

AWS Proof of Concept Project Report
Hsin-Fang Chiang, Dino Bektesevic,
2019-12-09

1 Background

In April 2019, LSST DM began a proof of concept project with the Amazon Web Services (AWS¹) and HTCondor² teams to explore whether a cloud deployment of the Data Release Production (DRP) is feasible. The execution plan is in DMTN-114. For this project AWS granted us credits to use their platform. The team met biweekly to discuss the progress and plans. In this document we report on the work carried out, lessons learned as well as the results of scaling and pricing tests carried out by the AWS PoC.

2 Approaches and strategies

AWS PoC's goal is to demonstrate that running Data Release Production (DRP) on AWS is possible. To be considered successful the following constraints needed to be met:

2.1 Progress in phases.

DMTN-114 proposed development phases that broadly outlined how progress should be made and when performance would be monitored. The outlined plan was to move the execution to AWS without any backend modifications and then gradually switch out each component to AWS.

DRP is planned to be executed on a shared POSIX-compliant filesystem. In the first phase such filesystem was replaced by Amazon Elastic File System (Amazon EFS). EFS provides a POSIX-compliant, scalable, elastic NFS file system designed to provide parallel shared access to thousands of AWS compute instances. The computational resource was provided by the Amazon Elastic Compute 2 (EC2) service. EC2 service provides secure, resizable compute capacity based on virtualization technologies. The data repository, including the datastore and the

¹<https://aws.amazon.com/>

²<https://research.cs.wisc.edu/htcondor/>

SQLite registry, was located on EFS that was mounted to each of the EC2 compute instances.

Simple Storage Service (S3) was set to replace EFS. S3 is an object storage that allows massive amounts of unstructured data where each object typically includes the data, related metadata and is identified by a globally unique identifier, to be stored and accessed from EC2 instances in a durable and highly scalable way.

Buckets, organizational units in which related objects are generally stored, enable easier access and privileges administration. Access, read, write, delete and other atomical units of action on the objects themselves can be allowed or forbidden at the account, bucket or individual object level. Logging is available for all actions on the Bucket level and/or at the individual object granularity. It is also possible to define and issue complex alert conditions on Bucket or object actions which can execute arbitrary actions or workflows, but S3 itself not POSIX compliant.

The SQLite registry was to be replaced by PostgreSQL, one of the most popular open source relational database systems available. The choice to go with PostgreSQL was based on the fact that it's a very popular and well supported open source software that suffers from no additional licensing fees usually associated with proprietary software. Relational Database Service (RDS) is the AWS cloud service that launches and configures databases with ease.

HTCondor is a well used, popular scheduler typically used in shared filesystem HPC type clusters. A condor_annex module allows HTCondor deployment on cloud resources. Instances allocated by condor_annex are added to a shared resource pool on which then jobs added in the job queue are scheduled on. Unused compute resources are automatically deallocated after some set time spent idling. Personal Condor Pool on a single EC2 instance was used before we tested the capability to launch new EC2 instances as the Condor workers and execute jobs. HTCondor has been successfully used by HEP Cloud project, scaling up to 60,000 cores in multi-region deployment on AWS, and by IceCube experiment, where they scaled up to 51,000 GPUs across multiple regions and multiple cloud providers (AWS, Azure, and Google Cloud).

Maybe add pegasus here, or move lods of this text out into a special section for AWS services under intro, or just punt all of this to glossary (but then not easy to read fluently)

This approach allowed us to debug and adapt more smoothly and always had a fallback option in integrating and testing new features and components.

2.2 Test with small tasks and dataset before scaling up.

Test with small tasks and dataset before scaling up. The canonical ci_hsc dataset and the CI workflow provide a minimal HSC test dataset and a representative DRP workflow and algorithms. The test workflow was always run as one of the first integration steps. For larger scale tests we used the RC2 dataset.

2.3 Use Gen. 3 Middleware.

One of the mandates in the AWS PoC was to use the Gen. 3 Middleware, which is designed to ease the DRP execution and automation compared to the previous Gen2 Middleware. Gen. 3 Middleware design made it more pliable to required changes to the codebase, compared to Gen. 2 Middleware, and, as it was still relatively early in its implementation, its simplicity helped move progress forward quicker than if Gen. 2 Middleware was adopted.

However, during the AWS PoC the Gen. 3 Middleware has been under active development. Backwards incompatible changes occur often and often interfere with additions to the code. Even after the code was accepted because there are no integration tests being run for the added code, changes often break functionality. Even if potentially possible additional integration tests would not be accepted contribution as they would double, or worse, the total time required to run the test suite.

