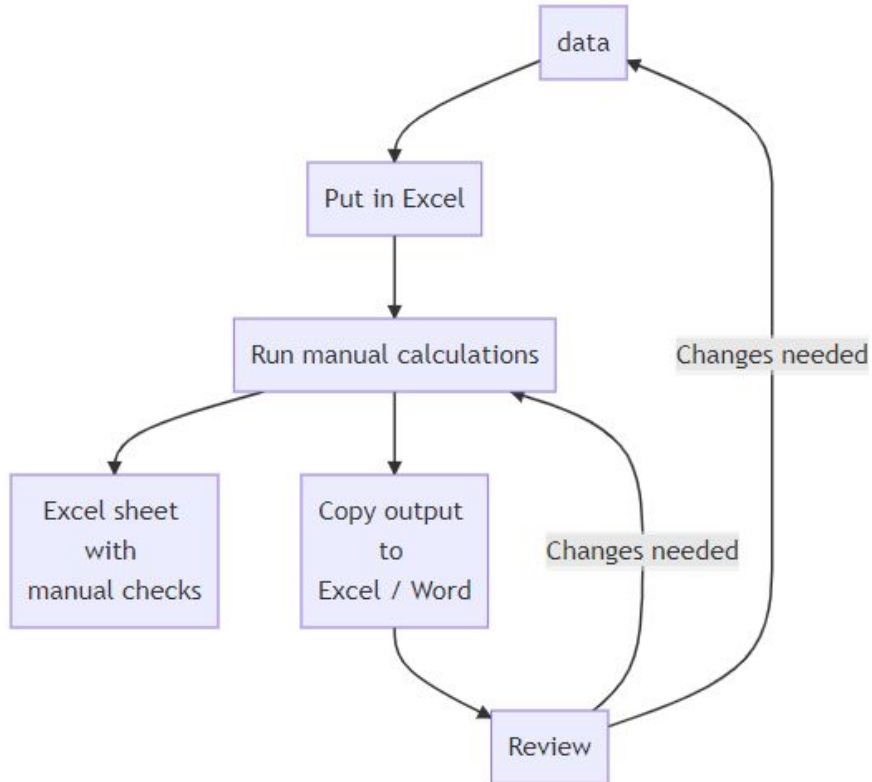


RAPping in the public sector

Shrividya Ravi
Ministry of Transport

Part 1: simple RAP

The ubiquity of bad processes



RAP as a paradigm

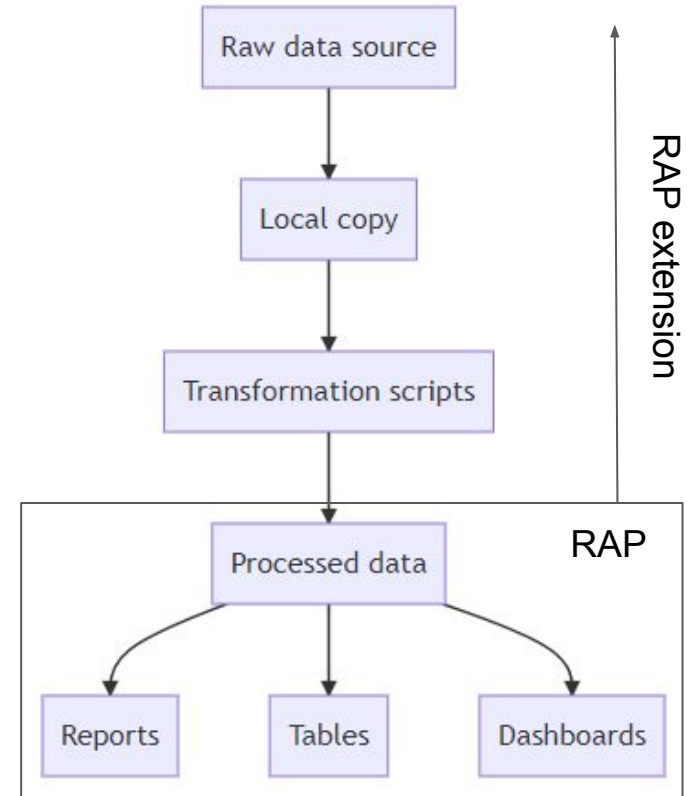
Reproducible Analytical Pipelines (RAPs) are automated statistical and analytical processes. They incorporate elements of software engineering best practice to ensure that the pipelines are reproducible, auditable, efficient, and high quality.

Practices include:

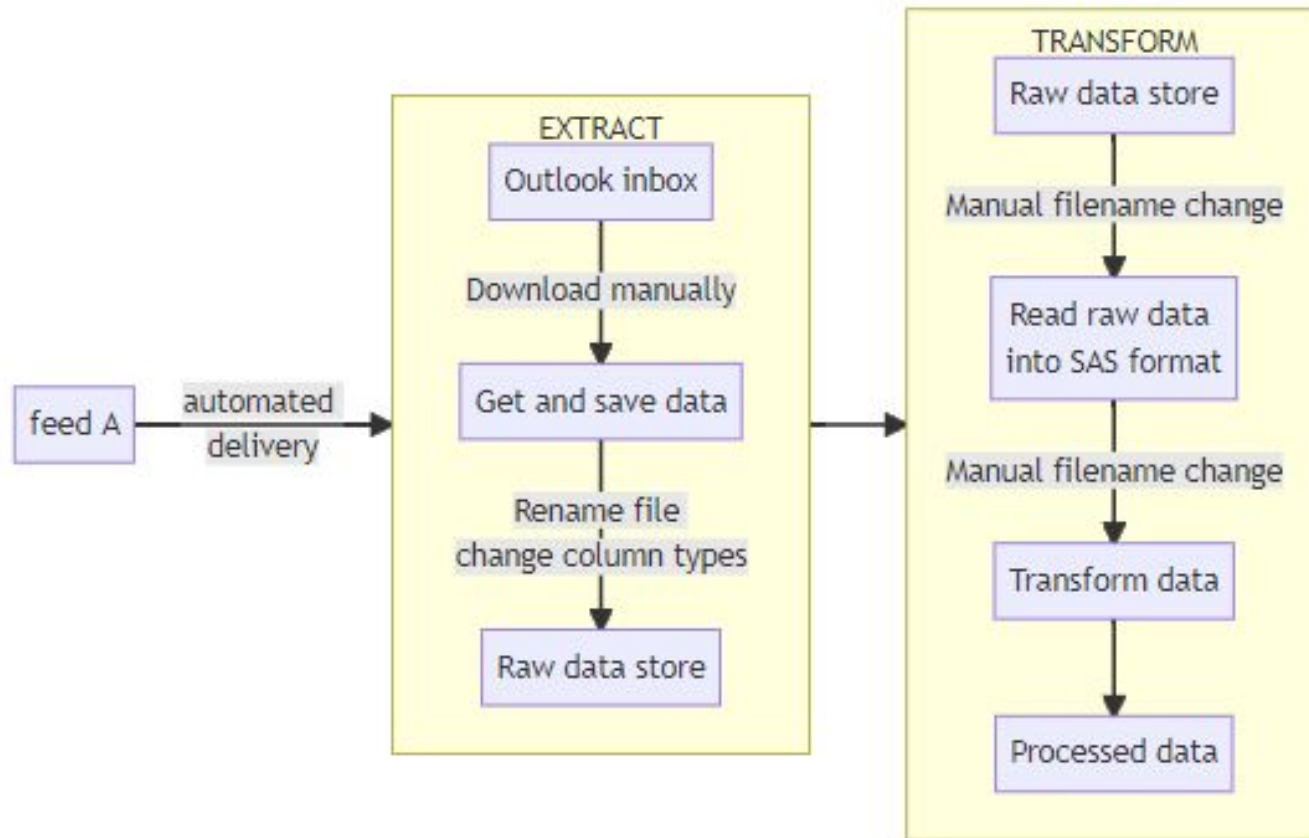
- substituting manual steps with code
- using modern, open source programming languages
- converting raw data to statistical output into pipelines / workflows
- using version control to keep records of development
- bringing in code review practices

