STScI | SPACE TELESCOPE SCIENCE INSTITUTE

EXPANDING THE FRONTIERS OF SPACE ASTRONOMY

# Server-Side Analytics

## Science Platforms

Arfon & Mike

# Some definitions

# Some definitions



MAST Portal &
MAST 'Classic'

Focus of this presentation

MAST APIs

# Common technologies, many implementers

# Highlights of what we're building

- Cloud-hosted copy of all HST public data
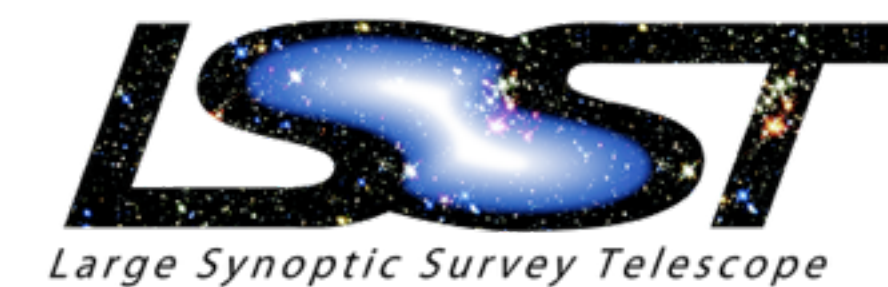
- (Live) JupyterLab environment with *some* compute/storage

- Collection of Docker containers installed with common tools

# Amazon Web Services: Public Dataset Program



- ~120TB public HST data for ACS, COS, STIS, WFC3, FGS
- Range of high-impact datasets
- Hosted in cloud - 'highly available'
- Enable new types of data analyses
- Hosted at no cost to STScI/NASA

- Data is hosted in an S3 region (for 'free')
- Conditions of inclusion in program: make the data useful:
  - An AMI with a demonstration of how to use the public dataset must be provided
  - AWS recover costs by making access to the data free from AWS services (EC2), making it cost effective for researchers to buy AWS computing time
  - Enables new types of analyses

