



Python
fwdays

The art of data engineering

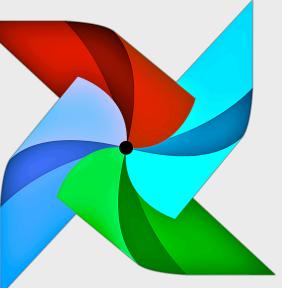
Andrii Soldatenko

TVTime

The art of Data engineering

14 December 2019,
Remote :)

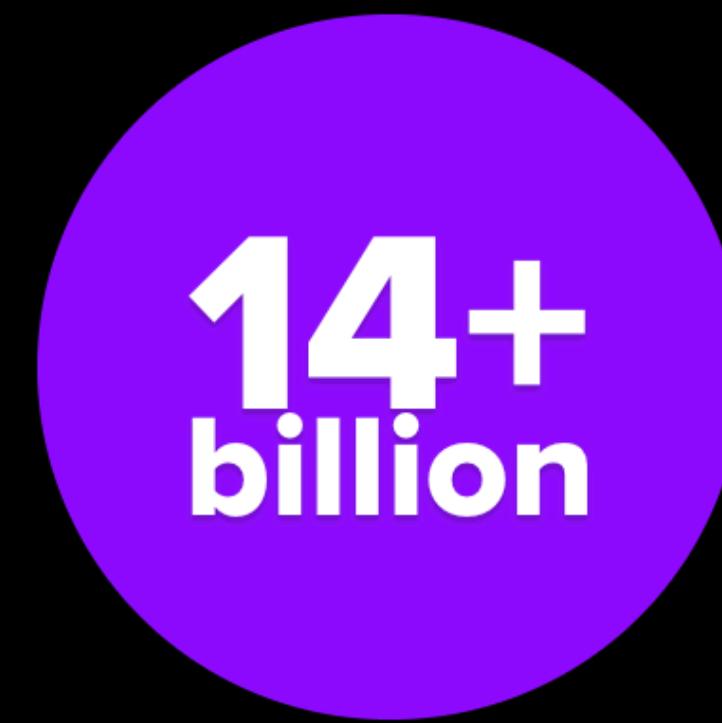
Andrii Soldatenko

- Senior Data Engineer TVTime
- Apache Airflow contributor 
- Public Speaker at many conferences, committee member of different PyCons
- blogger at <https://asoldatenko.com>





Registered
Users



Episodes Tracked
Across 90k+ Shows



Daily Active
Users



Years of
Tracking Data



Countries
Represented



Platforms
Covered

“...The Web was done
by amateurs”



- Alan Kay



P A R E N T A L
A D V I S O R Y
E X P L I C I T L Y R I C S

TL;DR



Mat Velloso

@matvelloso

Follow



Difference between machine learning
and AI:

