

DocuTrace
F20FC: Industrial Programming
Coursework 2

Sam Fay-Hunt — `sf52@hw.ac.uk`

December 5, 2020

Contents

1	Introduction	1
2	Requirements Checklist	1
3	Design Considerations	2
4	User Guide	3
5	Developer Guide	3
5.1	Summary of application flow	3
5.2	Analysis package	3
5.2.1	FileRead	3
5.2.2	DataCollector	4
5.2.3	ComputeData	4
5.2.4	Plots	4
5.3	Gui	4
5.3.1	main	4
5.3.2	Tab	4
5.3.3	Tasks	4
5.4	Utils	4
5.4.1	Exceptions	4
5.4.2	Logging	4
5.4.3	Tasks	4
5.4.4	Validation	4
6	Testing	4
7	Personal Development	5
8	Conclusions	5
A	References	6

The report should have between 10–15 pages and use the following format (if you need space for additional screenshots, put them into an appendix, not counting against the page limit, but don't rely on the screenshots in your discussion)

1 Introduction

The DocuTrace application is a moderate size, data-intensive application, its purpose is to analyse and display document tracking data from the website issuu.com. The website hosts a substantial number of documents, and provides anonymised usage statistics, the data is provided in the form of a sequence of individual JSON entries separated by new lines.

It is assumed the users of an application like DocuTrace would be someone with a high enough degree of technical competency to use simple Linux command line applications, the user would likely be a researcher (data science), or a business. The prior assumption leads to the assumption that the hardware running this application would be closer to server class than standard consumer hardware with higher CPU core counts and a lot more RAM. To try and compensate for the potential scale of the data a significantly large amount of RAM is not mandatory to run this application, but a pool of approximately 8GB of RAM should be installed on the system at a minimum when processing 3 million lines otherwise a significant performance penalty may be incurred.

DocuTrace was written in Python 3, it is intended to be run on Ubuntu 20.04, and has not been tested on other operating systems.

Comprehensive documentation has been generated, a list of dependencies, installation and run instructions are provided, note the recommended entrypoint of the application is via the “docutrace” shell script, this script will automatically configure an environment variable to specify the output directory of graph files generated.

2 Requirements Checklist

Here you should clearly show which requirements you have delivered and which you haven't.

Requirements marked with the Priority **Essential** are requirements explicitly requested in the task brief.

The priority for each requirement is encoded as follows:

- **Essential** - This priority indicates that this requirement must be implemented to satisfy basic functionality of the browser.
- **High** - A high priority indicates that this requirement is important for providing a good quality user experience.
- **Medium** - A medium priority requirement is nice to have but if it is missing it is acceptable.
- **Low** - Low priority indicates that this is unlikely to be fulfilled within the time frame of the project, a stretch goal at best.

Requirement	Description	Priority	Status
Runs on Linux	The application needs to run on an up-to-date Linux platform (Ubuntu 20.04)	Essential	Complete
Use Python 3	The core logic of the application should be implemented in Python 3	Essential	Complete
Views by country	Given a document UUID find from which countries the document has been viewed, display as histogram	Essential	Complete

Requirement	Description	Priority	Status
Views by continent	Given a document UUID find from which continent the document has been viewed, display as histogram	Essential	Complete
Views by browser	For an input file count the number of occurrences for each browser used to access each document in the file, display full and truncated strings as histogram	Essential	Complete
Reader profiles	For each user find the total time spent reading documents, display the top 10	Essential	Complete
Readers of document	For a given document identify all visitor UUIDs of readers of that document	Essential	Complete
Documents from reader	For a given reader identify all documents that reader has read	Essential	Complete
"Also likes" functionality	Using the readers of document and document from reader functionality return a list of documents sorted by a given sorting function as a parameter	Essential	Complete
"Also likes" graph	generate a graph that visualises the "Also likes" functionality	Essential	Complete
Simple GUI	Develop a simple GUI based on tkinter to display the statistical data and receive input parameters	Essential	Complete
Command-line usage	Provide a command line usage to test the application functionality in an automated way	Essential	Complete
Display query history in GUI	Extend the gui to display a history of queries that allow navigation to past queries	Medium	Incomplete
Provide documentation	Generate and host documentation for the application	High	Complete
Log exceptions and warnings	Log any exceptions encountered to be handled by journald	Low	Complete

3 Design Considerations

Here you should clearly state what you have done to your application to make it more usable and accessible.

Extensive documentation has been provided for other developers wishing to extend this codebase, additionally the `PEP8` style guide has been utilised to aid readability and help ensure consistent coding style.

In addition to the required parameters specified in the CLI requirement some secondary parameters have been included, this includes log level verbosity, an argument to limit the volume of data displayed in the GUI and command line, and a parameter to exit the application early (used when only a single task is needed).

The potentially long loading times when processing large datafiles motivated the inclusion of an animated loading bar to provide some feedback to the user.

A logger is used to handle debug, info, warning and exception messages. The logger sends messages to stdout so they can be managed by journald, this includes exceptions and traceback messages this is a common practice and lets system admins handle the log output as they see fit. The verbosity of

the log message can be set with the `-v` parameter, by default all messages of logging level WARNING and above are displayed.

4 User Guide

Use screen shots of the running application along with text descriptions to help you describe how to operate the application.

5 Developer Guide

Describe your application design and main areas of code in order to help another developer understand your work and how they might develop it. You may find it useful to supplement the text with code fragments.

