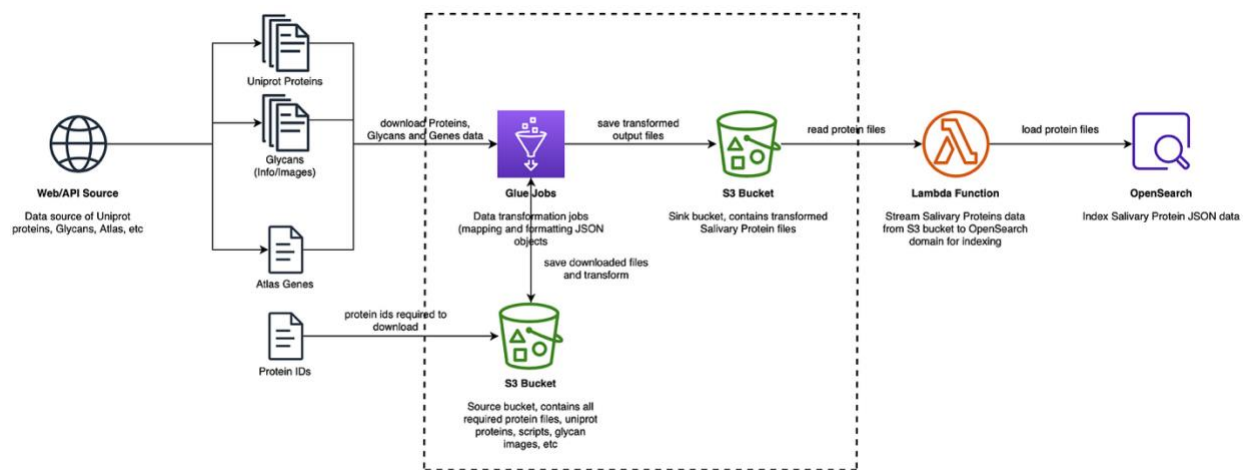# Salivary Proteins Data Pipeline using AWS Glue & OpenSearch

## Documentation

### AWS Architecture

The cloud architecture of Amazon Web Services utilizes several Glue Python jobs to perform various tasks. These tasks include fetching protein and glycan data along with glycan images from the Uniprot and Glygen websites, retrieving Human Atlas data through an API, and storing all of this information in S3 storage buckets. These files are subsequently employed for additional processing, where necessary salivary protein data is extracted using Glue's transformation capabilities. Lastly, the obtained salivary protein data is sent to OpenSearch using a Lambda function to facilitate the indexing process.
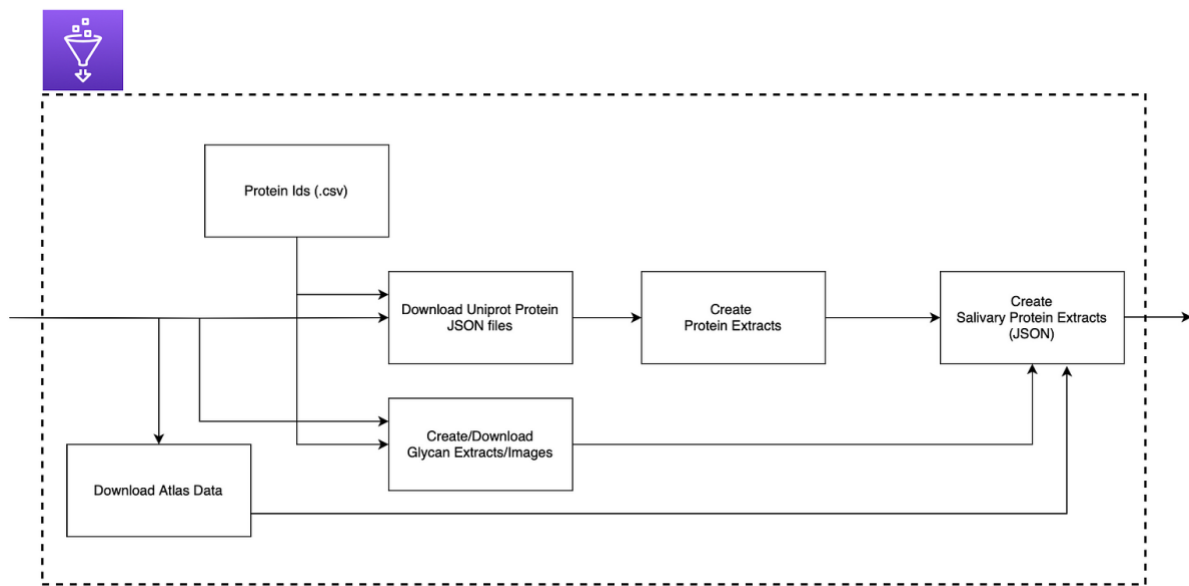