If it is written in Python, it's probably
machine learning

If it is written in PowerPoint, it's
probably AI

5:25 PM - 22 Nov 2018

1,186 Retweets 3,333 Likes



41

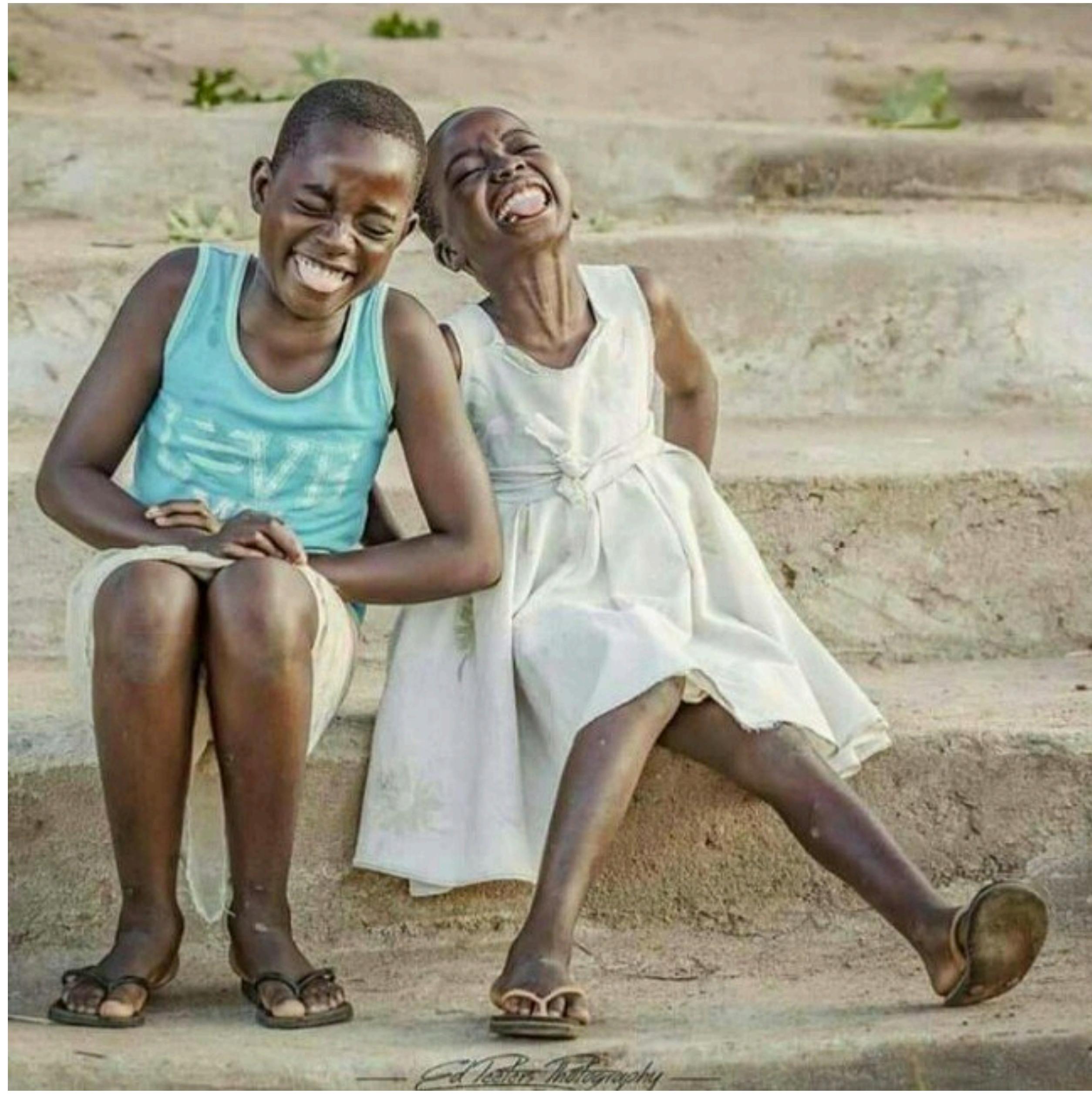


1.2K



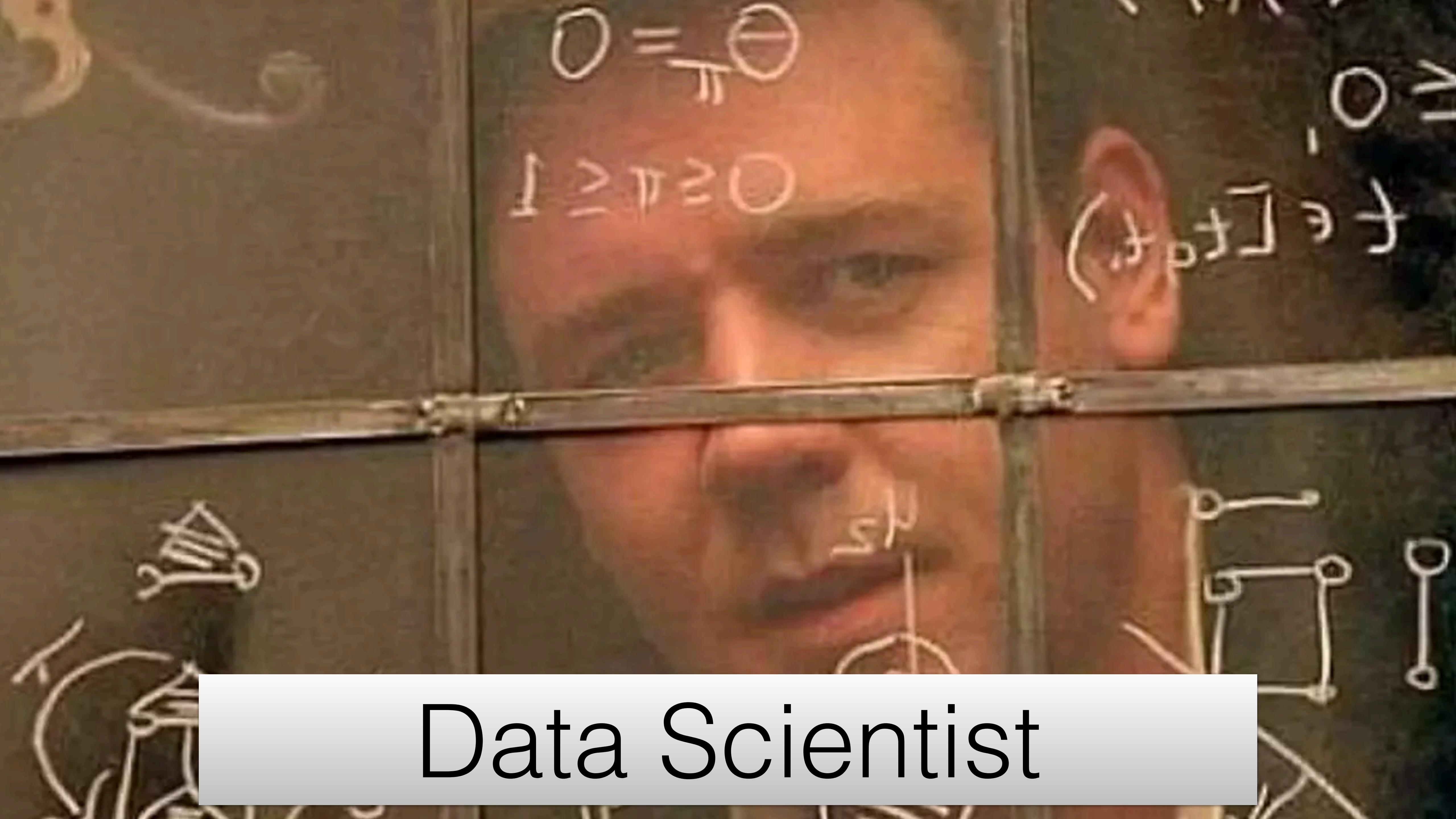
3.3K





— Ed Pastor Photography —



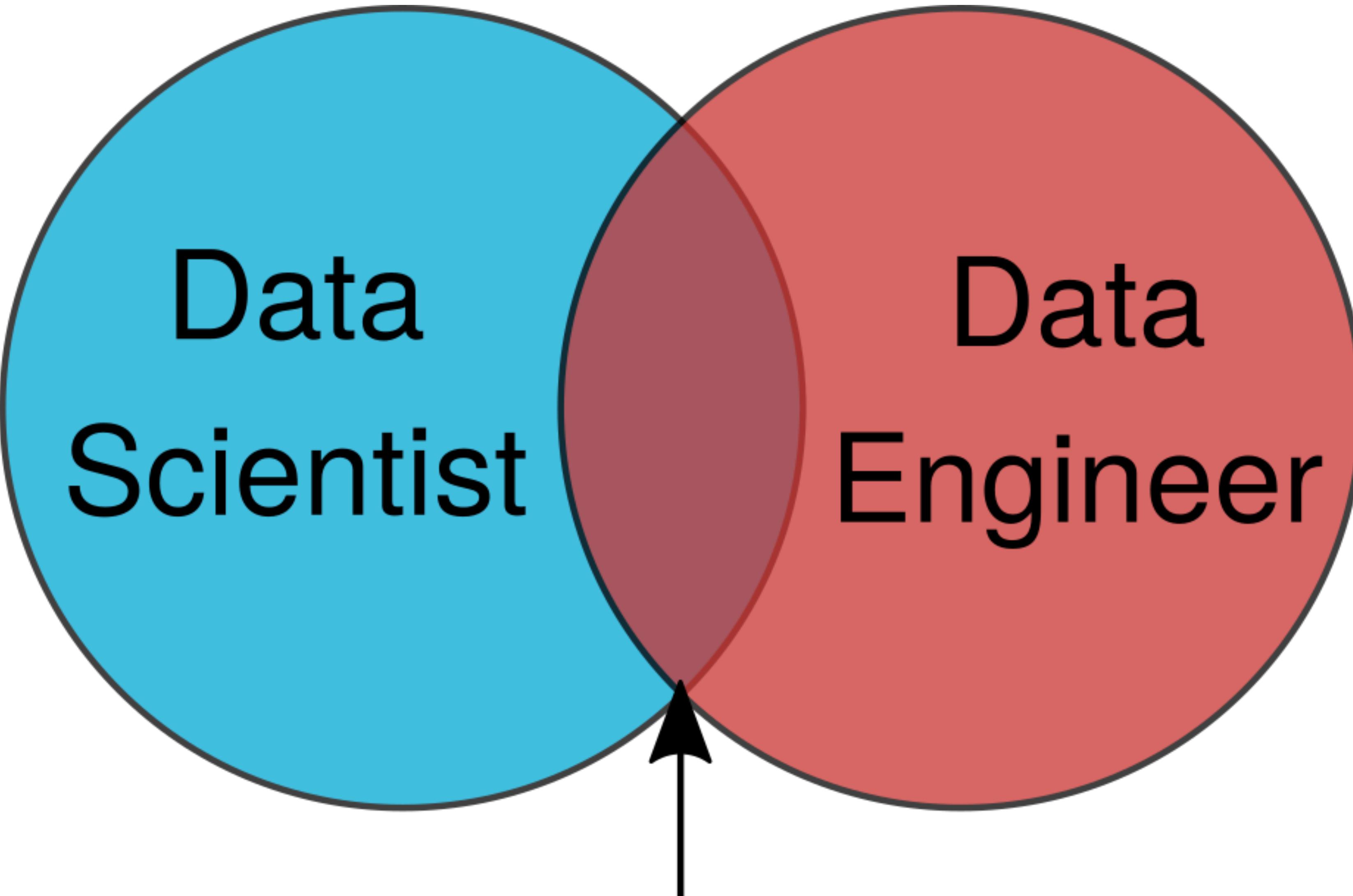
A photograph of a chalkboard covered in mathematical notation. At the top, there are two rows of equations: $O = \pi$ and $I > \pi > O$. Below these, there's a large, faint drawing of a figure standing next to a circle. To the right of the figure, there's a red chalk tray containing several pieces of chalk. The bottom half of the board features various geometric diagrams, including a compass rose-like figure on the left and several small rectangles with internal lines forming L-shapes on the right.

Data Scientist



Data Engineer

Data Engineering: A quick definition



Big Data

Twitter: 600 million tweets per day.

**Facebook: 600 terabytes of
incoming data each day, from 1.6
billion active users.**

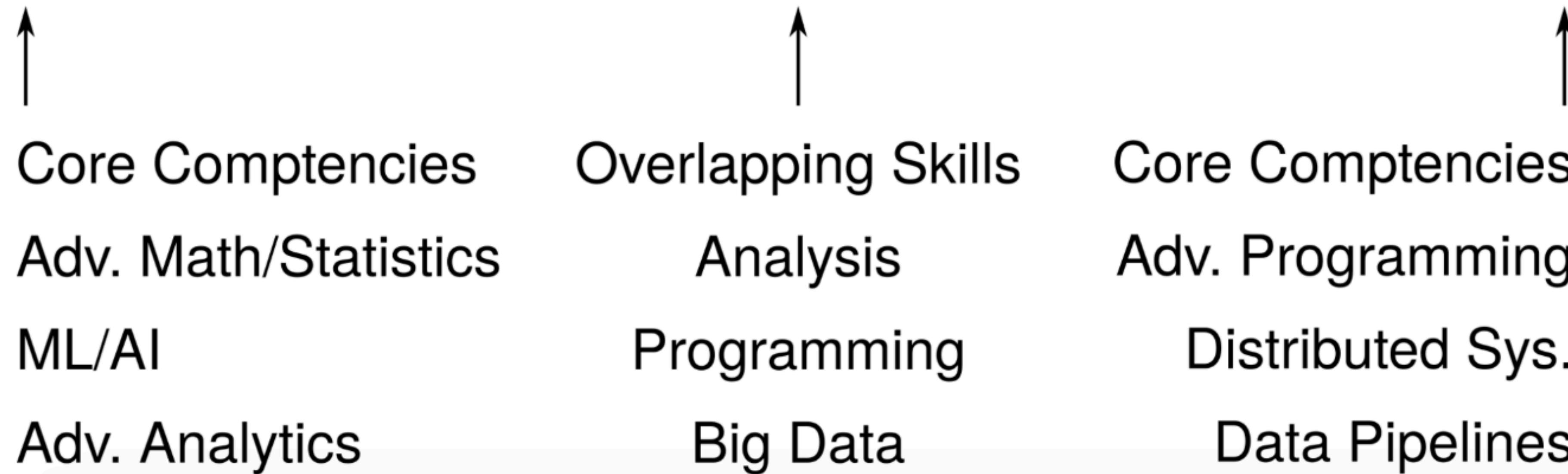
**Google: 3.5 billion search
queries per day.**

**Instagram: 52 million
new photos per day**

Data Scientist



Data Engineer





What are data
engineers focused on?

How to Become a Data Engineer? 🐚



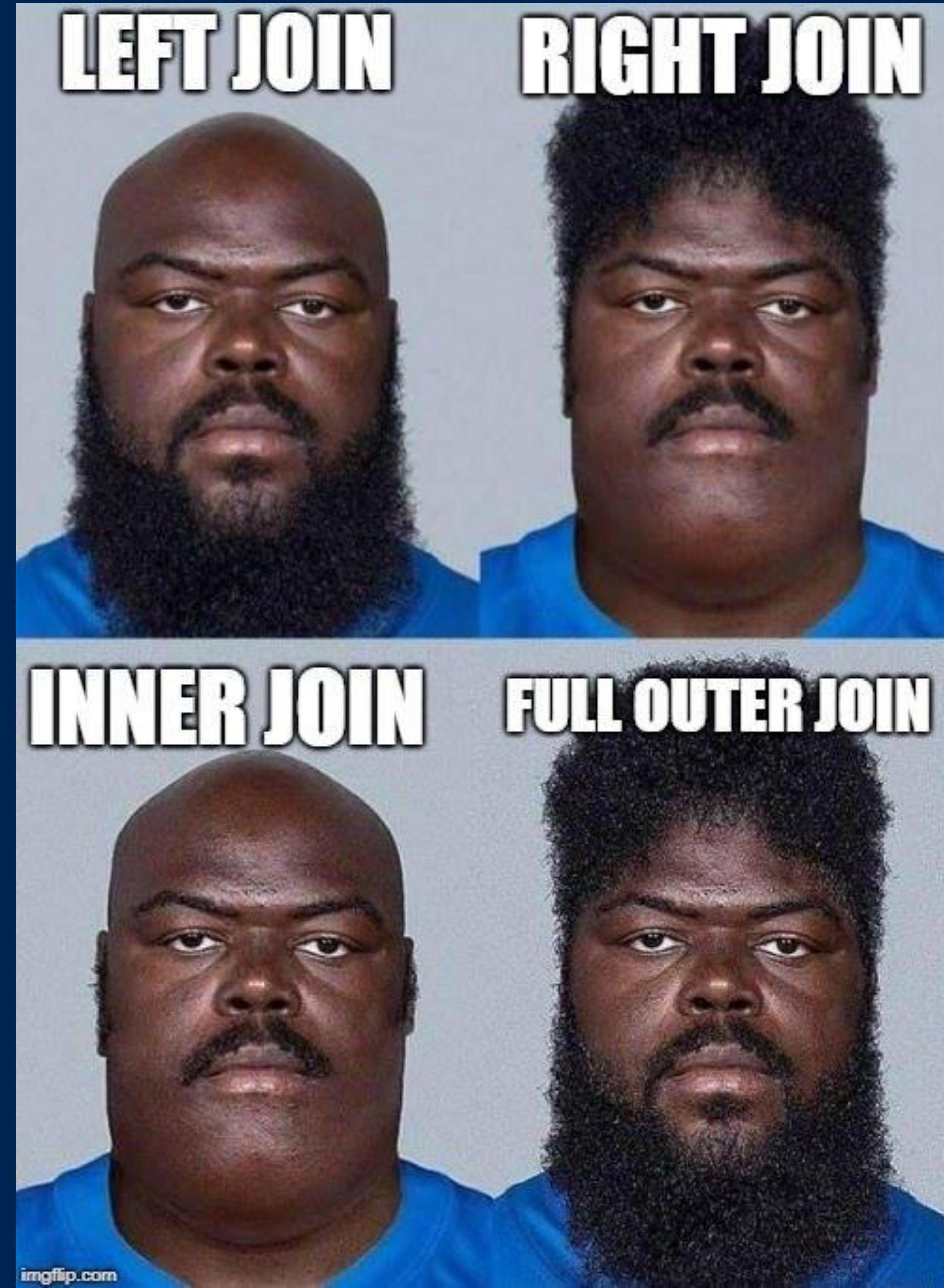
What skills do data engineers need?