Keeping this in mind, We decided not to always follow the bleeding edge version and updated only when unnecessary, such as important bug fixes. The stack versions should not have strong impacts in the AWS PoC conclusions. **This is too strong, the stack version can have a big effect on the performance, ergo on our conclusions. I would like to rephrase this.**

2.4 Focus on end-to-end execution.

Focus on end-to-end execution and leave optimization behind. Some potential optimizations and further investigations were identified throughout the PoC project but were not carried out. Ideas are described in Sec 6.

3 Architecture Design

The system design at the end of the PoC is shown in Figure 3. All components are hosted on the AWS platform.

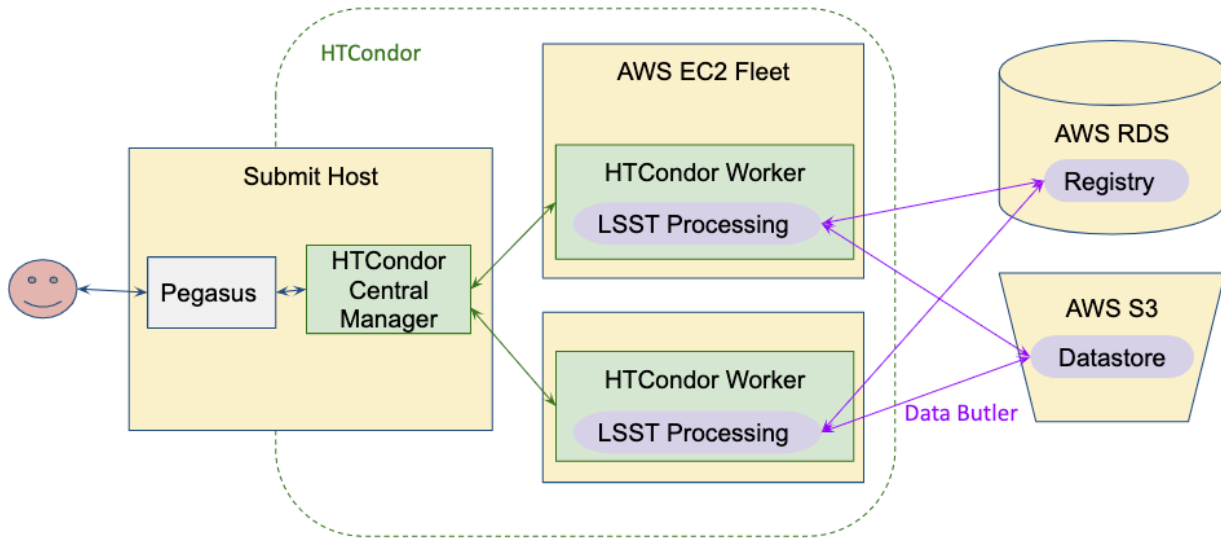


FIGURE 1: Diagram of DRP on AWS architecture.

Amazon Machine Image (AMI) bundle the Operating System (OS) with required data and configurations into a convenient image that can be launched on any EC2 instance type. We provide prebuilt AMIs containing LSST Stack, Pegasus, configured and running HTCondor in two flavors: master and the worker. These contain the complete environment required to run DRP workflows and scale them out with HTCondor. The user launches an on-demand instance using the master AMI, to which they log in and submit, scale out/down and run their workflows.

Compute resources are procured from the master node via condor_annex by targeting the worker AMI. Annex is a resource acquisition and external compute resource separation management unit. Each resource pool can have multiple annexes. Each annex manages its own lifecycle and generally speaking jobs will be submitted only to the annex they started and/or own. It is possible to request on-demand or spot instance fleets. On demand instances cost more but they can not be deallocated by AWS. Spot fleet instances tend to be significantly cheaper, but at any time AWS, given a 2 minute warning, can request the compute resources back and deallocate them.

The DRP, in Gen. 3 Middleware, is imagined as an execution of a Directed Acyclic Graph (DAG) known as Quantum Graph. A Quantum Graph link Tasks to datasets they consume and produce. Pegasus is used to submit the Quantum Graph to HTCondor as well as control, and monitor the workflow execution. Each LSST Pipeline Task quantum is resource-independent. Pegasus adds other necessary jobs, such as data transfer, to the executable workflow. HTCondor DAGMan is the workflow execution engine behind Pegasus and controls the processes.

The datastore is located in an S3 Bucket and follows the same hierarchical structure that POSIX datastore does. Consumed and produced datasets are read and written directly from S3 as bytes, whenever possible, and only downloaded to temporary files for objects that are not serializable.

The registry is a RDS PostgreSQL database. The RDS databases can be backed up into snapshots as well as exported to downloadable files on S3. Since the directory structure is preserved by the S3 datastore the entire data repository is trivially transferable between the cloud and a local filesystem.

