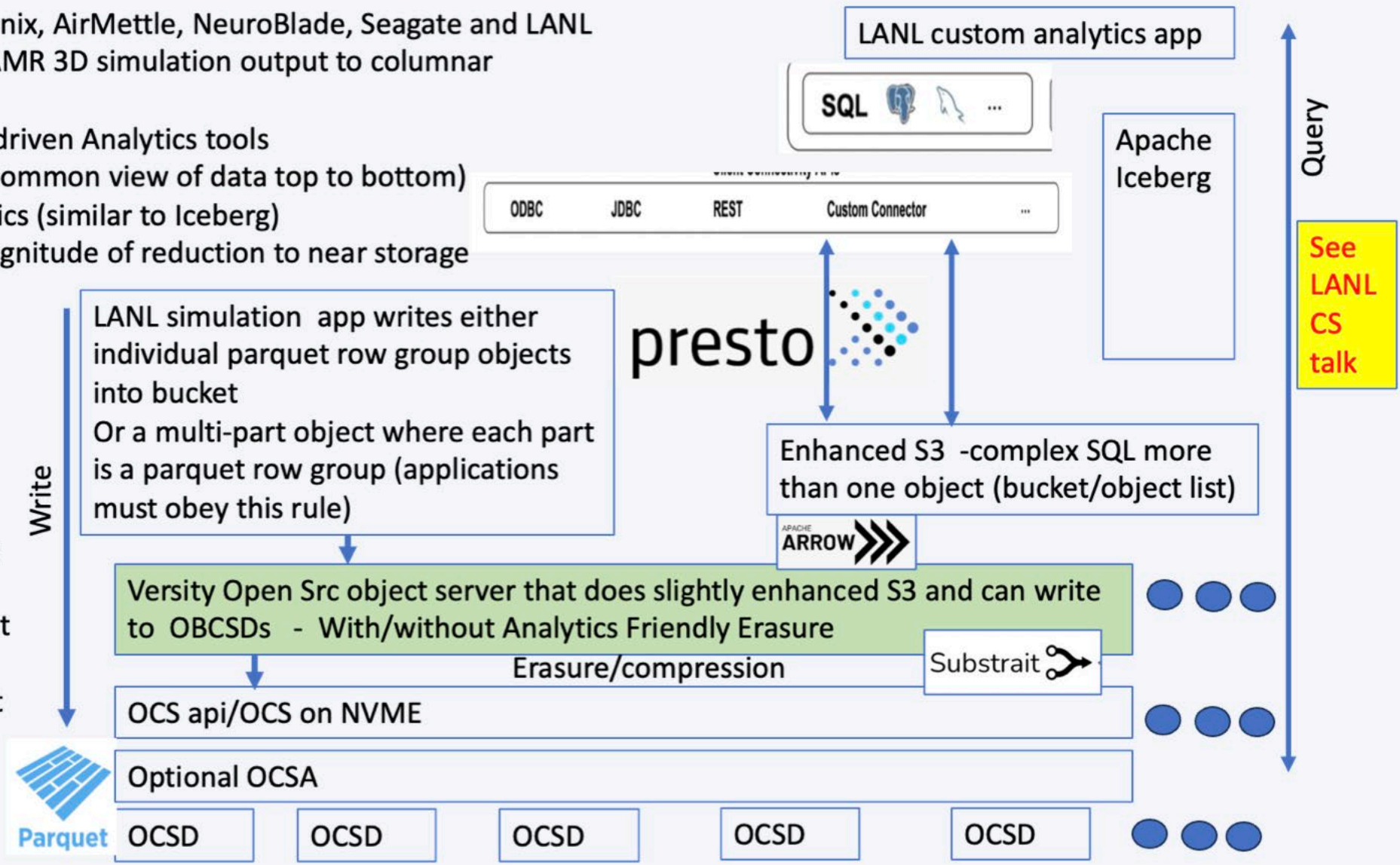Ben McClelland (Versity)
Gary Grider (LANL)
SDC 2023