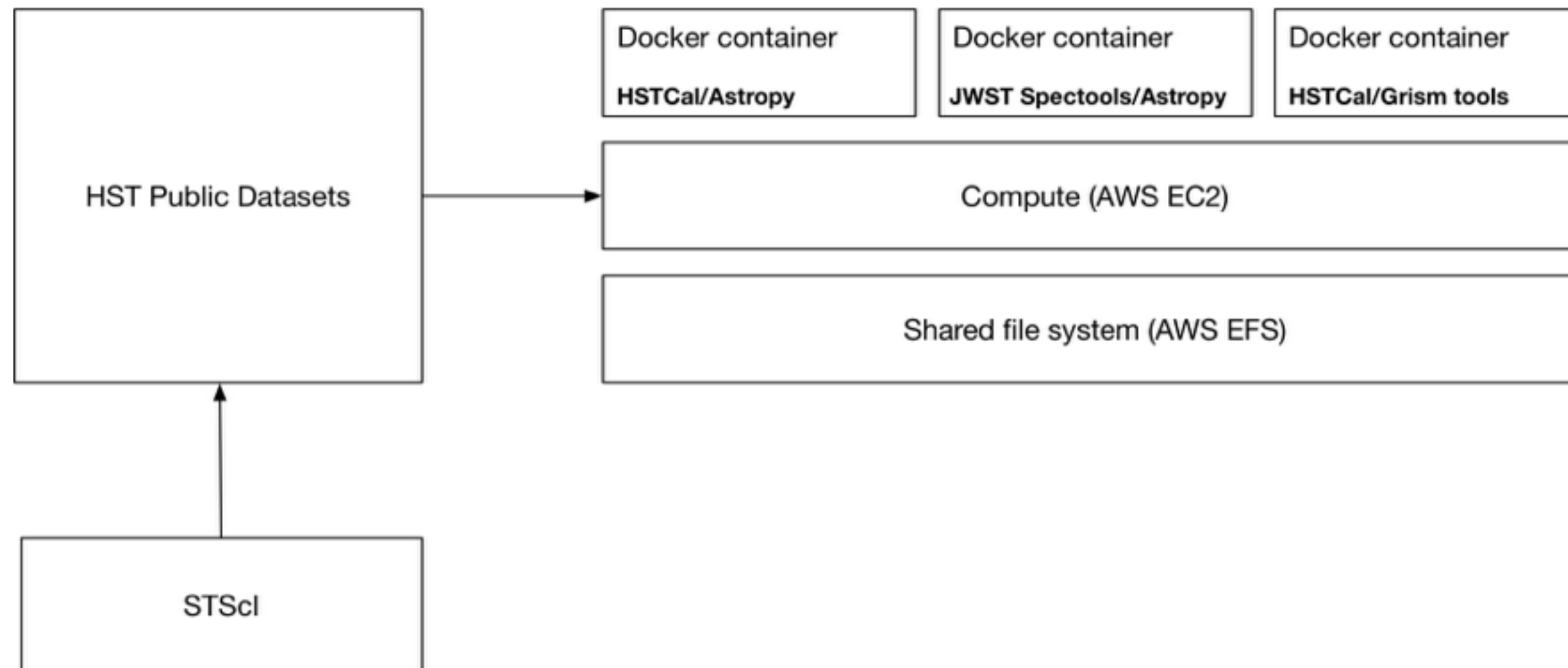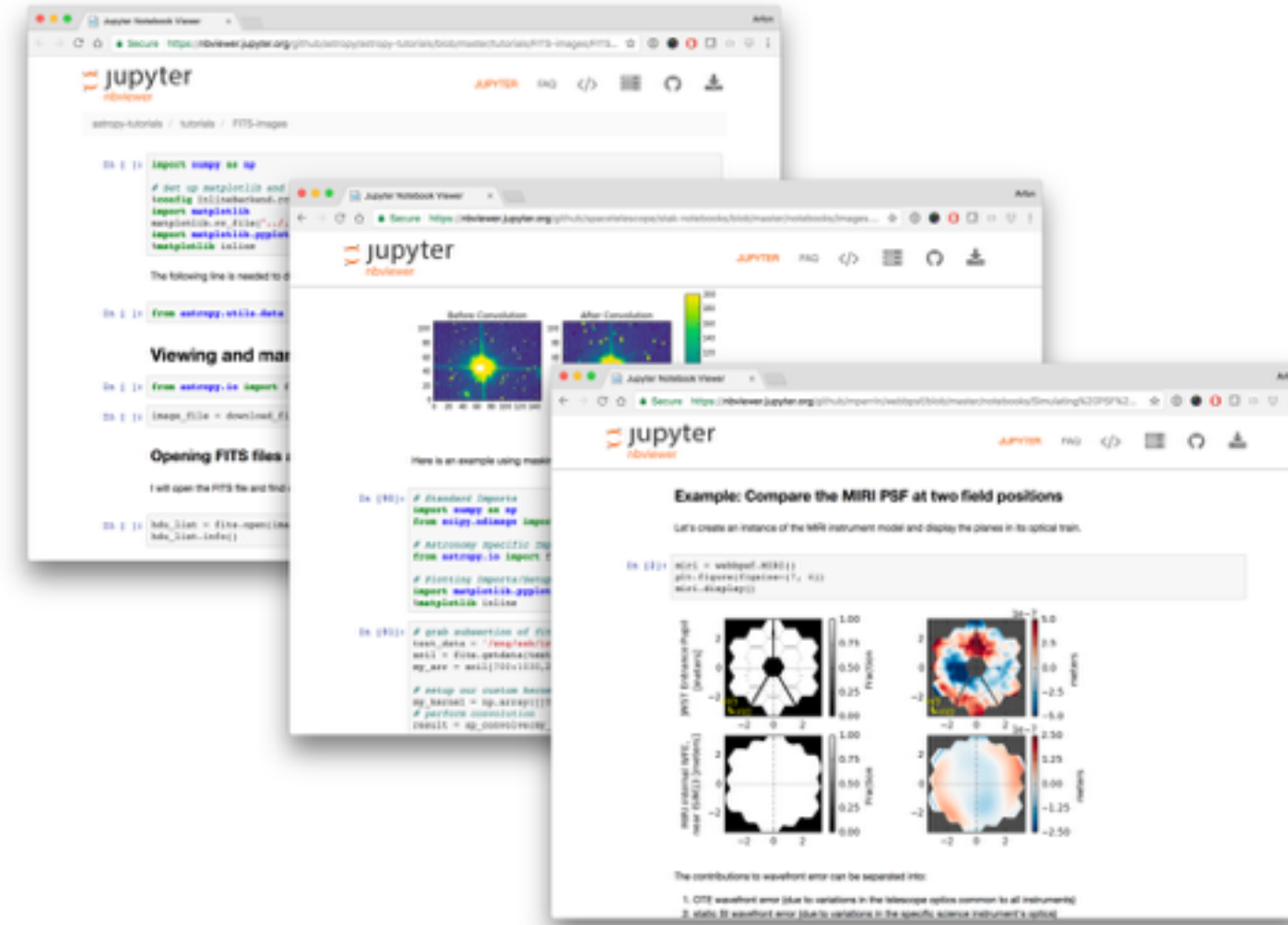# JupyterLab environment



- Interactive computing environment

- Where most development work is going from the core Jupyter team

- Works with community tools (e.g. Astropy)

STScI | SPACE TELESCOPE SCIENCE INSTITUTE

# Containers for different environments

# Technical details:

- JupyterHub, a multi-user Hub, spawns, manages, and proxies multiple instances of the single-user Jupyter notebook server.
- JupyterLab frontend provides notebook server, file management, and a terminal shell
- Using Docker to containerize science environments, allows a verified computing environment to be instantiated rapidly.
- Containers are versioned providing precise reproducibility of the research environment
- AWS computing resources scale with user load, providing good cost efficiency
- Container orchestration provides high availability, healing the cluster when there are hardware failures

# Core technical challenges

- Creating containers with pipeline/common software stacks
- Managing the cloud environment well:
  - User quotas (storage, compute etc.)
  - User storage (home directories), backups
  - Scalable, highly-available infrastructure (with cost caps/alerts)
- Relatively few large-scale JupyterHub deployments on AWS
- Inexperience of STScI with commercial cloud

STScI | SPACE TELESCOPE SCIENCE INSTITUTE

# Community coordination (i)



Iva++

# Community coordination (ii)

**Questions for the MUG (i)**

- What problems do you foresee with our approach?
  - Who is this useful for?
  - What are sensible defaults for a service like this?
  - Does staging data in commercial cloud confuse our value proposition as a NASA archive?
  - Is asking the community to use AWS a problem? How can we lower the barrier to entry?

# Questions for the MUG (ii)

- What extensions should we be thinking about?
  - Other missions (TESS?)
  - Other functionality? (e.g. batch processing)
  - Joint processing? (LSST/WFIRST/Euclid/JWST)
  - Accelerating JWST Early Release Science?
  - User support environment?
  - GO data delivery & archive proposals?
  - WFIRST (re)processing in the High Level Processing Partition