These new Data Butler backends were implemented during the PoC, see Sect 4. Files that are not managed by Data Butler are managed by the Condor File IO via Pegasus. These include the Data Butler configuration file, pipeline definition (Quantum Graph) files, and the log files. The master instance serves as the staging site for these files. Various network protocols transfer the files between the master and the workers; instances do not share a filesystem.

3.1 Alternative architecture designs

We discussed different architecture designs but did not pursue all of them due to time constraints.

One prominent idea was to use condor transfer to read and write files to S3, rather than relying on Data Butler to communicate directly with S3. In this design HTCondor controls all file transfers. The proposed implementation was external to the LSST code base. The discussed architecture implied that there is a local data repository from which data is read and processed into and data ingestion and data transfer was triggered from. This requires the development of many additional utilities (posix to S3 and vice-versa datastore copying, URL prediction and generation tools) which often mimic or outright duplicate the mechanisms already present in Gen. 3 Middleware, as well as describing all required and produced datasets

to HTCondor as well. These issues disfavored this architecture. The favorable potential benefits of this approach are potentials of bulk data transfer management optimizations.

Other discussed architecture proposed many individual output registries that would be collocated into a single output registry at the end of processing. The motivation was to avoid potential bottlenecks related to large scale parallel SQL transactions by bundling transactions of multiple jobs together while leaving individual jobs access to a fast local SQLite database. However, the registry depends on conflict resolution syntax to resolve insertions of complex datasets, not to mention the synchronization issues or the issues of dataset entires of chained Tasks would require a global registry synchronization lock in the middle of processing. While the local datastore solutions might be possible it's hard to imagine a decentralized, or local cache-like, registry solutions.

4 Building AWS support into the Data Butler and lessons learned

The Data Butler is the overarching data IO abstraction through which all LSST data access is mediated. Datasets are referred to by their unique IDs, or a set of identifying references, which are then resolved through an registry that matches the dataset IDs, or references, to the location, file format and the Python object type of the dataset. The system that persists, reads and potentially modifies the datasets is called the datastore. The Registry almost always backed by an SQL database and the Datastore is usually backed by a shared filesystem. A major focus of AWS POC was to implement, and investigate issues related to, an S3 backed Datastore and a PostgreSQL backed Registry.

At the time the AWS POC Group began generation 3 Butler, the latest implementation of the Data Butler, implemented PosixDatastore, a local or shared filesystem datastore, and a SqliteRegistry. OracleRegistry followed soon after LSST AWS POC group began work. Initially the focus was on implementing an S3 backed datastore called S3Datastore. The interface between AWS services and LSST Stack would be based on the official AWS SDK called boto3. In March 2019, Dino Bektesevic visited NOAO to work more closely with Tim Jenness, which proved to be instrumental in implementing the early versions of a new module in the 'daf_butler' called 's3utils', an 'S3Datastore' class, the PosixDatastore equivalent, and a set of appropriate unit tests that demonstrated its functionality and correctness. The unit tests for the datastore utilize the 'moto' library which mocks requests and responses sent to AWS services, so that no additional external infrastructure is required to use it. PostgreSQLRegistry

class was implemented partially during the visit and completed shortly after the visit. The initial implementation showcasing the required changes to the code was submitted as a Draft Pull Request PR-147.

The tentative implementation revealed issues with how the Data Butler treated Uniform Resource Identifiers, or URIs, which were, at the time, not being handled correctly, as per standards defined in FC-3986, by The Location class. After expansive discussions and an example re-implementation called S3Location to demonstrate the issues, in May 2019 Tim Jenness authored the 'ButlerURI' class (PR-167) resolving the issues. Major efforts were then invested into refining the newly added code to the level of production quality as well as updating the remaining Gen. 3 Butler to use the updated ButlerURI code instead. Every call to OS functionality had to be generalized to take a URI and from it determine the appropriate operation - a call to OS functionality, a AWS operation or something else. This led changes in Butler, Config, ButlerConfig and YAML Loader classes. These changes made the whole of Data Butler more general and pliable to future changes, such as adding support for other cloud providers.

Further integration of the S3 backend required a change to Formatter classes to enable data serialization and deserialization to and from bytes. Formatters present interfaces for reading and writing of Python objects to and from files. They are the mechanism underlying how Data Butler is capable of presenting data as science products in the form of Python objects, abstracting away the underlying file types. Modifications were made to JsonFormatter, PickleFormatter, YamlFormatter, PexConfigFormatter and the generic abstract class Formatter. This concluded the last of changes required for S3Datastore integration. After which Jenkins integration tests were run and the S3Datastore and supplemental code was merged to master branch of the 'daf_butler' repository in PR-179 (the associated Jira ticket is DM-13361). It became apparent that there are certain similarities that are shared between PosixDatastore and S3Datastore, similarities that would be shared by other future datastore implementations. To reduce code duplication the general datastore code was refactored and reorganized in PR-187 shortly after.