RAP to include data engineering

- RAP focuses on converting data from a commonly-managed data store into analytical outputs.
- In infrastructure-poor environments:
 - Data is often not in an accessible place
 - No automated process to transform it in a form fit for subsequent RAPping.
- **Concepts of RAP** need to be brought into the **data engineering domain**.



Sample legacy code



Sample legacy code

- Complexities in the code
 - Non-trivial to understand processing (e.g. unfamiliar paradigm)
 - Non-trivial business rules and algorithms
- Challenging data processing steps
 - Getting raw data from emails
 - Rolling window processing (or lack thereof)
 - Managing latent data

RAPping with legacy code?

- Two approaches:
 - Refactor completely to the RAP paradigm
 - Create interim pipelines that don't much change the code
- Interim pipelines:
 - Allow processing work to be done much faster
 - Reduce the mental overhead of managing manual processes
 - Give more breathing room for developing better solutions

Pipelining legacy code with snake oil

Why python?

- Python is a modern, multi-paradigm, evolving, open source programming language.
- Widely used across many domains - from web development to data science.
- Due to its breadth of use and popularity, there is an incredible ecosystem of packages.



Flask
web development,
one drop at a time

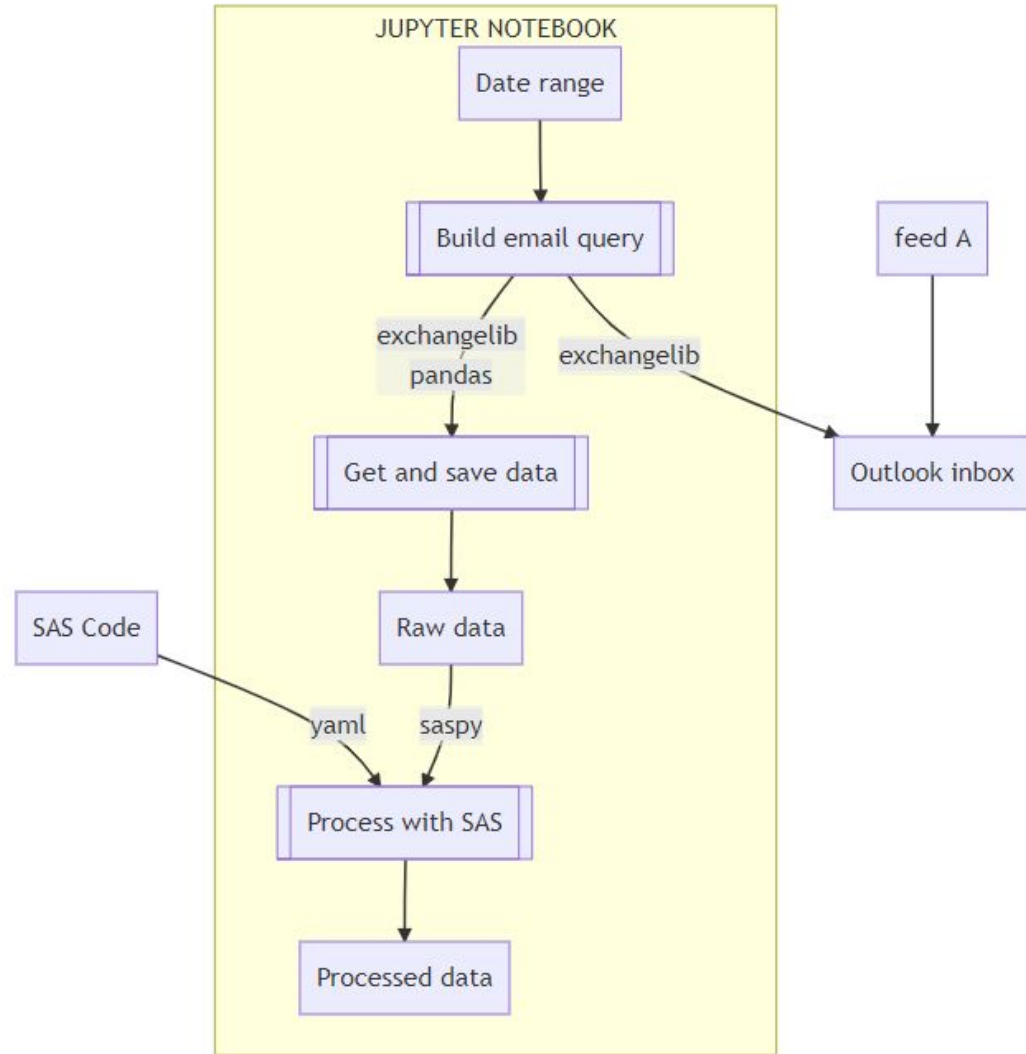


SciPy



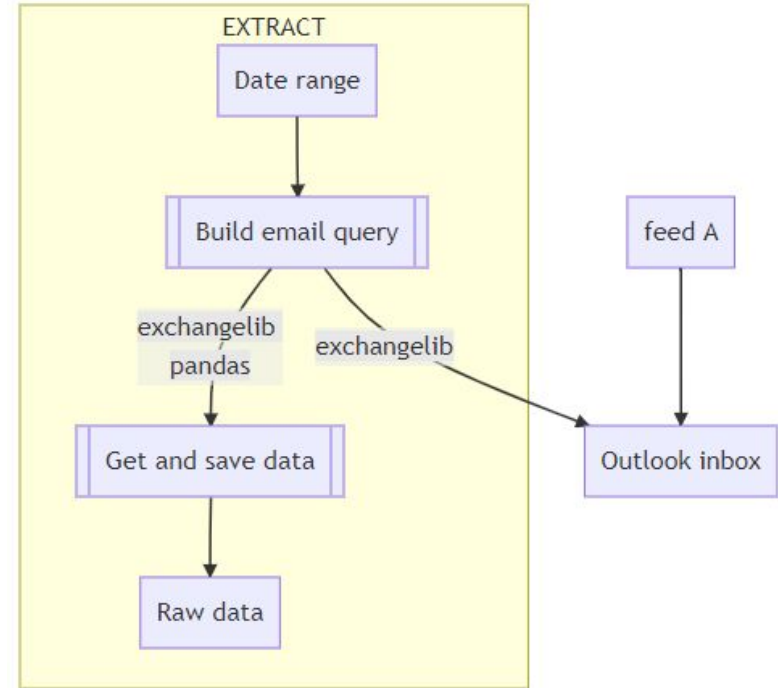
Python as glue

- Instead of manual adjustments steps, Python allows for programmatic “gluing”
- Some of the gluing is facilitated by packages, others just by basic Python functionality



Getting data with APIs - exchangelib

- Application programming interfaces (APIs) allows applications to communicate with each other.
- The Outlook email program has a rich API behind it called Exchange Web Services (EWS).
- Applications (like our RAP extract data application) can send EWS queries to push or pull data to Outlook objects like emails, contacts and calendars.



Code example - exchangelib

```
# Query to get weekly automated emails from POAL
# start, end are pre-defined date ranges
weekly_poal_query = figs_account.inbox.filter(
    datetime_received__range=(start, end),
    sender="BI-DWSupport@poal.co.nz",
    has_attachments=True)