# LANL use case of Versity Gateway - Object Computational Storage (OCS)

- Partnership: SK hynix, AirMettle, NeuroBlade, Seagate and LANL
- LANL grid-based AMR 3D simulation output to columnar Parquet
- Leverage Apache driven Analytics tools
- Leverage object (common view of data top to bottom)
- Extend S3 semantics (similar to Iceberg)
- Push orders of magnitude of reduction to near storage

- Big Data Analytics
  - You don't know what you are looking for, want compact model
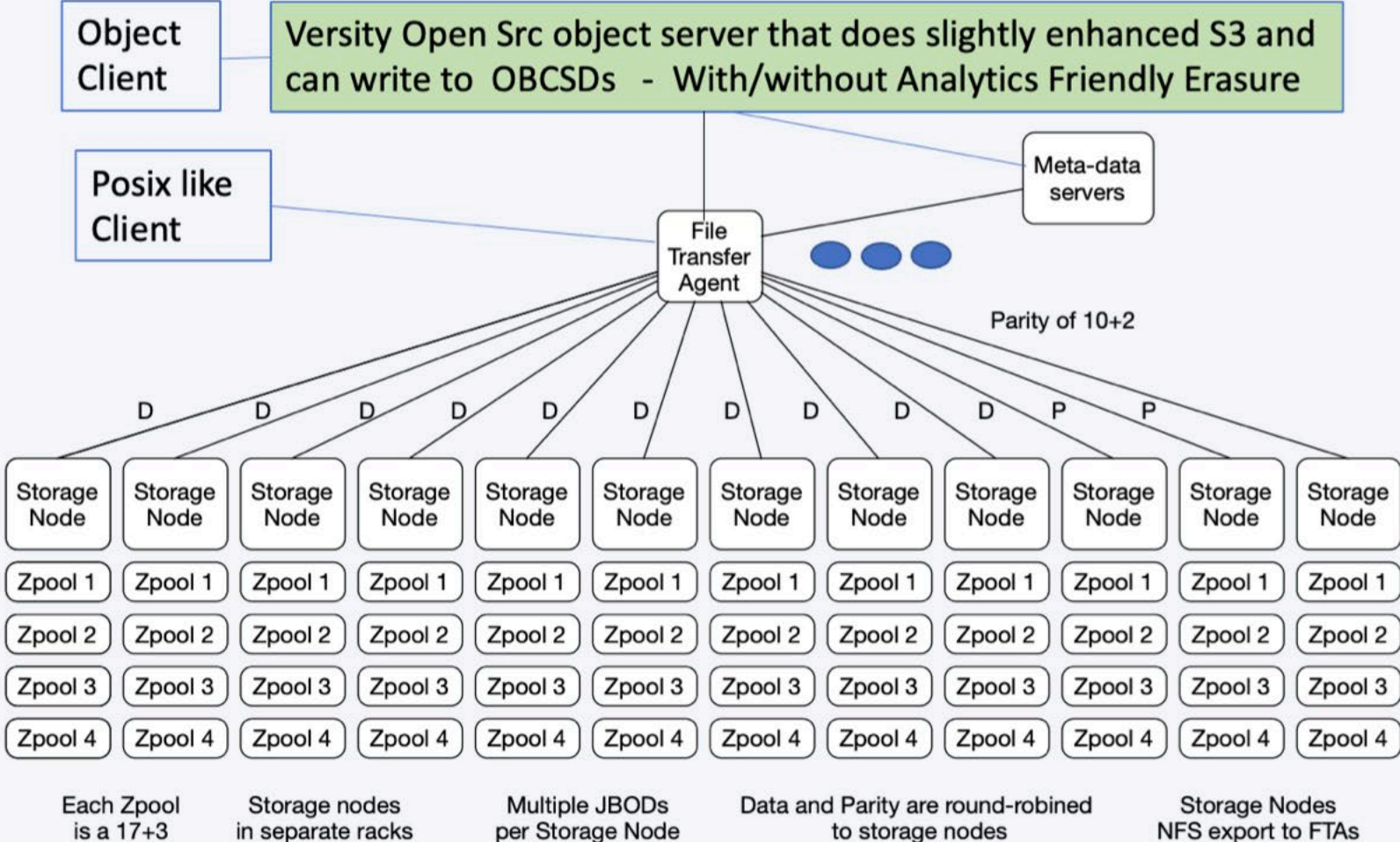  - You know what you are looking for but don't know where it is