PostgreSqlRegistry was not part of this PR. The initial implementation was based on OracleRegistry, due to the similarities between the two, but was re-implemented in terms of the generic SqlRegistry class in July. Problems were caused, for both Oracle and PostgreSQL, by the table naming conventions and additionally, for PostgreSQL, the table views did not conform to the assumptions made. In July the PostgreSqlRegistry was re-implemented in terms of the more general SqlRegistry and a new SQLAlchemy expressions compiler was written, so

that table views could be generated correctly. The policy for additional registry implementations was not to accept associated unit tests, as they are dependent on existing outside architecture, meant that checking whether it worked or not had to be based on manually executing one of the continuous integration tests such as `ci_hsc`. I migrated existing SQLite registries to PostgreSQL in July and August and made them available to the LSST AWS PoC group for testing. The code was merged into the master branch of Gen. 3 Butler in August with PR-161. A major issue was then discovered when issuing rollback statements during error recovery stemming from assumptions made when implementing how all of the current SQL registries handle errors during transactions. A stopgap solution, that works for all currently implemented registries, was implemented in PR-190 and a more complete solution was then implemented by Andy Salnikov in PR-196.

Outstanding issues are presented in terms of security and authorization when dealing with both `S3Datastore` and `PostgreSqlRegistry`, with `PostgreSqlRegistry` being especially sensitive to these issues. Security has received the outmost attention by the LSST AWS PoC group. Significant attention was paid to preserving the flexibility of the authentication in order to be able to incorporate external authenticators such as Oracle Wallets and AWS IAM Roles and Policies. There were several different iterations and improvements made to the authentication implementation (PR-189, PR-180 and PR-191) that resulted with the current implementation. An older Gen. 2 Butler module, `'db_auth'`, was re-implemented in Python by Kian-Tat Lim and added to Gen. 3 Butler so that the module would support basic file based authentication in absence of external authentication methods. Additional layers of security are achieved through EC2/S3/RDS interfaces by IP white/blacklisting, IAM, Policies etc. These policies can be very granular, affecting individually selected objects, Bucket-wide to placing all instances on the same, externally inaccessible, Virtual Private Network (VPN).

Adding the support for AWS into the Butler exercised almost the entirety of the Gen. 3 Data Butler. During the process many faults and unpredictable behaviors were discovered and solved. Many problems touched, and continue to exercise, the general Gen. 3 Data Butler implementation, as well as assumptions made during their implementation. Recounting the wide list of major improvements to the codebase, hopefully, reveals how productive this exercise has been in helping generalizing and strengthening the whole Gen3. Data Butler codebase.

4.1 PostgreSQL performance

The following described performance issue applies strictly to w_2019_38 builds and earlier. The QuantumGraph is created by issuing a very large SQL statement that, effectively, creates a cartesian product between many, of the 15-ish or so, tables in the registry. The results of the large query are then parsed in Python and a slew of many different, small, follow-up queries are issued. This diverse workflow presents a difficult challenge in DB optimization. The LSST DM team reports that on average it takes 0.5 to 1.5h to create QuantumGraphs on their infrastructure. In our tests we were unable to create even the simplest QuantumGraphs, even after doubling of the RDS instance resources, and were forced to abort multiple times after 30+h of execution. The lack of performance did not seem related to hardware limitations, since none of the performance metrics available showed heavy hardware loads nor did the problem go away even after significant increase of RDS resources.

PostgreSQL, and Oracle, regularly collect statistics on database objects used in queries. These statistics are then used to generate execution plans. These execution plans are then often-times cached. This in general results in performance gains. But because Butler generates all SQL dynamically there are no guarantees that a cached, or even pre-prepared execution plan can be executed. In PostgreSQL many of these optimizations then fail.

Specifically what seems to happen with PostgreSQL is that all the views materialize completely on disk. Even when the views are part of larger statements, in which the outer statements have strict constraints on them, the outer constraints do not seem to penetrate to the view statement. It is unclear, from within AWS PoC, whether Oracle DBMS does not suffer from this issue because it has a better SQL statement optimization engine or if this is something more directly tied to how PostgreSQL implements views.