# Go through all results from email query
for i, email in enumerate(weekly_poal_query):
    # loop through all attachment
    for attachment in email.attachments:

        # if attachment, then save file
        if isinstance(attachment, ex.FileAttachment):

            # Filename and data path
            filename = "weekly_poal_data" + date.today() + ".zip"
            local_path = os.path.join(data_path, "zip_files", filename)

            # Save attachment
            with open(local_path, 'wb') as f:
                f.write(attachment.content)
```

Running SAS through Python

There are two ways or running SAS outside the SAS program:

- With a SAS kernel in Jupyter
- Through SASPy

SASPy is officially supported by [SAS](#) and, available as an [open source package](#).
The library seems to be well-maintained and well-documented.

Running SAS through Python

At its core, SASPy is capable of **creating a SAS session and sending code to it for execution**, as well as returning the output (logs and other output) to the controlling Python script. Yet it is also a **powerful generator of SAS code**, which means that it **offers methods, objects, and syntax for use directly in idiomatic Python that it can then automatically convert to the appropriate SAS language** statements for execution. In most cases, SAS procedures or steps are mapped directly to Python methods as a one-to-one equivalent.

Value of running SAS through Python

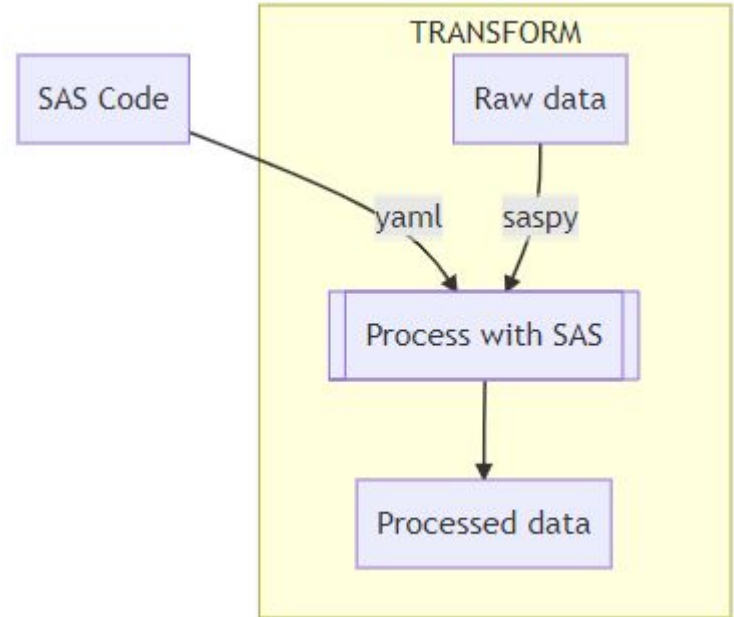
Incremental refactoring (pain management) principle

- When the (SAS) codebase is complex, large or both, it's convenient to instead just make it run easier.
- This approach allows to incremental refactoring - allowing hard-to-convert code to remain in SAS while moving easier code to Python.

Breaking up SAS code into chunks

To run existing SAS scripts, via Python, where parts need manual amendments needs an extra step. The easiest solution so far has three main steps:

- Breaking up a SAS script into yaml chunks for "immutable" components
- "Mutable" components are created in python
- The immutable and mutable are brought together with Python's f-strings



[yaml/pyyaml: Canonical source repository for PyYAML \(github.com\)](#)

[Python 3's f-Strings: An Improved String Formatting Syntax \(Guide\) – Real Python](#)

