

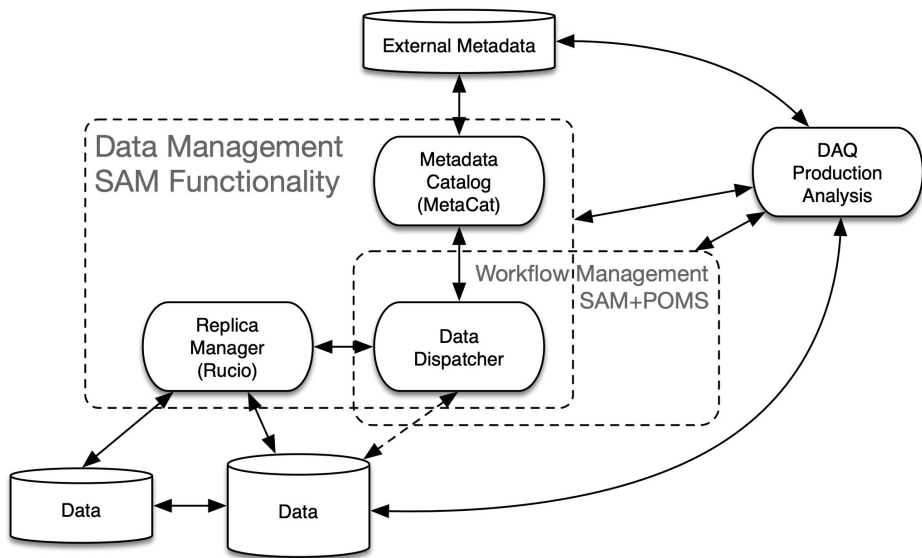


MetaCat and Data Dispatcher

Igor Mandrichenko, FNAL

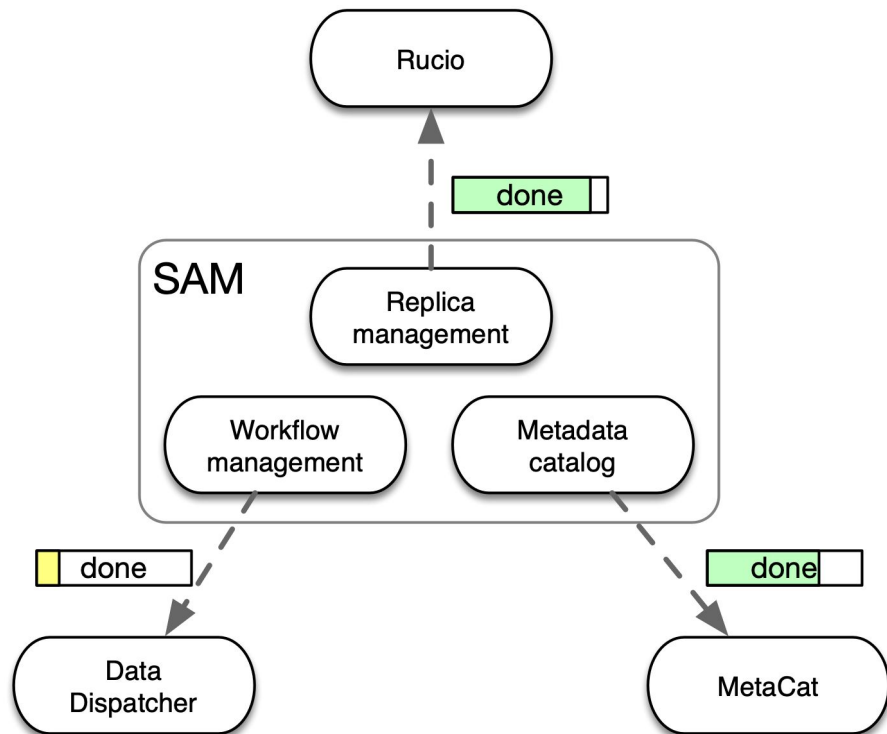
Essentials for DUNE Computing, December 14, 2021

