

Data Pipeline for Hyperspectral Camera in LDFZ

Naman J. Parikh

July 29, 2024

Design

Hyperspectral Camera

The primary data producer in this pipeline is the hyperspectral camera. The camera captures a line of pixels (horizontal line as it is set up in the LDFZ). In contrast to a standard camera which captures three values (red, green, and blue) for each pixel, the hyperspectral camera captures values at 371 wavelengths for each pixel. This allows us to determine temperature and emissivity of a blackbody imaged by the camera using blackbody physics. Data from captured images is stored by the Hyperspec III software running on a dedicated computer. Each image capture generates a folder that is named with the date and time of capture and contains the image files.

Pipeline Overview

The goal of this pipeline is to download and analyze data on a separate server for storage and efficiency while making the result of analysis available to the LDFZ user. To do this, we build two separate Kafka topics; one for raw data and one for the analysis result. We use OpenMSIStream [1] for streaming the files. The topics can be found on the brokers associated with `paradim01` as `hyperspec_LDFZ_data` for raw data and `hyperspec_LDFZ_result` for results.

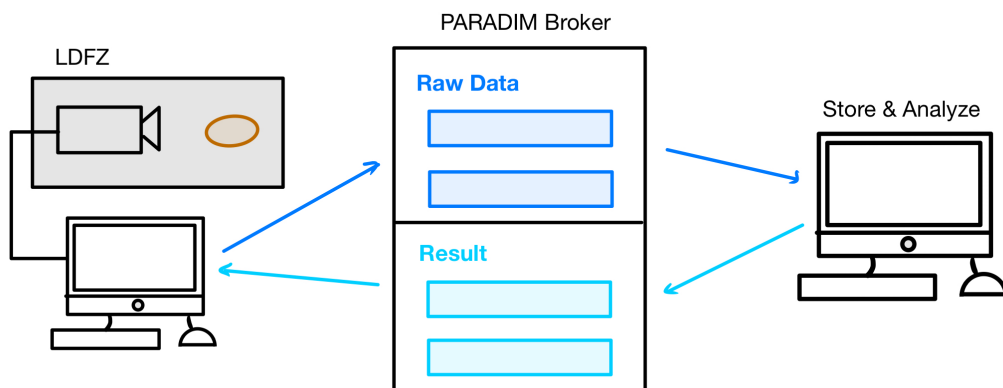


Figure 1: High-level design of the data pipeline

Camera: Data Producer and Result Processor

The camera's dedicated computer continuously runs a Python script, which is saved at `"C:\Users\Headwall\streaming.py"` and corresponds to `"StreamingScripts/hyperspec.py"` on my GitHub repository. This script prepares and spawns two threads.

1) The first runs a `DataFileUploadDirectory` producer watching `"C:\Headwall\sensor1\captured"`. This uploads image folders to `hyperspec_LDFZ_data` as they appear in the captured directory.

2) The second runs a `DataFileStreamProcessor` consuming files from `hyperspec_LDFZ_result`. The result files are assumed to be temperature arrays stored in NumPy (`".npy"`) files. The stream processor saves a copy of the file as well as a heatmap of the array, which represents a thermal gradient of the image.

Data Processor

The data processor script can be found at `"StreamingScripts/processor.py"` in my GitHub repository. This script runs a stream processor that tracks which files for each image have been received and triggers a pyrometry analysis when all requisite files are received. The pyrometry analysis can be found on my GitHub repository in the analysis function in `"StreamingScripts/temperature_analysis.py"`. To run this script on another machine, the directories for the broker config file and stream processor directory will need to be adjusted accordingly.

Additional Consumers

Girder consumer?

Usage and Troubleshooting

Dependencies

The scripts were written to be run with Python 3.9. Refer to the `openmsistream` documentation found [here](#) for information about setting up an environment for the package. Additionally, for my scripts and pyrometry analysis, the following packages are needed:

```
matplotlib
numpy
pathlib
scikit-learn
scipy
spectral
threading
time
tqdm
```

Alternatively, a JSON dump for the conda environment used is available on my GitHub at `"Documents/environment.json"`

Topic Configurations

Config.	Value
<code>retention.ms</code>	43200000
<code>max.message.bytes</code>	10485880

In other words, messages are retained in the topic for 12 hours and the maximum bytes per message is approximately 10 megabytes. The retention time can be adjusted, but it should be large enough to accomodate the time it takes to download large images. The maximum bytes for message is set to higher than default to accomodate the relatively large size of hyperspectral images.

Additionally, in the configuration file used by the data processor, under consumer, `max.poll.interval.ms` should be set to a large value (I have used 86400000, which is one day). This is because the stream processor will easily surpass the default value (30000) and be stopped when running pyrometry analysis on a large enough image.

Restarting the Data Processor

In the event that the data processor crashes and no image was being actively downloaded/processed at the time, I recommend clearing the `hyperspec_LDFZ_data` topic prior to restarting the processor. This will stop the processor from redownloading and analyzing images retained in the topic, which can take a substantial amount of time. The process to do this is as follows.

1. Run on `paradim01`

```
/opt/kafka/bin/kafka-configs.sh --bootstrap-server localhost:9092 --alter  
--entity-type topics --entity-name hyperspec_LDFZ_data --delete-config  
retention.ms
```

```
/opt/kafka/bin/kafka-configs.sh --bootstrap-server localhost:9092 --alter  
--entity-type topics --entity-name hyperspec_LDFZ_data --add-config  
retention.ms=1000
```

This will set the retention time to 1000 ms.

2. Wait a few minutes for the change to take effect and the messages to clear from the topic.
3. Reset the retention time by running the same commands from step 1 but replace 1000 with the desired retention time.

References

- [1] Margaret Eminizer, Sam Tabrisky, Amir Sharifzadeh, Christopher DiMarco, Jacob M. Diamond, K.t. Ramesh, Todd C. Hufnagel, Tyrel M. McQueen, and David Elbert. Openmsistream: A python package for facilitating integration of streaming data in diverse laboratory environments. *Journal of Open Source Software*, 8(83):4896, 2023.