



DAGSTER

Nick

Schrock

Founder, Elementl
@schrockn

“Our data is totally broken”

“Our data is totally broken”

- We don't know where our data comes from
- We don't know what it means
- We cannot reliably process and test it
- Our engineers don't want to deal with it
- It isn't “fun.” It isn't “sexy.”



DATA

Data Scientist: The Sexiest Job of the 21st Century

by [Thomas H. Davenport](#) and [D.J. Patil](#)

What do data scientists actually do on their day-to-day? Career (self.datascience)

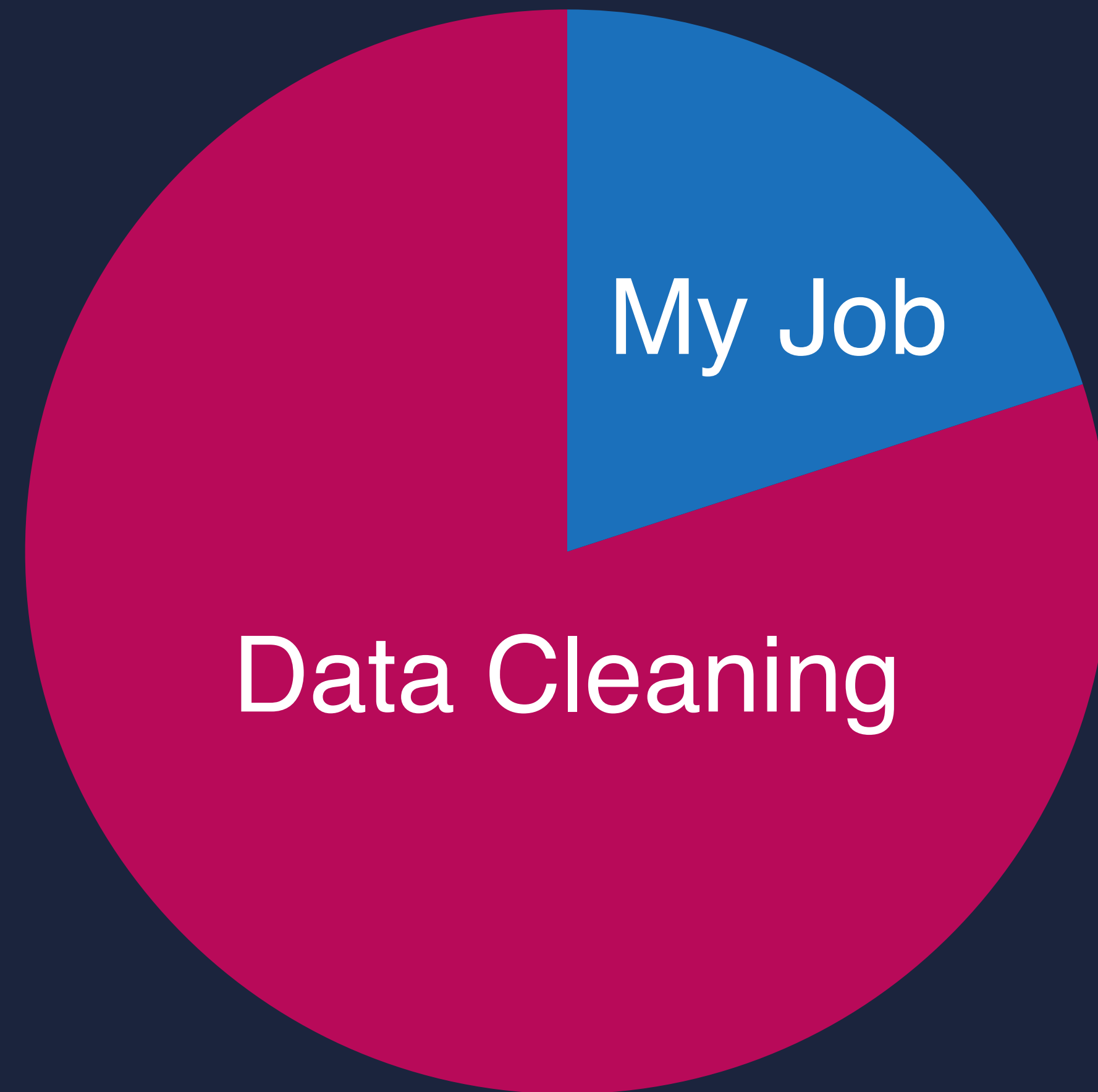


submitted 7 hours ago by [giokrist](#) to [r/datascience](#)

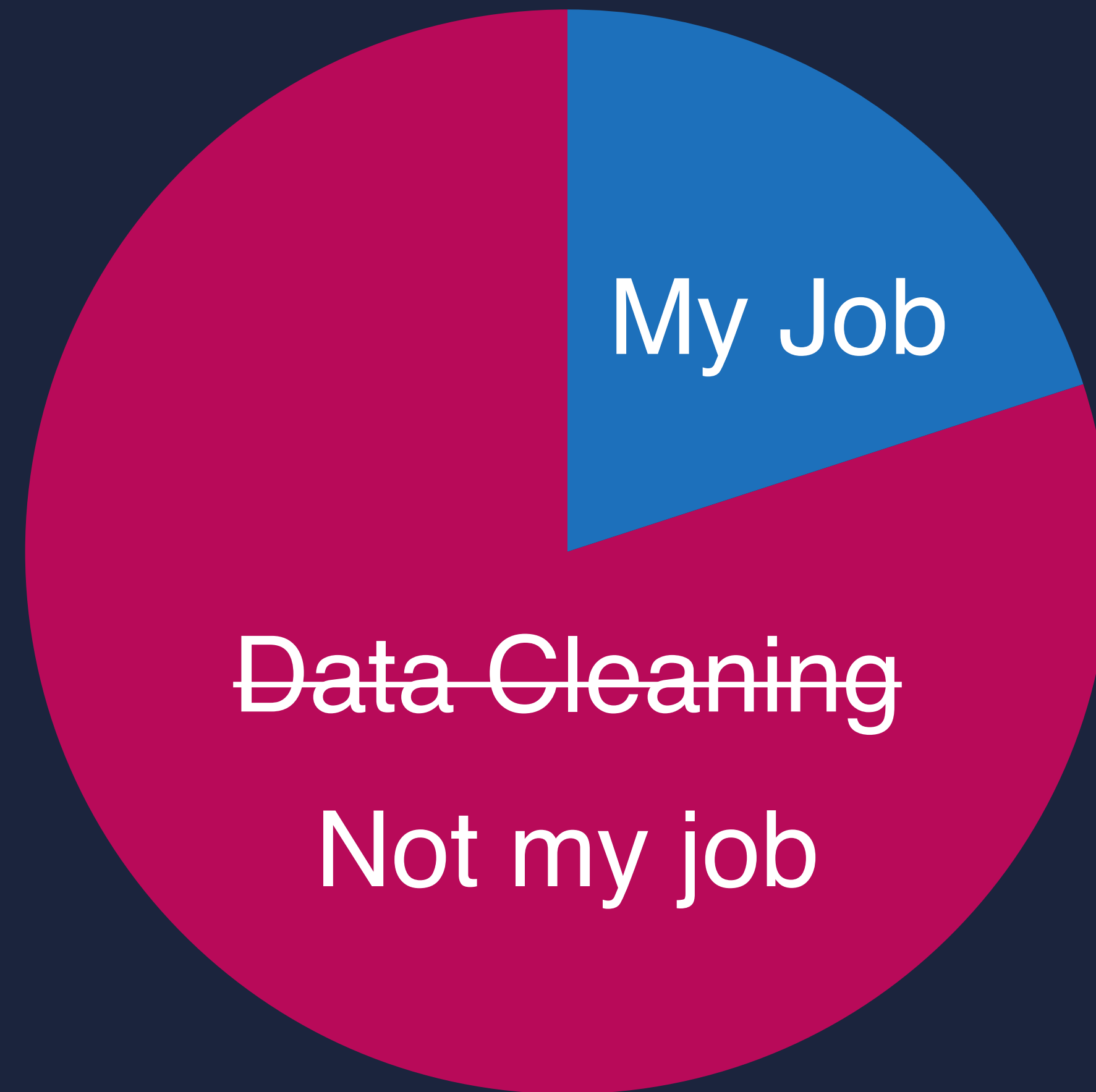
40 comments [share](#) [save](#) [hide](#) [give award](#) [report](#)