LEFT JOIN

INNER JOIN

RIGHT JOIN

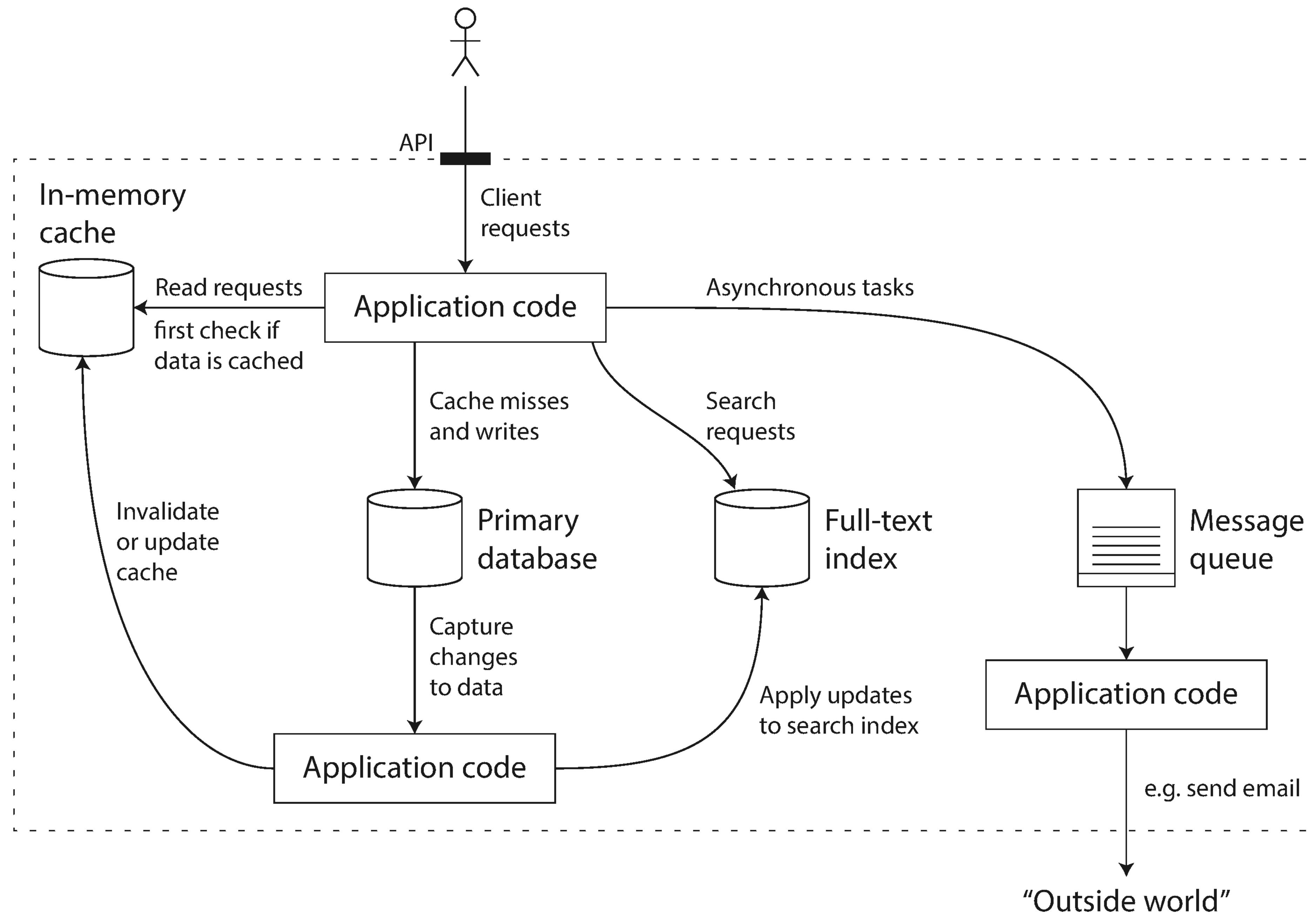
FULL OUTER JOIN







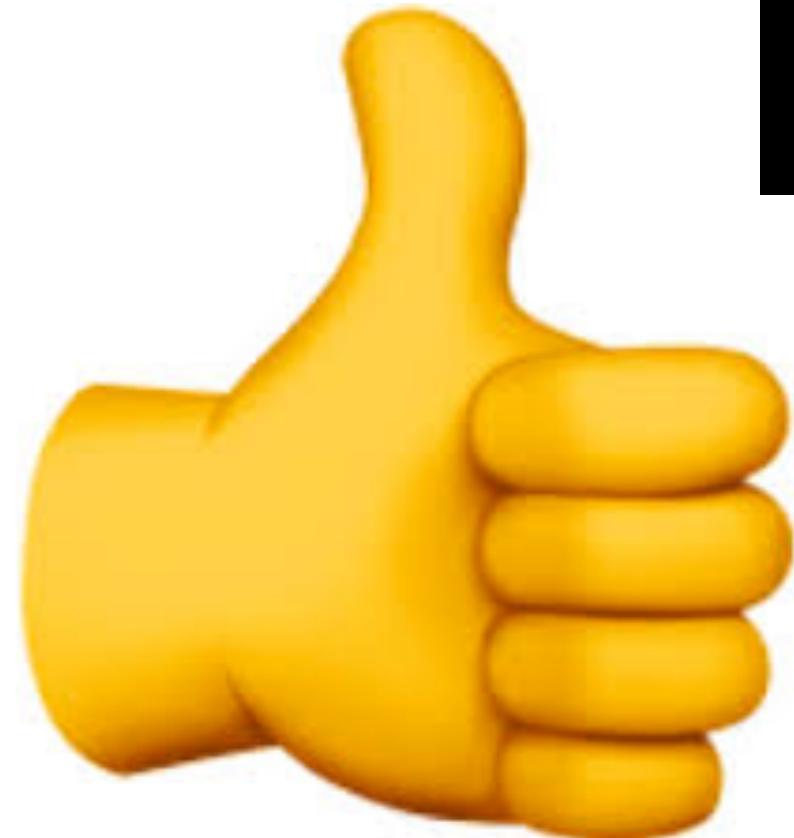
Possible architecture for a data system



Cron + bash script 🔥

#!/bin/bash

```
~root: env X="() { :;}; echo shellshock" /bin/sh -c "echo completed"  
> shellshock  
> completed
```



tes after the hour (0-59)
hour format (0-23).
ay of the month (1-31)

_____ 4. Month - Month of the year (1-12)
_____ 5. Weekday - Day of the week. (0-6, where 0 indicates Sun\$
_____ Command

```
# Web stats at https://  
1 */* * * perl /usr/lib/cgi-bin/awstats.pl -config=web -update >/dev/null  
1 */* * * perl /usr/lib/cgi-bin/awstats.pl -config=smtp -update >/dev/null  
# Backup LDAP at 03:00  
0 3 * * * bash /var/vmail/backup/backup_opendap.sh  
# Backup mysql at 03:30  
30 3 * * * bash /var/vmail/backup/backup_mysql.sh
```

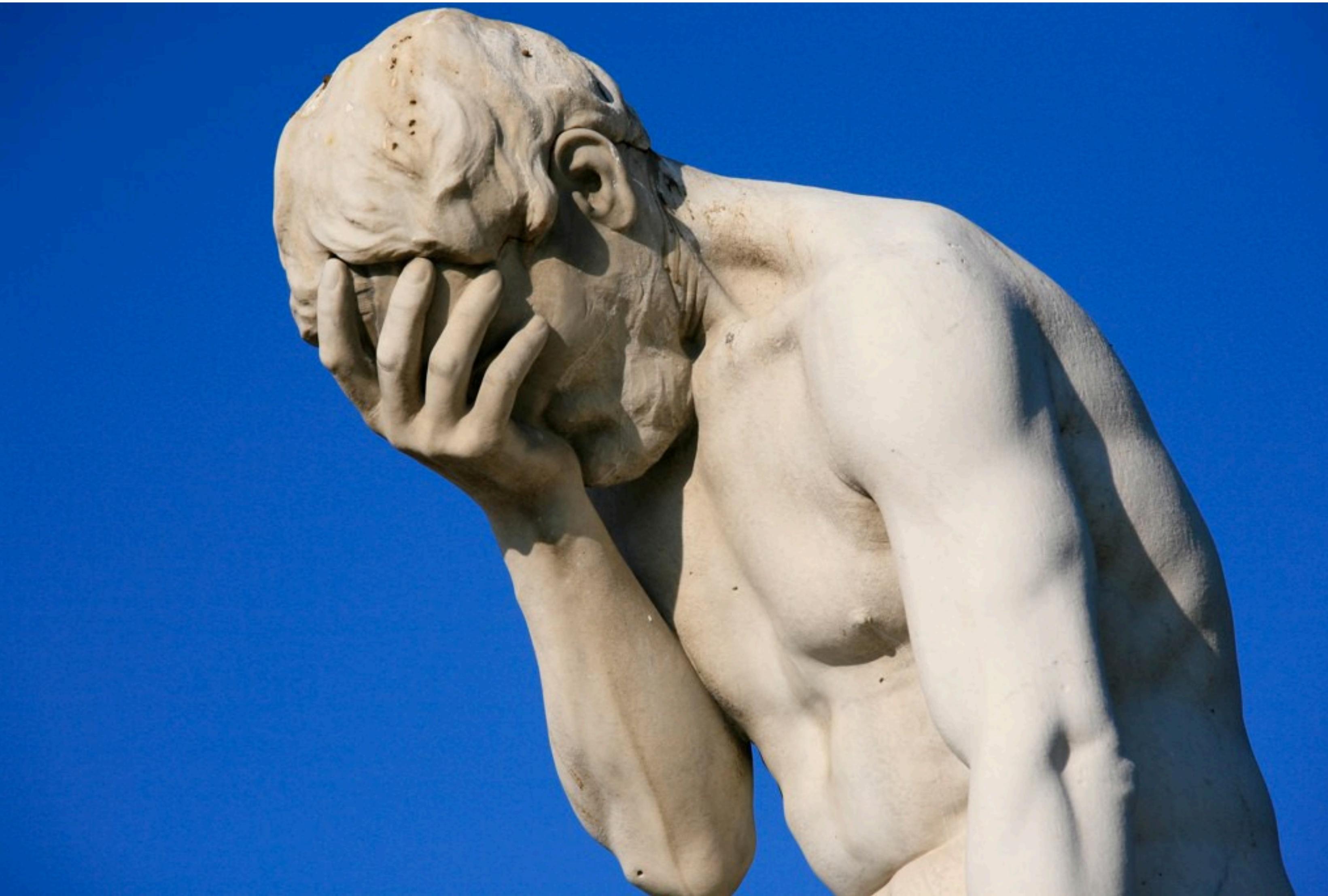
Simple data pipeline 🔥

TL; DR

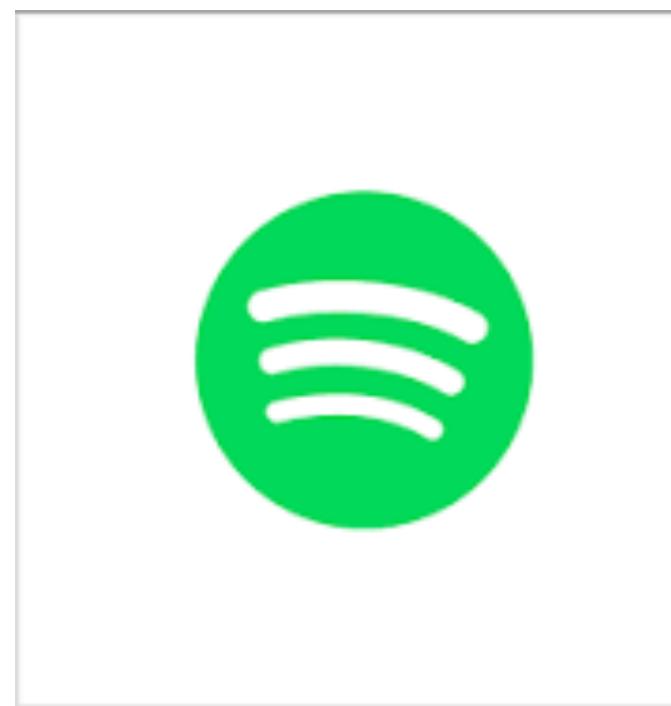
```
0 * * * * my_bash_pipeline.sh
```

```
#!/bin/bash
psql -U user -d database_name -c \
"COPY country FROM '/tmp/
country_data';"
```

