

# You do you.

How next-gen data platforms can stop weather and climate scientists from being software engineers and other perversions

**Theo McCaie** - Data Engineering Research Lead. UK Met Office and Joint Center for Excellence in Environmental Intelligence



80% of data science is...

...formatting

...data munging.

...cleaning

...wrangling

...validation

...etc

...discovery



# 80% of data science is data wrangling

---



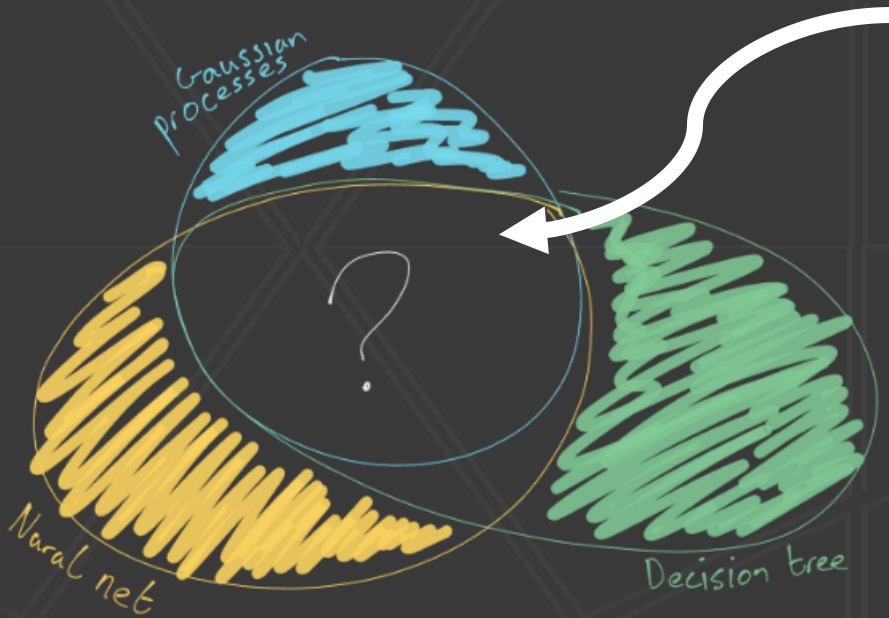
Data wrangling:

- Acknowledged as data science
- But oft distanfully
- But if you could make one part of this diagram 50% more efficient which would have most impact...



# More the same than different

---



Most machine learning approaches share many of the same foundations

- Access to clean, organised, data
- Pipelines to format data
- Compute infrastructure to run on
- Running many experiments e.g. hyperparameter tuning

This is the realm of Data Engineering



Better data engineering will be the most effective accelerant of data science.





# STAC

---



## SpatioTemporal Asset Catalogs:

- “a common language... more easily be indexed and discovered”
- “new code doesn't need to be written whenever a new data set or API is released.”
- Seeks to remove burdand of building unique pipelines, spur common tooling.





## SpatioTemporal Asset Catalog

# Welcome to the STAC Index!

Here you can find STAC Catalogs, Collections, APIs, Software and Tools.  
You can also add your own data and tools to the list.

You don't know STAC yet? Check it out at [stacspec.org](https://stacspec.org).

Catalogs

Collection Search

Ecosystem

### Recently added

#### Catalogs

[data.geo.admin.ch](https://data.geo.admin.ch)  

Data Catalog of the Swiss Federal Spatial Data Infrastructure

<https://data.geo.admin.ch/api/stac/v0.9/>

#### Ecosystem

EODAG 

EODAG is a CLI tool and a Python framework for searching, aggregating results and downloading EO data through a unified API regardless of the data provider. It can be run as STAC client



hello@informaticslab.co.uk



@informatics\_lab



informaticslab.co.uk

```
In [1]: import nc_ds_utils
import numpy as np
import iris
import datetime
from time import mktime
import os
```

STASH  
source

Batch by time

```
else:
    # Throw error
    if c < 0 or c > 9:
        raise ValueError
    # Converting
    if c <= 9:
        # Just co
```

```
In [8]: start, stop = next(iter(time_window_batch((1851, 1, 1, 0, 0), (1861, 1, 1, 0, 0), datetime.timedelta(hours=24*50))))
filter_files_for_dates(FILE, start, stop)
```

```
Out[8]: ['/data/cssp-china/mini-dataset-24-01-19/nc/hourly/apepda.pj51150.nc',
'/data/cssp-china/mini-dataset-24-01-19/nc/hourly/apepda.pj511f0.nc',
'/data/cssp-china/mini-dataset-24-01-19/nc/hourly/apepda.pj511g0.nc',
'/data/cssp-china/mini-dataset-24-01-19/nc/hourly/apepda.pj511h0.nc',
'/data/cssp-china/mini-dataset-24-01-19/nc/hourly/apepda.pj511i0.nc',
'/data/cssp-china/mini-dataset-24-01-19/nc/hourly/apepda.pj511j0.nc']
```

