

# Data Extraction HLD v0.1

Created by Neill Alexander, last modified about 3 hours ago

## Purpose of the Document

This document outlines a high level design for providing data export and data access from EDH, or an EDH-like environment on AWS.

## Background

### Export

The related CDS story is: <http://10.102.81.252/browse/CDSP-1949>

An export provides row level information about a declaration. No aggregation or analysis is done. The user can select the fields that he/she wants to include, some basic filtering (e.g. everything in February 2017), then run the export and download the results as csv.

### User Initiated

For a user initiated export, a user goes to a web page, does something to trigger the export, and the results are downloaded synchronously as a CSV file.

### Batch

Batch export is something that can be scheduled in advance, and will run without user interaction. In this scenario, we need a different delivery mechanism, since there will be no browser window sitting waiting to deliver the results. We can potentially support different delivery mechanisms e.g. UTM to CAF, email delivery (if not using EDH), etc.

### Access

The related CDS story is: <http://10.102.81.252/browse/CDSP-1925>

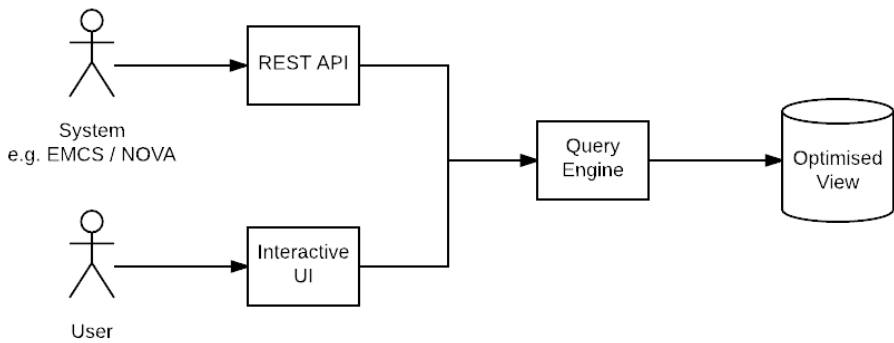
The specific requirements we have for this include:

- 1. NOVA Discovery Findings
- 2. EMCS Interface

For both of these, the existing solutions use the MSS Query Service to extract the data. We need to provide an equivalent service over the same data, callable by the NOVA and EMCS systems.

Moving to a higher level of abstraction, we need to be able to provide APIs which will return data from our Hadoop cluster, be that in EDH or AWS.

**Data Extraction = Data Export + Data Access**



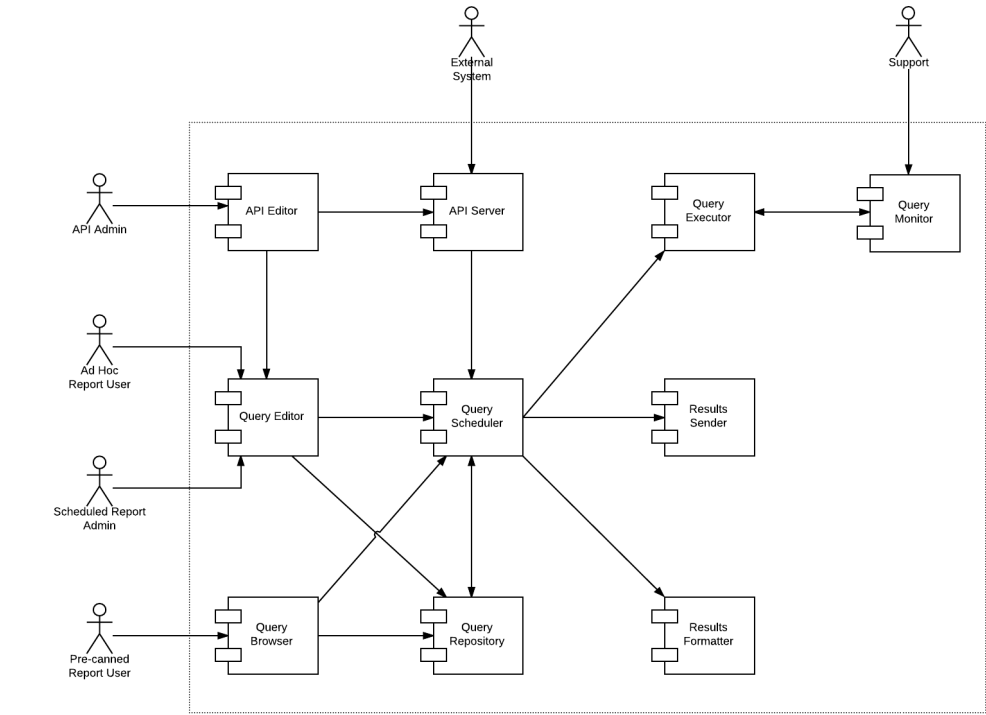
Whilst doing the initial investigation into the Data Export and Data Access requirements, and considering the various approaches we could take, it quickly became clear that a large degree of overlap exists between the two requirements. They both entail running queries against some kind of optimised view of the data. In reality, one optimised view will support both Data Export and Data Access.