Salivary Proteins Data Pipeline — AWS Glue Workflow

### AWS Glue Jobs workflow

AWS Glue Jobs simplifies the process of data transformation and ETL by automating many of the tasks involved, allowing users to focus on defining data transformations and extracting insights from their data.

In this workflow, once the pipeline is triggered, three jobs `download-create-protein-extracts`, `create-glycan extracts`, `download-glycan-images` run parallelly and once all are completed, the final job `create-salivary-protein-extracts` runs in the end to create salivary proteins JSON files.



Salivary Proteins Data Pipeline—AWS Glue Jobs Workflow

## Prepare AWS environment

Please note, the IAM roles and policies are managed by the AWS Administrator.

## 1. Create S3 buckets

Amazon Simple Storage Service (Amazon S3) is a scalable object storage service offered by Amazon Web Services (AWS). S3 provides secure, durable, and highly available storage for various types of data, making it suitable for a wide range of use cases.

| | | | | |
|---|---|---|---|---|
| ○ | hspw-dev-opensearch-upload | US East (Ohio) us-east-2 | Bucket and objects not public | August 15, 2023, 11:23:34 (UTC-04:00) |
| ○ | uniprot-proteins | US East (Ohio) us-east-2 | Bucket and objects not public | August 15, 2023, 16:25:33 (UTC-04:00) |
| ○ | proteins-reference | US East (Ohio) us-east-2 | Bucket and objects not public | August 15, 2023, 23:29:37 (UTC-04:00) |

S3 buckets—Ohio us-east-2 region

**proteins-reference**

In this bucket, we will keep all our reference files used in the data pipeline—

`protein_ids.csv`—All *protein ids* we need to process.
`rna_tissue_consensus.tsv`—Atlas *RNA tissue* reference data.
`normal_tissue.tsv`—Atlas *normal tissue* reference data.
`scripts/`—Folder to save all Glue Python job scripts.

**uniprot-proteins**

Here, all our protein files will be saved—

`uniprot_protein_files/`—Folder to download all uniprot protein files from uniprot website by Glue job `download-create-protein-extracts`.

`protein_extracts/`—Folder to save all protein extracts after downloading and transforming by Glue job `download-create-protein-extracts/create-protein-extracts`.

`glycan_extracts/`—Folder to save all glycan extracts after transforming glycan data by Glue job `create-glycan-extracts`.

`images/`—Folder to save all glycan images, downloading through glygen web API by Glue job `download-glycan-images`.

**hspw-dev-opensearch-upload**

This bucket contains all salivary protein extracts (output files from the data pipeline) which are ready to index in AWS OpenSearch.

`salivary-protein-extracts/`—Folder to save all salivary protein extracts created by job `create-salivary-protein-extracts`.