Data Management Big Picture



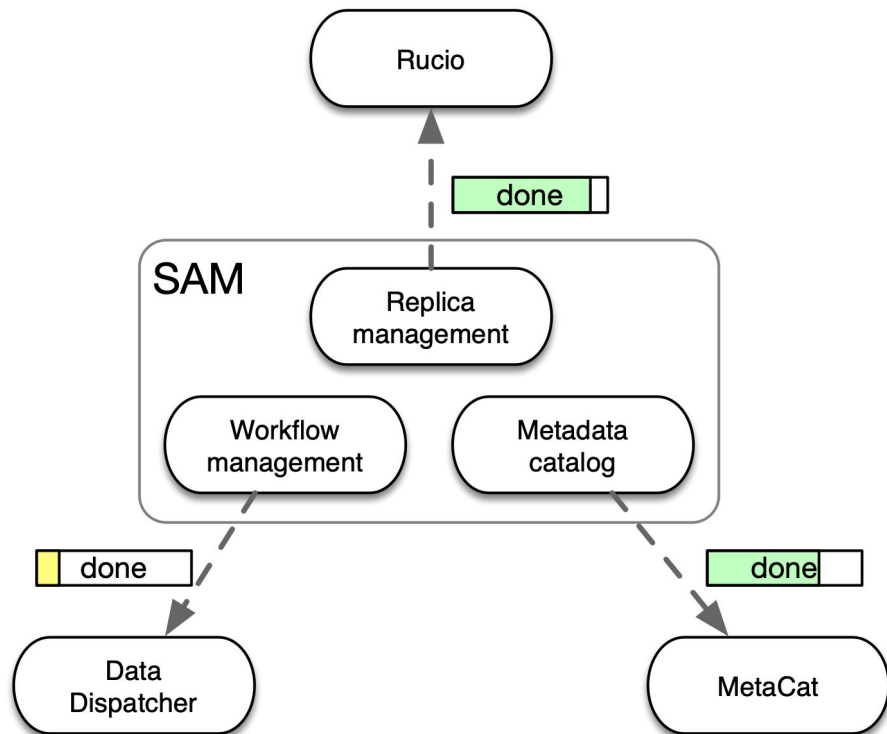
- Metadata catalog
 - What kind of information stored in the file ?
 - What files should be processed (analyzed) ?
- Replica Management
 - Where are copies of the files ?
- Workflow Management
 - File processing coordination
- Currently, SAM+POMS does all
- Goals:
 - Use Rucio as replica manager
 - Decompose SAM
 - Global WMS

Migration from SAM



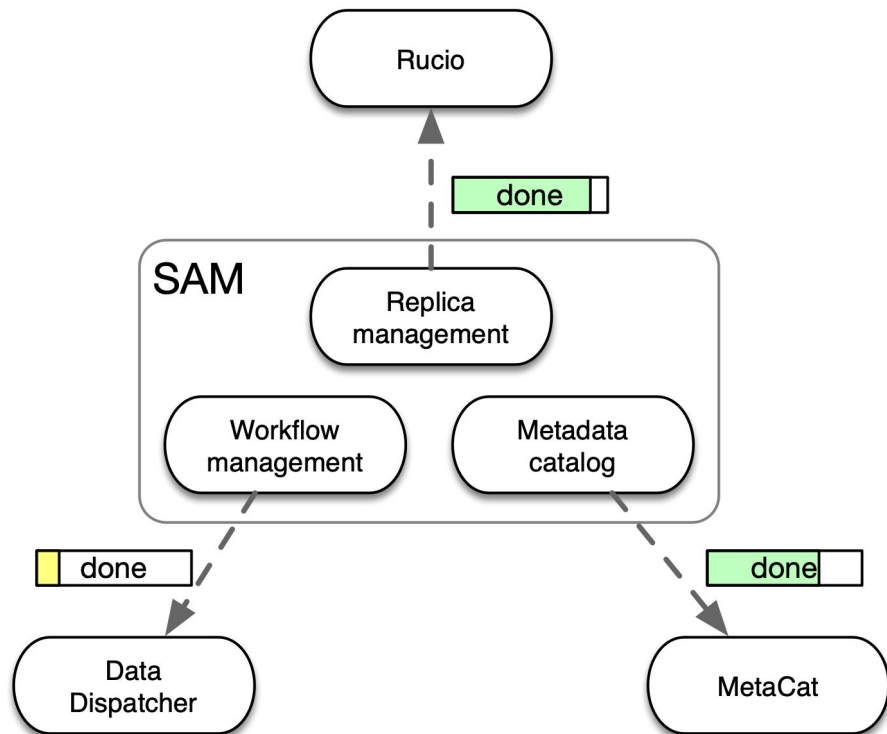
- Replica management
 - Almost done
 - Still uses SAM
 - Replica management delegated to Rucio
 - Virtual location