Breaking up SAS code into chunks - code

```
# open list of new data filenames after downloading from Outlook
files_list = open('..\sas_code\combine\files_list.sas').read()

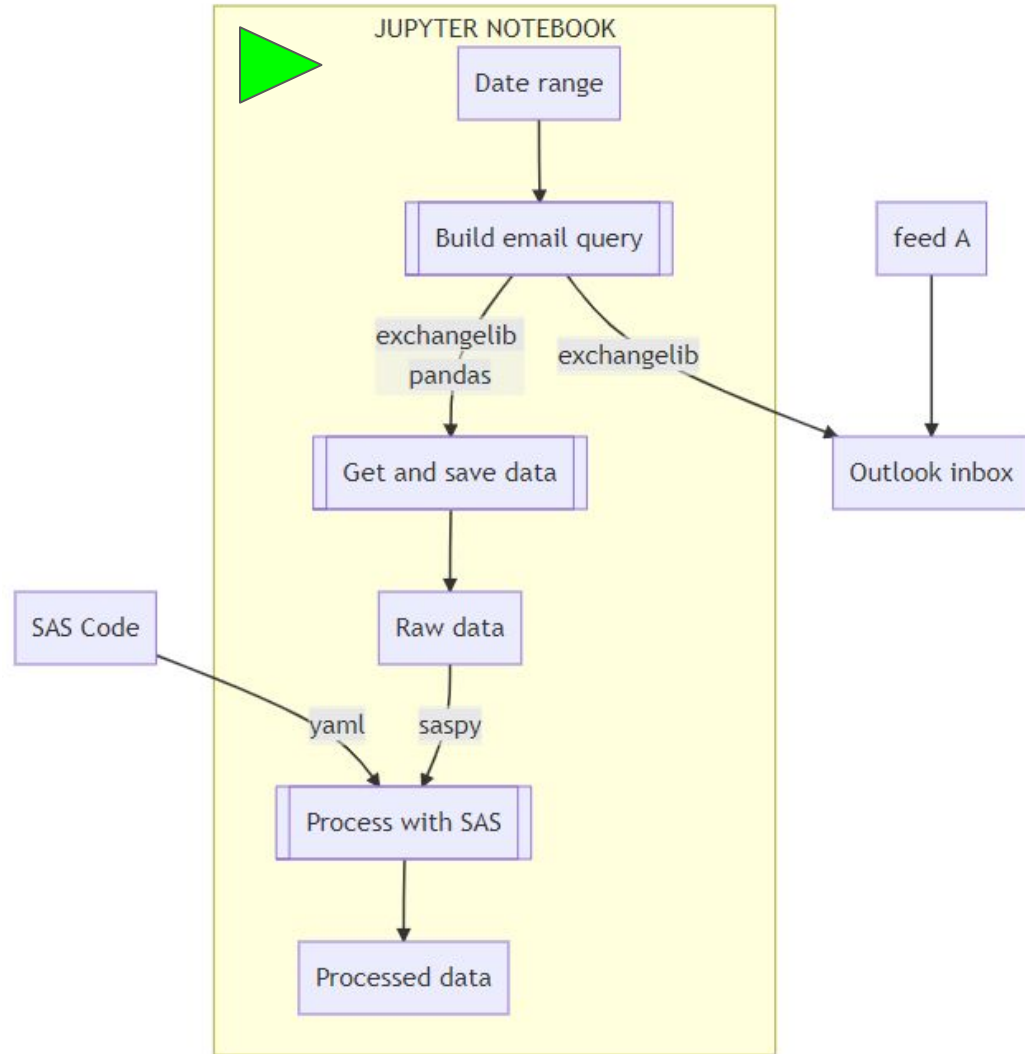
# open data processing code chunks
with open(r'..\sas_code\combine\poal_combine_sas.yaml') as file:
    sas_code = yaml.full_load(file)

# create code blocks by function
preamble = sas_code['preamble']
coarri_combine = sas_code['coarri_combine']
codeco_combine = sas_code['codeco_combine']

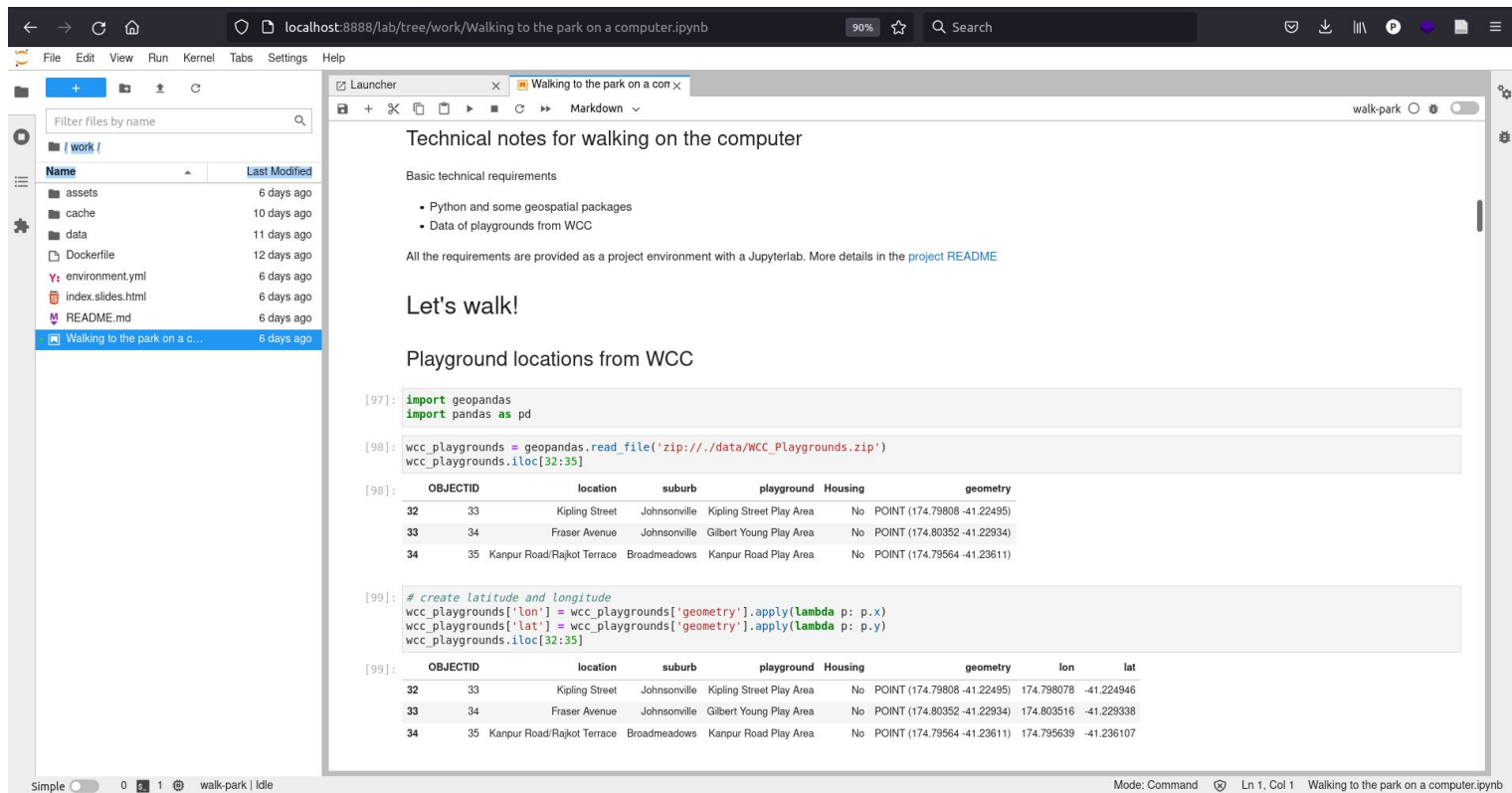
# create sas query - concat code blocs
sas_query = f"""{preamble}{files_list}{codeco_combine}"""
res = sas.submit(sas_query)
```