## Mapping the work so there is no overlap

Iris can not open the same HDF5 file (NetCDF) if the tasks without overlapping the files being accessed.

```
In [11]: TaskSet = namedtuple('TaskSet', ['index', 'start', 'stop', 'files'])
```

```
def create_task_set(i, start, stop):
    return TaskSet(i, start, stop, files)
```

```
windows = time_window_batch((1851, 1, 1, 0, 0), (1861, 1, 1, 0, 0),
                             step=datetime.timedelta(hours=24*50),
                             overlap=datetime.timedelta(hours=24*50))
```

```
tasks = [create_task_set(i, start, stop) for i, (start, stop) in enumerate(windows)]
```

```
In [12]: def do_overlap(task, other_tasks):
for other_task in other_tasks:
    for file in task.files:
        if file in other_task.files:
            return True
    else:
        return False
```

```
task_sets = []
for task in tasks:
    for task_set in task_sets:
        if not do_overlap(task, task_set):
            task_set.append(task)
            break
    else:
        task_sets.append([task])
len(task_sets)
```

Out[12]: 2

## Process the tasksets on distributed cluster

```
In [13]: Result = namedtuple('Result', ['success', 'value'])
```

```
def error_catch_rmean(start, stop):
    try:
        result = rolling_mean_at_point_for_date_range(start, stop)
        return Result(True, result)
    except Exception as e:
        return Result(False, e)
```

```
In [14]: DataSet = namedtuple('DataSet', ['data', 'times'])
```

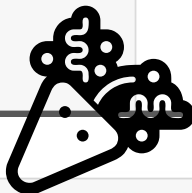
```
def to_datasets(i, result):
    if result.success:
        with dask.config.set(scheduler='distributed'):
            cube = result.value
            data = cube.data
            time_unit = cube.coord('time').units
            times_raw = cube.coord('time').times
            times = [datetime.datetime.strptime(t, time_unit) for t in times_raw]
        return DataSet(data, times)
    else:
        return result
```

```
In [15]: import distributed
from dask_kubernetes import KubeCluster
import dask.bag as db
```

```
In [16]: cluster = KubeCluster()
cluster.scale_up(20)
display(cluster)
```

```
VBox(children=(HTML(value='<h2>KubeCluster: 20 workers</h2>'),
                HTML(value='<pre>...</pre>'),
                ...))
```

```
In [19]: %%time
collections = []
for task_set in task_sets:
```

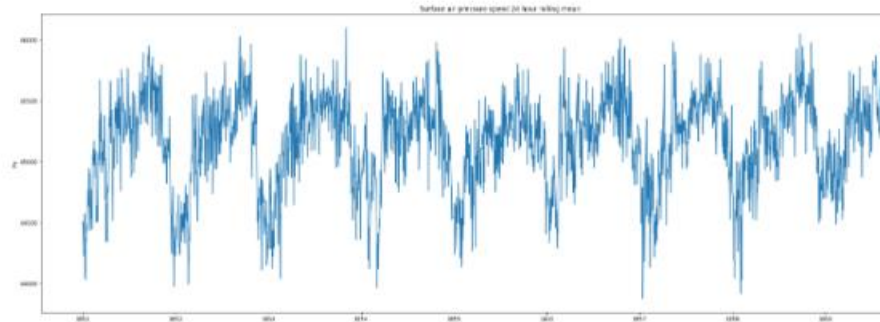


## Sort results and plot

```
In [18]: %%time
sorted_ds = sorted([ds for ds in chain(*collections) if isinstance(ds, DataSet)], key=lambda ds: ds.times)
datas, times = zip(*sorted_ds)
plt.title("Surface air pressure speed 24 hour rolling mean")
plt.ylabel("sample_cube.units")
plt.plot(list(chain(*times)), list(chain(*datas)))
```

CPU times: user 469 ms, sys: 4.55 ms, total: 474 ms  
Wall time: 467 ms

```
Out[18]: [<matplotlib.lines.Line2D at 0x7fe66aeb9210>]
```

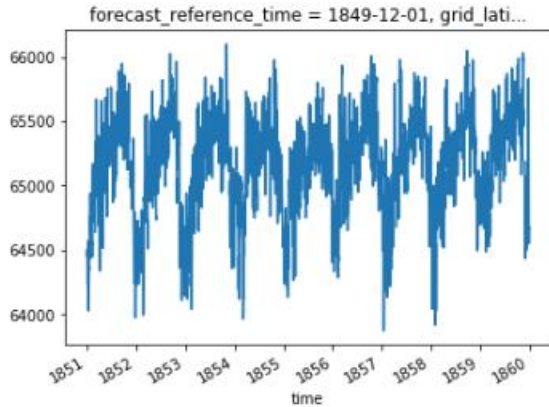




```
In [2]: import zarr
import xarray as xr
%matplotlib inline
```

```
In [12]: zarr_store = zarr.storage.ABSStore('cssp-china', prefix='zarr_hourly_1851-1859', account_name=account, account_key=key)
ds = xr.open_zarr(zarr_store)
ds.surface_air_pressure[:,100,100].rolling(time=24, center=False).mean().plot()
```