Migration from SAM



- Metadata catalog
 - MetaCat
 - Development mostly done
 - Volume tests:
 - ProtoDUNE SAM ~10M files
 - NOvA SAM ~200M files, 5.5B name/value pairs
 - External data access test - ProtoDUNE Runs DB
 - Need more feedback !

Migration from SAM



- Data Dispatcher *as piece of SAM*:
 - Requirements and functionality understood
 - Prototype exists
- Interaction with global WMS
 - *Unclear*

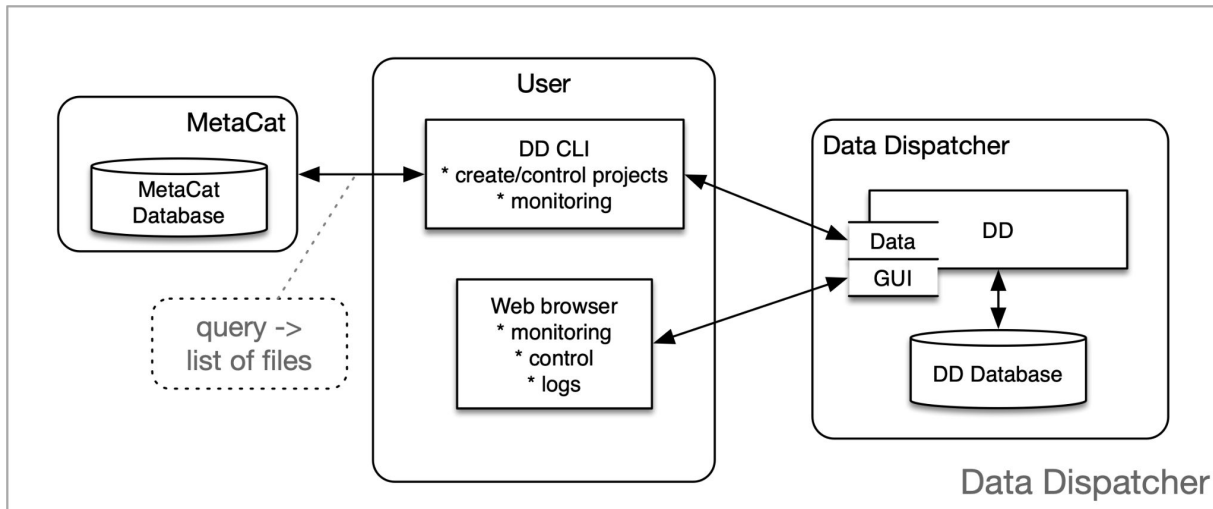
Data Dispatcher Functionality (SAM Station)

- Coordinate workflow as a stream of *projects*
- Project is a list of *files* to process with given executable
- Dispatch file *replicas*, as they become available, to *workers*
 - I am a working on this Project, give me next file
 - File processed successfully
 - File processing failed
- Optimization - replica to CPU proximity, etc.
- Monitoring, logging, control, report generation

Data Dispatcher and MetaCat

- Select “interesting” files
 - MetaCat query
- Create project with Data Dispatcher
- Wait for project completion
 - Monitor project status/progress
- Project final status
 - Logs
 - Reprocess some files ?