The performance penalty of the behaviour is debilitating. Bellow is an abbreviated query plan for the simplest case of `select * from visit_detector_patch_join limit 4;` that suffers from the described issue.

```

1 Limit .... (actual time=12732.766..12732.774 rows=4 loops=1)
2  ->Unique .... (actual time=12732.765..12732.770 rows=4 loops=1)
3    ->Sort .... (actual time=12732.764..12732.767 rows=6 loops=1)
4      Sort Key: ....
5      Sort Method: external merge  Disk: 209008kB
6      ->Hash Join .... (actual time=763.846..1524.504 rows=4642107 loops=1)
7        Hash Cond: ....
8        ->Seq Scan on .... (actual time=0.008..16.289 rows=206960 loops=1)
9        ->Hash .... (actual time=763.723..763.724 rows=3259107 loops=1)
10         Buckets: 65536 Batches: 64 Memory Usage: 3635kB
11         -> Seq Scan on patch_skypix_join .... (actual time=0.006..240.128 rows=3259107 loops=1)

```

```
12 Planning time: 0.323 ms
13 Execution time: 12759.391 ms
```

Note that the total materialized size of the view on disk is 200MB and that it took 12 seconds to retrieve only 4 results! There is no clear solution to this problem from within the Middleware codebase.

The solution is however possible by manually redefining the views as materialized views instead, adding triggers that recreate the views on any insert statement to the underlying tables that make the view, and then adding indexing onto the materialized views. This reduces the time it takes to create a QuantumGraph to the point where it's comparable to that reported by the DM team. The query plan for the same SQL select statement as described above now looks like:

```
1 Limit .... (actual time=0.009..0.010 rows=4 loops=1)
2 ->Seq Scan on visit_detector_patch_join .... (actual time=0.008..0.008 rows=4 loops=1)
3 Planning time: 0.053 ms
4 Execution time: 0.021 ms
```

Recently the SQL schema and datamodel for the registries was reworked and completely re-implemented. We have not yet tested the performance of the new registry schema in PostgreSQL. The issues described here will optimize the set of specific queries, at the cost of a performance hit for all inserts into visits, patches..., but an important take away lesson here should be that because all of the SQL is generated completely dynamically guaranteeing execution plan stability and performance is difficult. Not all DBMSs optimize SQL statements equally well. Tying the CI and testing to only one database doesn't sufficiently test the code performance and because the SQL is generated completely dynamically to fix potential poor performance requires schema modifications to be made external to the Middleware code. Because some of the performance loss can be debilitating, it is possible that certain DBMSs are, in this implementation, tied into implementing unoptimal solutions just to be able to run DRP.

5 Results of the tract-sized runs

After successful execution with the `ci_hsc` dataset, we scaled up the run to one full tract of the HSC-RC2 dataset (DM-11345). The full HSC-RC2 input repository contains 108108 objects and totals ~1.5TB, including 432 raw visits in 3 tracts and ~0.7TB of calibration data. In this project, we targeted tract=9615 which was executed with the Oracle backend on the NCSA

Task	Count
Initialization	1
IsrTasks	6787
CharacterizeImageTasks	6787
CalibrateTasks	6787
MakeWarpTasks	4580
CompareWarpAssembleCoaddTasks	395
DetectCoaddSourcesTasks	395
MergeDetectionsTasks	79
DeblendCoaddSourcesSingleTasks	395
MeasureMergedCoaddSourcesTasks	395
MergeMeasurementsTasks	79
ForcedPhotCoaddTasks	395
Total	27075

TABLE 1: Task breakdown of the HSC-RC2 tract=9615 workflow

cluster in July 2019 as the S2019 milestone of the Gen3 team; see DM-19915. In terms of raw inputs, tract=9615 contribute around 26%, or ~ 0.2 TB, of the raw data in the HSC-RC2 dataset. We ignored patch 28 and 72 due to a coaddition pipeline issue as reported in DM-20695. A Butler repo was first made on NCSA's GPFS with a sqlite registry, and then transferred to the S3 bucket and the RDS instance. All tract-sized runs reported in this DMTN used LSST Software Stack release `w_2019_38`. The version of HTCondor was 8.9.3.

The workflow contains 1 initialization job and 27074 regular PipelineTask jobs. Science configurations from <https://github.com/lstt-dm/gen3-hsc-rc2> were used to generate a Quantum Graph. We transformed the Quanta into jobs in the Pegasus format with one-to-one mapping. The breakdown of the tasks in the workflow is in Table 1.

Generally speaking, there are two types of jobs: small-memory and large-memory jobs. Small memory jobs take less than 4GB per jobs, and large memory jobs can take up to ~ 30 GB per job. For simplicity, we consider all MakeWarpTask, CompareWarpAssembleCoaddTask, DeblendCoaddSourcesSingleTask, and MeasureMergedCoaddSourcesTask as large-memory jobs and for all of them set the required amount of memory to be ~ 30 GB as their job requirements, hence HTCondor only match them to machines with sufficient memory. AWS instances come with different flavors and the `r` family provides memory optimized instances with ~ 8 GB per core. The worker instances are configured to be HTCondor partitionable slots which dynamically splits resources and creates new slots to suit the jobs.