I'm sorry Cron, I've met
Airflow :)



Data Pipelins: Open source ToolZZZ



OS workflow tools:

	Luigi	Airflow	Pinball
repo	https://github.com/spotify/luigi	https://github.com/airbnb/airflow	https://github.com/pinterest/pinball
docs	http://luigi.readthedocs.org	https://airflow.readthedocs.org	none
github forks	1786	3494	58
github stars	10520	10072	922
github watchers	510	626	126
commits in last 30 days	3700 commits	5618 commits	133 commits
<u>architecture</u>			
web dashboard	not really, minimal	very nice	yes
code/dsl	code	code	python dict + python code
files/datasets	yes, targets	not really, as special tasks	?
calendar scheduling	no, use cron	yes, LocalScheduler	yes
backfill jobs	yes	yes	?
persists state	kindof	yes, to db	yes, to db
tracks history	yes	yes, in db	yes, in db
code shipping	no	yes, pickle	workflow is shipped using pickle, jobs are not?
priorities	yes	yes	?
parallelism	yes, workers, threads per workers	yes, workers	?
control parallelism	yes, resources	yes, pools	?
cross-dag deps	yes, using targets	yes, using sensors	yes
finds new deployed tasks	no	yes	?
executes dag	no, have to create special sink task	yes	yes
multiple dags	no, just one	yes, also several dag instances (dagruns)	yes
<u>scheduler/workers</u>			

https://docs.google.com/spreadsheets/d/1wfKibn4yS7AshPzY8ch2Dq_VmQVKnGSXGj08A2u1YQI/edit#gid=0

Airflow



Airflow DAGs Data Profiling ▾ Browse ▾ Admin ▾ Docs ▾ About ▾

DAGs

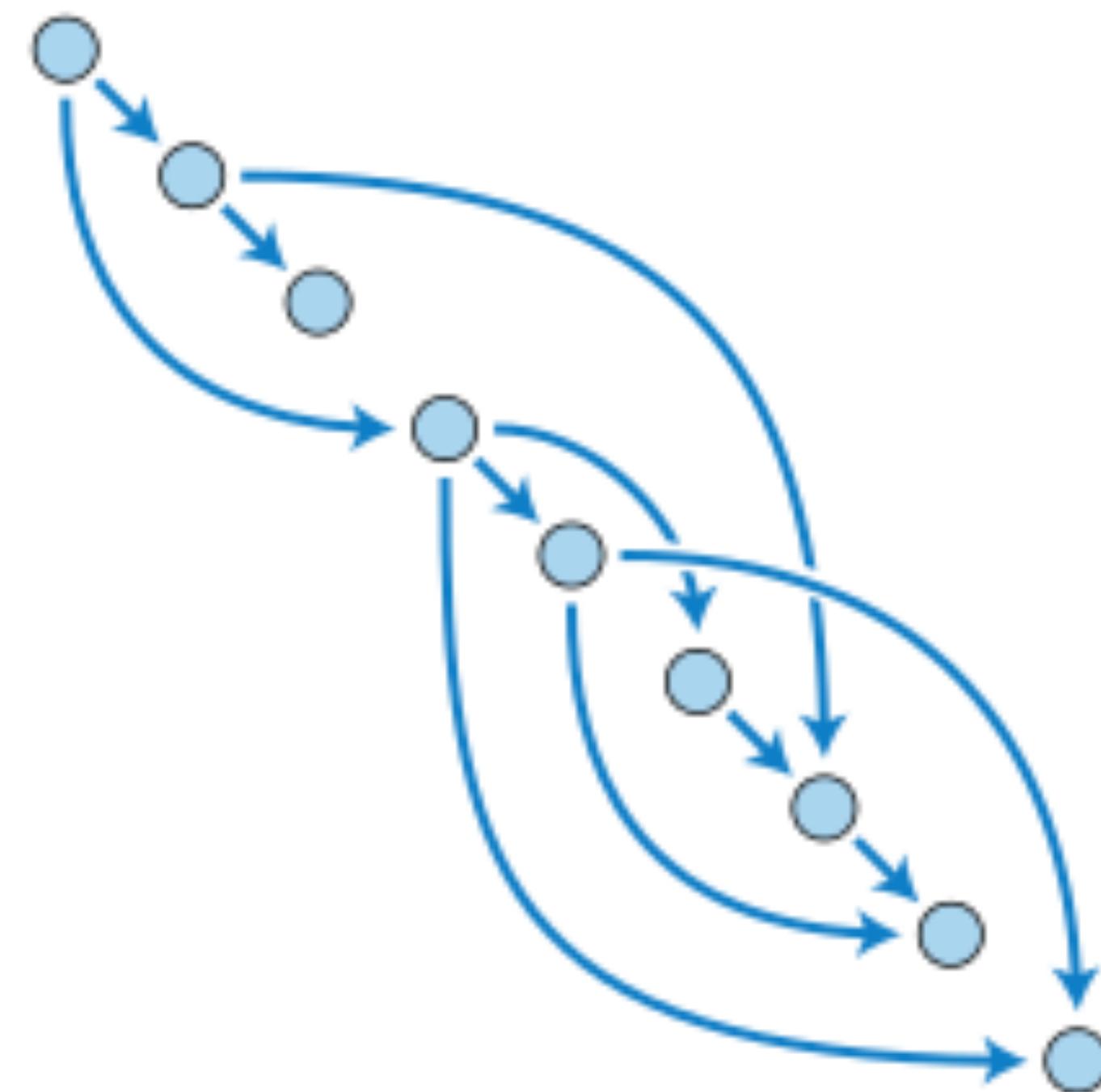
Search:

		DAG	Schedule	Owner	Recent Tasks	Last Run
			1 day, 0:00:00	airflow		2018-11-22 00:00
				airflow		
			0 12 * * *	airflow		2018-08-26 12:00
			1 day, 0:00:00	airflow		2018-11-22 00:00
				airflow		
				airflow		
				airflow		
			1 day, 0:00:00	airflow		2018-11-22 00:00
			1 day, 0:00:00	airflow		
			1 day, 0:00:00	airflow		2018-11-22 00:00