The main differences between Data Export and Data Access are:

- 1. The method of initiating the query
- 2. The method of delivering the query results.

For simplicity, we propose to combine both sets of requirements into a single **Data Extraction** capability.

## Component View



### Query Editor

#### Responsibilities

- Provide a way for user to build and execute queries
- 2 views:
  - Simple
    - Allow users to build query visually
    - Select fields
    - Simple filters
    - Simple joins
  - Advanced
    - User can enter SQL
    - We should still make sure they don't try to do anything crazy
      - Perhaps parse the SQL and disallow certain potentially destructive constructs?
- User can choose to execute the query straight away
- Can also create / update / delete queries
  - Though, for delete, will need to take into account whether or not it is used in an API (see below)
- Can use the UI to set up a scheduled execution of the query
  - This would require allowing different delivery mechanisms, since wouldn't be interactive
    - Delivery to CAF directory
    - Email?

#### Visualization

*i* NB: All UI visualizations are to communicate the general idea only. They are not UX designs to be slavishly adhered to.

#### Simple View

Query Editor

http://10.11.12.128/queryEditor

Simple

Advanced

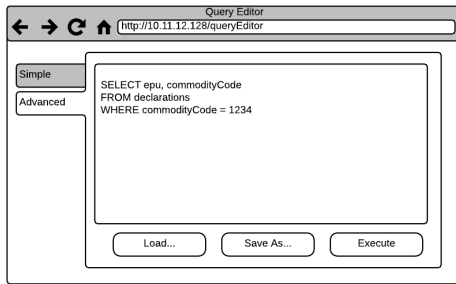
Field	Display	Filter Operation	Filter value
EPU	<input checked="" type="checkbox"/>		
Commodity Code	<input checked="" type="checkbox"/>	Equals	1234

Load...

Save As...

Execute

#### Advanced View



## Query Browser

### Responsibilities

- Provides a view into the repository - can see all queries
- Can link to the editor from here
  - They might end up being the same physical component i.e. a web app to allow user to view / update etc queries

### Visualization

Query	Description	Owner	Created	Last Modified
Unicorns 1	A list of everyone who imported Unicorns	nalexander	2017-10-01 12:23	2017-11-02 14:33
Butter Exporters	Everyone who exported butter (by month)	jlabbs	2017-11-01 08:21	2017-11-02 17:33

## Query Repository

### Responsibilities

- This is where we store the queries
- Could be HBase or HBase + Phoenix
  - For online access needs to be quick
- Could be a RDBMS
- Provides a very simple CRUD interface

## Query Scheduler

### Responsibilities

- When to run a query
- Responsible for calling the executor at the correct time
  - Can be 'now'
- Will need it's own persistent store
  - Can we spin up a DB for this?
  - Or just use HBase + Phoenix?
- Will require a front end view and a way to CRUD schedules
- Called by the query editor to run an ad-hoc query on demand

### Visualization

Query	Description	Owner	Schedule
Unicorns 1	A list of everyone who imported Unicorns	nalexander	Daily at 12:59
Butter Exporters	Everyone who exported butter (by month)	jlabbs	On Demand

## Query Executor

### Responsibilities

- Query executor will in reality be a fairly simple component
  - It will execute the SQL against Impala and return the result
- Assumption that Impala will be used to run these queries
  - It's designed specifically for this kind of thing
  - Gives ultimate flexibility in terms of what can be queried and how queries can be put together
  - Performance should not be an issue. Even 7 years worth of post-Brexit declarations is just about Big Data...
    - This will be validated by running some queries against production
  - And if over time it doesn't prove fast enough, we can provide further optimized layers for reports to run against, but keep the same architecture
    - Could also have custom executors for specific queries that need optimization
    - e.g. say there is an API published that is being called at very high volumes, we could provide a custom implementation as a Spark job, or something else, still referenced by name, and still with the same parameters.