## 2. Create Glue Python jobs

All Python scripts—

| ☐ | Job name | ▽ | Type | Last modified | ▼ | AWS Glue version | ▽ |
|---|---|---|---|---|---|---|---|
| ☐ | create-salivary-protein-extracts | | Python shell | 8/14/2023, 2:49:19 PM | | | |
| ☐ | download-glycan-images | | Python shell | 8/14/2023, 2:44:05 PM | | | |
| ☐ | create-glycan-extracts | | Python shell | 8/14/2023, 2:42:15 PM | | | |
| ☐ | create-protein-extracts | | Python shell | 8/14/2023, 2:38:48 PM | | | |
| ☐ | download-create-protein-extracts.py | | Python shell | 8/14/2023, 2:35:45 PM | | | |

AWS Glue Python Jobs

**download-create-protein-extracts** — Job to retrieve *protein ids* from the S3 location `proteins-reference/protein_ids.csv`, simultaneously download corresponding UniProt protein files from the UniProt website using multi-threading and batch processing to expedite the procedure. After each batch, initiate a subsequent task named `create-protein-extracts` passing *batch id* and *protein ids* as job parameters. This task involves generating protein extracts and storing them within the `uniprot-proteins/uniprot_protein_files/` directory. The process incorporates a retry mechanism that reschedules the job at a consistent interval to manage potential errors caused by concurrent job executions.

| | Run status | ▽ | Retry | ▽ | Start time | ▼ | End time | ▽ | Duration | ▽ | Capacity | ▽ | Worker type | ▽ | Glue version | ▽ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ⦿ | ✓ Succeeded | | 0 | | 08/14/2023 15:28:41 | | 08/14/2023 16:05:35 | | 36 m 47 s | | 0.0625 DPUs | | - | | 3.0 | |

download-create-protein-extracts—Duration ~ 30 m

**create-protein-extracts** — This job is triggered by `download-create-protein-extracts` and takes protein ids as job parameters. It creates protein extracts using required protein data and subsequently stores them in the `uniprot-proteins/protein_extracts/` directory.

| | Run status | ▽ | Retry | ▽ | Start time | ▼ | End time | ▽ | Duration | ▽ | Capacity | ▽ | Worker type | ▽ | Glue version | ▽ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ⦿ | ✓ Succeeded | | 0 | | 08/14/2023 16:04:25 | | 08/14/2023 16:04:51 | | 20 s | | 0.0625 DPUs | | - | | 3.0 | |
| ○ | ✓ Succeeded | | 0 | | 08/14/2023 16:03:22 | | 08/14/2023 16:03:48 | | 20 s | | 0.0625 DPUs | | - | | 3.0 | |
| ○ | ✓ Succeeded | | 0 | | 08/14/2023 16:02:18 | | 08/14/2023 16:02:50 | | 25 s | | 0.0625 DPUs | | - | | 3.0 | |
| ○ | ✓ Succeeded | | 0 | | 08/14/2023 16:01:15 | | 08/14/2023 16:01:41 | | 20 s | | 0.0625 DPUs | | - | | 3.0 | |
| ○ | ✓ Succeeded | | 0 | | 08/14/2023 16:00:12 | | 08/14/2023 16:00:41 | | 22 s | | 0.0625 DPUs | | - | | 3.0 | |
| ○ | ✓ Succeeded | | 0 | | 08/14/2023 15:59:09 | | 08/14/2023 15:59:34 | | 20 s | | 0.0625 DPUs | | - | | 3.0 | |

create-protein-extracts—batch runs—Duration ~ 20 s (30 m)

**create-glycan-extracts** — This job involves reading *protein ids*, retrieving glycan data from the glycogen API, and then generating glycan extracts using the necessary glycan information. These extracts are subsequently stored in the `uniprot-proteins/glycan_extracts/` directory.

| | Run status | ▽ | Retry | ▽ | Start time | ▼ | End time | ▽ | Duration | ▽ | Capacity | ▽ | Worker type | ▽ | Glue version | ▽ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ● | ⊘ Succeeded | | 0 | | 08/14/2023 15:28:41 | | 08/14/2023 16:12:56 | | 44 m 7 s | | 0.0625 DPUs | | - | | 3.0 | |

create-glycan-extracts—Duration ~ 40 m

**download-glycan-images**—Job to download all glycan images and save them to `uniprot-proteins/images/` directory reading *protein ids* list.