crosspost

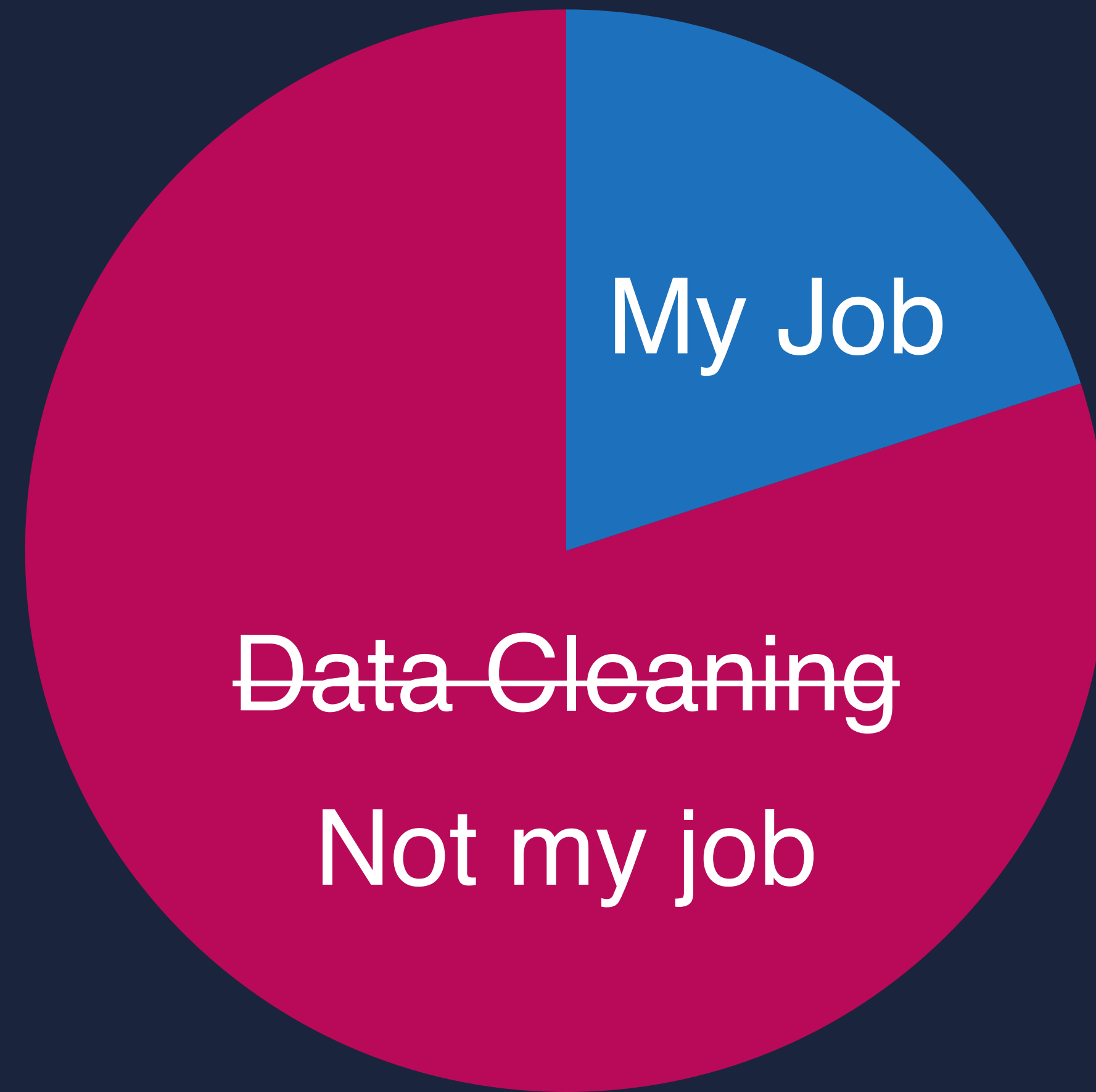
WHAT THEY SAY



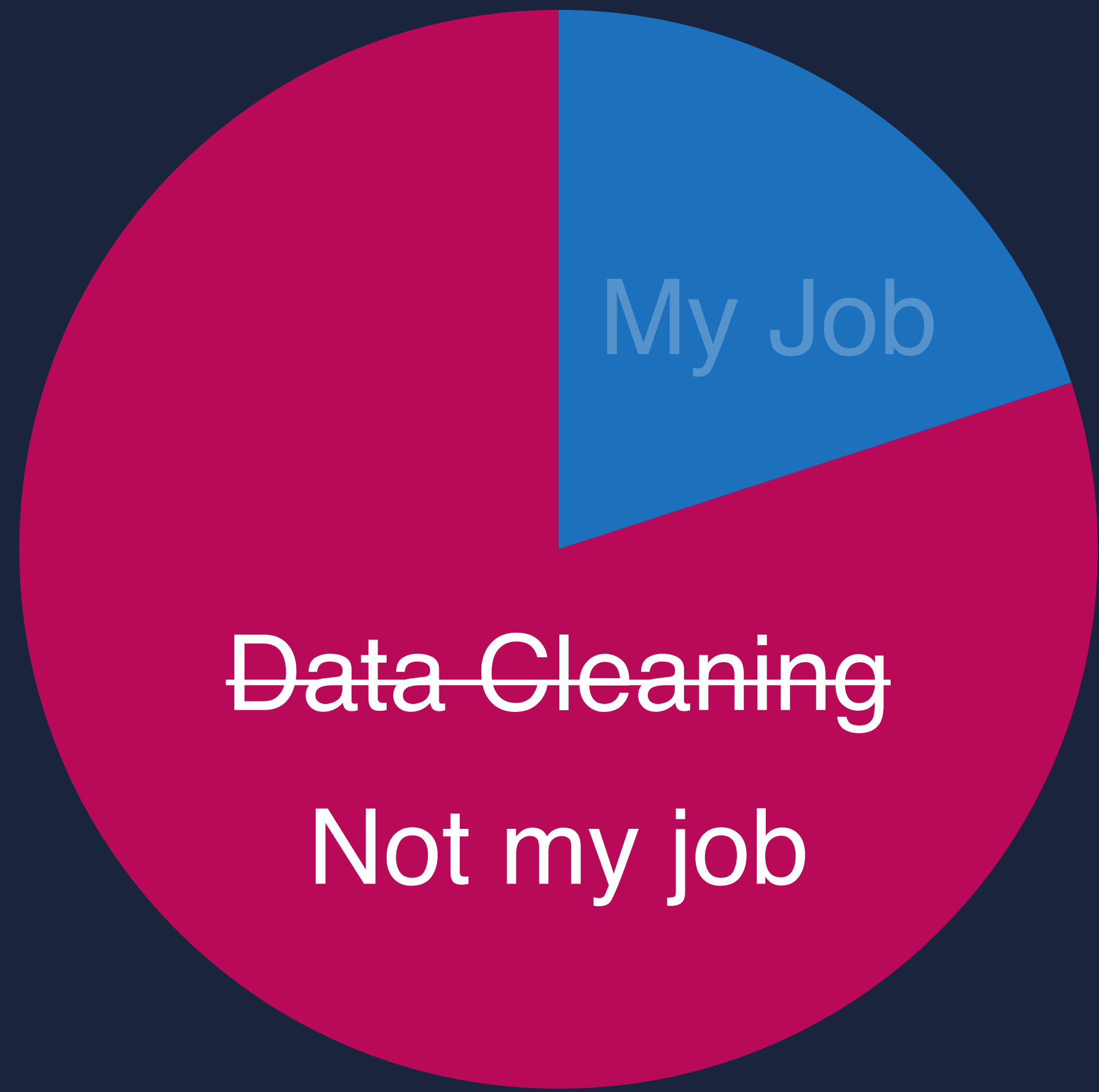
WHAT THEY MEAN



WHAT THEY MEAN



WHAT THEY MEAN



- Rolling their own infrastructure
- Repeated work
- Maintaining unreliable processes

FAILURE IS THE NORM

Big data strategies disappoint with 85 percent failure rate

BY JAMES WALKER NOV 23, 2017 IN TECHNOLOGY

[LISTEN](#) | [PRINT](#)

Big data projects aren't delivering the transformations companies had anticipated. 85% of all big data strategies are failing due to a combination of challenges, including the problems presented by legacy technologies and pre-existing corporate biases.

Engineers: I don't want to touch it.

Data scientist: I waste most of my time.

Business Leader: Failure is the norm.

