

# TERRA-REF Pipeline: A Containerization Story

Max Burnette

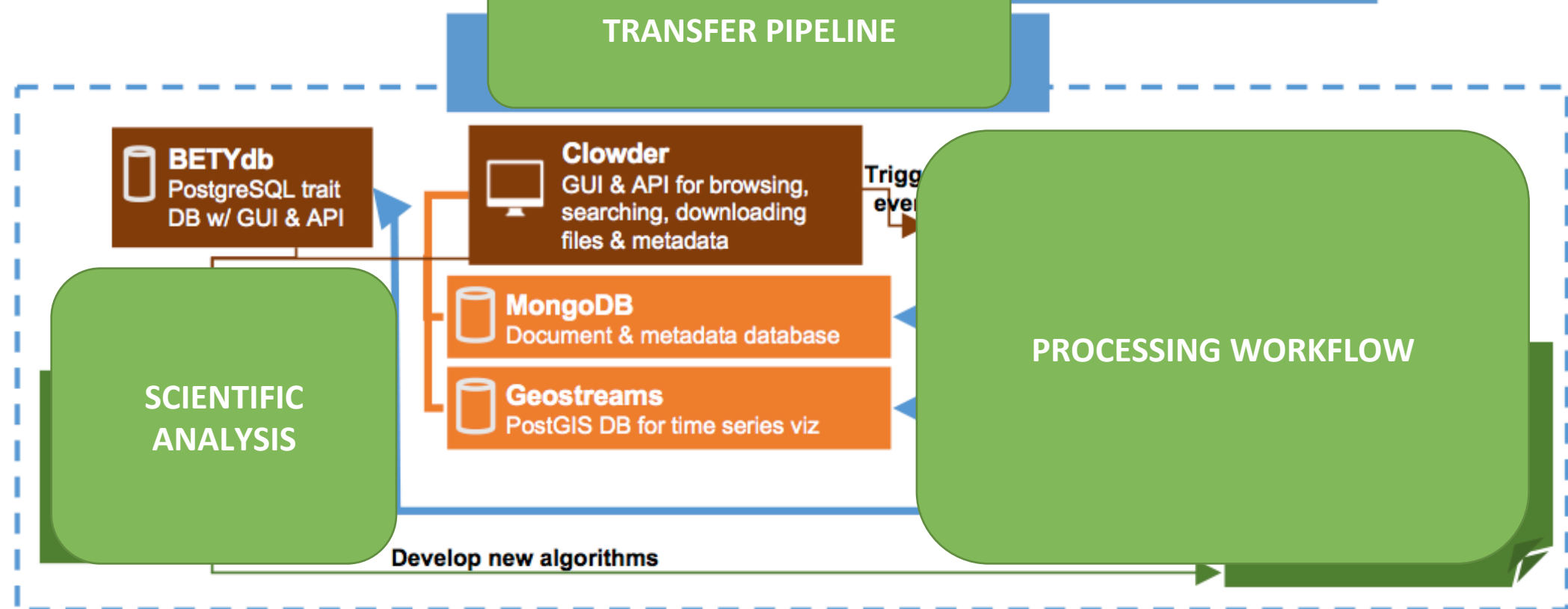
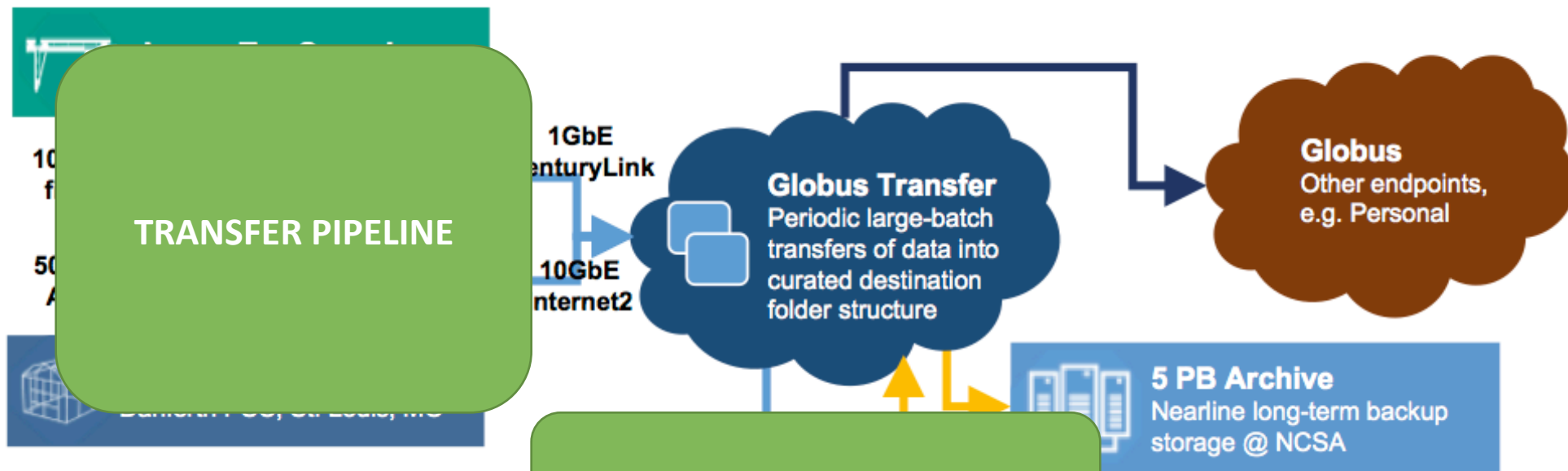
August 14, 2017