**OCS Value Prop**

**Query**

**See LANL CS talk**

LANL custom analytics app

SQL ...

| ODBC | JDBC | REST | Custom Connector | ... |

Apache Iceberg

**presto**

LANL simulation app writes either individual parquet row group objects into bucket
Or a multi-part object where each part is a parquet row group (applications must obey this rule)

**Write**

Enhanced S3 -complex SQL more than one object (bucket/object list)

APACHE ARROW >>>

Versity Open Src object server that does slightly enhanced S3 and can write to OBCSDs - With/without Analytics Friendly Erasure

Substrait

Erasure/compression

OCS api/OCS on NVME

Optional OCSA

**Parquet**

| OCSD | OCSD | OCSD | OCSD | OCSD |

**Versity**

# LANL use case of Versity Gateway
# Object Access to Science Campaign Data

- **Leverage massive 2 tier erasured science campaign storage data lake from posix like users and from object users.**

- **Data is already in immutable files**

- **Hundreds of gigabyes/sec**

- **Stable for years**

- **Disk cool storage and even erasured tape in the future.**

Object Client

Versity Open Src object server that does slightly enhanced S3 and can write to OBCSDs - With/without Analytics Friendly Erasure

Posix like Client

Meta-data servers

File Transfer Agent

Parity of 10+2

| D | D | D | D | D | D | D | D | D | D | P | P |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Storage Node | Storage Node | Storage Node | Storage Node | Storage Node | Storage Node | Storage Node | Storage Node | Storage Node | Storage Node | Storage Node | Storage Node |
| Zpool 1 | Zpool 1 | Zpool 1 | Zpool 1 | Zpool 1 | Zpool 1 | Zpool 1 | Zpool 1 | Zpool 1 | Zpool 1 | Zpool 1 | Zpool 1 |
| Zpool 2 | Zpool 2 | Zpool 2 | Zpool 2 | Zpool 2 | Zpool 2 | Zpool 2 | Zpool 2 | Zpool 2 | Zpool 2 | Zpool 2 | Zpool 2 |
| Zpool 3 | Zpool 3 | Zpool 3 | Zpool 3 | Zpool 3 | Zpool 3 | Zpool 3 | Zpool 3 | Zpool 3 | Zpool 3 | Zpool 3 | Zpool 3 |
| Zpool 4 | Zpool 4 | Zpool 4 | Zpool 4 | Zpool 4 | Zpool 4 | Zpool 4 | Zpool 4 | Zpool 4 | Zpool 4 | Zpool 4 | Zpool 4 |

Each Zpool is a 17+3    Storage nodes in separate racks    Multiple JBODs per Storage Node    Data and Parity are round-robined to storage nodes    Storage Nodes NFS export to FTAs