The provided documentation includes a definition of expected parameter types for all functions and methods. Additionally primitive type stubs have been added to the function and method signatures to improve readability. More comprehensive type information could be added in the future by defining custom type stubs for all the classes, and union types.

5.1 Summary of application flow

This application can be broken down into 4 parts: Analysis, Gui, Utils packages and the `main.py` script; `main.py` is the primary entrypoint of the application, Analysis module handles all the data processing needs of the application, Util provides convenience functions like logging, exceptions, validation and also a task selection script, and finally the Gui module handles all functionality related to the graphical user interface.

The main entrypoint (`DocuTrace/main.py`) is a script to define the arguments, checks the validity of the filepath (a thread is started to begin processing the data right away) and task identifier from the command line, finally it starts the logic that parses the task and runs the given task once the data has finished processing.

Once processing is finished the task is launched, this will open the GUI. An instance of the `ComputeData` class is instantiated using the data obtained by the `DataCollector` class, it is then passed to the gui, this class handles all interactions between the gui and the data.

Once the gui opens the selected task should be visible on the screen, with the data visible, at this point new parameters for the various tasks can be provided to view other datapoints related to this task, or the number of instances displayed can be modified, depending on the task.

5.2 Analysis package

5.2.1 FileRead

The `FileRead` module handles all interaction with the filesystem with regards to opening and closing the file. The core function used to read from the file is the `stream_file_chunks()` function, this function creates an iterator and lazily reads the file on demand. There are 3 classes within this module, `ParseFile`, `JsonProcessContextManager`, and `JsonParseProcess`. A file is actually read by the `parse_file()` method of the `ParseFile` class, when concurrency is enabled it uses the `JsonProcessContextManager` to add chunks of the JSON data file to a `JoinableQueue` [1] managed by the context manager, see Fig. 1. The context manager starts all the processes, and maintains 2 queues: the aforementioned queue that receives JSON chunks, and a queue to receive processed data back from each process known as the `feedback_queue`.

The `JsonParseProcess` subclasses python's `Process` [1] class to get around the global interpreter lock, it takes chunks from the JSON queue, processes the chunk and then adds a `DataCollector` class to the feedback queue until all items have been processed. There is substantial room for optimisation in this class, `deepcopy` is utilised liberally as a workaround to prevent race conditions, however this adds substantial overhead to the processing time (almost double) for collection of data.

```

1         with JsonProcessContextManager(data_collector, max_workers) as jtcn:
2             for chunk in self.file_iter:
3                 jtcn.enqueue(chunk)
4                 logger.debug('Finished queueing chunks')

```

Figure 1: Context manager to start Processes and enqueue chunks from the JSON file, inside the ParseFile class.

When the context manager main thread is finished enqueueing JSON chunks it starts a thread to deque and merge the feedback queue back into a single instance of the DataCollector class.

5.2.2 DataCollector

The DataCollector module features 3 convenience classes (ReadingData, BrowserData, and DocLocation) to help with internal data representation in the DataCollector class. The convenience classes all overload the `__add__()` method which is used when the `merge_dict()` function is used, additionally ReadingData and BrowserData overload `__eq__()` and `__lt__()` methods which when combined with the `@total_ordering` decorator from `functools` [2] allows all comparison operations to be performed.

The constructor of the DataCollector class itself has parameters to disable collection of a specific metric if desired, this can be useful to speed up testing, the path parameter is important; it is passed to an instance of ParseFile (see Section 5.2.1). During instantiation the constructor builds a list of methods that will be used by the ParseFile class to process the JSON data once the `gather_data()` method is called.

5.2.3 ComputeData

The ComputeData class contains several helper functions including some sorting functions and functions to get country and continent names from an Alpha2 country code.

5.2.4 Plots

5.3 Gui

5.3.1 main

5.3.2 Tab

5.3.3 Tasks

5.4 Utils

5.4.1 Exceptions

5.4.2 Logging

5.4.3 Tasks

5.4.4 Validation

6 Testing

Show the results for testing all cases and prove that the outputs are what are expected. Preferably, use unit testing to test core functionality of the implementation. If certain conditions cause erroneous results or the application to crash then report these honestly.

7 Personal Development

A short discussion on lessons learnt from the feedback given on CW1 and a discussion how you integrated this feedback into CW2. Cover both coding and report writing, possibly more (project management, preparing for interview style questions etc). Lessons learnt from the experience of CW1: Started out using test driven development

Feedback from CW1: — code — Less code duplication too much global state limited input validation no custom exceptions not restrictive enough access modifiers

— report — intro: should cover short spec cover goals cover env

dev section: should discuss class dependencies should discuss method interfaces := params/return, assumptions of args

conclusion: should discuss adv lang features

8 Conclusions

Reflect on what you are most proud of in the application and what you'd have liked to have done differently. You should reflect on the produced software, and compare software development in scripting vs. systems languages. Most proud of: - Concurrency during file reading - Structure of the code, good decoupling front and back end

Do differently: - Leave more time to work on the gui - Use a different library for the gui - subclass the datacollector class to break it into smaller tasks

A References

References

- [1] (). “Multiprocessing — Process-based parallelism — Python 3.8.6 documentation,” [Online]. Available: <https://docs.python.org/3.8/library/multiprocessing.html?highlight=process#module-multiprocessing> (visited on 12/05/2020).
- [2] (). “Functools — Higher-order functions and operations on callable objects — Python 3.9.1rc1 documentation,” [Online]. Available: <https://docs.python.org/3/library/functools.html> (visited on 12/05/2020).