[terraref.org](http://terraref.org)

# outline

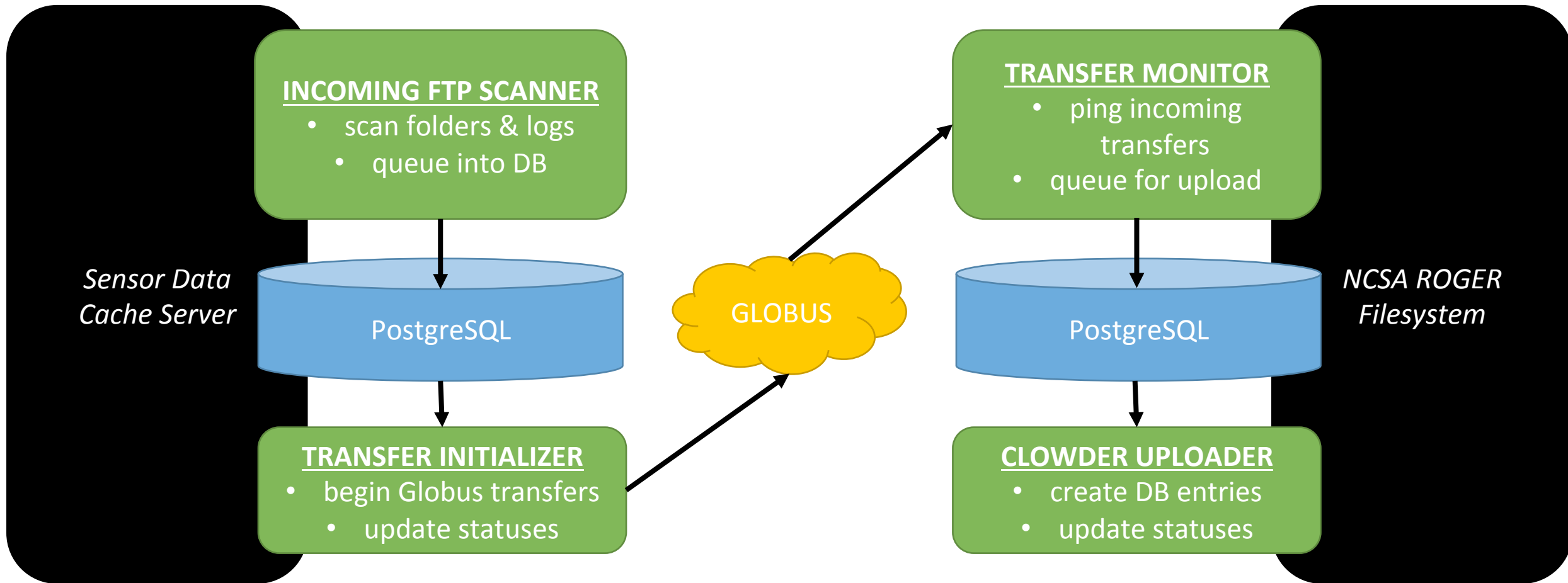
- Brief project overview
- Containers I - **transfer pipeline**
- Containers II - **processing workflow**
- Containers III - **scientific analysis**
- Q/A?