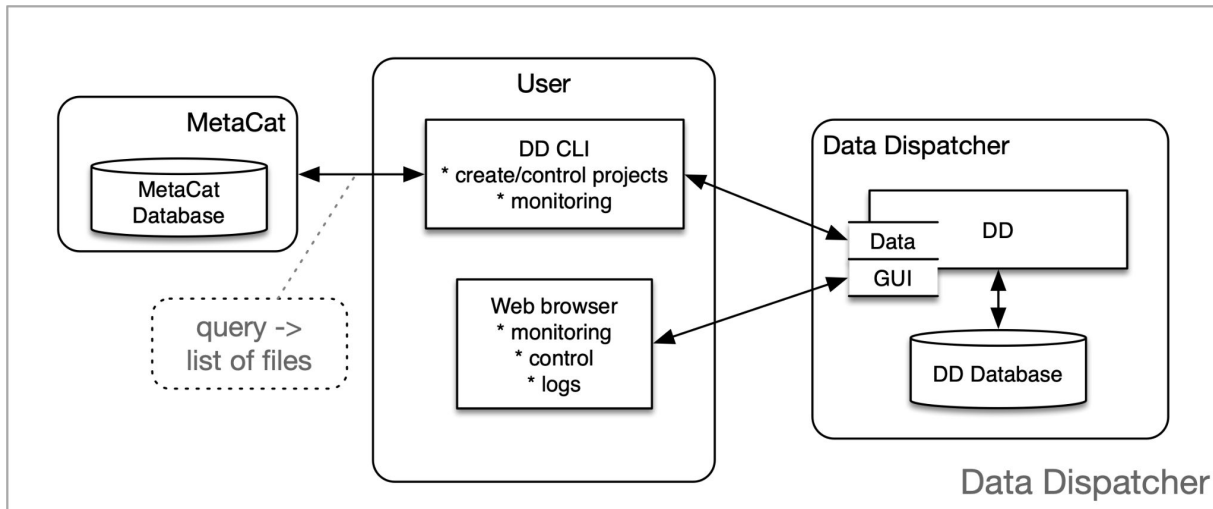
Data Dispatcher and MetaCat



DD = Data Dispatcher

- Query MetaCat for files based on the user's criteria
- Create project with given set of files to process
- DD will run the project

Data Dispatcher and MetaCat



- Query MetaCat for files based on the user's criteria
- Create project with given set of files to process
- DD will run the project

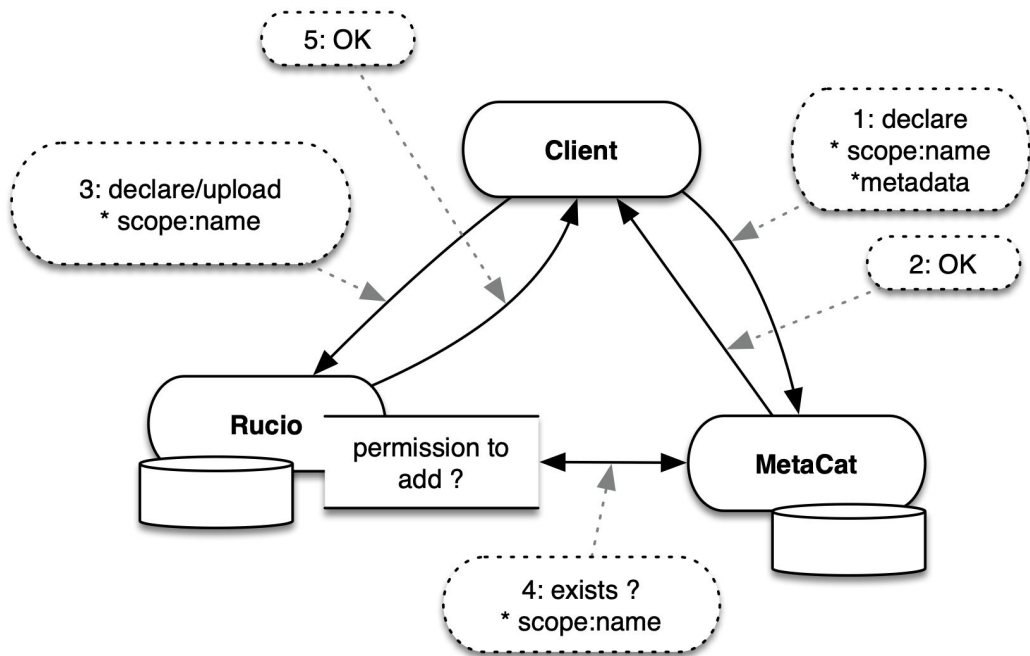
Options:



More SAM-like



Rucio/MetaCat Synchronization

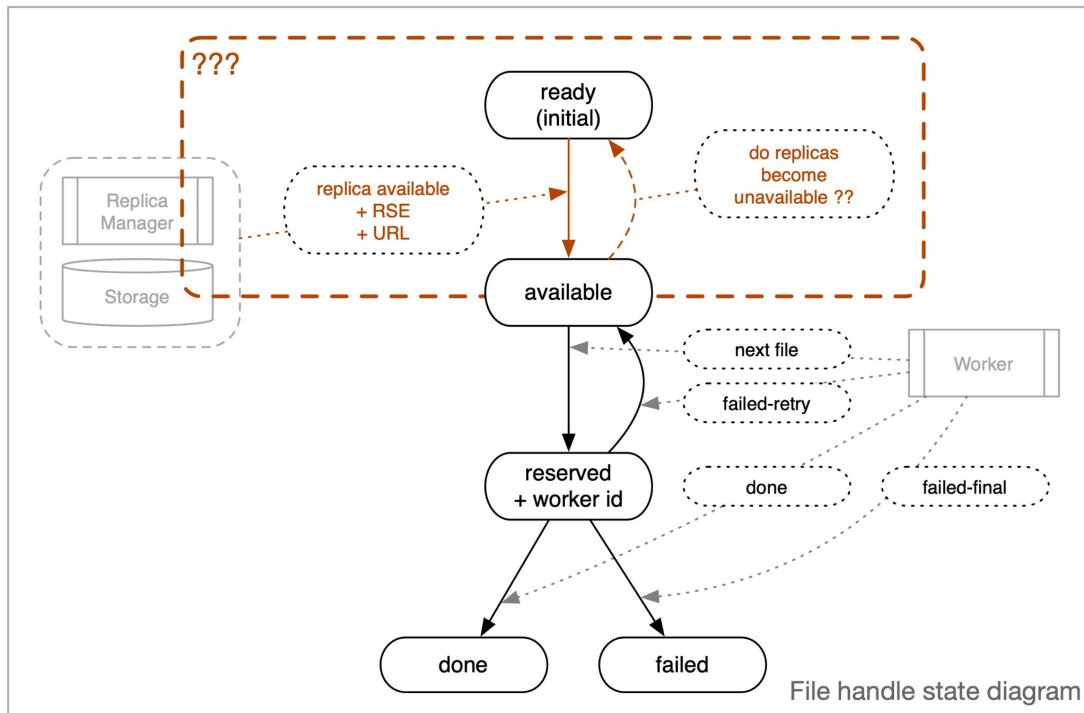


Implementation:

- Customizable permission-to-add-file Rucio module
- Permission is granted only if the file is already known to MetaCat
- Rucio/metadata catalog interface mechanism
 - Rather trivial
 - Written, to be tested