One-click linear pipelines

- Put code of pipeline steps into cells of a Jupyter notebook
- Run all steps with “one click”
- Jupyter notebooks can print outputs, checks, visualisations in a clean IDE interface



Jupyterlab IDE and Jupyter notebooks



The screenshot displays the JupyterLab IDE interface. On the left is a file explorer showing a directory structure with files like `assets`, `cache`, `data`, `Dockerfile`, `environment.yml`, `index.slides.html`, `README.md`, and `Walking to the park on a computer.ipynb`. The main area shows a notebook titled "Walking to the park on a computer.ipynb". The notebook content includes:

- Technical notes for walking on the computer
- Basic technical requirements
 - Python and some geospatial packages
 - Data of playgrounds from WCC
- All the requirements are provided as a project environment with a Jupyterlab. More details in the [project README](#)
- Let's walk!
- Playground locations from WCC

The notebook shows two code cells. The first cell imports `geopandas` and `pandas` as `pd`. The second cell reads a file and displays a table of playground locations:

```
[97]: import geopandas
import pandas as pd

[98]: wcc_playgrounds = geopandas.read_file('zip://./data/WCC_Playgrounds.zip')
wcc_playgrounds.iloc[32:35]
```

OBJECTID	location	suburb	playground	Housing	geometry	
32	33	Kipling Street	Johnsonville	Kipling Street Play Area	No	POINT (174.79808 -41.22495)
33	34	Fraser Avenue	Johnsonville	Gilbert Young Play Area	No	POINT (174.80352 -41.22934)
34	35	Kanpur Road/Rajkot Terrace	Broadmeadows	Kanpur Road Play Area	No	POINT (174.79564 -41.23611)

The third cell creates latitude and longitude columns:

```
[99]: # create latitude and longitude
wcc_playgrounds['lon'] = wcc_playgrounds['geometry'].apply(lambda p: p.x)
wcc_playgrounds['lat'] = wcc_playgrounds['geometry'].apply(lambda p: p.y)
wcc_playgrounds.iloc[32:35]
```

The final output shows the table with `lon` and `lat` columns added:

OBJECTID	location	suburb	playground	Housing	geometry	lon	lat	
32	33	Kipling Street	Johnsonville	Kipling Street Play Area	No	POINT (174.79808 -41.22495)	174.798078	-41.224946
33	34	Fraser Avenue	Johnsonville	Gilbert Young Play Area	No	POINT (174.80352 -41.22934)	174.803516	-41.229338
34	35	Kanpur Road/Rajkot Terrace	Broadmeadows	Kanpur Road Play Area	No	POINT (174.79564 -41.23611)	174.795639	-41.236107

The bottom status bar shows "Simple" mode, "0" files, "1" notebook, "walk-park | Idle", and "Mode: Command | Ln 1, Col 1 | Walking to the park on a computer.ipynb".

Jupyterlab IDE

The screenshot displays the Jupyterlab IDE interface. The top bar shows the file path: `localhost8888/lab/tree/work/utills/load_data.py`. The left sidebar contains a file explorer with a search bar and a list of files: `corrections...` (4 days ago), `load_data.py` (2 minutes ago), `s3_read_wri...` (6 days ago), and `transform.py` (31 minutes ago). The main area is divided into two panes. The left pane shows the code editor for `load_data.py`, which contains Python code for cleaning and mapping trip data. The right pane shows the console output, displaying logs and the execution of various data processing steps.

```
119 # Add 'sa' zones to trips OD
120 #PAH_tractiv,
121 'PAH_purpose', 'trleave', 'trarriv', 'jstaddno',
122 'traddno', 'PAH_trmode', 'trip_n']
123 trips.columns = ['person_key', 'hh_key', 'tripday', 'trip_n',
124 'Year', 'samo', 'person', 'pump', 'tst', 'tet',
125 'ozone', 'dzone', 'mode', 'seq']
126 log.info("Trips data created!")
127 return(trips)
128
129 def clean_trips(trips):
130     # map other to car
131     def mapper(inp):
132         return ("other": "car").get(inp, inp)
133     trips["mode"] = trips["mode"].apply(mapper)
134
135     # Remove NA-purpose trips
136     log.info("Removing {} trips (out of {}) with NA purpose".format(
137         sum(trips['pump']=='NA')==sum(trips['pump'].isna()), len(trips))
138     )
139     trips = trips[~((trips['pump']=='NA')|(trips['pump'].isna()))]
140
141     # Remove NA modes
142     log.info("Removing {} trips (out of {}) with NA mode".format(
143         sum(trips['mode']=='NA')==sum(trips['mode'].isna()), len(trips))
144     )
145     trips = trips[~((trips['mode']=='NA')|(trips['mode'].isna()))]
146
147     # add tripday to uid so that days are unique
148     trips["pid"] = ["{}({}person_key)-{}day)".format(person_key, day) for person_key, day in zip(trips.person_key,
149     trips.tripday)]
150     trips["hid"] = ["{}({}hh_key)-{}day)".format(hh_key, day) for hh_key, day in zip(trips.hh_key, trips.tripday)]
151
152     return(trips)
153
154 def remap_zones(clean_trips, sa_mapping):
155     # Add 'sa' zones to trips OD
156     log.info("Mapping zones")
157     clean_trips.ozone = clean_trips.ozone.map(sa_mapping).map(int)
158     clean_trips.dzone = clean_trips.dzone.map(sa_mapping).map(int)
159     return clean_trips
160
161
162 def expand_diary_by_sampling_days(trips, sample, expanded_idx, target_idx):
163     expanded_sample = trans.expand_optimised(
164         trips,
165         persons,
166         expanded_idx=expanded_idx,
167         target_idx=target_idx,
168
```