## Query Monitor

Responsibilities

- Keeps track of running queries
- Allows to see the history of queries that have been run
  - i.e. audit
- Also useful for support i.e. why did my query not run

Visualization

Query Monitor				
Query	Started	Started By	Finished	Status
Unicorns 1	2017-10-01 23:30	system	2017-10-01 23:45	COMPLETE
Unicorns 2	2017-10-02 23:30	system		IN PROGRESS
Butter Exporters	2017-11-01 10:30	jlabbs	2017-11-01 10:31	FAILED

Results Formatter

Responsibilities

- Format results to desired output format
- Format can be specified dynamically
  - Start simple and add more formatters

Results Sender

Responsibilities

- Deliver the results to the user via there preferred method
- Import for the asynchronous exports

Results Repository

Responsibilities

- Store the results as they were sent out to the user
- Could be as simple as writing files to HDFS
  - Job to package them up on a weekly basis to avoid the 'small files' problem

API Editor

Responsibilities

- Once we have a saved query, can use the API Editor to make it available via REST
- Would need checks to make sure that once an API has been published, the query can't be changed
  - Wouldn't want someone to accidentally change a query and break an API
  - Think the Query Editor should have that kind of control
    - i.e. a persisted query has metadata, including whether or not it has been published as an API
    - Once published, can't be revoked? Can only publish a new version?
      - /extract/[QUERY\_ID]/[VERSION]?param1=blah&param2=hello
- Possible to unpublish an API
  - Can provide stats on how often it has been called etc to check if it is still being used
- Define the permission required to connect via the API here

Visualization

List APIs

API Editor							
API Path	Query	Version	State	Created	Last Called	Total Calls	Calls Per Day
Unicorns	Unicorns 1	1	PUBLISHED	2017-11-01 10:30	2017-11-05 14:10	120	10
/butter	Butter Exporters	0	DRAFT	n/a	n/a	0	0

Edit APIs

API Editor

http://10.11.12.128/apis/edit?api=1

Query

Path

Version

Butter Exporters

/butter

3

Publish New Version

Parameters

Name	Description	Validation Regex	Mandatory
epu	The processing unit	lw(4)	<input checked="" type="checkbox"/>
edate	Entry Date	^[0-9]{4}-[0-9]{2}-[0-9]{2}\$	<input type="checkbox"/>

History

Version	State	Number of Hits	Last Modified	Owner
2	PUBLISHED	250	2017-11-01 12:13	tscott
1	REMOVED	10	2017-11-01 12:12	tscott



Example Query

```
{
  "correlationId": "query_123",
  "querySql": "SELECT ...",
  "userId": "jsmith"
}
```

Example Response

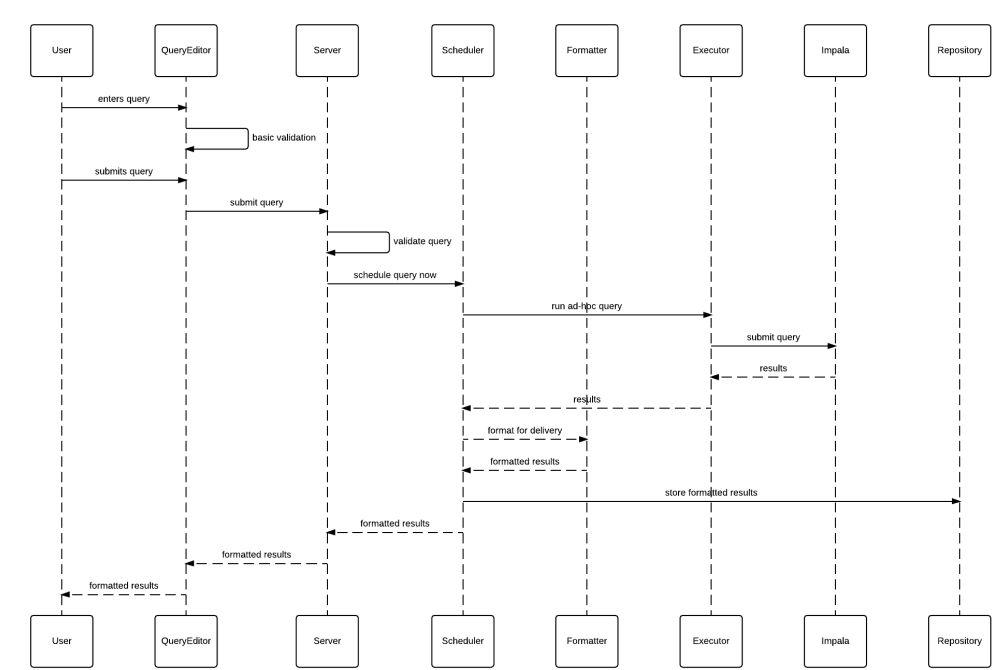
```
{
  "correlationId": "query_123",
  "response": [
    ["2017-01-01", "ID1", "User 1", "Notes"],
    ["2017-01-01", "ID2", "User 2", "Notes"]
  ]
}
```

For this to work, all requesting components would need to listen to all responses, and ignore those they aren't interested in (i.e. no consumer groups for the 'query\_response' topic).