Run ID	Workflow wall time	Cumulative job wall time	Pipetask Min (sec)	Pipetask Max (sec)	Pipetask Mean (sec)	Pipetask Total (sec)
20191026T041828	28.4 hrs	61 days, 10 hrs	17.025	5936.038	195.465	5292217.997
20191121T015100	11.7 hrs	65 days, 8 hrs	18.272	5852.514	207.691	5623244.556
20191127T192022	8.7 hrs	62 days, 16 hrs	17.861	6243.819	199.636	5405141.464
20191127T192345	10.0 hrs	62 days, 23 hrs	19.297	6300.657	200.601	5431273.315

TABLE 2: Run summary

The submit host is an AWS on-demand instance, typically m5.large or larger. Spot fleets are requested after the Pegasus workflow start. Typically m4 or m5 instances are used for the single frame processing or other small-memory jobs, and r4 instances are used for large-memory jobs. After the workflow finishes, remaining running Spot instances may be terminated on the AWS console. Besides the 27075 pipetask invocations, Pegasus added 2712 data transfer jobs and one directory creation job. The total output size from the tract=9615 workflow is ~4.1 TB with 74360 objects.

5.1 Notes from the successful runs

Details of all runs are summarized in DM-21817, and in the following we summarize the successful runs only. Table 2 lists the runtime as reported by Pegasus tools.

In the first successful run 20191026T041828+0000, a fleet of 40 m5.xlarge instances were used for single frame processing and then a fleet 50 r4.2xlarge memory optimized instances for the rest. A m5.large on-demand instance served as the master. The single frame processing part finished in 4 hours; coadd and beyond took 16 hours. In this run, the memory requirement of the large-memory jobs was slightly higher than half of a r4.2xlarge, resulting in instance resources not fully used. This run spanned two billing days.

In the repeated run 20191121T015100+0000, the master was also a m5.large on-demand Instance. Fleets containing 75 m4.xlarge instances and 50 r4.2xlarge were launched. The memory requirement of the large-memory jobs was adjusted so that two such jobs can run on a r4.2xlarge simultaneously. Due to the larger fleet, the whole workflow finished within 12 hours in one bill day.

We then ran two of the same workflow graphs 20191127T192022+0000 and 20191127T192345+0000 simultaneously, simulating a larger input size. The master was a m5.2xlarge on-demand in-

stance. A Spot fleet containing 150 m5.xlarge instances ran the single frame processing for the first three hours, and a fleet of 150 r4.2xlarge instances ran the rest of the workflow. 600 jobs ran simultaneously during single frame processing. Larger fleets were used to help finishing the workflows in a shorter wallclock time. This run spanned two billing days.

5.2 Cost Analysis

The main components of the charges come from (a) EC2 Spot instances, (b) EC2 on-demand instances, (c) S3 storage, (d) RDS, and (e) others. The numbers reported here are from the AWS Cost Explorer tools.

1. **EC2 Spot instances.** EC2 Spot instances are the workers that execute the processing jobs, so this essentially is the cost of the compute power and scales with the compute resources needed to accomplish the processing campaign. The exact pricing for Spot instances varies based on supply and demand of the overall EC2 capacity. For instances used in our test runs, Spot instances cost around 20-25% of the on-demand instance price. Charges continue as long as the instances are up and running, even if no jobs are assigned to the instances. One hour of Spot instance typically costs \$0.045 for m5.xlarge and \$0.08 for r4.2xlarge. Different mix of instance types could affect the performance as well as the total cost. Including the instance idle time, we paid ~ \$0.035 per computing hour in average. We could pay less with better workload control and instance lifecycle control.
2. **On-demand EC2 instance.** We use an on-demand EC2 instance to serve as the submit host and the central workflow manager because we do not want it terminated by AWS. The current price of m5.xlarge on-demand instances is \$0.192 per hour.
3. **RDS.** Throughout our test runs we used a db.m5.xlarge instance, which has 4 vCPU and 16 GB of memory, to host the Butler Registry of the HSC-RC2 repository. We are aware this DB instance is more powerful than we usually need but we keep it running. The charge is therefore proportional to the span of time. For simplicity we count all RDS charge during the runs towards the cost of the runs, which is an overestimate because we also host other small database instances for testing purposes. For example, a db.t2.micro has been running alongside to host a Butler Registry for the ci_hsc repo which costs \$1.3 per day.
4. **S3.** The charge of S3 is dominated by the data storage cost, which is \$0.023 per GB

per month for the first 50TB. This means it costs \sim \\$1.2 per day storing the input repo (\sim 1.5TB) alone, and \sim \\$3.1 per day storing one set of the tract=9615 workflow outputs (\sim 4.1TB). Note that outbound data transfer is not free at AWS. At the rate of \\$0.09 per GB, transferring one tract=9615 output dataset out of AWS would cost \sim \\$370. We also incur charges per quantity of requests, which cost \sim \\$4.3 for each run of the tract=9615 workflow.