Console output:

```
[36]: clean_trips = trips.head(1000).pipe(data.create_trips).pipe(data.clean_trips, sa_mapping)
27-Jul-21 03:46:59 - load_data - INFO - Indexing legs
27-Jul-21 03:46:59 - load_data - INFO - Simplifying legs as trips
27-Jul-21 03:46:59 - load_data - INFO - Dropping columns and renaming to PAH schema
27-Jul-21 03:46:59 - load_data - INFO - Trips data created!
Removing 0 trips (out of 965) with NA purpose
Removing 0 trips (out of 965) with NA mode

[37]: reload(data)
[37]: <module 'utils.load_data' from '/home/jovyan/work/utills/load_data.py'>

[38]: clean_trips = trips.head(1000).pipe(data.create_trips).pipe(data.clean_trips, sa_mapping)
27-Jul-21 03:47:52 - load_data - INFO - Indexing legs
27-Jul-21 03:47:52 - load_data - INFO - Simplifying legs as trips
27-Jul-21 03:47:52 - load_data - INFO - Dropping columns and renaming to PAH schema
27-Jul-21 03:47:52 - load_data - INFO - Trips data created!
27-Jul-21 03:47:52 - load_data - INFO - Removing 0 trips (out of 965) with NA purpose
27-Jul-21 03:47:52 - load_data - INFO - Removing 0 trips (out of 965) with NA mode
27-Jul-21 03:47:52 - load_data - INFO - Mapping zones

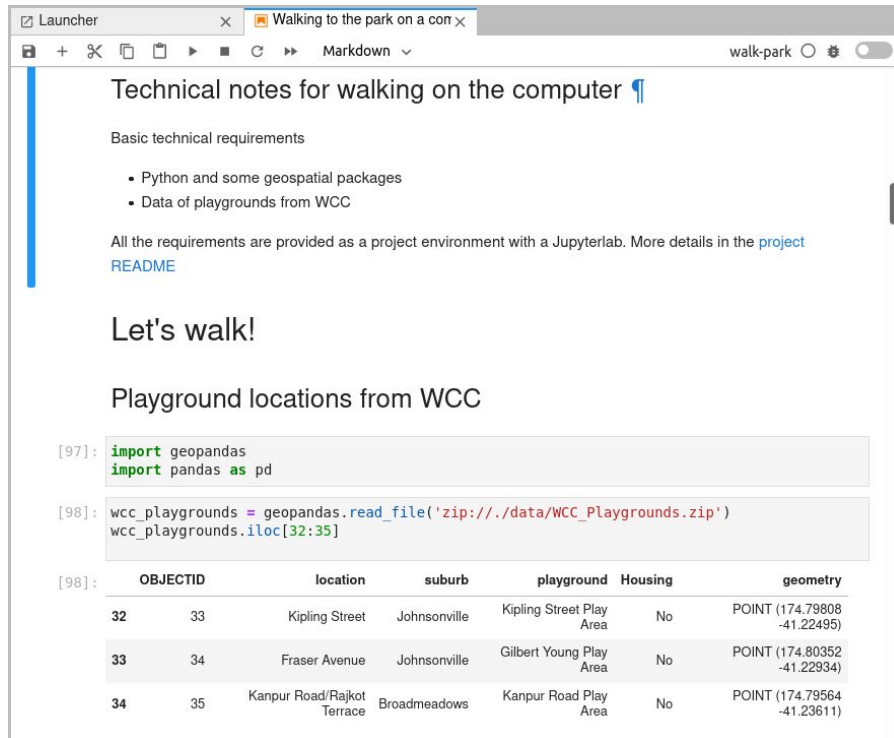
[39]: trips = data.read_hts_legs(tr_path)
27-Jul-21 03:48:40 - load_data - INFO - Loading HTS trips data
27-Jul-21 03:48:47 - load_data - INFO - HTS trips loaded as MATSIM legs

[40]: clean_trips = trips.pipe(data.create_trips).pipe(data.clean_trips, sa_mapping)

[ ]:
```

Jupyter to bind, execute and document

- Jupyter notebook cells can be executed in serial with the 'Run All Cells'
- Cells can be either code or markdown.
- Documenting process next to code is easy and can be supported by table outputs, graphs and even external images for explanation.



Launcher x Walking to the park on a computer x

Basic technical requirements

- Python and some geospatial packages
- Data of playgrounds from WCC

All the requirements are provided as a project environment with a Jupyterlab. More details in the [project README](#)

Let's walk!

Playground locations from WCC

```
[97]: import geopandas
import pandas as pd
```

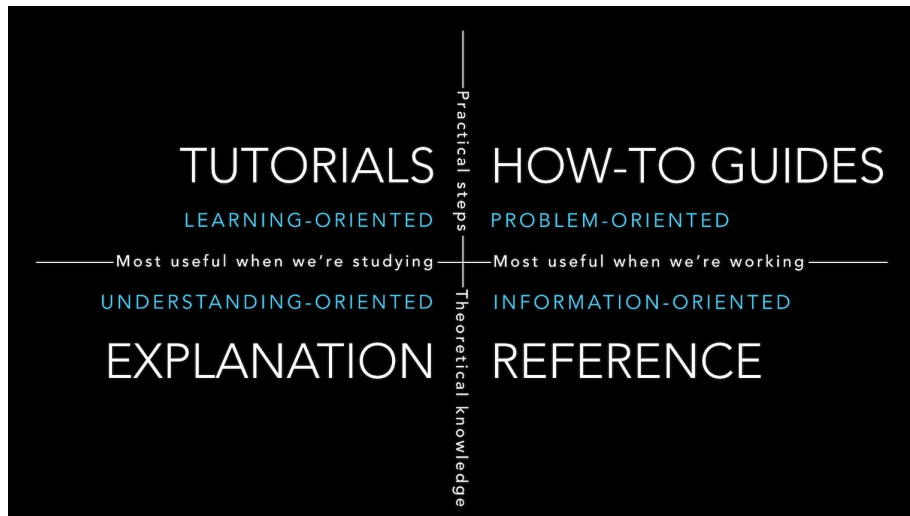
```
[98]: wcc_playgrounds = geopandas.read_file('zip:///data/WCC_Playgrounds.zip')
wcc_playgrounds.iloc[32:35]
```

	OBJECTID	location	suburb	playground	Housing	geometry
32	33	Kipling Street	Johnsonville	Kipling Street Play Area	No	POINT (174.79808 -41.22495)
33	34	Fraser Avenue	Johnsonville	Gilbert Young Play Area	No	POINT (174.80352 -41.22934)
34	35	Kanpur Road/Rajkot Terrace	Broadmeadows	Kanpur Road Play Area	No	POINT (174.79564 -41.23611)

Part 2: extensions

Good code *and* documentation

- Documenting alongside code helps explain the particular steps in the pipeline
- But, good documentation exists for several use-cases
- A ‘website’ that contains all the relevant documentation explaining process, setup, tutorials and how-tos can be helpful.