# TERRA-REF project goals

- **transfer pipeline**
  - move data from Arizona to Illinois efficiently
- **processing workflow**
  - handle large variety of sensor data in standardized scalable way
- **scientific analysis**
  - make tools, software & data easily accessible without large downloads

# containers I - transfer pipeline



# containers I - transfer pipeline

## THE GOOD

- **version/environment control**
  - avoid dependency desynchronization
  - hardware can change with minimal impact
  - same setup between development/testing/production
- **minimize downtime**
  - install & deploy container with software updates and new libraries with effectively zero switchover time
  - administration on deployment side is easy

INCOMING FTP SCANNER  
TRANSFER INITIALIZER  
TRANSFER MONITOR  
CLOUDER UPLOADER

# containers I - transfer pipeline

## THE BAD

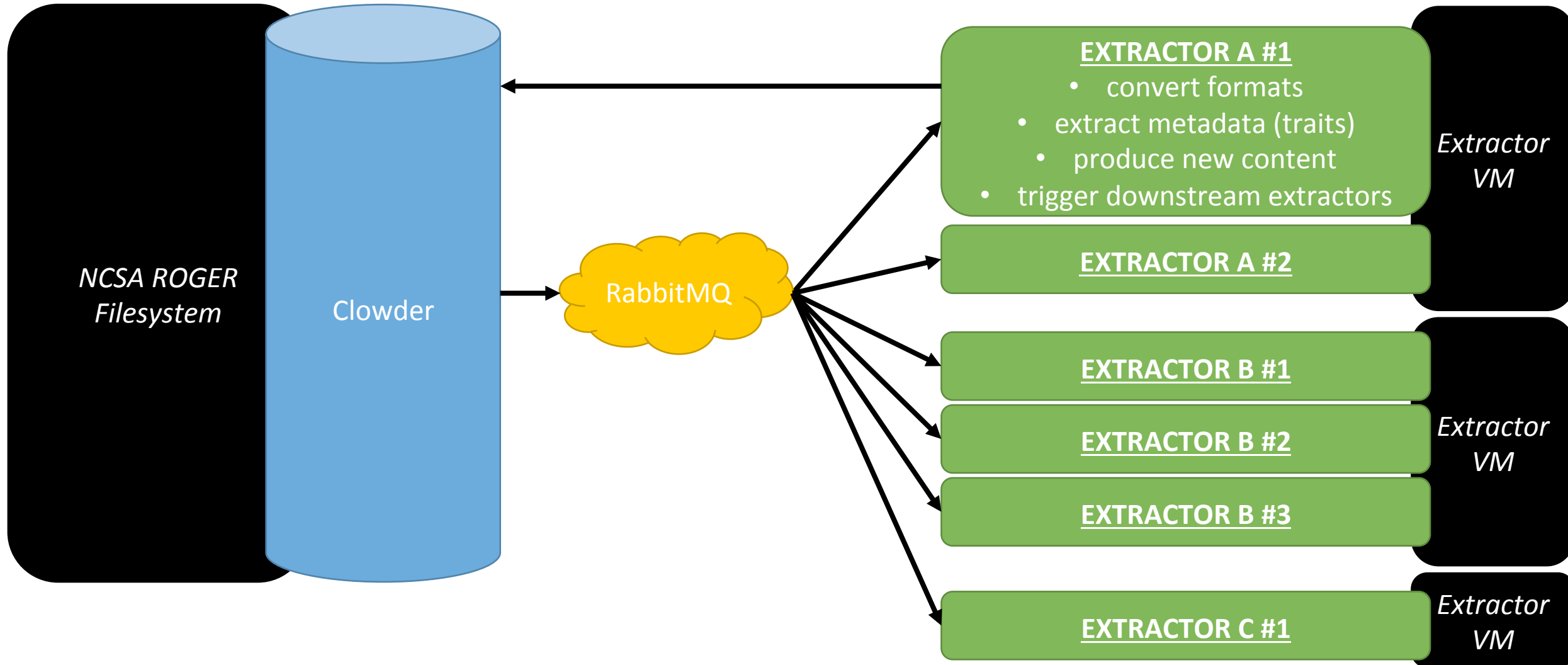
- **one more deployment step**
  - hotfixes and such become more strict

## THE UGLY

- **permission control on mounted storage**
  - complex user/group rules need special accommodation
  - users with permission can run container and lose permission
  - can require manhandling of UIDs and GIDs of user

INCOMING FTP SCANNER  
TRANSFER INITIALIZER  
TRANSFER MONITOR  
CLOUDER UPLOADER

# containers II - processing workflow





# containers II - processing workflow

## THE GOOD

- **highly scalable**
  - dynamically spin up or down full stacks instantly
  - separate temp space and logging\*

### EXTRACTOR A #1

- convert formats
- extract metadata (traits)
- produce new content
- trigger downstream extractors