| | Run status | ▽ | Retry | ▽ | Start time | ▼ | End time | ▽ | Duration | ▽ | Capacity | ▽ | Worker type | ▽ | Glue version | ▽ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ● | ⊘ Succeeded | | 0 | | 08/14/2023 15:28:41 | | 08/14/2023 16:06:52 | | 38 m 4 s | | 0.0625 DPUs | | - | | 3.0 | |

download-glycan-image—Duration ~ 30 m

**create-salivary-protein-extracts**—The final job is to read the list of *protein ids* and create salivary protein extracts. These extracts encompass essential information obtained from both protein and glycan extracts. Subsequently, these files are saved within the `hspw-dev-opensearch-upload/salivary-protein-extracts/` directory. This action will then activate a Lambda function to facilitate their transfer to the OpenSearch domain index.

**create-salivary-protein-extracts**                    Last modified on 8/14/2023, 2:49:19 PM    Actions ▼    Save    **Run**

Script | Job details | Runs | Data quality New | Schedules | Version Control

**Job runs (1/2)** Info          Last updated (UTC)  August 15, 2023 at 13:31:36  ⟳  View details  Stop job run        **Table View**  Card View

Q Filter job runs by property                                                          ‹ 1 › ⚙

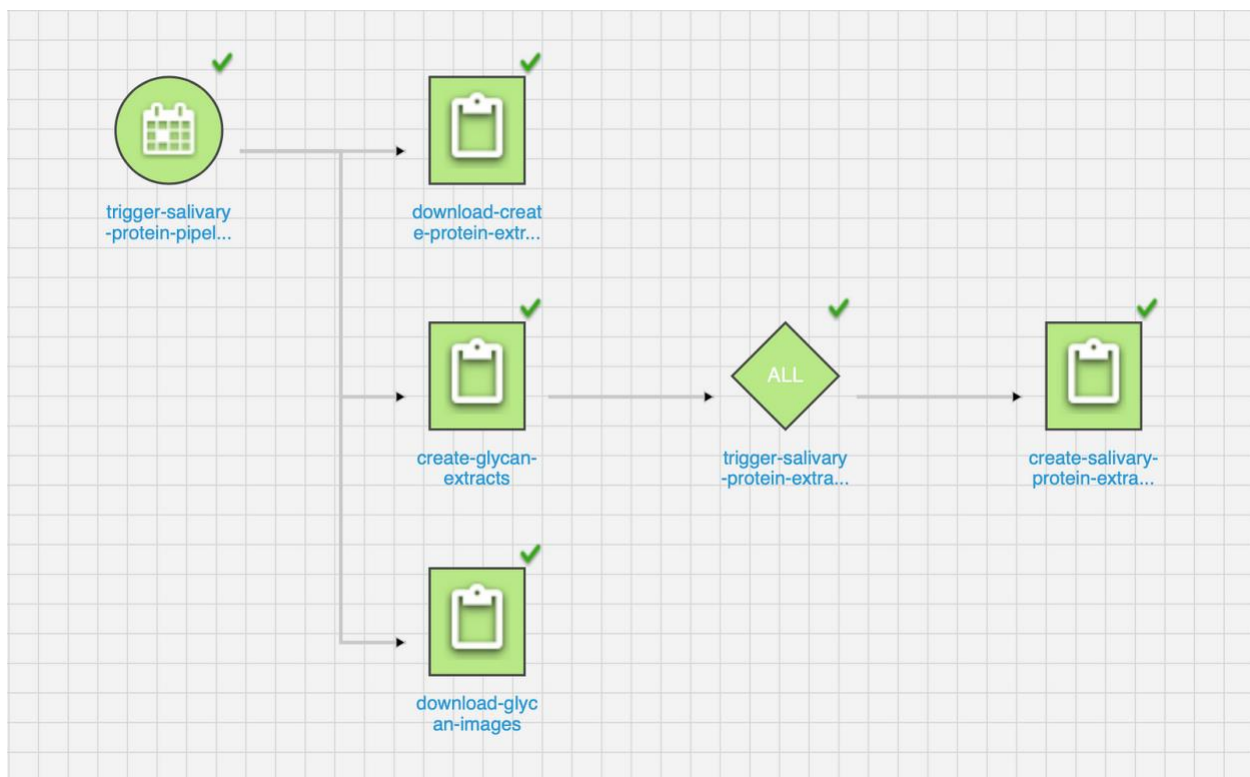| | Run status | ▽ | Retry | ▽ | Start time | ▼ | End time | ▽ | Duration | ▽ | Capacity | ▽ | Worker type | ▽ | Glue version | ▽ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ● | ⊘ Succeeded | | 0 | | 08/14/2023 16:13:26 | | 08/14/2023 16:42:53 | | 29 m 21 s | | 0.0625 DPUs | | - | | 3.0 | |