- Other charges.** The majority of other charges come from the Elastic Block Store (EBS) that provides storage for use with EC2. This includes SSD-backed volumes and snapshots, both of which are priced per size and time. Other charges are relatively small, such as a Business Support plan to get help from AWS engineers, and CloudWatch for additional monitoring information. In our accounting here this also includes charges from instances used in work independent of the workflow execution, so this should also be seen as a upper bound.

Table 3 provides an overview of the cost breakdown for each run. This includes all charges incurred on the billing days during which the runs were done, so this can be seen as an upper bound. We also have not optimized the usage; more discussions are in Sect 6.

Category	Runs		
	20191026T041828	20191121T015100	20191127T192022 20191127T192345
EC2 on-demand instances for master	2.88	1.54	2.05
EC2 Spot instances	94.69	58.55	52.94
RDS	26.38	13.82	13.82
S3 (including other storage cost)	17.29	13.48	21.48
Others/mostly EC2	18.38	7.12	6.37
Total	159.62	94.51	96.66

TABLE 3: Cost breakdown for each run

5.3 Cost projection

The raw input exposures in the HSC-RC2 tract=9615 workflow contains \sim 0.2TB from 112 visits. Compared to one night of LSST data in full operations, this is only \sim 11% in the number of visits. In one of our tests, we doubled the input size to explore how it scales up. It scales roughly linearly in cost and total compute resources. In processing wallclock time, it does not take longer because more EC2 instances can be deployed. For one tract of DRP test workflow the cost was around \\$95. If such scaling relationship holds for larger amount of data, 1000 visits

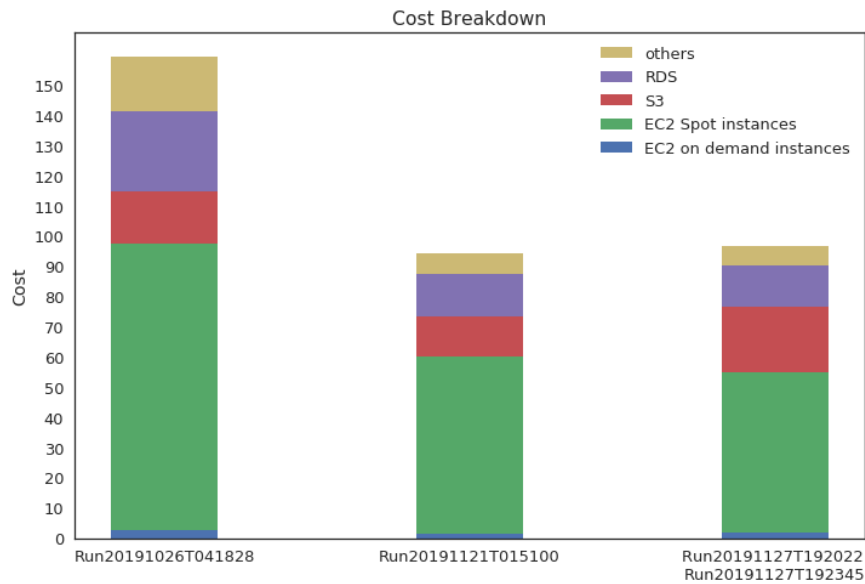


FIGURE 2: Cost breakdown for each run

will cost around \$850 to process.

6 Potential Improvements and more lessons learned

In this session we describe issues we have encountered during the execution and ideas on how to improve the code or the instance and workflow management. We discuss both intermittent failures that we understand and expect to occasionally encounter even in production, as well as higher level design or tooling improvements.

Failures can occur due to non-pipeline issues such as underlying infrastructure. The fault rate may be small but as we scale up we start to encounter some. Some examples are listed below; most seem transient.

1. Database connection timeout. Attempting to connect to the RDS instance failed.

```
sqlalchemy.exc.OperationalError: (psycopg2.OperationalError) could not connect to server: Connection timed out
```

2. After a file was added to a S3 bucket and during ingestion in the Butler registry, S3 reported a file does not exist. This will be fixed in DM-22201.

```
FileNotFoundError: File at 's3://hsc-rc2-test1-hfc/repo-w38/hfc30/srcMatchFull/1316/srcMatchFull_1316_24_HSC_17.fits' does not exist; note that paths to ingest are assumed to be relative to self.root unless they are absolute.
```

3. S3 read timeout before science processing started in a job.

```
botocore.exceptions.ReadTimeoutError: Read timeout on endpoint URL: "None"
```

4. Out-of-memory while running jobs. For the same pipeline and input data, this is reproducible. But we may not always have accurate memory usage prediction for any input data before running the jobs. We can configure HTCondor to increase the memory requirement in the retry. However sometimes OOM crashed the instances and appeared as a network issue which is not desirable behaviour.