2009: UI development is awful

- I spend 80% of my time fighting the browser



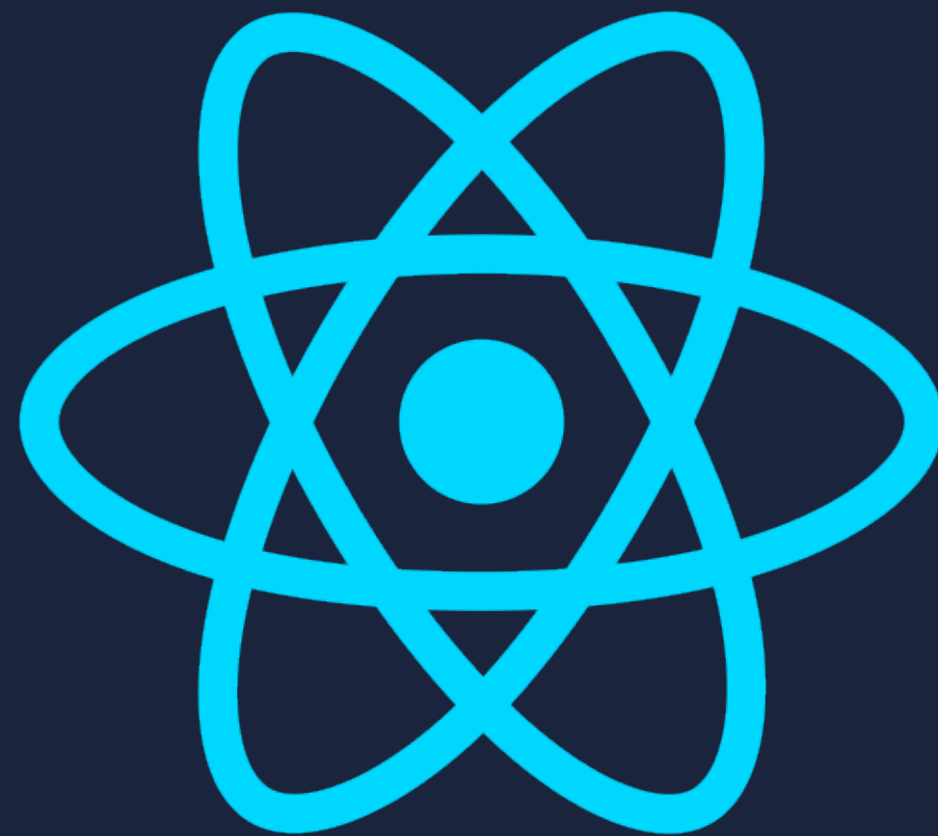
2009: UI development is awful

- I spend 80% of my time fighting the browser
- We can't change our UI—there's no testing
- It breaks all the time.
- Our engineers don't want to touch it

2019: A (UI) world transformed

Browsers did get better.

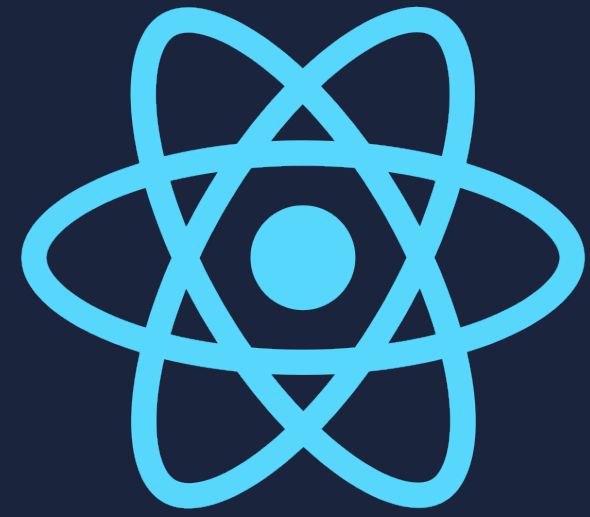
But it was the **software abstractions** that proved decisive.



Scripts → Full applications

React acknowledged complexity

It respected the discipline



React → Frontend Applications



Dagster → Data Applications

PRINCIPLES

- ▶ Solves a real problem
- ▶ Incremental adoption path
- ▶ Preserve tools that work
- ▶ Immediate value and productivity gains

WHAT WOULD YOU SAY

YOU DO HERE?



Graphs of functional computations
that produce and consume data assets



```
> pip install dagster
```

🐍 pagerank.py ×



🐍 pagerank.py

1

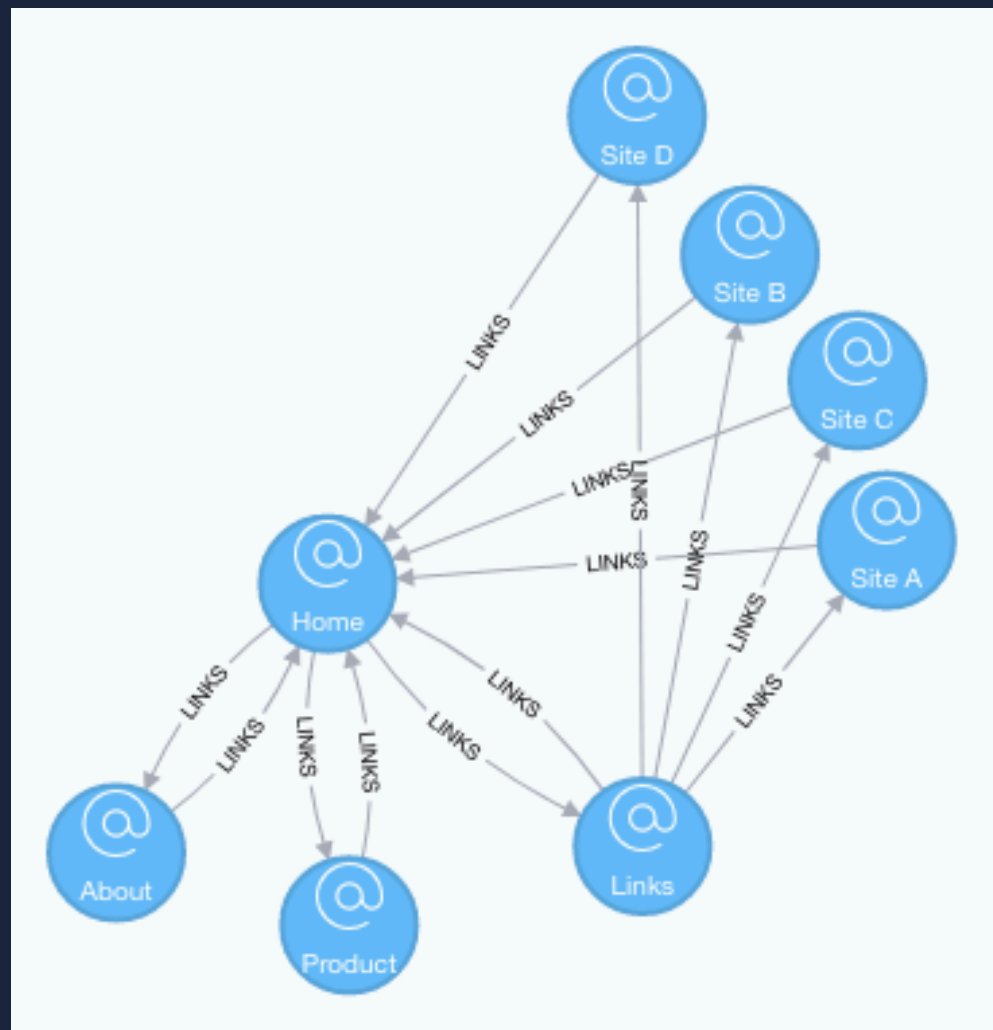


```
(dagster-3.7.1) [schrockn@mbp ~/code/dagster/examples/pyspark_pagerank (schrockn/screencast-notebooks)]$ dagit -f pagerank.py -n define_pipeline
```

```
}
```

DAGSTER CONCEPTS

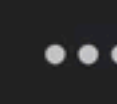
- Solid: A unit of functional computation
- Pipeline: A DAG of solids



Page Rank

on





original.py ▸ ...

```
82
83
84 if __name__ == "__main__":
85     if len(sys.argv) != 3:
86         print("Usage: pagerank <file> <iterations>", file=sys.stderr)
87         sys.exit(-1)
88
89     print(
90         "WARN: This is a naive implementation of PageRank and is given as "
91         + "an example!\nPlease refer to PageRank implementation provided "
92         + "by graphx",
93         file=sys.stderr,
94     )
95
96     # Initialize the spark context.
97     execute_pagerank(sys.argv[1], int(sys.argv[2]))
98
```





No pipeline selected

Select a pipeline in the navbar

DAGSTER CONCEPTS

- Solid
 - Inputs: Inputs are the data
 - Config: Config modifies how data is computed
- Pipeline

```

@solid(
    inputs=[InputDefinition('path', Path)],
    config_field=Field(Dict({'iterations': Field(Int)})),
)
def execute_pagerank(context, path):
    spark = SparkSession.builder.appName("PythonPageRank").getOrCreate()
    lines = spark.read.text(path).rdd.map(lambda r: r[0])
    links = lines.map(parseNeighbors).distinct().groupByKey().cache()
    ranks = links.map(lambda url_neighbors: (url_neighbors[0], 1.0))
    iterations = context.solid_config['iterations']

```



```

@solid(inputs=[InputDefinition('path', Path)])
def load_pagerank_data(_context, path):
    spark = SparkSession.builder.appName("PythonPageRank").getOrCreate()
    lines = spark.read.text(path).rdd.map(lambda r: r[0])
    return lines.map(parseNeighbors)

@solid(
    inputs=[InputDefinition('links')],
    config_field=Field(Dict({'iterations': Field(Int)})),
)
def execute_pagerank(context, links):
    spark = SparkSession.builder.appName("PythonPageRank").getOrCreate()
    cached_links = links.distinct().groupByKey().cache()
    ranks = cached_links.map(lambda url_neighbors: (url_neighbors[0], 1.0))
    iterations = context.solid_config['iterations']