Versity

**What We Do:**

# Manage Large Unstructured Data Collections at Low Cost

- Software-Defined Storage Platform
- Mass Storage & Large Archive

Versity

# Global Installed Base

**AMERICAS**

**EMEA**

**APJ**

## Total Data

2 Exabytes Data

80 Sites
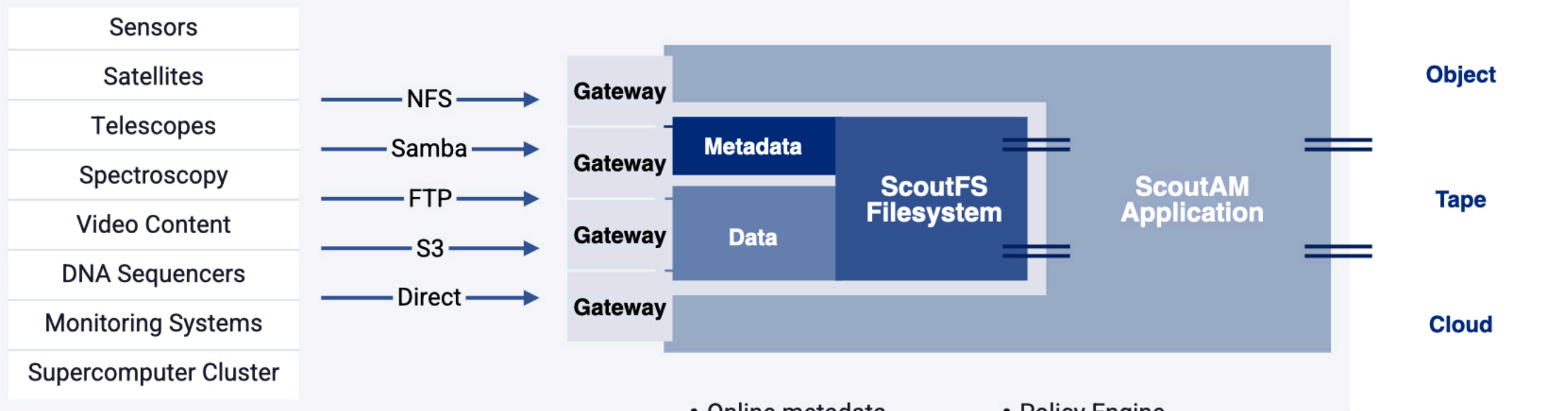
Versity

# Versity Gateway

- Scalable S3 Service
- High Performance
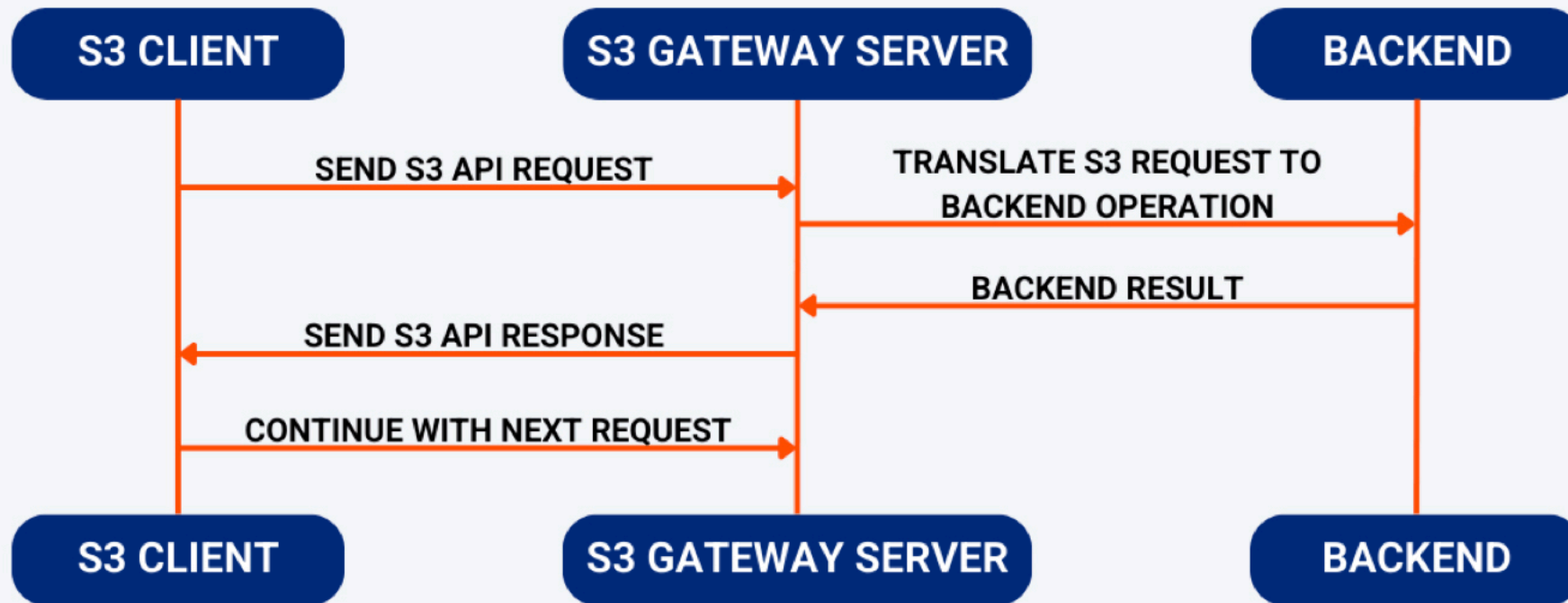- Modular Backend Support
- Flexible Open Source Licensing (Apache 2.0)