5. Launching HTCondor Annex workers failed with connectivity check collector issues.

```
Connectivity check found wrong collector (f5fc15573ffb9c93 vs a006066e73c412da).
```

6. Dataset ConflictingDefinitionError. In rare occasions we observed dataset conflict errors from the registry without obvious reasons such as duplicate collection names or retries. DM-21201 has refactored the code to robustify such transactions. It could also be related to Spot instances getting terminated; see next section on failure recovery.

More generally, improvements in system design and tooling rise to prominence. We discuss some ideas below.

1. Better job failure recovery strategies. Our jobs write directly to the S3 bucket and the RDS instance via the Data Butler. If a job fails in a state that partial outputs are written but the job does not fully finish, recovery is not trivial. Such failed partial writes are handled on the registry side the DBMS but on the datastore side they are handled by wrapping failable transactions in a transaction context manager. The context manager takes a callable that then proceeds with the cleanup. For S3 datastore, at the time of writing, the utility functions performing the cleanup have not yet been written (in progress). The larger issue with job failures is the lack of DRP workflow failure handling. There is no way yet to reliably pause the DAG execution, overwrite files from a failed job, or other ways to handle such scenarios.
2. Container based software stack. We have found it tricky to handle the LSST stack installation, dependencies, and environments to be used together with other software. One possible way to avoid the headaches is to use docker based stack releases. This may also ensure consistency of software on the master and the workers more easily, as well as improve shareability.

3. Better cluster management tooling. Our current operational approach requires manually deploy suitable types and sizes of fleets based on our understanding of the overall workload. Strategies on the instance choices and timing of requests therefore affect the cost, and manual adjustments are usually needed to reduce cost. Annex can remove idle instances out of the condor pool but may not terminate the instances until the lease expires. Also once instances drop out of the pool they can't be added back easily. Tooling become essential for the operations. For example we may use scripts to auto-scale the Annex condor pool based on the current demand for resources, or based on predicting resource demands from the Quantum Graph. Such situations can also be helped by utilizing some of other AWS services, such as Simple Workflow service (SWF) or AWS Batch service. This of course has the potential to bind many of the find solutions strictly to AWS therefore any such requirements need to be handled with care.
4. End-to-end CI. This should include all operational components to do an end-to-end run. This includes Butler repo generation, registry generation, Quantum Graph generation, job composition, workflow translation, job execution, and so on. Many of the components were in the development phase and workarounds were used during the PoC. For example a native Gen3 ingestion was not available so a Gen2-to-Gen3 conversion was needed, and in fact many of the registries were migrated to RDS partially manually by executing bespoke scripts made specifically for such purpose that are not generally portable cross backwards-incompatible changes to the Middleware code. As we put together the pieces, the absence of to do so in an automatic fashion became a key burden. In retrospect we probably should have invested more time to automate the end-to-end workflow, even with workarounds, but the efforts were plagued by the lack of canonical ingest workflow. Additionally, in conversations with DM members, concerns were raised about having overly-inclusive CI tests. Adding S3 datastore CI tests to the default CI suite, while welcomed, would double the total CI execution time.
5. Credential handling. Security concerns received special attention in the AWS PoC. Natively none of the Middleware code requires the credentials to exist on the instance in plain text, even though they support that option. Authentication can be performed via Identity and Access Management Roles to both S3 and RDS services. IAM Role authentication is more secure than other authentication methods, but require understanding of how AWS security and permissions work and sometimes require long-ish setup. For example for RDS IAM access an account admin needs to create the IAM Role then DB admin needs to log in to the DB interactively and create a role with `rds_iam` permissions and then also grant usage privileges on the DB schema, tables etc. While we have tried

an IAM setup in the AWS PoC we often stored credentials onto the master and worker instances in `/.lsst/db-auth.yaml`, `/.aws/credentials` or in environmental variables which is not the best practice of handling the access. Utilities and tooling that help speed up such setup would be highly desirable.

6. Robustify and give better error messages. It has been observed that sometimes the error messages could be misleading.

7 Summary

In this PoC project we have demonstrated the feasibility of LSST DRP data processing on the cloud. We implemented AWS backends in the LSST Gen 3 middleware, allowing processing entirely on the AWS platform using AWS S3 object store, Postgres database, and HTCondor software. We analyzed cost usage in our test execution, and estimated cost for larger processing campaigns. We also showcased our progress in a live demonstration in the LSST Project Community Workshop, as well as a hands-on tutorial in the Petabytes to Science Workshop. Ideas of improvements necessary for larger-scale production are identified and discussed.

A References

References

[DMTN-114], Lim, K.T., Guy, L., Chiang, H.F., 2019, *LSST + Amazon Web Services Proof of Concept*, DMTN-114, URL <http://dmtn-114.lsst.io>

B Acronyms