create-salivary-protein-extracts—Duration ~ 30 m

## 3. Create a Glue workflow

Finally, create the Glue workflow to build the data pipeline and execute all jobs with on-demand/scheduled triggers.

An AWS Glue workflow is a sequence of interconnected Glue jobs and triggers that automate and orchestrate data processing tasks. This workflow simplifies the process of managing, monitoring, and executing various ETL (Extract, Transform, Load) tasks across your data sources.



AWS Glue Workflow

**salivary-protein-data-pipeline**—Glue workflow to download and create protein extracts, glycan extracts and images, and salivary protein extracts.

**trigger-salivary-protein-pipeline**—Trigger `salivary-protein-data-pipeline` scheduled monthly on 1st day of every month. It also triggers other parallel

jobs—`download-create-protein-extracts`, `create-glycan extracts`, and `download-glycan-images`.

**`trigger-salivary-protein-extracts`**—Trigger `create-salivary-protein-extracts` job, after all previously triggered jobs are executed successfully.



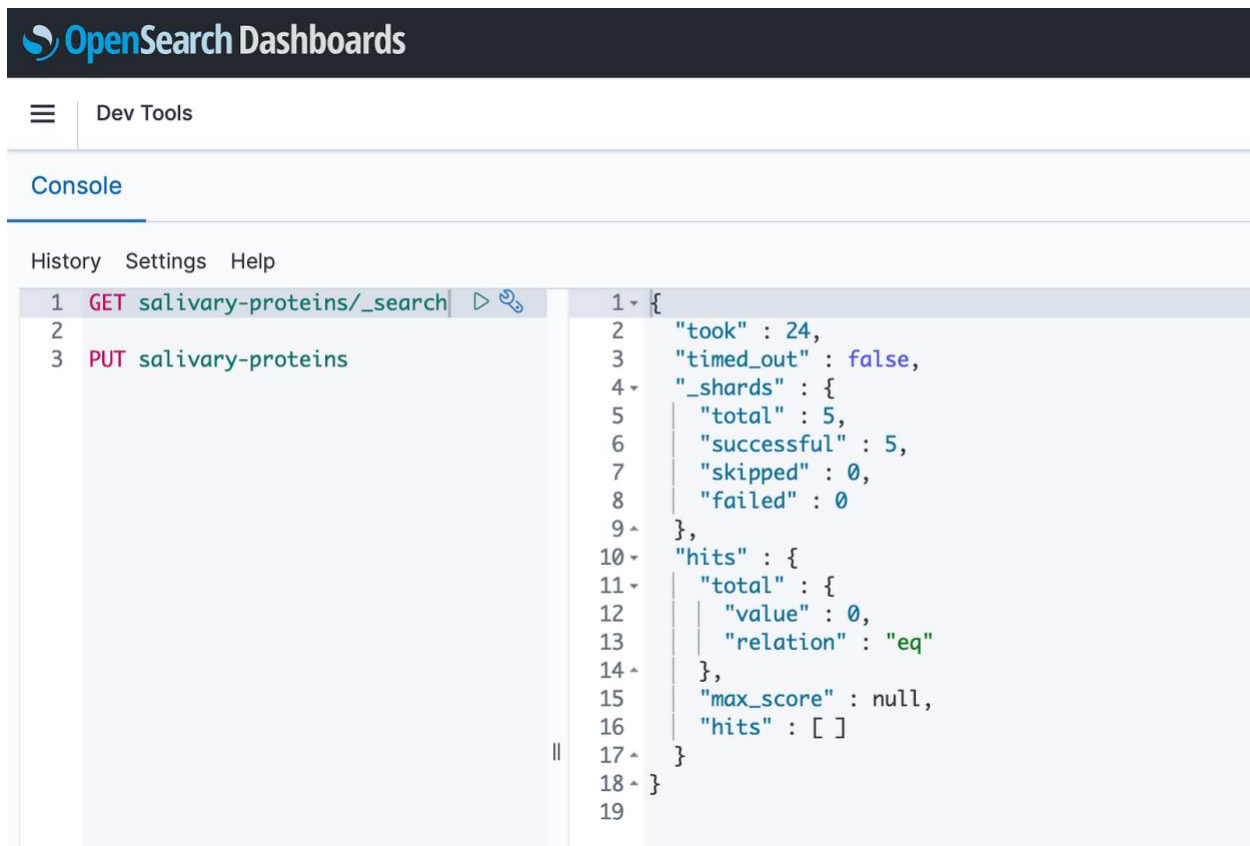Workflow Run Details—Duration ~ 1 hr

## 4. Create an OpenSearch domain

Amazon OpenSearch (formerly known as Amazon Elasticsearch Service) is a managed service offered by Amazon Web Services (AWS) that provides a scalable and reliable solution for deploying and operating the open-source Elasticsearch and Kibana software for search and analytics.



AWS OpenSearch Domain

OpenSearch domain indexes all salivary protein data and make it available end-to-end for search queries. In our case, accessible by the website.

**hspw-dev2**—OpenSearch service domain to ingest salivary protein data from the S3 bucket and index them.
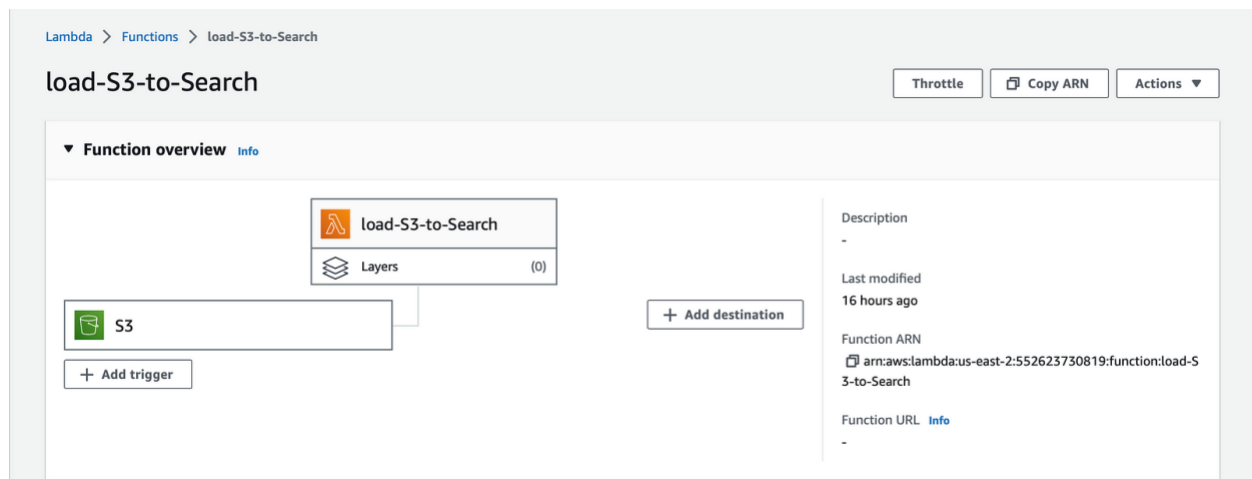