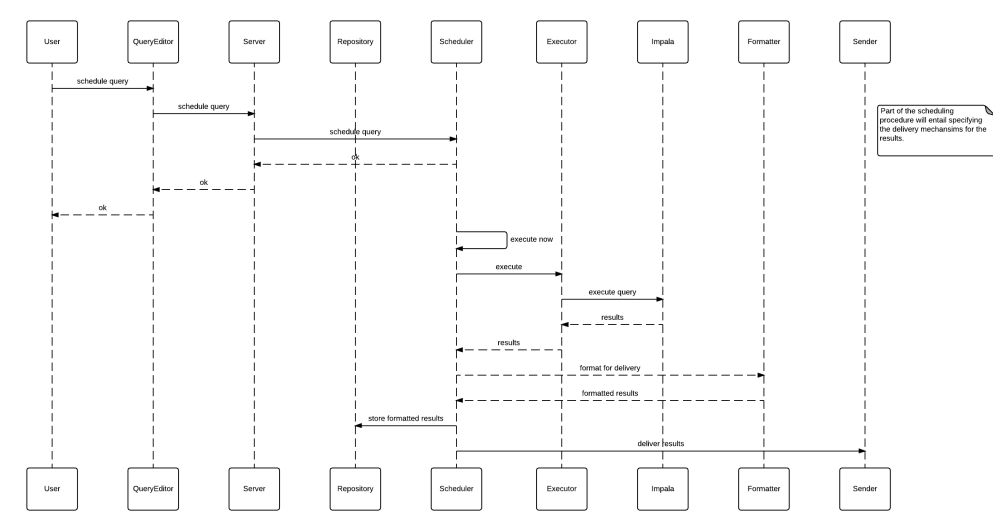
Sequence View

Following are some sequence diagrams to illustrate how the various components could interact to support the requirements.

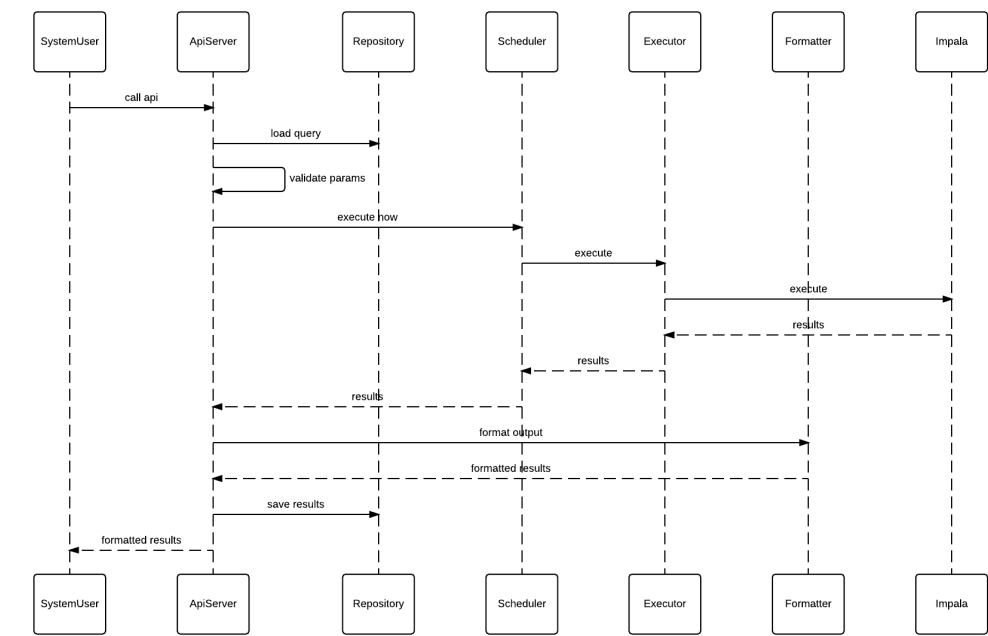
User Initiated Export



Scheduled Export



Data Access



### Delivery Mechanisms

These will matter particularly with scheduled data export. For user-initiated export the hope would be that the query runs quickly enough that the user can simply wait for it to finish and then save a file to their local machine, via browser download. However, since we will need to support different delivery methods for the scheduled exports, we could also support them for user-initiated exports.