# Scale-Out Archive Manager: ScoutAM

## A Modern Data Management Platform



Sensors
Satellites
Telescopes
Spectroscopy
Video Content
DNA Sequencers
Monitoring Systems
Supercomputer Cluster

NFS
Samba
FTP
S3
Direct

Gateway

Metadata

Data

ScoutFS Filesystem

ScoutAM Application

Object

Tape

Cloud

- Online metadata
- Indexed attributes
- Scalable namespace
- Policy Engine
- Scheduling
- Parallel Data Movement

Versity

# Versity Gateway

# Modular Backend

Current Backends:

- POSIX - filesystem with xattr support
- ScoutFS - (https://github.com/versity/scoutfs)

Open to collaboration on new backends

Versity

# ScoutFS Optimized Backend

## Optimized Writes

- Multi-part upload segments are written once to the underlying storage then
- Parts of data upload are combined into a single file with a system call.
- **Eliminates one full read/write cycle**
- Speeds up the process and reduces overall upload time.

## S3 Glacier Mode

- Cold storage feature for data archiving and long-term data retention.
- Data may be stored on local tape systems using Versity's ScoutAM platform.
- Manage storage costs more effectively by accessing low cost storage.



Versity

# Project Layout



📁 auth

📁 backend

📁 cmd/versitygw

📁 integration

📁 s3api

📁 s3err

📁 s3event

📁 s3log

📁 s3response

Versity

# Project Layout



→ Multi-Tenant
IAM account integrations

Extensible to support new authentication systems

Currently Supported:
Local accounts

auth
backend
cmd/versitygw
integration
s3api
s3err
s3event
s3log
s3response

Versity

# Multi-Tenant

root account defined with cli flags/env vars
admin/user accounts stored in IAM service

|  | See All Buckets | Create New Buckets | Create New Users | Assign Bucket Ownershi | See Only Owned Buckets |
|---|---|---|---|---|---|
| **root** | X | X | X | X | |
| **admin** | | X | X | X | X |
| **user** | | | | | X |

Versity

# Project Layout



auth

backend →

cmd/versitygw

integration

s3api

s3err

s3event

s3log

s3response

## Storage Backends

Extensible to support new storage systems

Currently Supported:
POSIX
ScoutFS

Versity

# Backend

- Embed backend.BackendUnsupported
- Implement backend.Backend methods as needed to satisfy interface
- Any unimplemented will return s3err.ErrNotImplemented back to client

```go
package mystorage

type MyStorage struct {
    backend.BackendUnsupported
}
```

```go
type Backend interface {
    // bucket operations
    ListBuckets(_ context.Context, owner string, isRoot bool) (s3response.ListAllMyBucketsResult, error)
    …

    // multipart operations
    CreateMultipartUpload(context.Context, *s3.CreateMultipartUploadInput) (*s3.CreateMultipartUploadOutput, error)
    …

    // standard object operations
    PutObject(context.Context, *s3.PutObjectInput) (string, error)
    …

    // special case object operations
    RestoreObject(context.Context, *s3.RestoreObjectInput) error
    …

    // object tags operations
    GetTags(_ context.Context, bucket, object string) (map[string]string, error)
    …
}
```