```

```
@solid(inputs=[InputDefinition('path', Path)])
def load_pagerank_data(context, path):
    # ...

@solid(
    inputs=[InputDefinition('links')],
    config_field=Field(Dict({'iterations': Field(Int)})),
)
def execute_pagerank(context, links):
    # ...
```

```
@solid(inputs=[InputDefinition('path', Path)])
def load_pagerank_data(context, path):
    # ...

@solid(
    inputs=[InputDefinition('links')],
    config_field=Field(Dict({'iterations': Field(Int)})),
)
def execute_pagerank(context, links):
    # ...
```



```
def define_pipeline():
    return PipelineDefinition(
        name='pagerank',
        solids=[load_pagerank_data, execute_pagerank],
        dependencies={
            'execute_pagerank': {
                'links': DependencyDefinition('load_pagerank_data')
            }
        },
    )
```

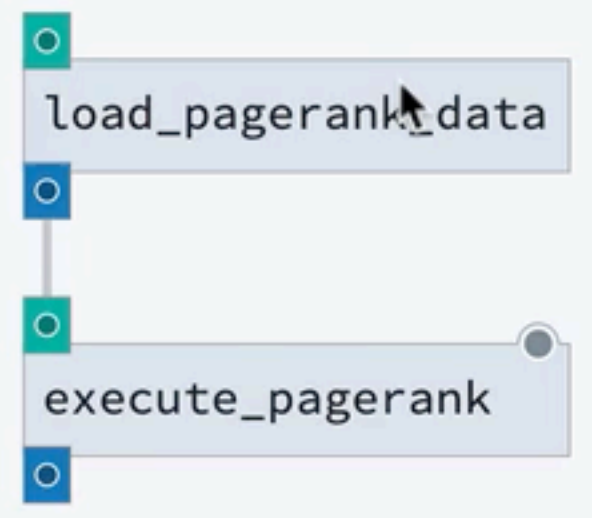
Select a Solid... ▾ Filter...

Pipeline
pagerank

DESCRIPTION ▾

CONTEXTS ▾

```
default  
  
/* A configuration dictionary with typed  
fields */  
{  
  log_level?: log_level  
}
```



DAGSTER CONCEPTS

- Solid
 - Inputs & Config
- Pipeline
- Dependencies

Before:

```
@solid(inputs=[InputDefinition('path', Path)])
def load_pagerank_data(_context, path):
    # Initialize the spark context.
    spark = SparkSession.builder.appName("PythonPageRank").getOrCreate()

    # two urls per line with space in between)
    lines = spark.read.text(path).rdd.map(lambda r: r[0])
```

After:

```
@solid(inputs=[InputDefinition('path', Path)])
def load_pagerank_data(context, path):
    # two urls per line with space in between)
    lines = context.resources.spark.read.text(path).rdd.map(lambda r: r[0])
```

```
from dagster_pyspark import spark_session_resource
```

```
def define_pyspark_pagerank():  
    return PipelineDefinition(  
        name='pyspark_pagerank',  
        context_definitions={  
            'local': PipelineContextDefinition(  
                resources={'spark': spark_session_resource}  
            )  
        },  
    )
```



pagerank

Explore Execute Runs

0.3.5

Select a Solid... Filter...

Info Types

Pipeline pagerank

DESCRIPTION

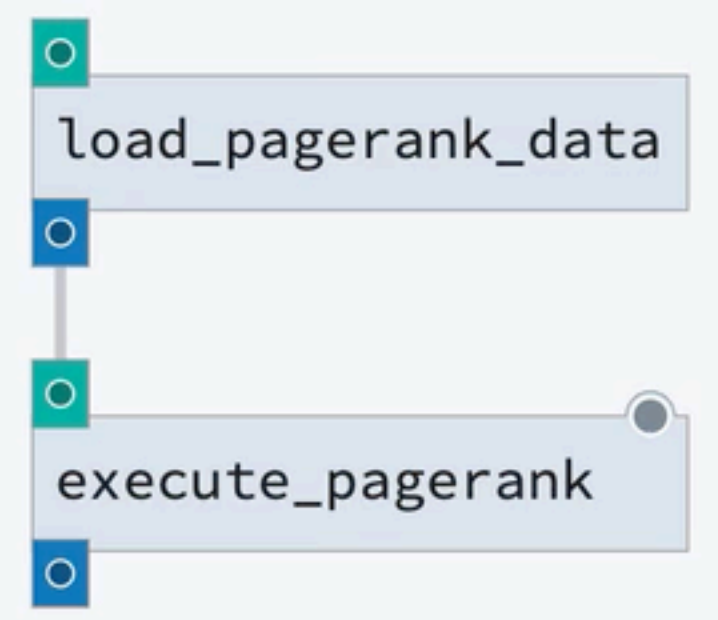
CONTEXTS

local

```
/* A configuration dictionary with typed fields */
{
  log_level?: log_level
}
```

spark

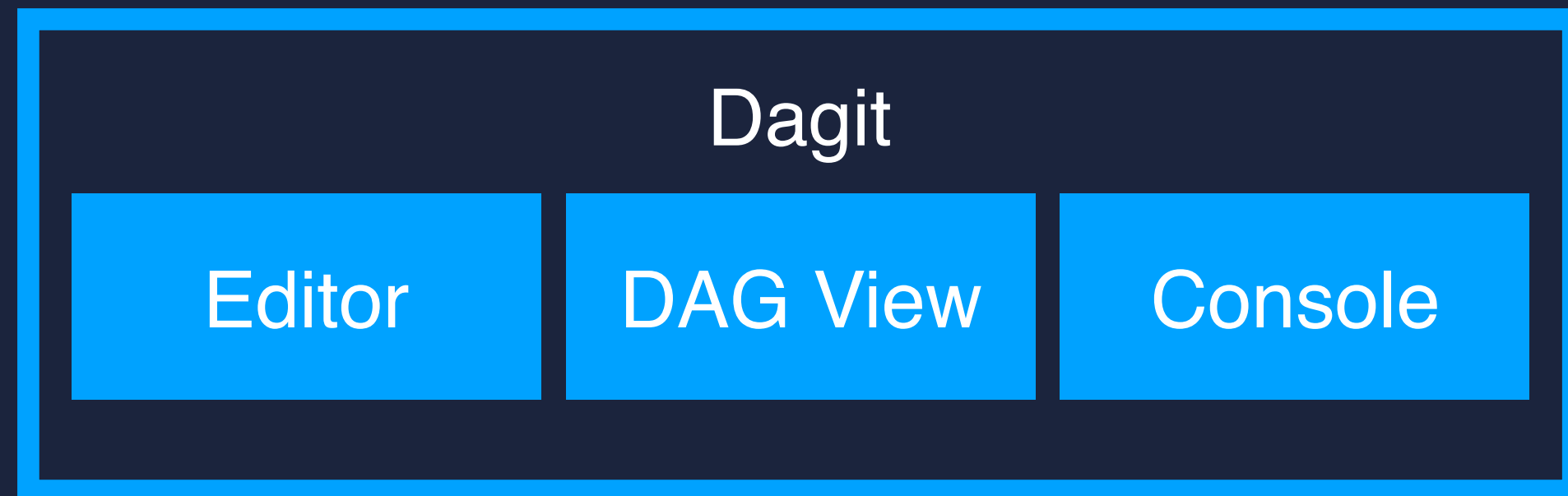
```
/* A configuration dictionary with
{
  spark_conf?: {
    spark?: {
      app?: {
        /* Application Properties:
```



DAGSTER CONCEPTS

- Solid
 - Inputs & Config
- Pipeline
- Dependencies
- Context
 - Logging: Structured Logging
 - Resources: Connections, Services, Etc

Beautiful, High-Quality Tools



API



Python library

Dagster Libraries and Integrations

PySpark

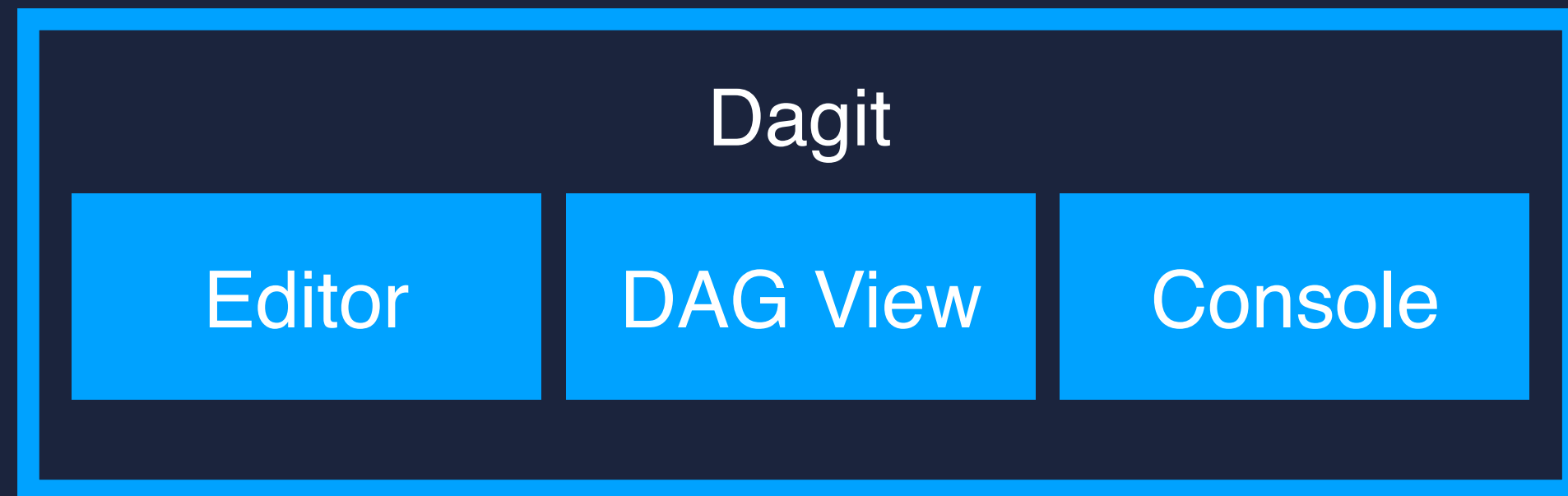
Graph of Functional Computation

API

- Queryable and Introspectable
- Operable
 - Executable and Configurable
- Monitorable
 - Logging and Live Subscriptions