OpenSearch Dashboard

**salivary-proteins**—OpenSearch index which contains all salivary protein data.

## 5. Deploy a Lambda function

Now, we need a Lambda function to transfer salivary protein data from the S3 bucket to AWS OpenSearch.

AWS Lambda is a serverless computing service provided by Amazon Web Services (AWS) that allows you to run code without provisioning or managing servers. With Lambda, you can execute your code in response to various events and triggers, enabling you to build scalable and cost-effective applications.

AWS Lambda Function

**load-S3-to-Search**—Lambda function is triggered when any salivary protein extract file is uploaded to the S3 location `hspw-dev-opensearch-upload/salivary-protein-extracts/` which then sends data to OpenSearch index `salivary-proteins`.

## CloudWatch Logs

Amazon CloudWatch Logs is a service provided by Amazon Web Services (AWS) that allows you to monitor, store, and analyze log data from various AWS resources and applications. It enables you to collect and centralize log data for better visibility, troubleshooting, and compliance purposes.



CloudWatchLog Groups

```
16    CLOUDWATCH LOGS -
17
18    -- Download uniprot protein files
19    Downloading Uniprot protein file https://www.uniprot.org/uniprotkb/000206.json ...
20    000206.json downloaded successfully and saved to uniprot-proteins/source_files
21    Processing of protein id 000206 completed successfully.
22
23    Total protein ids processed: 1/3
24    Downloaded: 1   Failed: 0
25
26    -- Create protein extracts
27    Reading Uniprot protein file 000206.json ...
28    protein_extract_000206.json created successfully and saved to uniprot-proteins/protein-extracts
29    000206.json moved to uniprot-proteins/source_files/completed
30    Processing of protein id 000206 completed successfully.
31
32    Total protein ids processed: 1/3
33    Completed: 1   Failed: 0
34
35    -- Create glycan extracts
36    Getting Glycan protein data from https://api.glygen.org/protein/detail/000206 ...
37    Getting Glycan protein mass from https://api.glygen.org/glycan/detail/ ...
38    Downloading Glycan image from https://api.glygen.org/glycan/image/ ...
39    ___ downloaded successfully and saved to uniprot-proteins/glycan-extracts/images
40    glycan_extract_000206.json created successfully and saved to uniprot-proteins/glycan-extracts
41    Processing of protein id 000206 completed successfully.
42
43    Total protein ids processed: 1/3
44    Completed: 1   Failed: 0
45
46    -- Create salivary protein extracts
47    Reading protein_extract_000206.json ...
48    Reading glycan_extract_000206.json ...
49    Creating salivary protein extract ...
50    salivary_protein_extract_000206.json created successfully and saved to hspw-dev-opensearch-upload/salivary-proteins-extracts
51    Processing of protein id 000206 completed successfully.
52
53    Total protein ids processed: 1/3
54    Completed: 1   Failed: 0
55
56    -- Lambda S3 to OpenSearch data transfer
57    Pushed salivary protein data (Protein Id 000206)- to opensearch index successfully.
```

CloudWatch Logs

## Full Code at GitHub

You can get the full code in the JCVenterInstitute GitHub [repository](#).

**[HSPW-V3/awsjobs at main · JCVenterInstitute/HSPW-V3](#)**
*Contribute to JCVenterInstitute/HSPW-V3 development by creating an account on GitHub.github.com*