Versity

# ScoutFS/Glacier Example

ScoutFS/ScoutAM - Archiving Filesystem
- POSIX filesystem
- Supports offline files where data only resides on tape/mass storage
- Automatic recall when data requested

When in Glacier Emulation Mode,
scoutfs backend enables following:

HEAD offline object returns
storage-class: GLACIER
x-amz-restore: (transition state)

GET of offline file returns
Invalid Object State

Restore Object
triggers backend data recall

```go
type Backend interface {
    …

    // special case object operations
    RestoreObject(context.Context, *s3.RestoreObjectInput) error
    …

}
```

Versity

# Project Layout



auth

backend

cmd/versitygw → CLI command support

integration

s3api

s3err

s3event

s3log

s3response

Versity

# Project Layout

📁 auth

📁 backend

📁 cmd/versitygw

📁 integration ———→ Integration Tests

📁 s3api

📁 s3err

📁 s3event

📁 s3log

📁 s3response

GitHub CI runs test suite against POSIX for PRs

Manually run test suite against any storage backend

Versity

# Project Layout

| | |
|---|---|
| 📁 | auth |
| 📁 | backend |
| 📁 | cmd/versitygw |
| 📁 | integration |
| 📁 | s3api |
| 📁 | s3err |
| 📁 | s3event |
| 📁 | s3log |
| 📁 | s3response |

