

NeuroData's Open Data Cloud Ecosystem

Introduction

- Neuroscience imaging produces datasets too large for traditional approaches
- We have deployed a comprehensive suite of tools to support
 - Hosting hundreds of datasets
 - Fast visualization
- Scale to petabytes
- Difficult for end users and maintenance
- We explore different methods to manage data
- AWS Open Data bucket for hosting public data
- Data stored in flat file system
- Flywheel for managing derivatives and metadata
- Flat file system is faster and less expensive
- Need to build additional missing services



https://neurodata.io/ndcloud/



- Neuroglancer: visualization of image volumes, annotations, segmentations, skeletons, meshes¹
- BOSS: hosting datasets using cloud infrastructure²
- ndweb: web application which allows custom interactions with BOSS³ - channel selection, TIFF cutouts, custom neuroglancer links
- Ndex: exchange data w/ neurodata BOSS programmatically



2: <u>https://api.boss.neurodata.io</u> 3: <u>https://ndwebtools.neurodata.io</u>

Neuroglancer deployment https://zbrain.viz.neurodata.io

BOSS

Design

- Built for MICrONS project¹
- Traffic goes through endpoint(s)
- Mix of microservices and monolith,
- making it complex 1: https://www.iarpa.gov/index.php/research-programs/microns

RDS, MySQL Instance ata model objects & Read Cache / Write Buffer

Cons

- \$\$\$
- Central point of failure
- Requires cache to be performant
- (Re)deployments are *tricky*
- Downsampling could be faster
- Limited LIMS system
- No place to put derived data (meshes, skeletons, computations)
- Difficult for end users to use
- Bugs are difficult to track down due to interconnected services
- Requires login
- Coordinates have to be positive

Static storage for image volumes

Neuroglancer precomputed format (created by Jeremy Maitin-Shepard) is a flat file storage specification that has several advantages, and some disadvantages, over a centralized system such as BOSS. Specification is open source and available at the following link: https://github.com/google/neuroglancer/b

x 644, y 286, z 61 x 469, y 276 2e+1 X 2zbrain_atlas 13 X +

13: Forebrain: Diencephalon

ob/master/src/neuroglancer/datasource/p ecomputed/README.md

Pros

- Flat file format is flexible (hosting)
- Custom chunk size per dataset
- Direct access to data
- No endpoint
- Fewer things can go wrong
- \$



- Cache Database (bounding box lookups)
- LIMS
- Infrastructure

Benjamin Falk, Joshua T. Vogelstein | Johns Hopkins University, Kavli Neuroscience Discovery Institute



Pros

- Provides everything
- Authentication
- Scales reasonably well
- Downsample inexpensive
- Timeseries support
- Database indexing of
- segmentation locations for
- easy lookup



Lacking

- Authentication
- No timeseries support

Benchmarks



- Precomputed can be as fast as BOSS

Conclusions and Future Directions

- complications of running the BOSS.
- datasets requires minimal infrastructure.
- resources and infrastructure, which has yet to be finalized.
- working with precomputed volumes.
- organization, and pipeline orchestration.
- 1: https://aws.amazon.com/opendata/
- 2: https://github.com/seung-lab/cloud-volume
- 3. https://github.com/seung-lab/igneous/





We do not observe a large benefit of cache for BOSS in this test

- Cube size has a strong effect on transfer speed for both BOSS and precomputed

 Neuroglancer precomputed format allows hosting of large image volumes using distributed storage provided by cloud providers, without the

• Availability of AWS Open Data¹ and cheap storage pricing, hosting these

• Pipelines for downsampling, registration, statistical analysis and computer vision can operate on these precomputed volumes, but they take additional

• Cloud-Volume² and Igneous³ provide excellent command line utilities for

• Commercial solutions (e.g. Flywheel) could provide the needed LIMS frontend,