In all cases, we should store all data export files for a period of time in HDFS (it is a big data system after all, so let's store all the data!).

### CAF

Copy file to a shared folder. We could use UTM to achieve this. The user would then know to expect the file in a specific location using a specific file format.

### Email

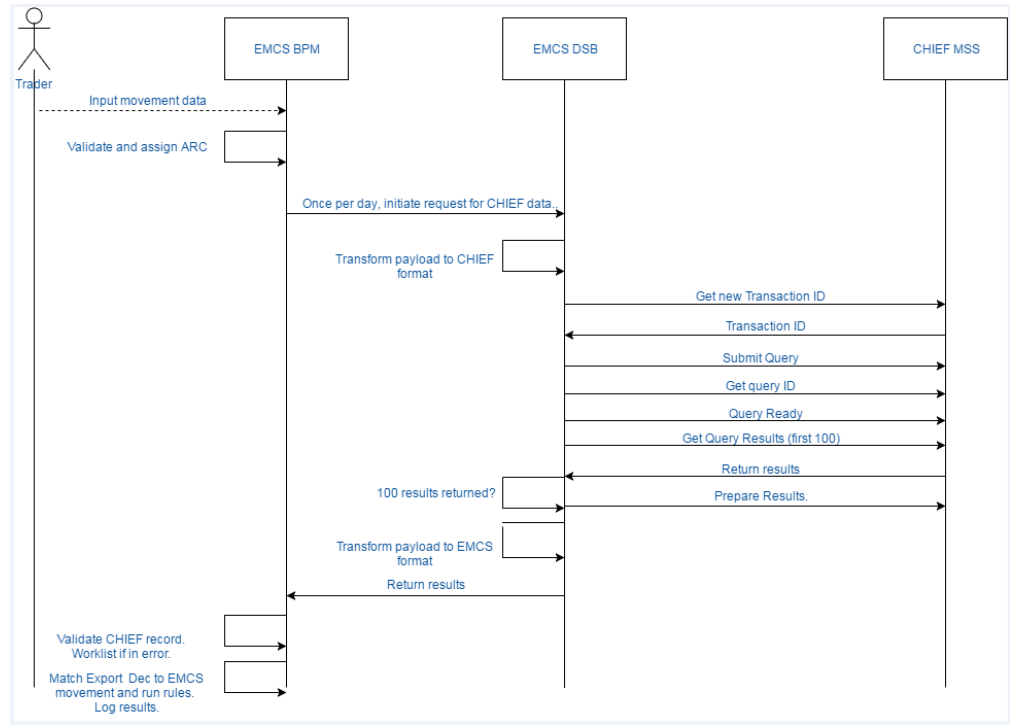
Since the data export is internal only (for now) we could probably use email to send files. Or, if we can't send the files, send a link which would allow a user to download the file, providing they have the correct authorisation to access it.

### Download File

For user-initiated export, the browser could block and download the file once the results have been returned. We could also provide a link via the Query Monitor page which would allow the user to download the results - clicking on the link would require us to retrieve the file from HDFS and return to the user.

### HTTP

For the data access API, we could return a synchronous or asynchronous response. If synchronous, the call simply blocks until results are available. If asynchronous, we would need to return a token to the calling system which would allow it to get the results at a later date. This kind of mechanism is used currently for the CHIEF MSS query API, with the transaction ID:



(reproduced from EMCS Interface)

### Delivery Formats

Applicable to both data access (via the API) and export (via the UI / scheduler).

#### JSON

Something like:

```
{
  "queryId": "1234_xyz",
  "fields": {
    "fieldA": "fieldA_value",
    "fieldB": "fieldB_value"
  }
}
```

#### XML

Something like:

```
<?xml version="1.0" encoding="UTF-8"?>
<queryResults>
  <queryId>1234_xyz</queryId>
  <fields>
    <field name="fieldA">fieldA_value</field>
    <field name="fieldB">fieldB_value</field>
  </fields>
</queryResults>
```

#### CSV

Something like:

```
fieldA, fieldB
fieldA_value, fieldB_value
```

#### MSS Legacy (maybe)


If existing users of the CHIEF MSS Query interface don't have the capacity to change the way they process the returned data, then we may need to replicate the output format, which is described at [EMCS Interface](#) as: "Data is returned in a MSS Enterprise specific format. (EMCS DSB transforms the data from the format output by MSS to a format that is used by EMCS BPM.)"

Ideally we would not need to do this.



## Security Considerations

- Use Kerberos / RBAC to control access to declarations.
- Could introduce an API Gateway to control access to the API Server

 [Like](#) Be the first to like this

No labels