Dagster: a platform for building tools

Beautiful, High-Quality Tools



API



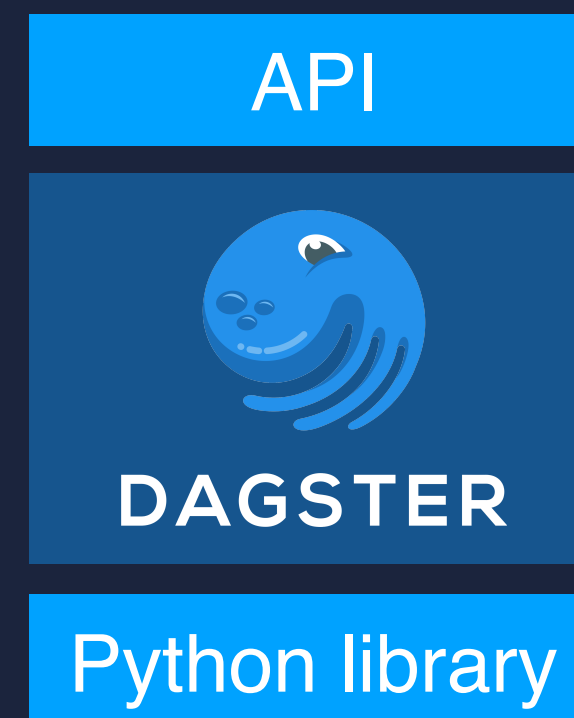
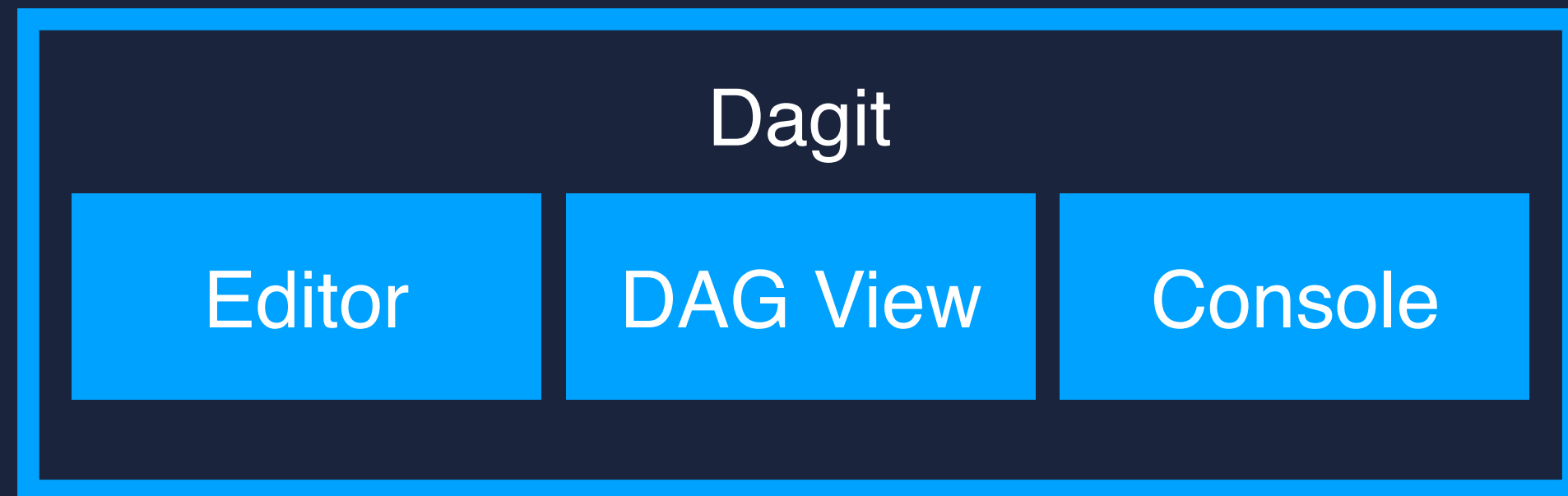
DAGSTER

Python library

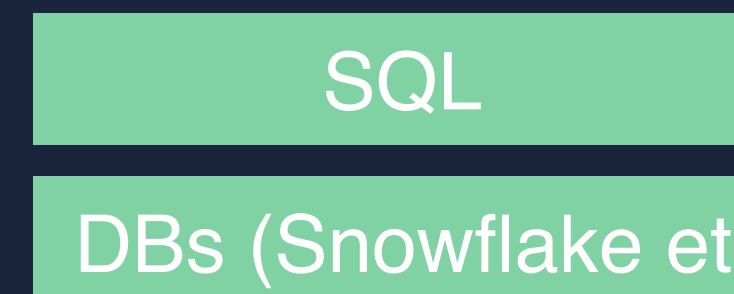
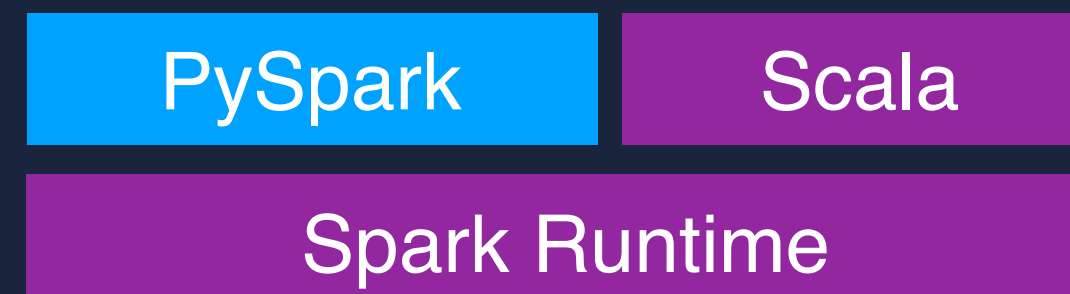
Dagster Libraries and Integrations

PySpark

Beautiful, High-Quality Tools



Dagster Libraries and Integrations



```
event_ingest = SparkSolidDefinition(
    name='event_ingest',
    main_class='io.dagster.events.EventPipeline',
    description='Ingest events from JSON to Parquet',
)

snowflake_load = SnowflakeLoadSolidDefinition(
    'snowflake_load', table='events'
)

return PipelineDefinition(
    name='event_ingest_pipeline',
    solids=[download_from_s3, gunzipper, event_ingest, snowflake_load],
    dependencies={
```

No pipeline selected

Select a pipeline in the navbar

- Open Source, Python Library
- Multi-lingual integration
- Beautiful Tooling

The Jupyter logo consists of a central orange smiley face with two dark grey dots above and below it. The word "jupyter" is written in a dark grey, lowercase, sans-serif font across the middle of the smiley face.

jupyter

O'REILLY®

jupytercon

Brought to you by NumFOCUS Foundation
and O'Reilly Media Inc.

jupytercon.com
#JupyterCon

I don't like notebooks.

Joel Grus

*Allen Institute for
Artificial Intelligence*

Beyond Interactive: Notebook Innovation at Netflix



Netflix Technology Blog [Follow](#)

Aug 16, 2018 · 13 min read

By [Michelle Ufford](#), [M Pacer](#), [Matthew Seal](#), and [Kyle Kelley](#)

Notebooks have rapidly grown in popularity among data scientists to become

README.md



build **passing**  codecov **90%** docs **passing**  highlight binder  cli binder code style black

papermill is a tool for parameterizing, executing, and analyzing Jupyter Notebooks.

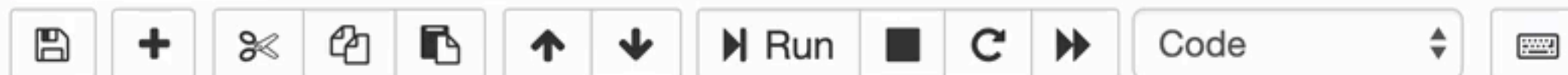
Papermill lets you:

- **parameterize** notebooks
- **execute** notebooks

No pipeline selected

Select a pipeline in the navbar





In [18]:

```
import dagstermill as dm
from event_pipeline_demo.repository import define_repo
dm.register_repository(define_repo())
import pandas as pd
```

In [19]:

parameters ✕

```
df = pd.read_csv('table_to_df.csv')
```

In [20]:

```
df.head()
```

Out[20]:

	latitude	longitude	timestamp
0	46.18396	6.10237	2019-01-02 15:59:59
1	37.74615	-25.66689	2019-01-02 15:59:57
2	-30.60106	-71.19901	2019-01-02 15:59:54
3	35.85000	117.70000	2019-01-02 15:59:52
4	6.03333	37.55000	2019-01-02 15:59:47

In [21]:

```
import os
import geopandas as geo
import matplotlib.pyplot as plt
import pandas as pd
from matplotlib.backends.backend_agg import FigureCanvasAgg as FigureCanvas
```

Select a Solid...