Out[12]: [`<matplotlib.lines.Line2D at 0x7f70f056e950>`]



# Analysis Ready Data

---

- Analysis Ready data is a dataset that is trivial to load and analyse
- ARCO (Analysis Ready Cloud Optimized) takes this to the next step ensuring that the data is efficiently to use on public cloud infrastructure.
- ARCO applies to a whole dataset (many thousands of millions of files/blobs/chunks)
- Data provider does most of the cleaning, amalgamating, indexing, etc
- Some example formats: Zarr, TileDB



# Pangeo

---

- Open community
- Open source software
- Open source infrastructure
- Have working groups on ML amongst others

*"Pangeo is first and foremost a community of people working collaboratively to develop software and infrastructure to enable Big Data geoscience research."*

<https://pangeo.io/>



hello@informaticslab.co.uk



@informatics\_lab

<https://discourse.pangeo.io/>



informaticslab.co.uk

The image shows a JupyterLab environment with several open notebooks. The left notebook, titled 'Table of Contents', contains code for setting up a Kubernetes cluster and connecting to a distributed system. The right notebook, titled 'Table of Contents', displays a dashboard with two charts: a 'Task Status' chart showing a distribution of tasks over time, and a 'Table of Contents' chart showing a list of tasks with their durations.

**Table of Contents (Left Notebook)**

- 1. Create and connect to a distributed cluster
- 2. Load data into xerney from an intake catalog
- 3. Create a distributed job
- 4. Run a distributed job
- 5. Monitor the job
- 6. Clean up the cluster

**Table of Contents (Right Notebook)**

- 1. Create a distributed job
- 2. Run a distributed job
- 3. Monitor the job
- 4. Clean up the cluster

**Task Status Chart (Right Notebook)**

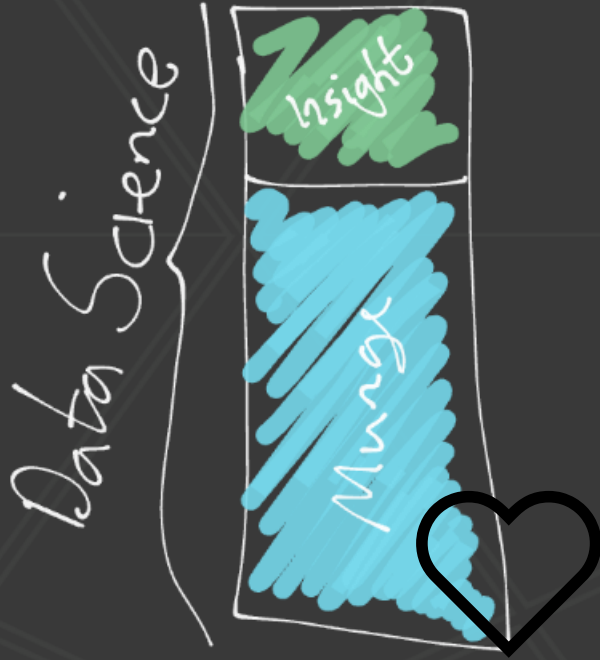
Task Name	Start Time	End Time	Duration
Task 1	10:00	10:05	5m
Task 2	10:05	10:10	5m
Task 3	10:10	10:15	5m
Task 4	10:15	10:20	5m
Task 5	10:20	10:25	5m
Task 6	10:25	10:30	5m
Task 7	10:30	10:35	5m
Task 8	10:35	10:40	5m
Task 9	10:40	10:45	5m
Task 10	10:45	10:50	5m

**Table of Contents Chart (Right Notebook)**

Task Name	Start Time	End Time	Duration
Task 1	10:00	10:05	5m
Task 2	10:05	10:10	5m
Task 3	10:10	10:15	5m
Task 4	10:15	10:20	5m
Task 5	10:20	10:25	5m
Task 6	10:25	10:30	5m
Task 7	10:30	10:35	5m
Task 8	10:35	10:40	5m
Task 9	10:40	10:45	5m
Task 10	10:45	10:50	5m

# Conclusions

---



- Give love to data wrangling
- Embrace and engage with data engineering efforts
- This will accelerate the data science revolution faster than anything else





# Thank you!

---

## Links:

- This presentation: <https://bit.ly/3ePaJLE>
- Me: Tho McCaie [theo.mccaie@informaticslab.co.uk](mailto:theo.mccaie@informaticslab.co.uk)
- Pangeo: <https://pangeo.io/>
- STAC: <https://stacindex.org/>
- [Blog about Analysis Ready Data](#)
- More work from my team: <https://medium.com/informatics-lab>

