

ButlerWG Meeting 2017-08-31

Meeting began 31 Aug 2017 @ 10am; ended 11:59am Project Time.

Attending

- [Michelle Gower](#)
- [Brian Van Klaveren](#)
- [Tim Jenness](#)
- [Jim Bosch](#)
- [Unknown User \(pschella\)](#)
- [Russell Owen](#)
- [Te-Wei Tsai](#)
- [Simon Krughoff](#)

Next meeting: 05 Sep 2017 @ 11:30am Project Time.

Examination of Draft Use Cases

DRP

- ☐ [Unknown User \(pschella\)](#) to merge some DRP use cases that are almost identical.
- We addressed the open questions on the DRP use cases.
 - Is memory mapping important? We talked about whether Butler has a need to give fine-grained control to the developer over the form of data access. In particular, whether state should be cached (such as a file handle) or if Butler should be stateless. It was felt that stateless is preferred. The memory mapping argument was not really compelling although it was thought that there was a useful concept for being able to retrieve multiple postage stamps from a single PVI if this would allow for efficiencies internally for memory mapping subsets without having to load the entire PVI into memory. This led to the retirement of the DRP use case but a DAX one is to be added. This is effectively a "vectorized get".
 - ☐ [Brian Van Klaveren](#) to add a use case for retrieving multiple postage stamps from a single PVI efficiently.
 - Do we need a transactional put? We felt that there is a need for a single put() to be atomic in the sense that if the put involves writing a file to disk and adding an entry to a database, then if the database insert fails the file should be deleted and an exception thrown. We felt there was no need for a put of multiple objects to be transactional. [Russell Owen](#) asked if put always blocks and the answer is yes.
 - ☐ [Unknown User \(pschella\)](#) to augment DRP22 to include atomic put example.

ARCH

- [Jim Bosch](#) asked if ARCH2 was the only use case requiring a single put writing to multiple output repositories in different formats. If that is the case then this could be deprioritized as you can always run the same processing twice with different output repositories.
- There was some discussion of Object Stores. IN2P3 drove Swift being supported by the current butler and S3 is being driven by SQuaRE although there are no S3 use cases currently in the spreadsheet.
 - ☒ [Simon Krughoff](#) to investigate the SQuaRE use case for object stores.

OPS

- [Michelle Gower](#) needs to re-read them following edits made by [Tim Jenness](#) to ensure that she agrees with the new wording.
- We talked about Provenance tracking. NCSA are going to propose as a baseline that DESDM provenance tracking be used. This relies on obtaining provenance from the workflow DAG and not from the tasks handling the pixel processing. There are cases where this provenance can not be completely accurate (for example when making coadds) and LSST have to decide whether there is a requirement for SuperTask to be generating additional provenance information. A distinction was made between provenance tracking data Contributing to a data product and that tracking the data from which the output was Derived From (the latter example being an algorithm that drops the two outlier images from a stack: you can't give that algorithm the 8 images that you used last time because then it would remove 2 and only use 6).

"France" (contributed by [Dominique Boutigny](#))

- ☐ [Michelle Gower](#) to read the new versions of the IN2P3 operations use cases and ensure they are consistent with the OPS vision.

COMM

- We discussed the "Major Questions" listed by [Simon Krughoff](#) at the bottom of the spreadsheet. It was made clear that the batch processing system will not be able to query the EFD and that any metadata will have to be either stored in file headers (which are created on the fly when raw images are accessed) or, for multi-valued EFD data covering a range of times, some mechanism must be available to query the EFD and write the results into files that can be handed to the batch system. Should the butler know how to read those EFD files?
- [Michelle Gower](#) reported that the Commissioning Cluster and Batch Processing system will have the same interface. She has forwarded on the questions relating to the size of the raw data cache at the base facility and the size of the compute system at NCSA during commissioning.

AP

- [Russell Owen](#) has started to work on the AP use cases and will finish transferring them from confluence by early next week.

Use Cases Future

- [Michelle Gower](#) asked if the enclave columns in the spreadsheet are useful. The consensus was that they are not and should be removed.
- [Michelle Gower](#) also noted that we are inconsistent in our usage of Actor titles in the use case and provided some standard names. Everyone was asked to go through their use cases and extract the Actor and write it down in the Actor tab so we can see how much duplication we have and can normalize the names.
- [Simon Krughoff](#) agreed to write a table of contexts (batch processing, local desktop, laptop on plane) and what data services they will have access to.