Filter...

[Info](#) [Types](#)

Pipeline

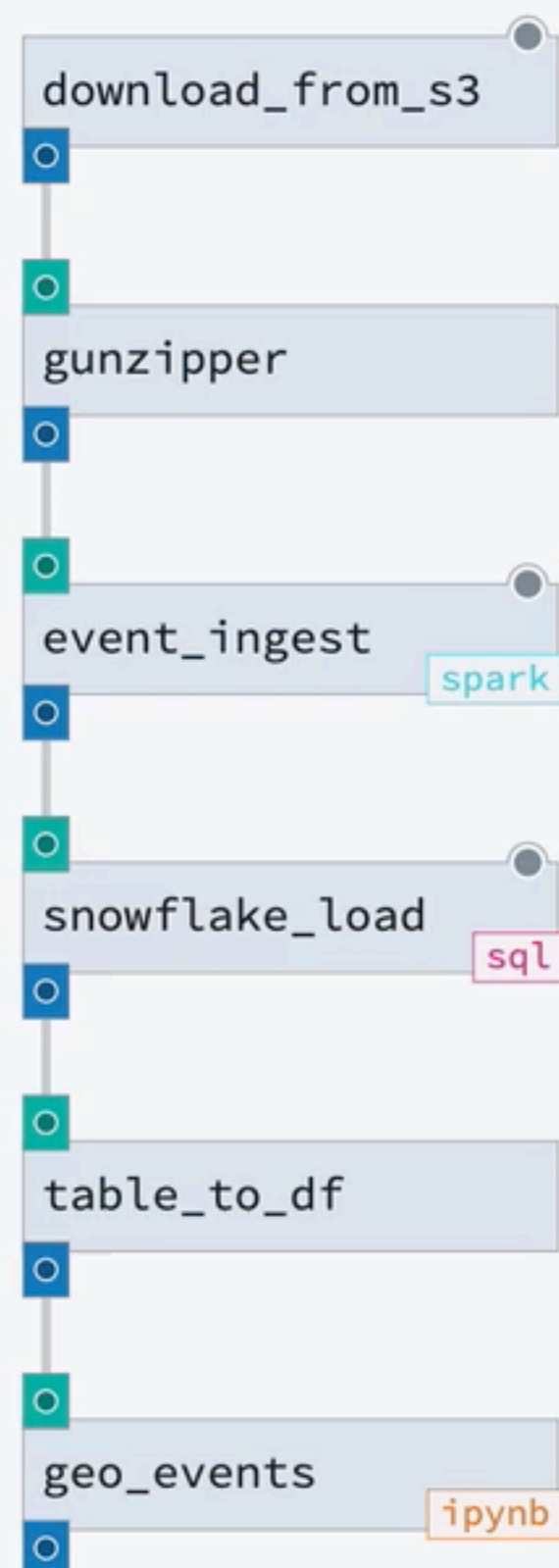
event_ingest_pipeline

DESCRIPTION

CONTEXTS

default

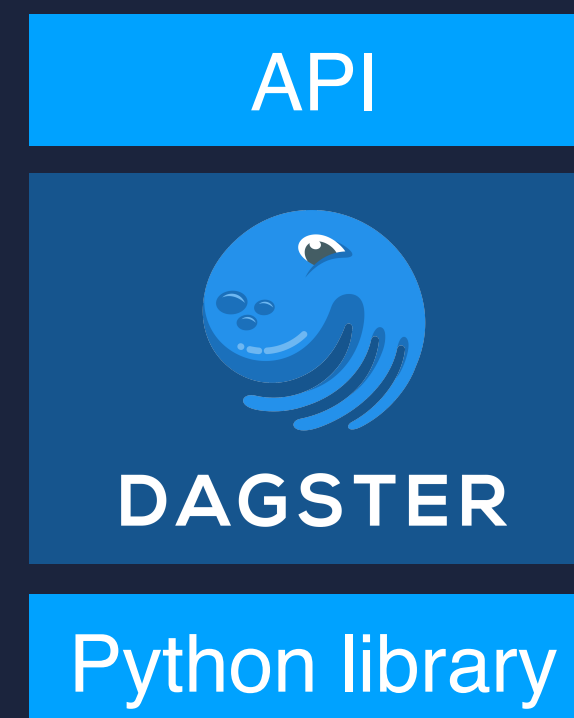
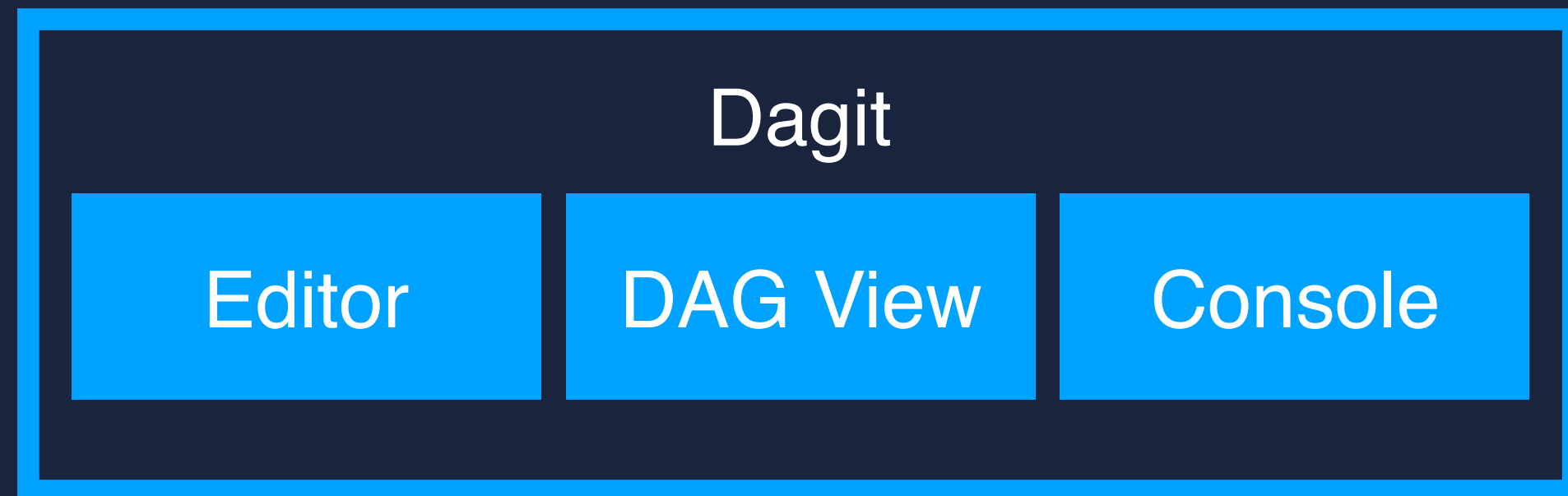
```
/* A configuration dictionary with typed
   fields */
{
  log_level?: log_level
}
```



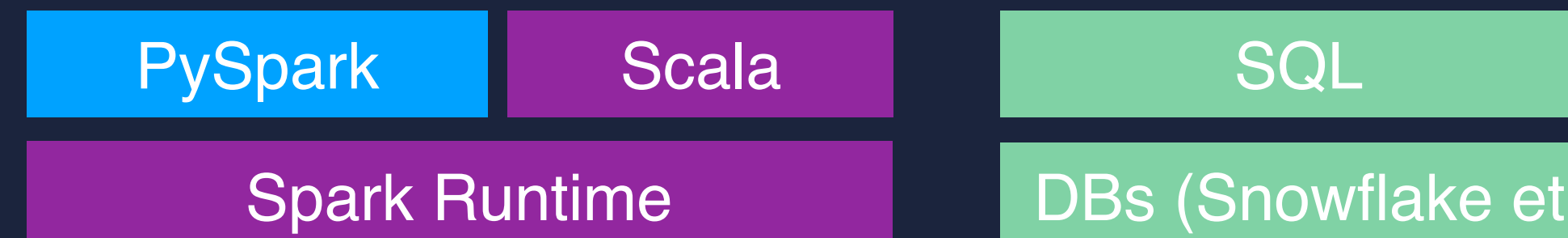
(dagster-3.7.1) [schrockn@mbp ~/code/dagster (schrockn/airflowize *)]\$