## THE UGLY

- **port exposure to outside world**
  - APIs and other HTTP services require additional forwarding
- **making sure logs aren't lost**
  - some container environments (Docker) will happily overwrite

# containers II - processing workflow

## THE BAD

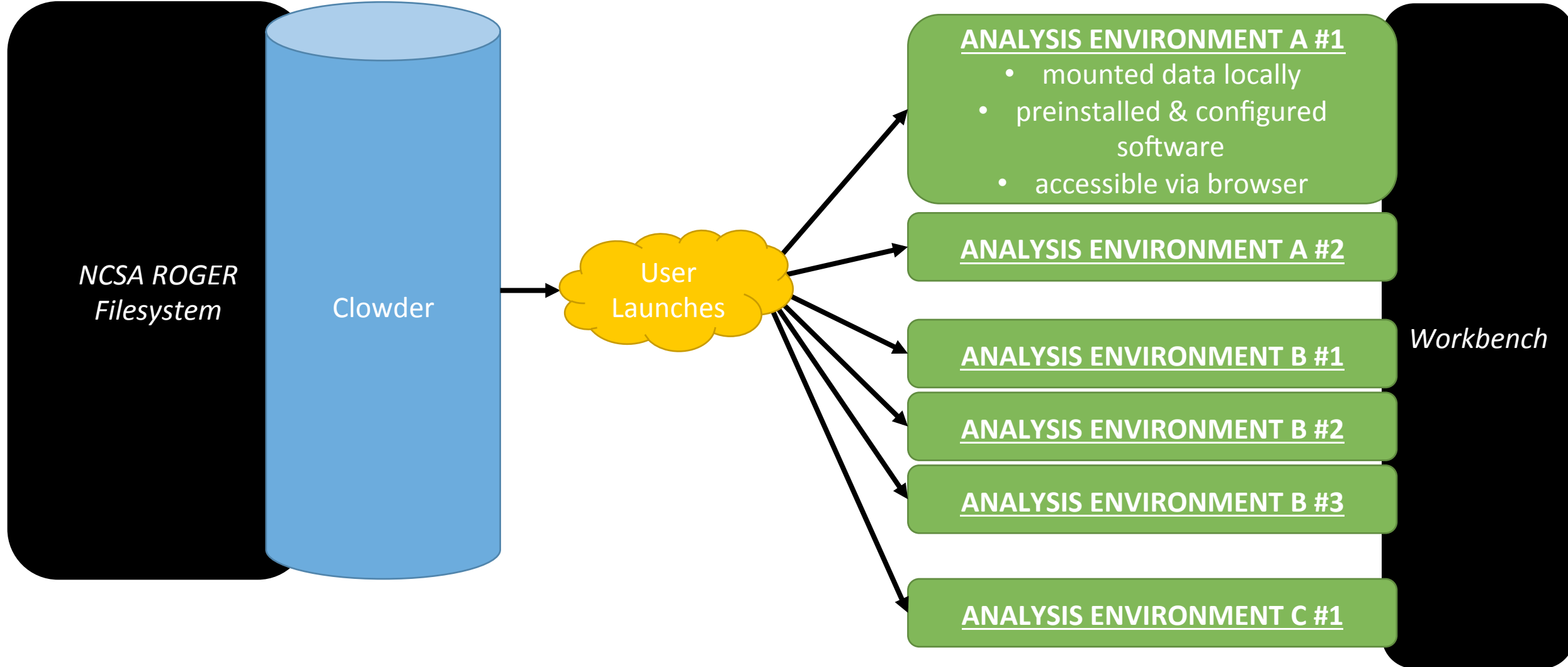
- **recreating external file structures**

- some code written against expected paths
- even mounted directories must replicate those, or account for them
- “local” files will not appear local unless correctly configured

### EXTRACTOR A #1

- convert formats
- extract metadata (traits)
- produce new content
- trigger downstream extractors

# containers III - scientific analysis



# containers III - scientific analysis

## THE GOOD

- **easy to use**
  - no installation or tracking versions
  - identical environments across users
- **sharable**
  - simply sharing URLs and credentials is sufficient

## THE BAD

- **ephemeral by nature**
  - Expectations of saving data between sessions?
  - How to move data in and out of a remote container?

### ANALYSIS ENVIRONMENT A #1

- mounted data locally
- preinstalled & configured software
- accessible via browser

# Q/A?

- Brief project overview
- Containers I - transfer pipeline
- Containers II - processing workflow
- Containers III - scientific analysis
- “huh?”