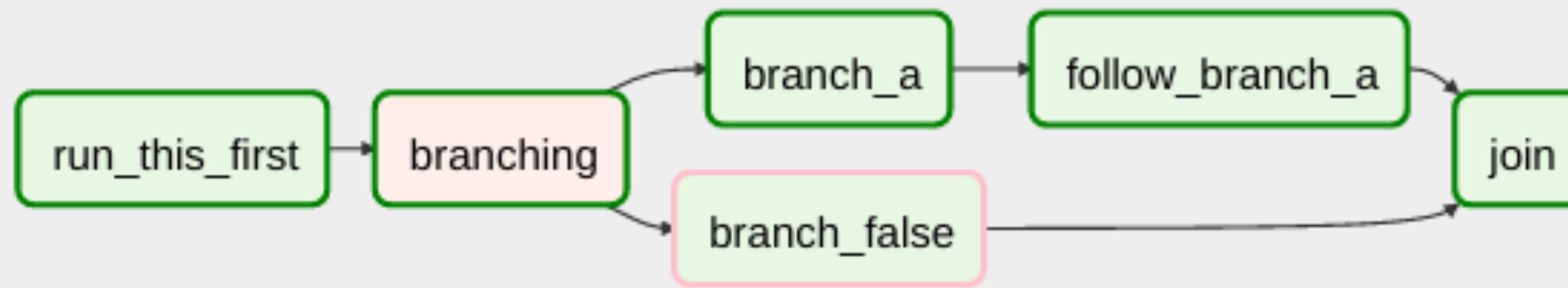
Airflow Concepts



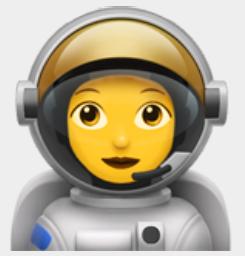
Airflow Core Ideas: DAGS



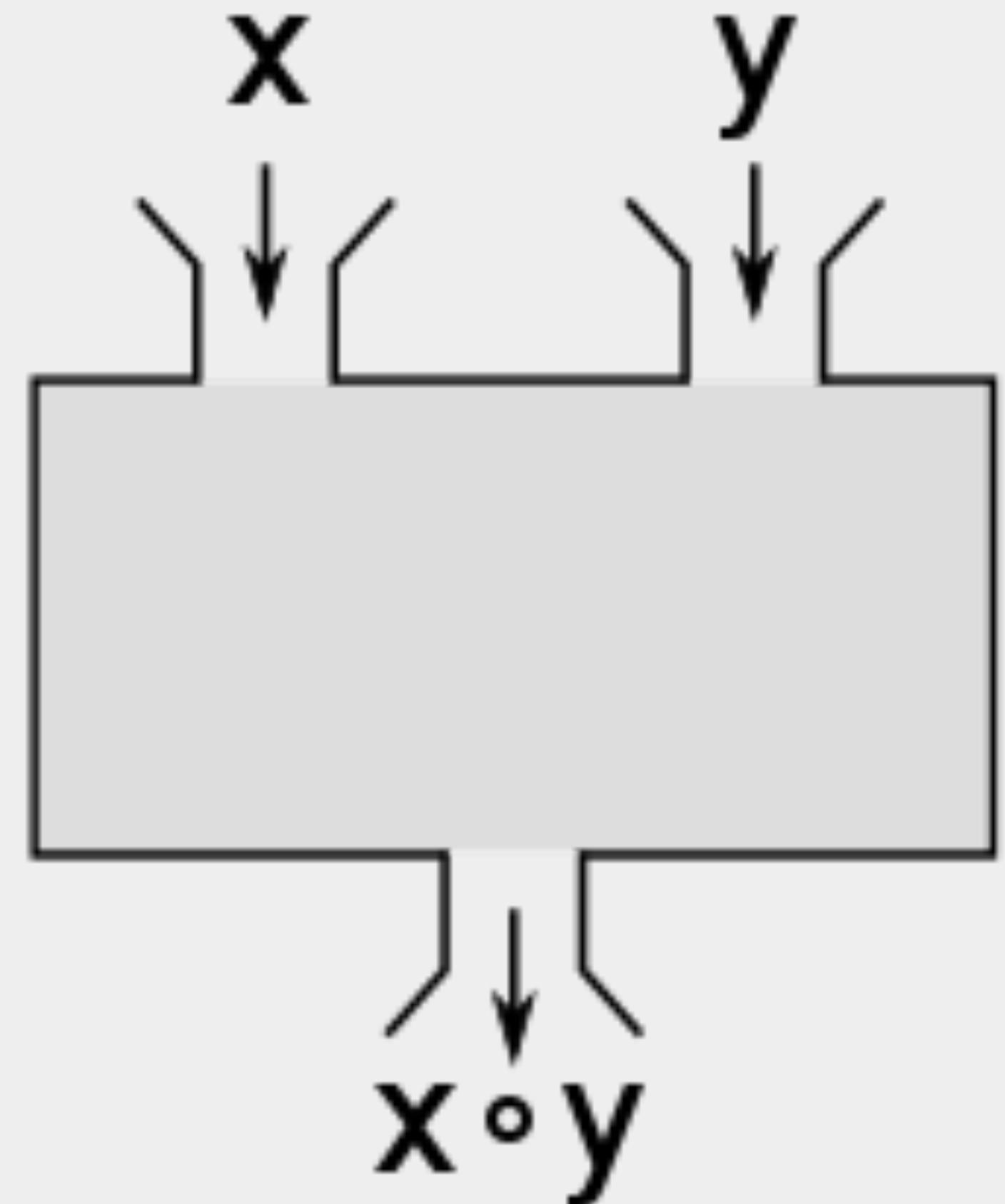
Airflow Core Ideas



Airflow Core Ideas: Operators

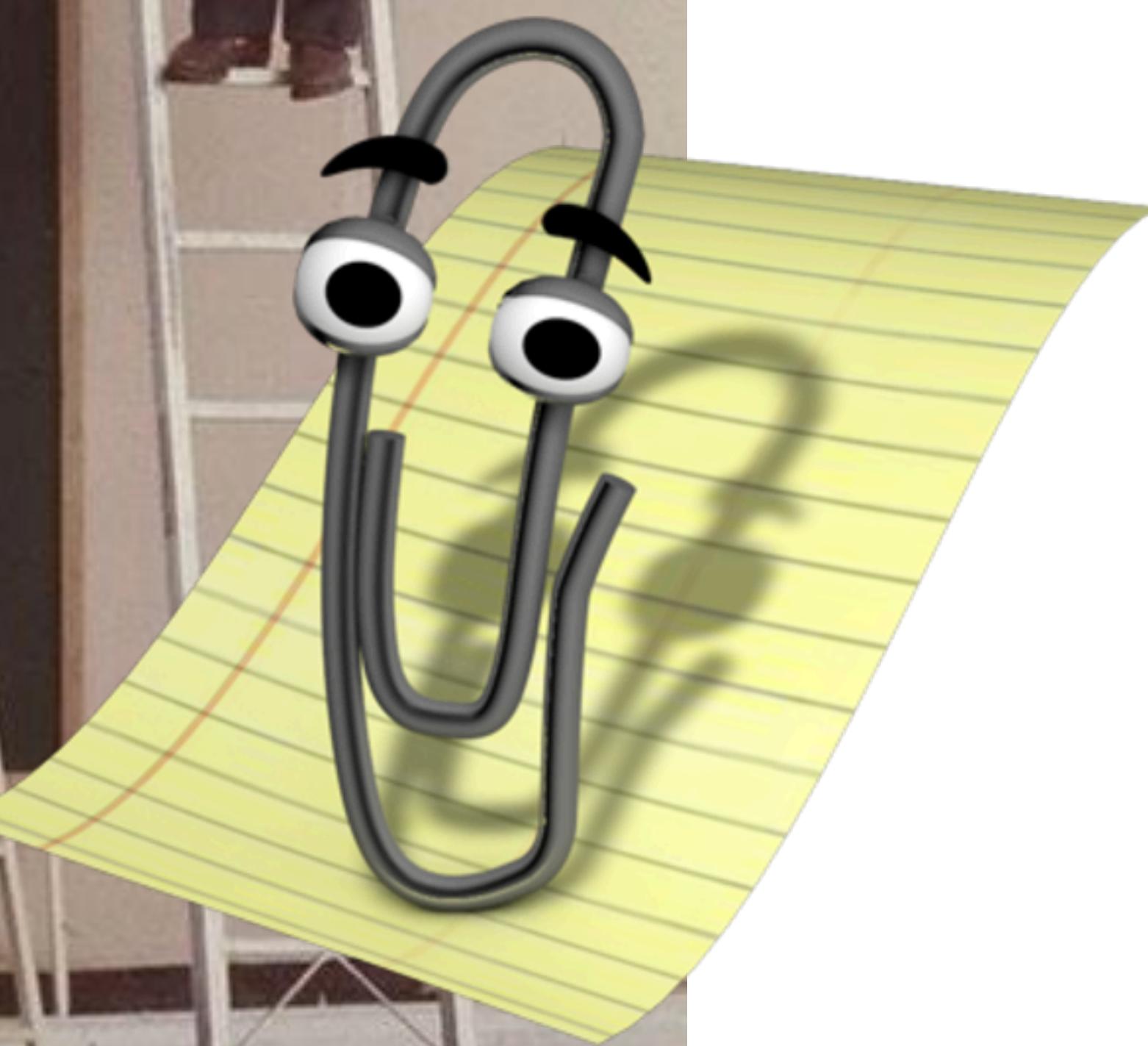


- BashOperator
- PythonOperator
- MySqlOperator
- S3FileTransformOperator ... you get the idea



Return to real world

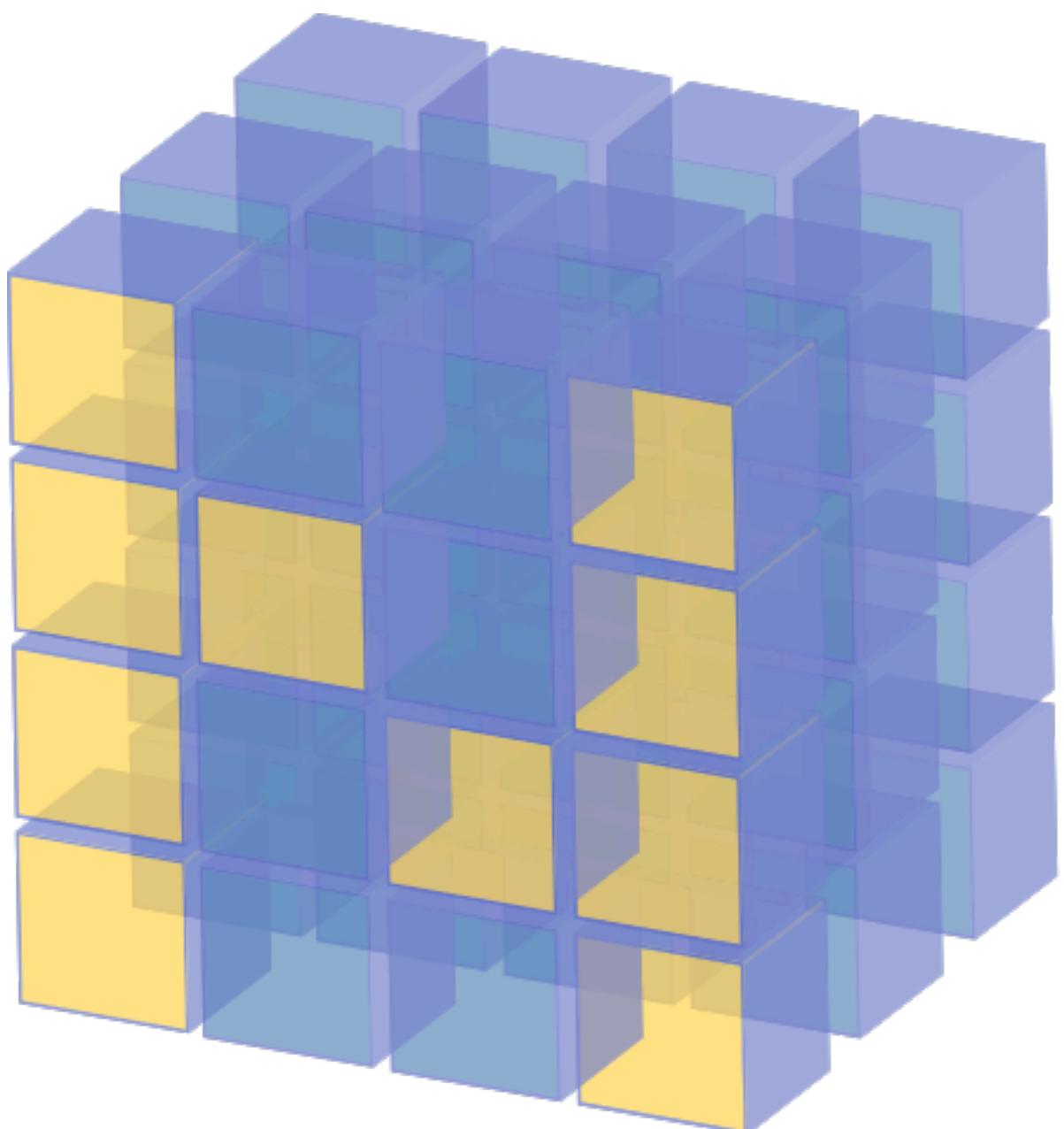




ML Dependency hell:



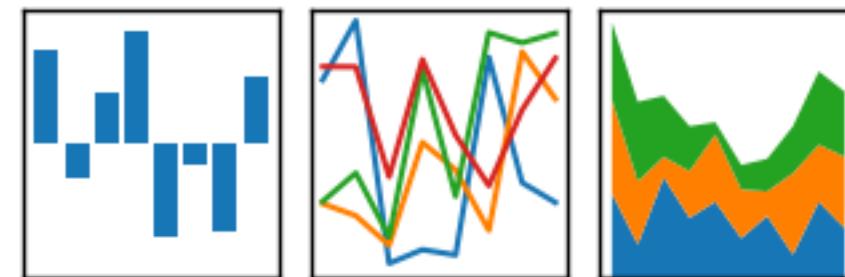
SciPy



NumPy

pandas

$$y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$$



spaCy

The python NLTK logo includes the Python logo (a yellow and blue snake) followed by the text "python™ Natural Language Analyses with NLTK".