→ ## Frontend API Handlers

Go Fiber framework for HTTP(s) handlers
Supporting S3 and Admin API

Versity

# Project Layout

S3 errors contains all errors that clients would expect from AWS S3

Non-s3err types returned from backend treated as Internal Server Error

📁 auth

📁 backend

📁 cmd/versitygw

📁 integration

📁 s3api

📁 s3err → S3 error response types

📁 s3event

📁 s3log

📁 s3response

Versity

# Project Layout

Same event structure as AWS

Extensible event integrations, currently supported:
- NATS
- Kafka

📁 auth

📁 backend

📁 cmd/versitygw

📁 integration

📁 s3api

📁 s3err

📁 s3event   ⟶   S3 Event Notifications

📁 s3log

📁 s3response

**Versity**

# Project Layout

Same event structure as AWS

Extensible log integrations,
currently supported:
- file
- webhook

📁 auth

📁 backend

📁 cmd/versitygw

📁 integration

📁 s3api

📁 s3err

📁 s3event

📁 s3log → **S3 Server Access Logs**

📁 s3response

Versity

# Project Layout

📁 auth

📁 backend

📁 cmd/versitygw

📁 integration

📁 s3api

📁 s3err

📁 s3event

📁 s3log

📁 s3response

AWS defines specific xml structure for some requests/responses

For these requests, use well specified struct field members and xml struct tags
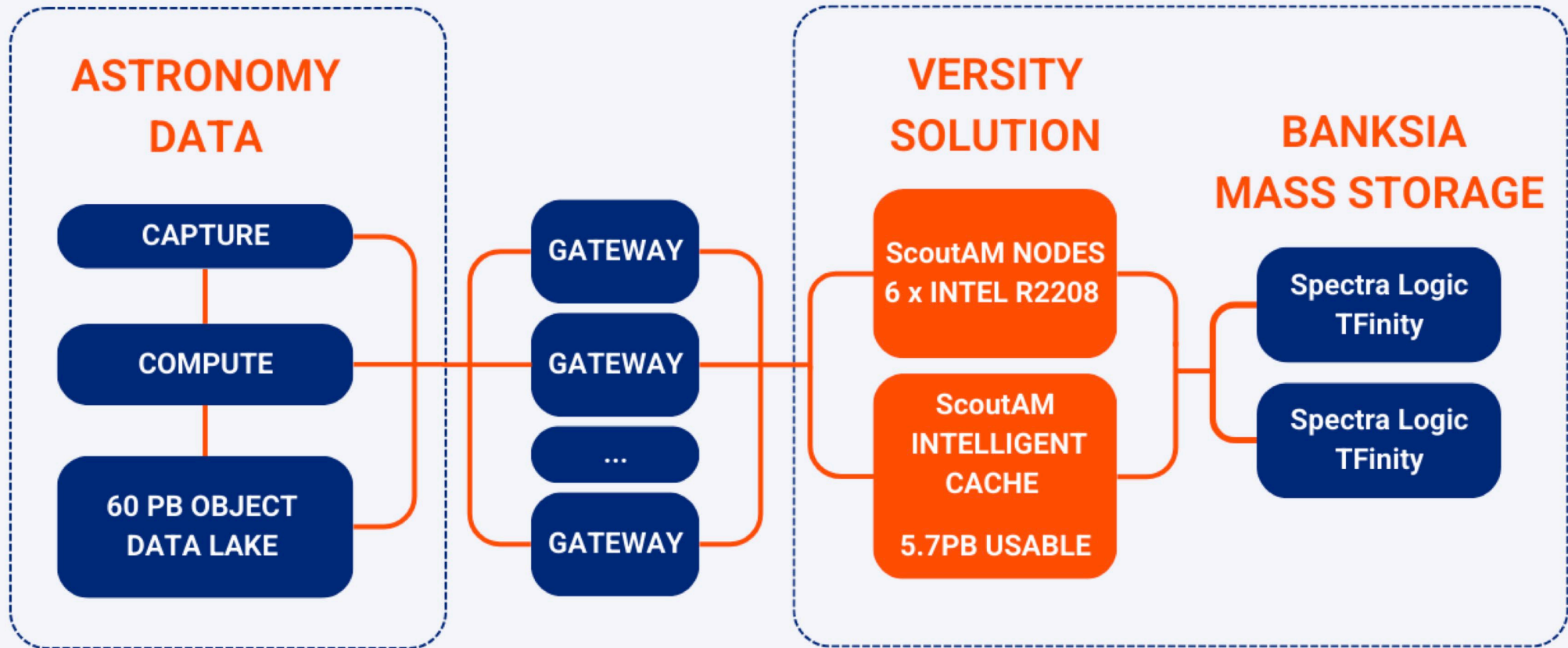
→ S3 Response formatting

Versity

# Gateway enables S3 workloads for Pawsey's Large Scale Archive

- Versity's ScoutAM managing tape-based archive, Banksia
  - 150 PB tape
  - 34 tape drives
  - 2 tape libraries
  - 5 PB cache
- Incoming data is S3 via gateway
- Cluster load balanced

# Pawsey Solution Architecture



**ASTRONOMY DATA**
- CAPTURE
- COMPUTE
- 60 PB OBJECT DATA LAKE

GATEWAY
GATEWAY
...
GATEWAY

**VERSITY SOLUTION**
- ScoutAM NODES 6 x INTEL R2208
- ScoutAM INTELLIGENT CACHE 5.7PB USABLE

**BANKSIA MASS STORAGE**
- Spectra Logic TFinity
- Spectra Logic TFinity

Versity

# Developed in the Open on GitHub

https://github.com/versity/versitygw

- Bugs/Features tracked in GitHub Issues
- Documentation in project wiki
- CI using GitHub Actions

Versity

# Thank You

info@versity.com

@versitysoftware

https://github.com/versity/versitygw

Versity