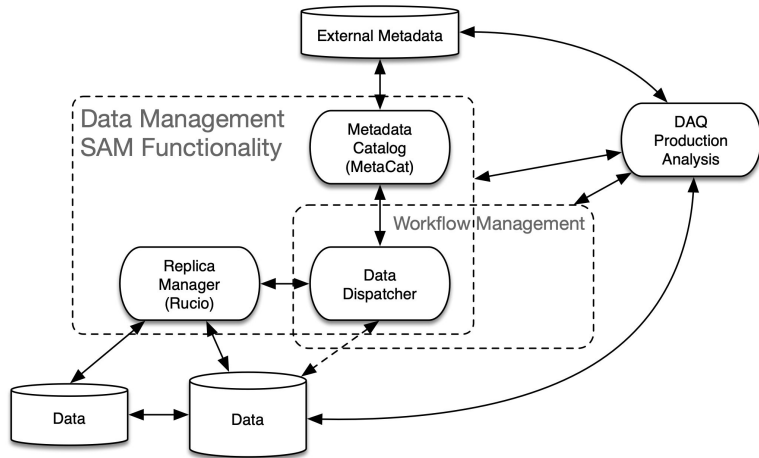
Goal: Make sure all files known to Rucio have metadata in MetaCat

Data Dispatcher and Replica Manager



- Dispatch *replicas* to workers as replicas become available
- Availability notification mechanism options:
 - Poll Rucio ?
 - Subscription/notification mechanism ?
- Data Dispatcher prototype has interface to update replica availability information
 - RSE, URL

Data Dispatcher and Global WMS

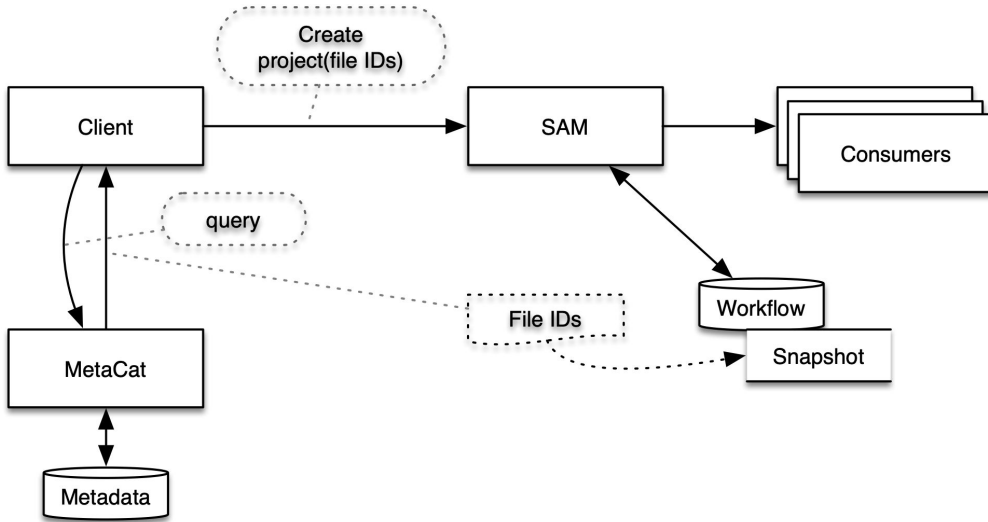


Beyond SAM decomposition - unclear:

- DD/WMS scope separation
- Functionality separation
- Architecture
 - DD is a separate component
 - Or part of WMS
- Communication
 - If separate, how does it communicate with WMS - has to be defined

Further Data Dispatcher development needs these to be cleared

Option: keep using SAM Station as is



1. Client sends query to MetaCat
 - -> List of file IDs
2. Client contacts SAMWeb to create a new snapshot with the list of file IDs
3. SAM creates snapshot
4. SAM creates project for the new snapshot

Minor changes in SAM: add function to create a project from a list of files (file ids)