TensorFlow

AWS: the good, the bad and the ugly

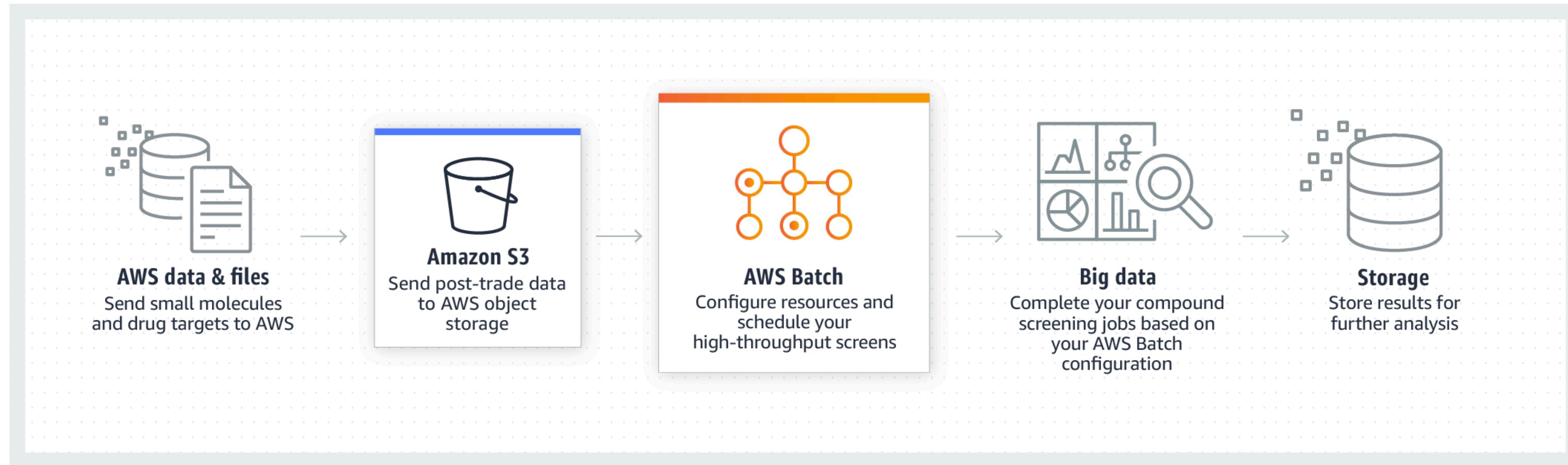


ML Dependency hell:



<https://www.instagram.com/natgeo/>

Case #1 Amazon Batch:



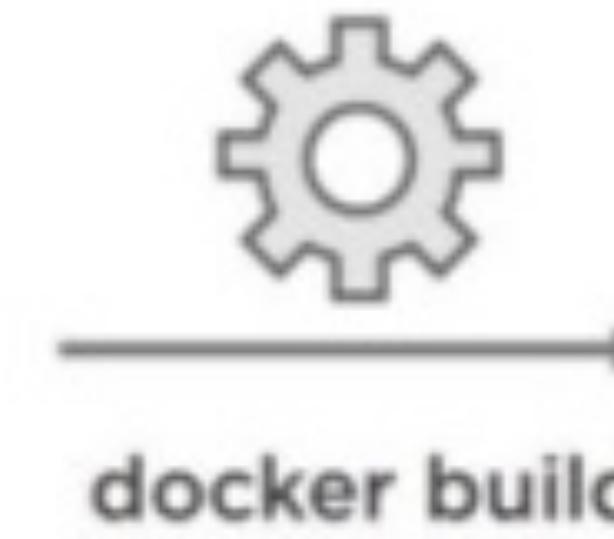
AWS Batch and Airflow



AWS Batch and Airflow



Dockerfile



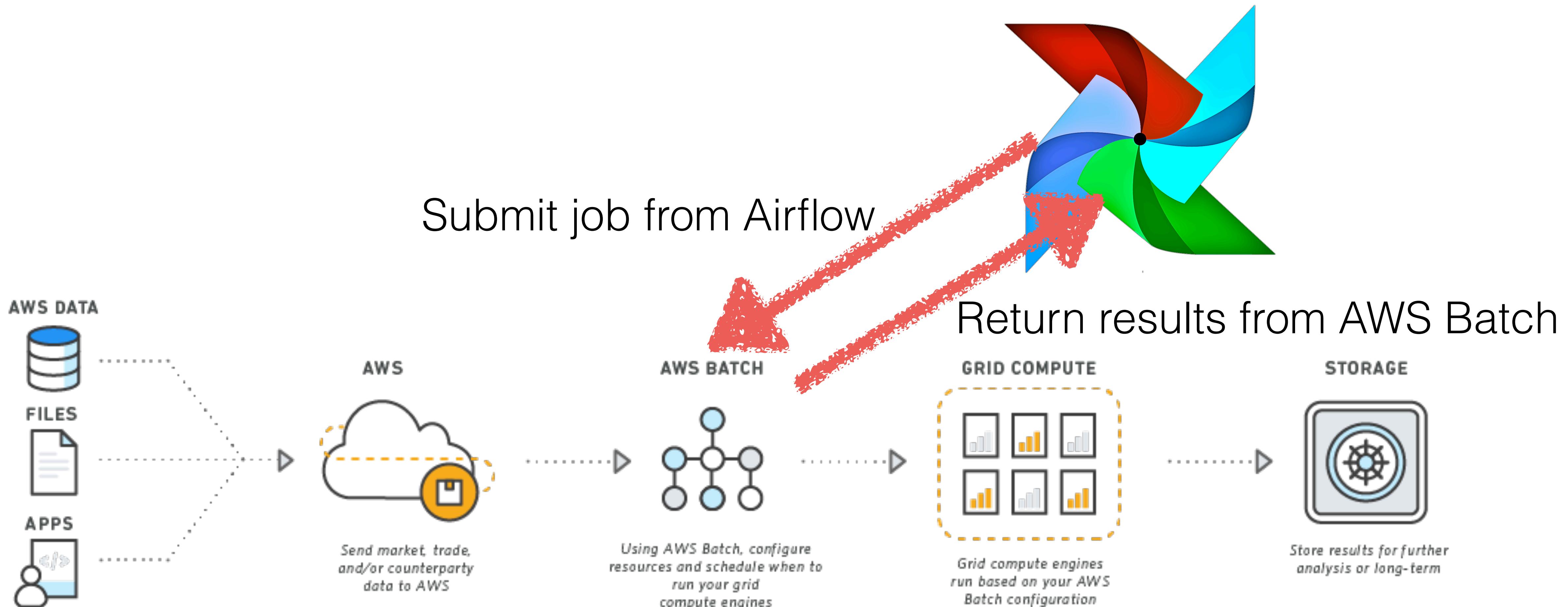
Docker Image

AWS Batch and Airflow: Example

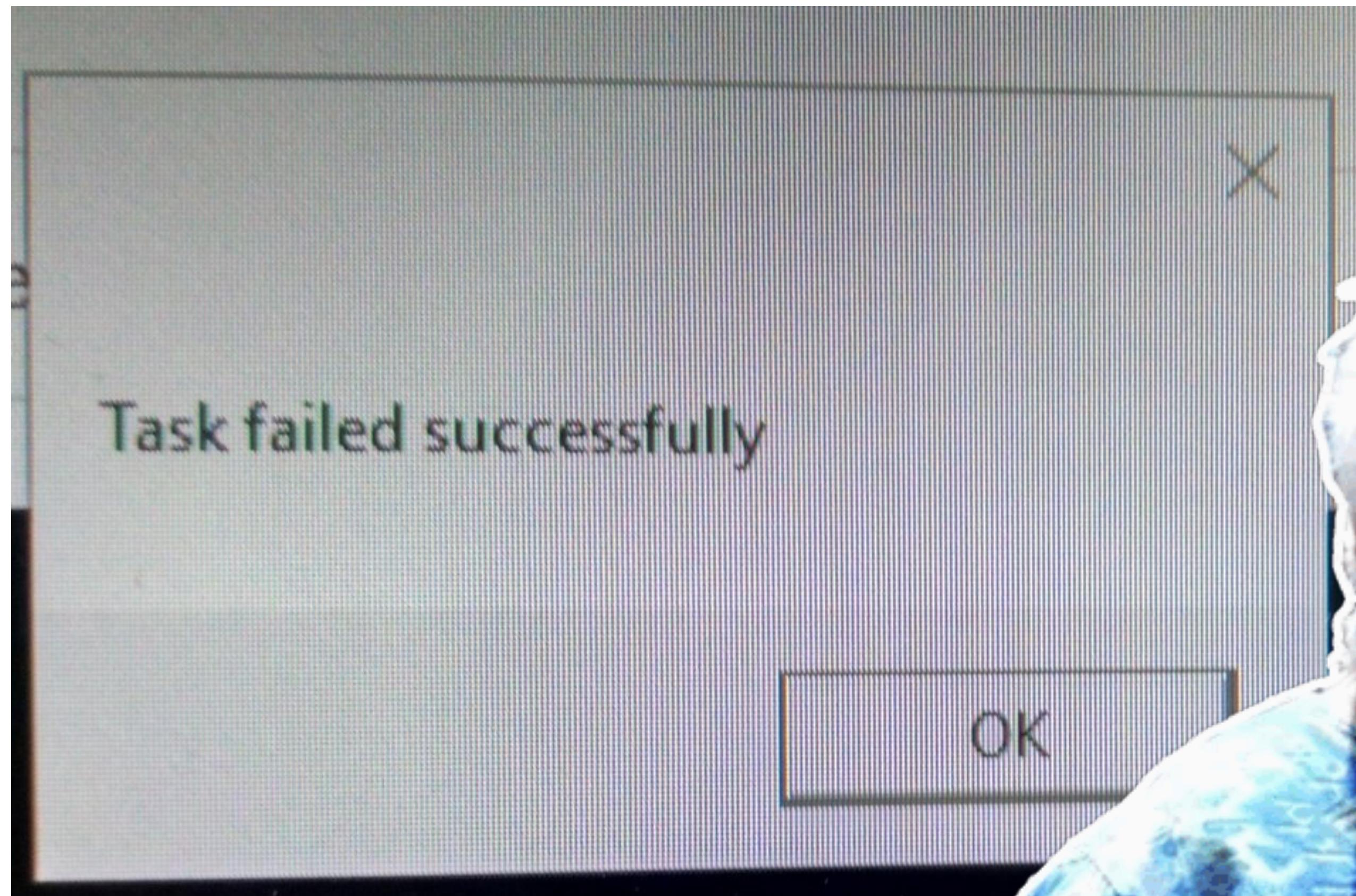
```
FROM amazonlinux:latest
RUN yum -y groupinstall 'Development Tools'
...
WORKDIR /scratch
RUN git clone https://github.com/cjlin1/
libmf.git /tmp/libmf && cd /tmp/libmf && make
ENTRYPOINT [ "/usr/local/bin/entrypoint.sh" ]
```

LIBMF is a library for large-scale sparse matrix factorization. For the optimization problem it solves and the overall framework, please refer to [3].