}

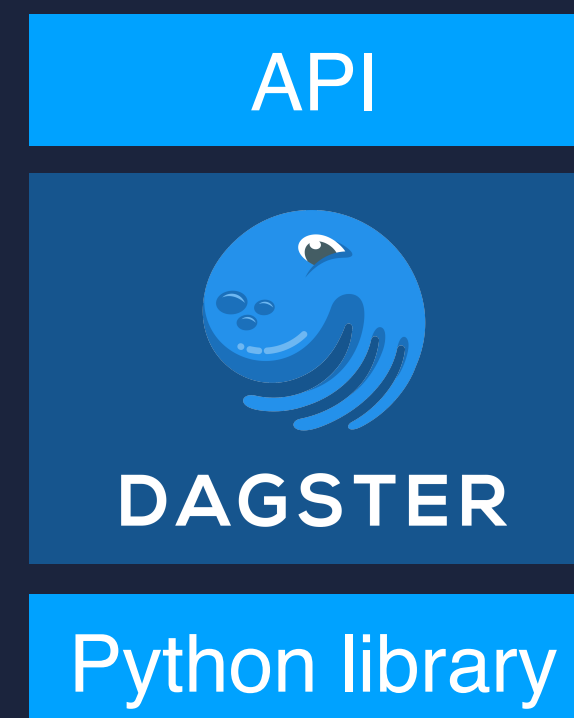
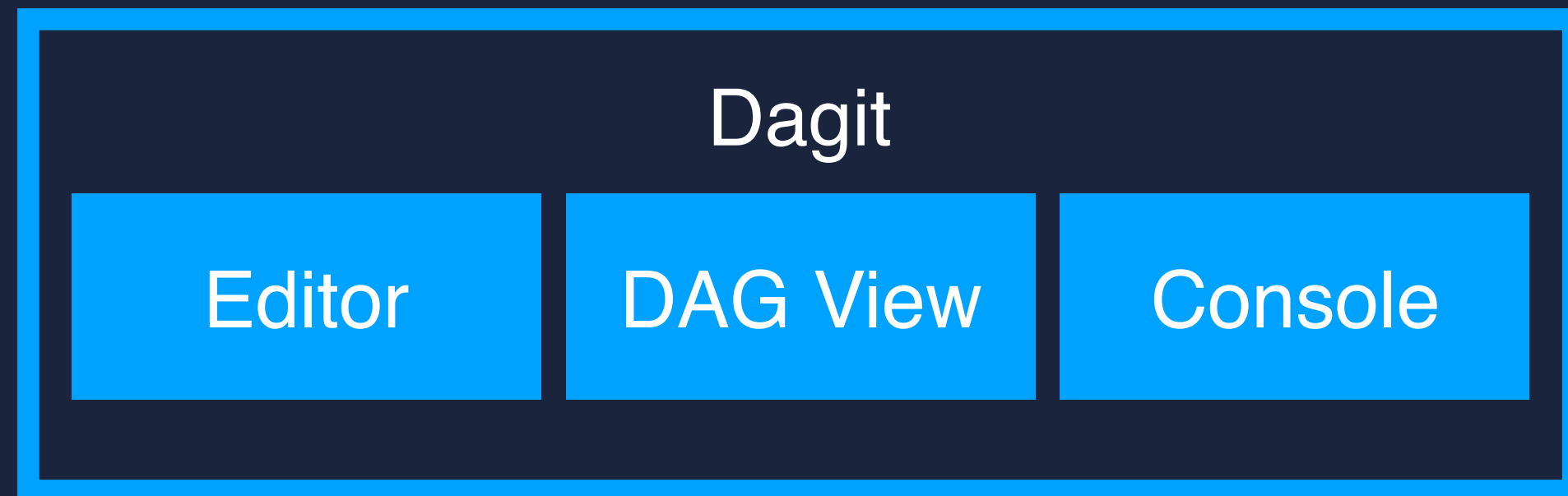
Beautiful, High-Quality Tools