Git(hub/lab) Pages with bookdown for documentation

Technical Monty

1 Motivation

- 1.1 Technical setup

2 Introduction

3 PopulationSim methodology overview

- 3.1 Inputs
- 3.2 Objective
- 3.3 Algorithms
- 3.4 Running PopulationSim
- 3.5 Current configurations

4 Freight

- 4.1 General notes about freight
- 4.2 Approaches
- 4.3 New Zealand specific resources

References

Published with bookdown

Technical Monty - documentation of Monty ABM

Ministry of Transport (NZ)

Chapter 1 Motivation

This bookdown project was motivated by straggling documentation across Gitlab and Github and, the need for us to be able to present Work In Progress reports of the Monty project. It is anticipated that technical understanding of the ABM's setup, components and features will be written up here. There is no other design at present but as the documentation builds up, we will set up some principles around style, intended audience, structure etc. (To come!)

1.1 Technical setup

This documentation is written in **Markdown**. You can use anything that Pandoc's Markdown supports, e.g., a math equation $a^2 + b^2 = c^2$.

The **bookdown** package can be installed from CRAN or Github:

```
install.packages("bookdown")  
# or the development version  
# devtools::install_github("rstudio/bookdown")
```

[Home](#) | [Bookdown](#)

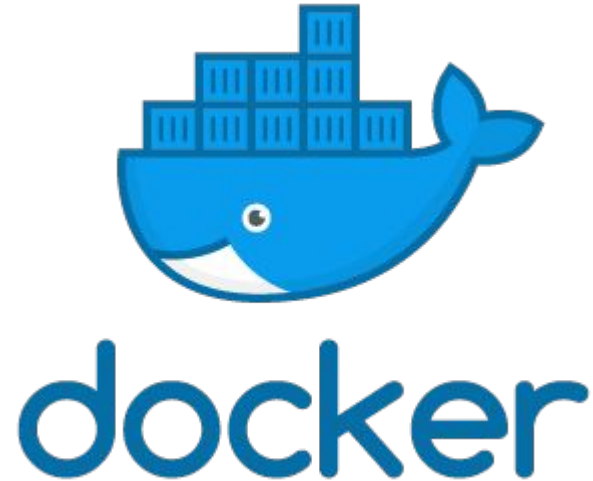
[GitLab Pages](#) | [GitLab](#)

[GitHub Pages](#) | Websites for you and your projects, hosted directly from your GitHub repository. Just edit, push, and your changes are live.

Creating applications with docker

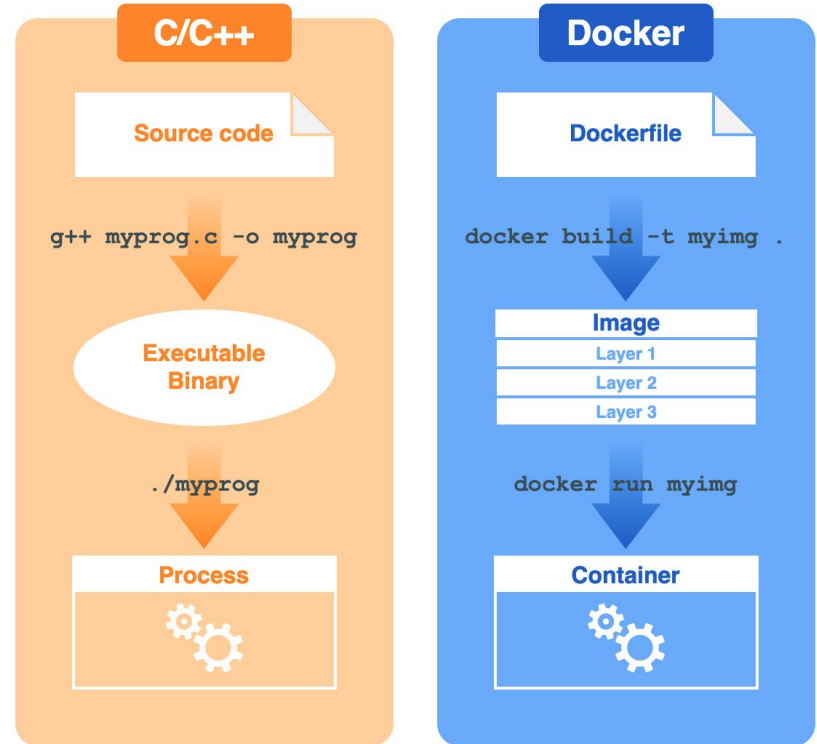


- Images are
 - described by a Dockerfile
 - can be built on demand
 - stored in a common registry e.g. docker, gitlab, Amazon ECR
- Containers can run:
 - IDEs
 - Scripts
 - Services (like a database)



Creating applications with docker

- Like compiling source code for a programming language, creating a container also starts with a plain text file (Dockerfile)
- Similar to using a compiled binary file to launch a program, the image is then run to create a container instance.

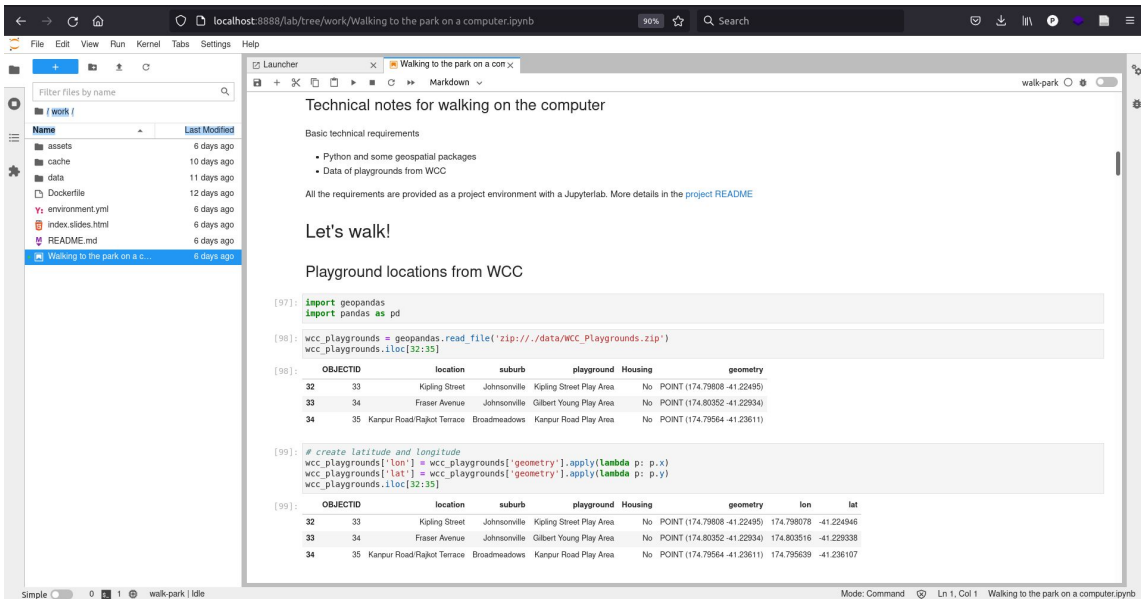


© Benjamin D. Evans, 2020. This work is released under the CC-BY 4.0 license <https://creativecommons.org/licenses/by/4.0/>