AWS Batch and Airflow



AWS Batch and Airflow



**22G volume mounted by
default , ~ 10Gb**

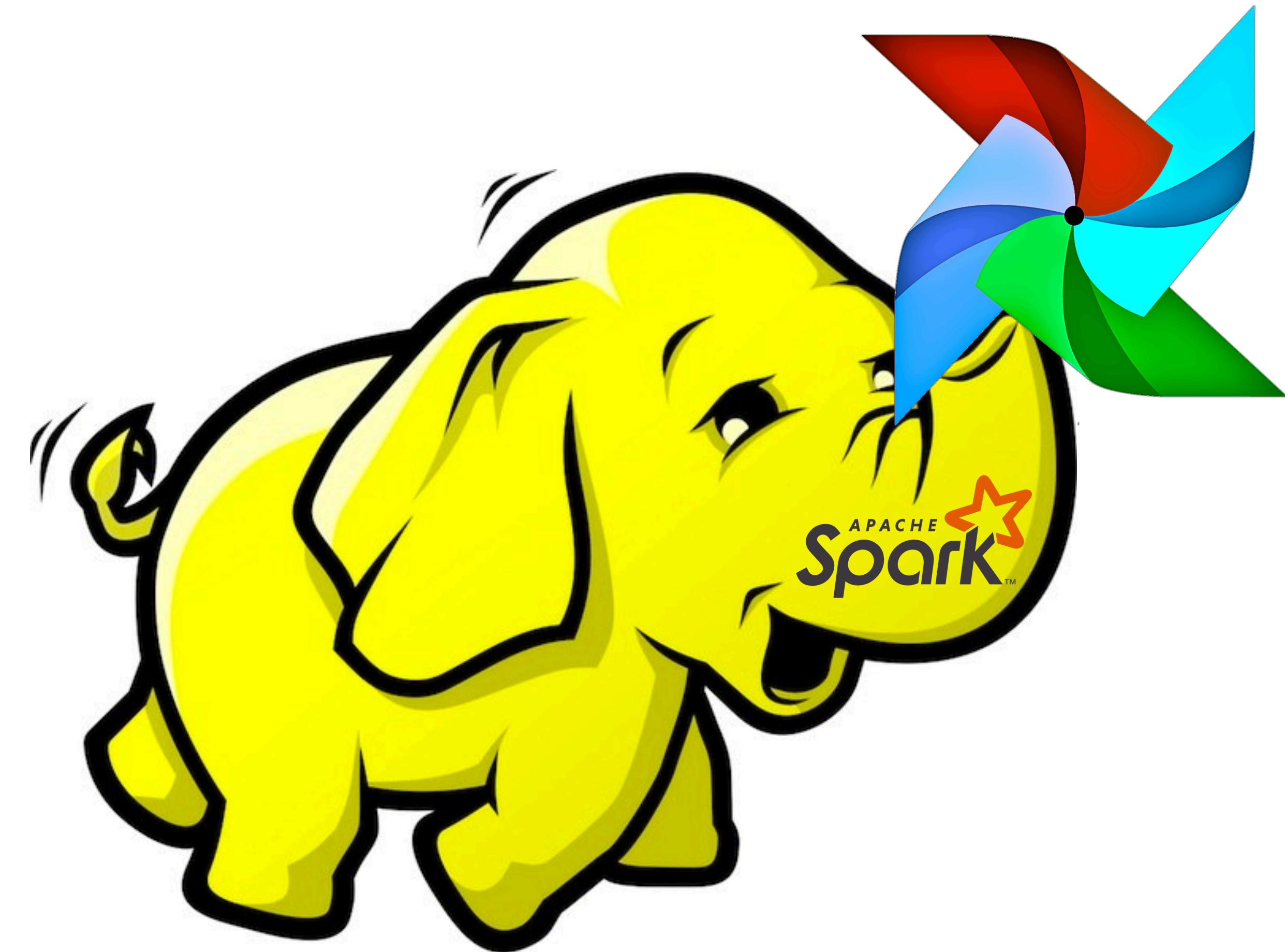


<https://forums.aws.amazon.com/thread.jspa?threadID=250705>

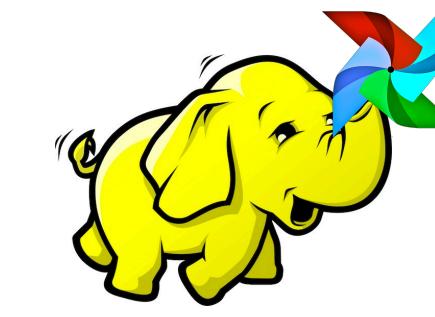
Solution =)

- You need to create EC2 based on ECS optimized AMIs
 - Attach another volume separate from your root volume, you will want to modify the /etc/fstab
 - Create image (AMI) Or use (ami-XXXXXX with extra 20Gb space)
 - Create AWS Batch compute environment based on AMI
-
- useful link <https://aws.amazon.com/blogs/compute/building-high-throughput-genomic-batch-workflows-on-aws-batch-layer-part-3-of-4/>

Case #2 EMR Spark and Airflow



EMR and Airflow



`airflow/contrib/operators/emr_add_steps_operator.py`
`airflow/contrib/operators/emr_create_job_flow_operator.py`
`airflow/contrib/operators/emr_terminate_job_flow_operator.py`

Debug EMR cluster:

Clone Terminate AWS CLI export

Cluster: [REDACTED] 2018-11-27 Terminated Terminated by user request

Summary Application history Monitoring Hardware Events Steps Configurations Bootstrap actions

YARN applications > application_1543401061957_0010 (Spark) C

Jobs Stages Executors

Executors (20)

Filter: Filter executors 20 executors (all loaded) C

Executor ID	Address	Status	RDD blocks	Storage memory	On-heap storage memory	Off-heap storage memory	Disk used	Cores	Active tasks	Failed tasks	Complete tasks
driver	ip-10-96-156-62.eu-west-1.compute.internal:36549	Active	0	0.0 B / 6.2 GB	0.0 B / 6.2 GB	0.0 B / 0.0 B	0.0 B	0	0	0	0
1	ip-10-9-48-160.eu-west-1.compute.internal:35521	Active	0	0.0 B / 30.7 GB	0.0 B / 30.7 GB	0.0 B / 0.0 B	0.0 B	16	0	0	1021
2	ip-10-96-156-62.eu-west-1.compute.internal:33469	Active	0	0.0 B / 30.7 GB	0.0 B / 30.7 GB	0.0 B / 0.0 B	0.0 B	16	0	0	3555
3	ip-10-9-48-160.eu-west-1.compute.internal:42199	Active	0	0.0 B / 30.7	0.0 B / 30.7	0.0 B / 0.0 B	0.0 B	16	0	0	1037

[View logs](#)

[View logs](#)

[View logs](#)

[View logs](#)

As a final step
load data to warehouse service:



Amazon **Redshift**



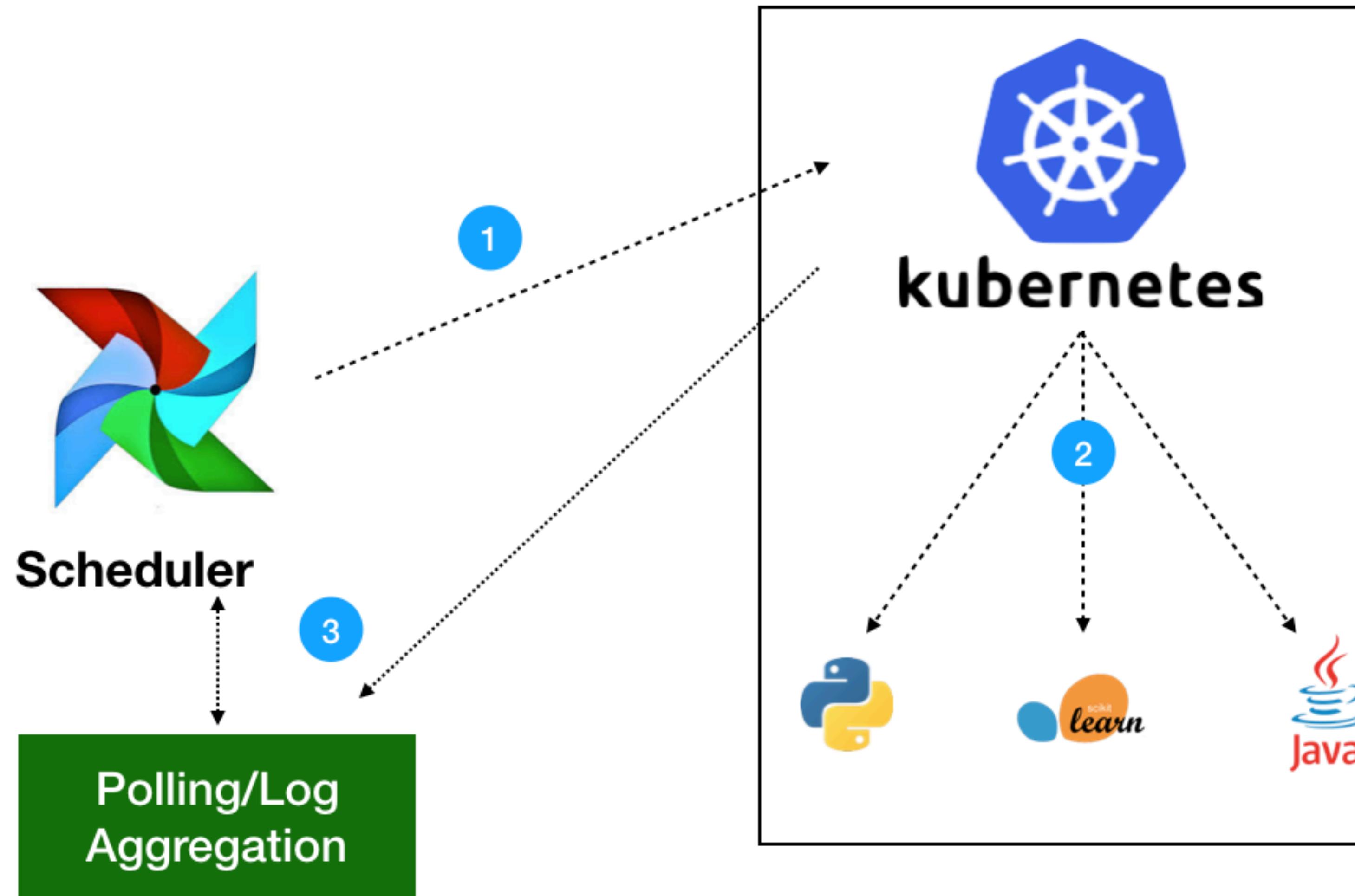
As a final step
load data to warehouse service:



`COPY INTO table (colA, colB)
from (SELECT (colA, colB)
FROM @s3_stage_parquet)
ON_ERROR = CONTINUE`



Case #3 Kubernetes Pod Executor and Airflow



CI/CD pipelines



CI/CD pipelines

- unit tests/integration tests 
- code quality checks 
- automatic deployment after merge into master 
- Pull Request code review flow 

unit tests/integration tests



```
@pytest.mark.parametrize('dag_file', DAG_FILES)
def test_import_dag_files(dag_file):
    """Import dag files and check for DAG."""
    module_name, _ = os.path.splitext(dag_file)
    module_path = os.path.join(DAG_PATH, dag_file)
    mod_spec = importlib.util.spec_from_file_location(module_name, module_path)
    module = importlib.util.module_from_spec(mod_spec)
    mod_spec.loader.exec_module(module)
    assert any(
        isinstance(var, af_models.DAG)
        for var in vars(module).values() )
```

Data quality checkers

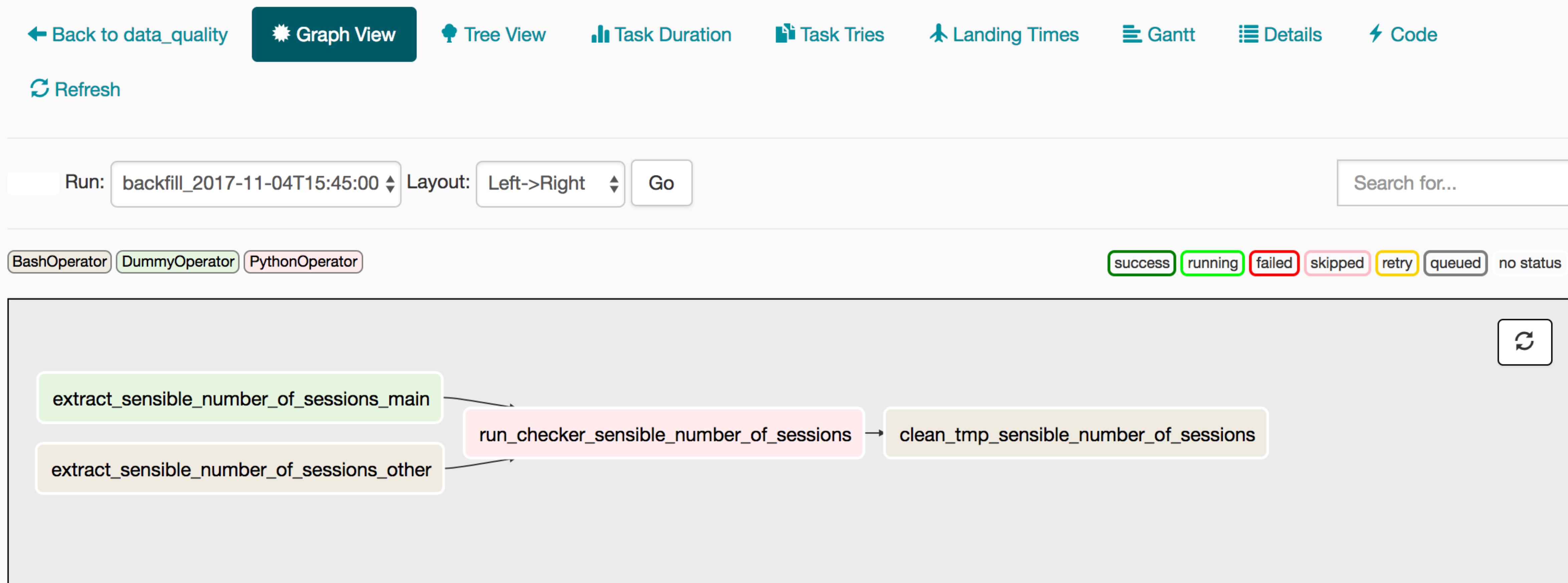
```
COPY DATA to NEW_TABLE  
APPLY DATA QUALITY CHECKS IF ERROR -> SEND ALERTS EXIT(1)  
BACKUP ORIGINAL TABLE  
RENAME NEW_TABLE to ORIGINAL_TABLE
```

Data quality checkers



SUBDAG: data_quality.sensible_number_of_sessions

schedule: */15 * * * *

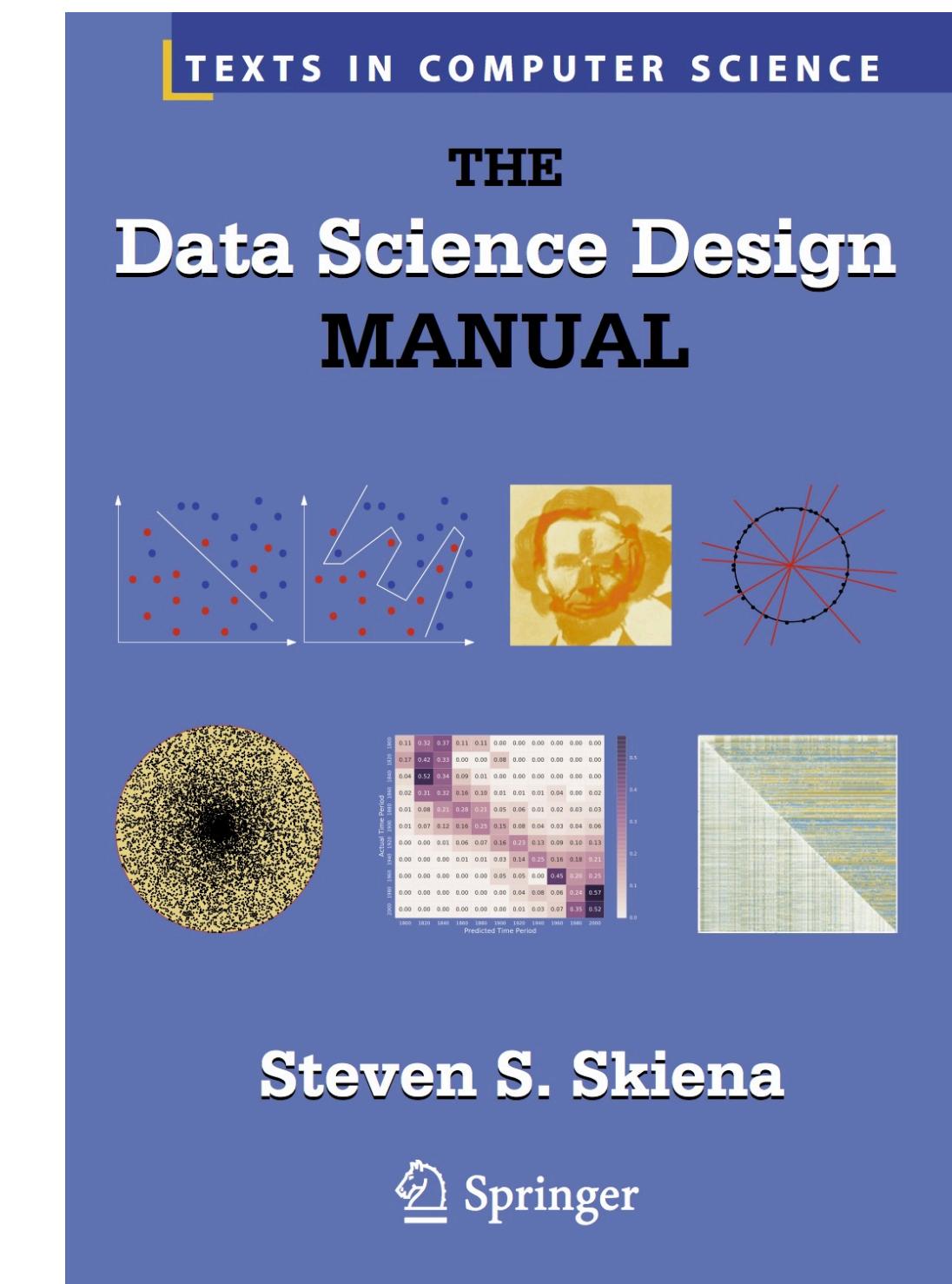
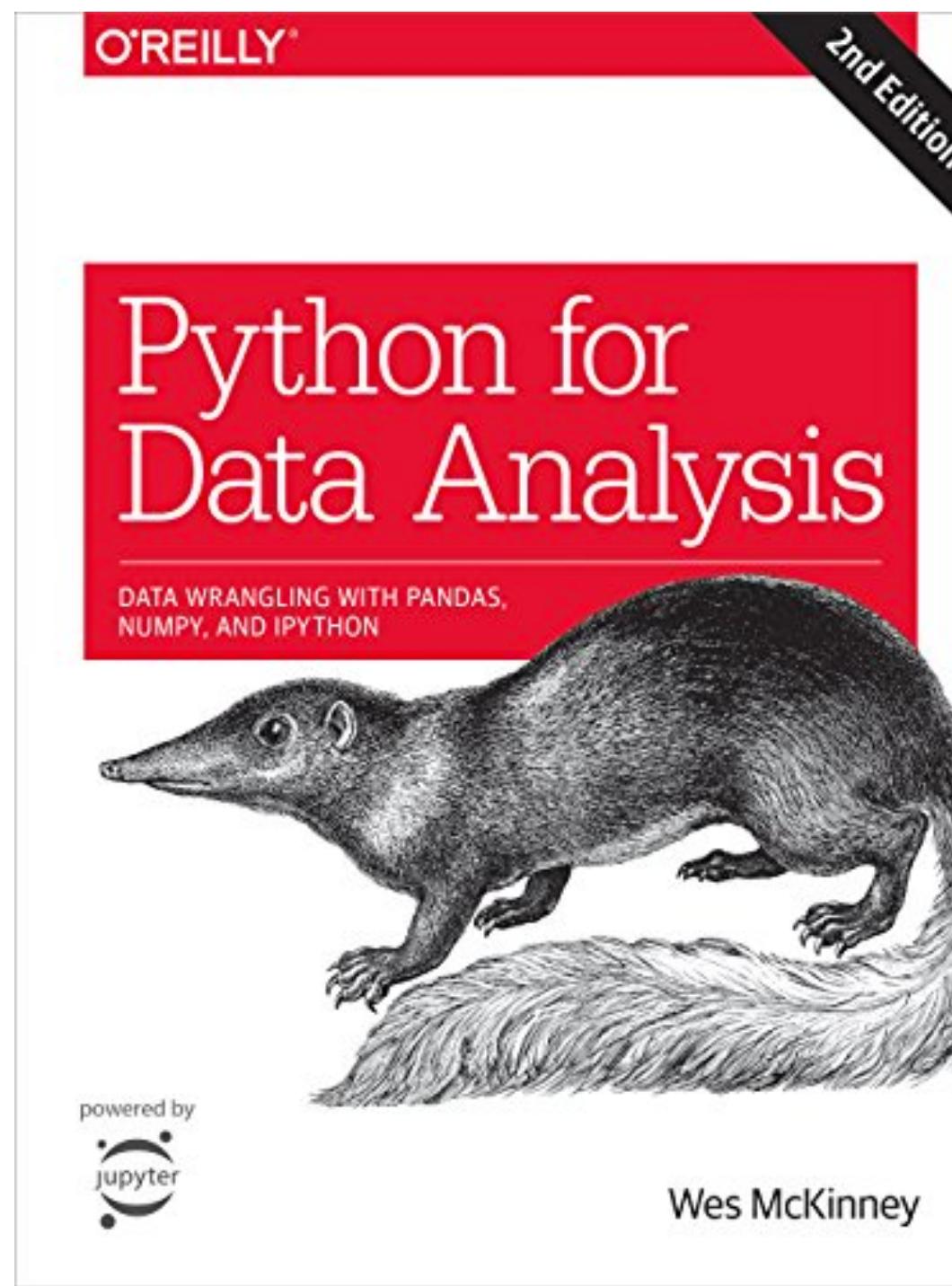


Learn more

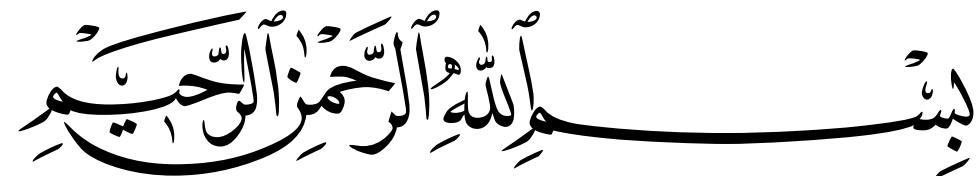


Learn more

- Awesome Data Engineering ([igorbarinov/awesome-data-engineering](https://github.com/igorbarinov/awesome-data-engineering))
- Awesome Apache airflow ([jghoman/awesome-apache-airflow](https://github.com/jghoman/awesome-apache-airflow))



Conclusion:

- Data engineering new trend with old roots 
- Data intensive apps is fun 
- ETL code base is just yet another source code 
- You must know better than others you data storages 

Thank You



Thank You

andrii.soldatenko@gmail.com

Questions