Dockerfile to container running an IDE

```
Dockerfile x
Dockerfile
1 # Choose your desired base image
2 FROM jupyter/minimal-notebook:latest
3
4 # name your environment and choose python 3.x version
5 ARG conda_env=walk-park
6
7 # Create the conda environment within the image
8 COPY environment.yml .
9 RUN conda env create -f environment.yml
10
11 # create Python 3.x environment and link it to jupyter
12 RUN $CONDA_DIR/envs/${conda_env}/bin/python -m ipykernel install --user --name=${conda_env} && \
13     fix-permissions $CONDA_DIR && \
14     fix-permissions /home/$NB_USER
15
16 # prepend conda environment to path
17 ENV PATH $CONDA_DIR/envs/${conda_env}/bin:$PATH
18
19 # if you want this environment to be the default one, uncomment the following line:
20 ENV CONDA_DEFAULT_ENV ${conda_env}
21
```

```
docker run -d -p 8888:8888 \
-v ${PWD}:/home/jovyan/work \
-e JUPYTER_ENABLE_LAB=yes \
-e CHOWN_HOME_OPTS='-R' \
-e CHOWN_HOME=yes \
maptime-test \
start.sh jupyter lab --LabApp.token=''
```



Technical notes for walking on the computer

Basic technical requirements

- Python and some geospatial packages
- Data of playgrounds from WCC

All the requirements are provided as a project environment with a Jupyterlab. More details in the [project README](#)

Let's walk!

Playground locations from WCC

```
[97]: import geopandas
import pandas as pd

[98]: wcc_playgrounds = geopandas.read_file('zip:///data/WCC_Playgrounds.zip')
wcc_playgrounds.iloc[32:35]
```

OBJECTID	location	suburb	playground	Housing	geometry
32	33	Kipling Street	Johnsonville	Kipling Street Play Area	No POINT (174.79808 -41.22495)
33	34	Fraser Avenue	Johnsonville	Gilbert Young Play Area	No POINT (174.80352 -41.22934)
34	35	Karupur Road/Rajkot Terrace	Broadmeadows	Karupur Road Play Area	No POINT (174.79564 -41.23611)

```
[99]: # create latitude and longitude
wcc_playgrounds['lon'] = wcc_playgrounds['geometry'].apply(lambda p: p.x)
wcc_playgrounds['lat'] = wcc_playgrounds['geometry'].apply(lambda p: p.y)
wcc_playgrounds.iloc[32:35]
```

OBJECTID	location	suburb	playground	Housing	geometry	lon	lat
32	33	Kipling Street	Johnsonville	Kipling Street Play Area	No POINT (174.79808 -41.22495)	174.798078	-41.224948
33	34	Fraser Avenue	Johnsonville	Gilbert Young Play Area	No POINT (174.80352 -41.22934)	174.803516	-41.229338
34	35	Karupur Road/Rajkot Terrace	Broadmeadows	Karupur Road Play Area	No POINT (174.79564 -41.23611)	174.795639	-41.236107

Dockerfile to container running a script

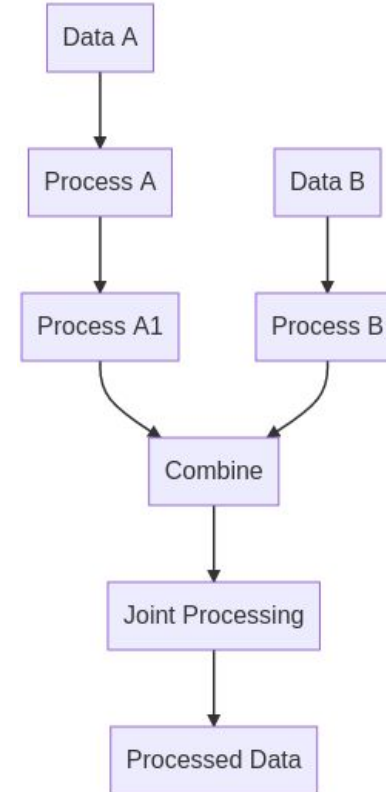
- Run a script as a self-contained application
- Pass in arguments to docker run command

```
docker run -v ${PWD}:/tmp rail-etl \
--credentials_file='/tmp/credentials.yaml' \
--date_start='2021-06-14' \
--date_end='2021-07-15' \
--output_path='/tmp'
```

```
Dockerfile
1 FROM continuumio/miniconda3
2
3 WORKDIR /app
4
5 # Create the conda environment within the image
6 COPY environment.yml .
7 RUN conda env create -f environment.yml
8
9 # Make RUN commands use the new environment:
10 SHELL ["conda", "run", "-n", "rail-etl", "/bin/bash", "-c"]
11
12 # change? create a subfolder with app?
13 COPY etl /app
14
15 # # The code to run when container is started:
16 ENTRYPOINT ["conda", "run", "--no-capture-output", "-n", "rail-etl", "python", "etl_main.py"]
17
```

Linear to DAG workflows

- Directed Acyclic Graphs are very common flow structures in data processing
- Until outputs interact, task execution can be managed as parallel steps




Orchestrating DAGs

Two main approaches

- **The “Data Scientist” route:** local, often within codebase
 - kedro - Python
 - targets (formerly drake) - R
 - GNU make
 - Popper
- **The “Data Engineer” route:** cloud-deployable, scalable, schedulable
 - Airflow
 - Luigi
 - AWS Step Functions

Basic container-native workflows with popper

- A container-native task automation engine
- Runs on distinct container engines
- Orchestration framework
- Simple YAML files to describe workflows



```
1 steps:
2 # download CSV file with data on global CO2 emissions
3 - id: download
4   uses: docker://byrnedo/alpine-curl:0.1.8
5   args: [-L, git.io/JUcRU, -o, global.csv]
6
7 # obtain the transpose of the global CO2 emissions table
8 - id: get-transpose
9   uses: docker://getpopper/csvtool:2.4
10  args: [transpose, global.csv, -o, global_transposed.csv]
```


Other topics being explored

- Parquet / feather for Persistent Staging Areas (PSAs)
- SQLite for indices or data needing appends
- Rolling window processing with Parquet

Thank you!