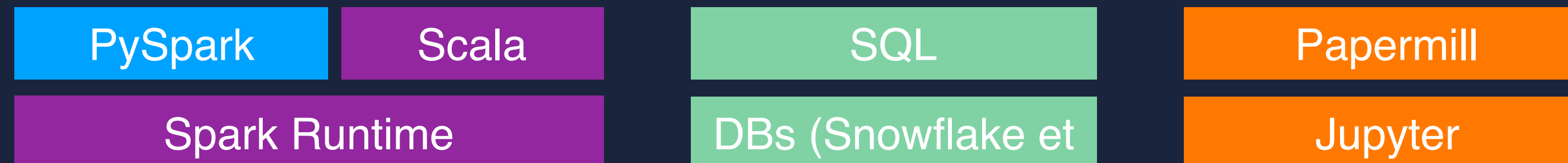
Dagster Libraries and Integrations



Beautiful, High-Quality Tools



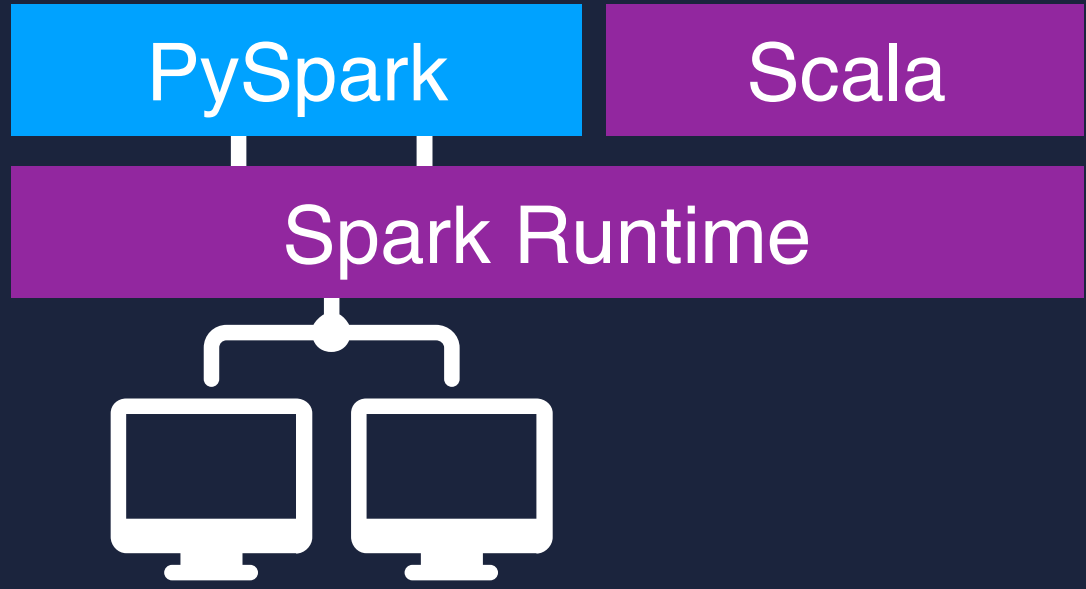
Dagster Libraries and Integrations



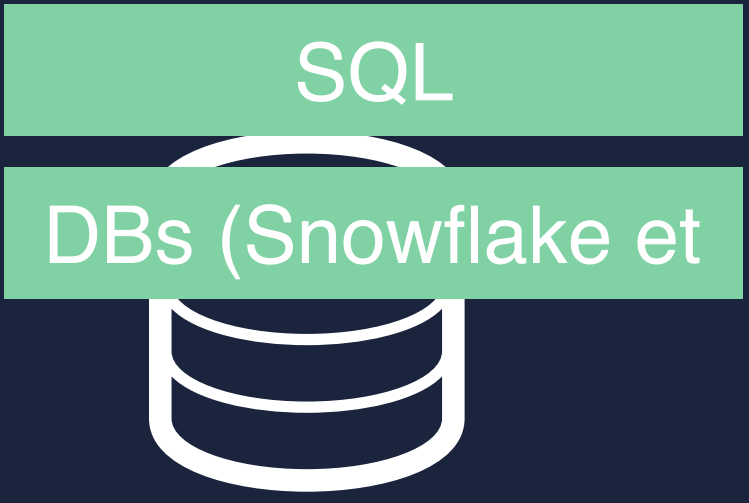


Python library

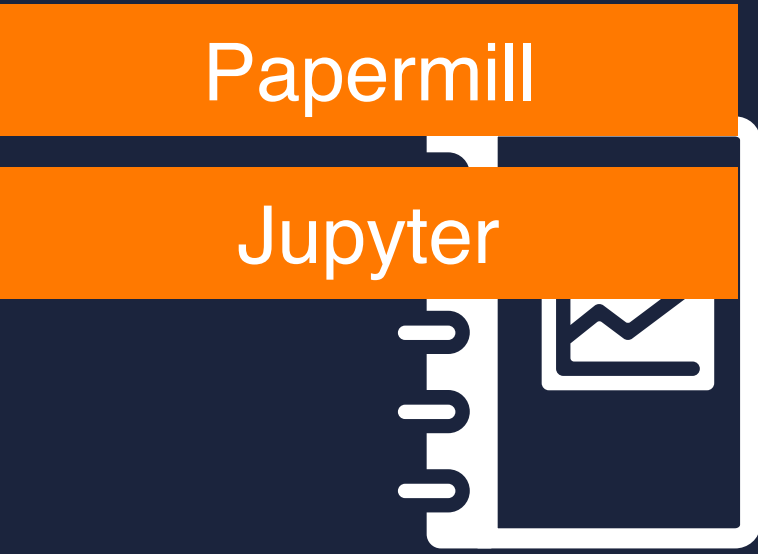
Dagster Libraries and Integrations



Data Engineering

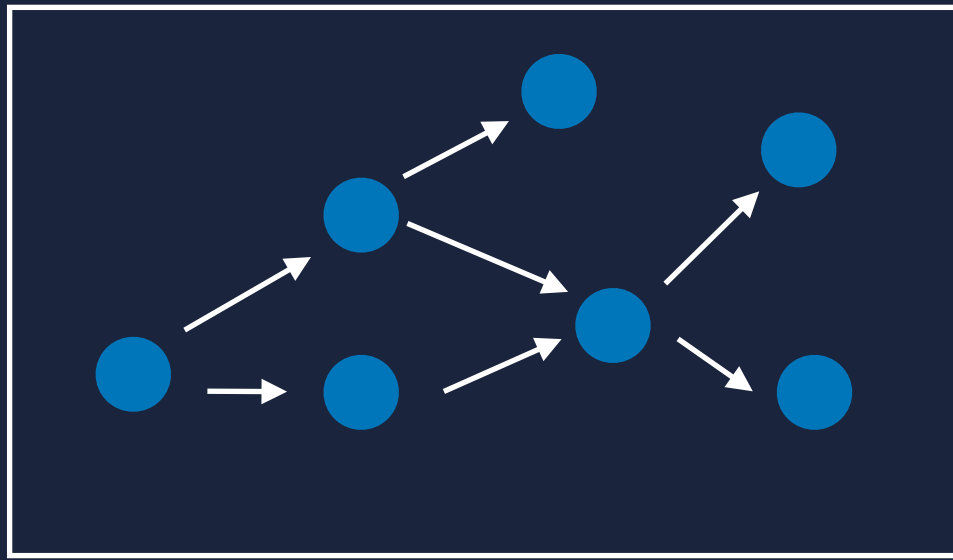


Analysts

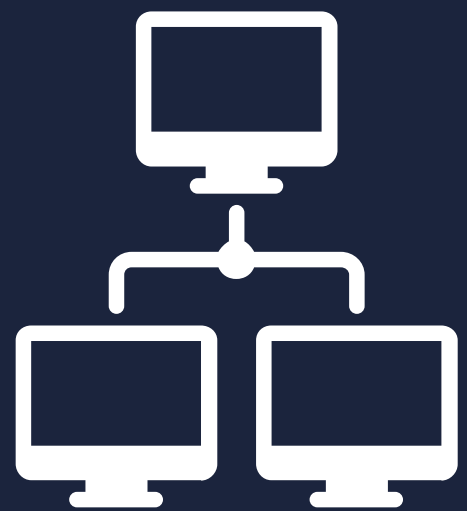


Data Science

Data Application: A Graph of Computations



Dagster Libraries and Integrations



Data Engineering



Analysts



Data Science

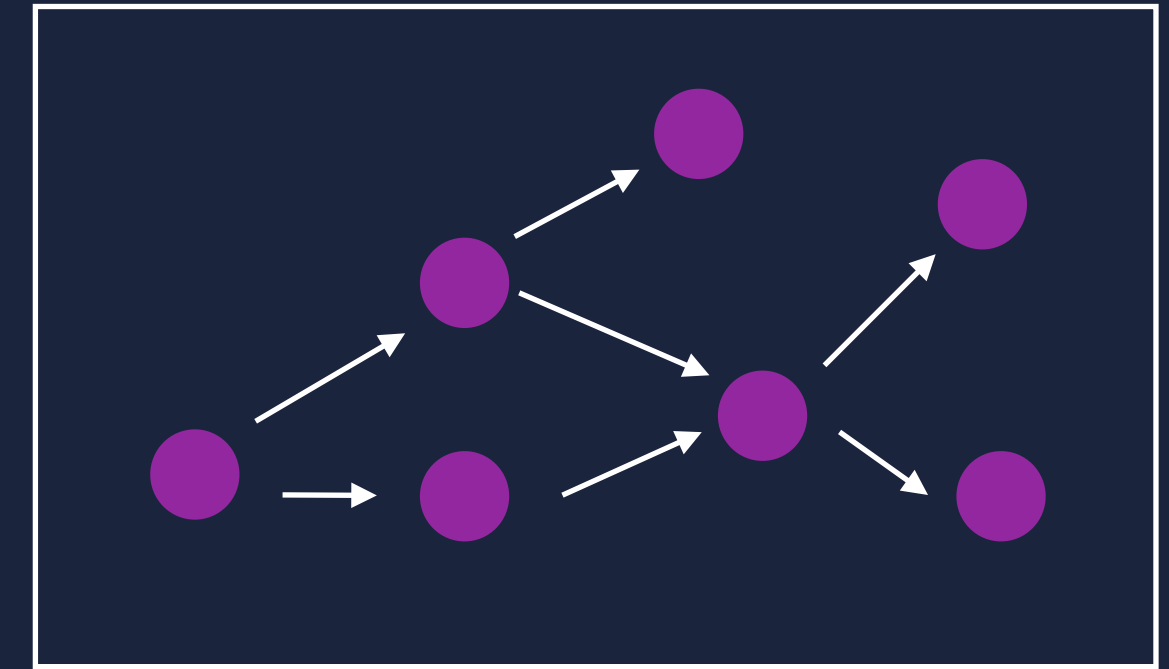


API

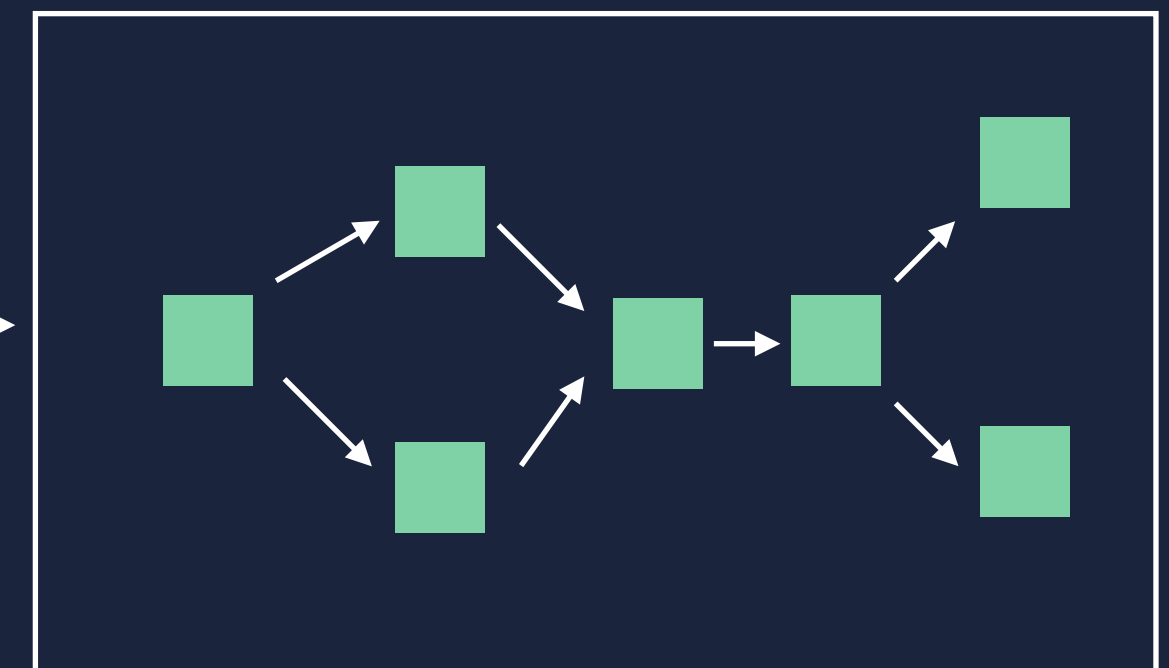
dagit

dagster cli

Local Executor



Airflow



dagster-

```
1 from dagster.utils import load_yaml_from_globs
2 from dagster_airflow.factory import make_airflow_dag
3 from event_pipeline_demo.pipelines import define_event_ingest_pipeline
4
5 PATH_TO_CONFIG_YML = '/Users/schrockn/event-pipeline-demo/environments/*.yaml'
6
7 dag, steps = make_airflow_dag(
8     define_event_ingest_pipeline(),
9     env_config=load_yaml_from_globs(PATH_TO_CONFIG_YML),
10    dag_id='event_ingest_pipeline',
11    dag_description='A demo Airflow DAG corresponding to the event ingest pipeline.',
12    dag_kwargs=None,
13    op_kwargs=None,
14 )
15
```


2. bash

(dagster-3.6.7) [schrockn@mbp ~/code/dagster (schrockn/airflowize *)]\$ airflow webserver

DATA ENGINEERING

- An emerging discipline
- At an inflection point

Scripts  Data Applications

ELEGANT PROGRAMMING

MODEL

NEW, BEAUTIFUL TOOLING

FLEXIBLE AND INCREMENTAL

FLEXIBLE AND INCREMENTAL

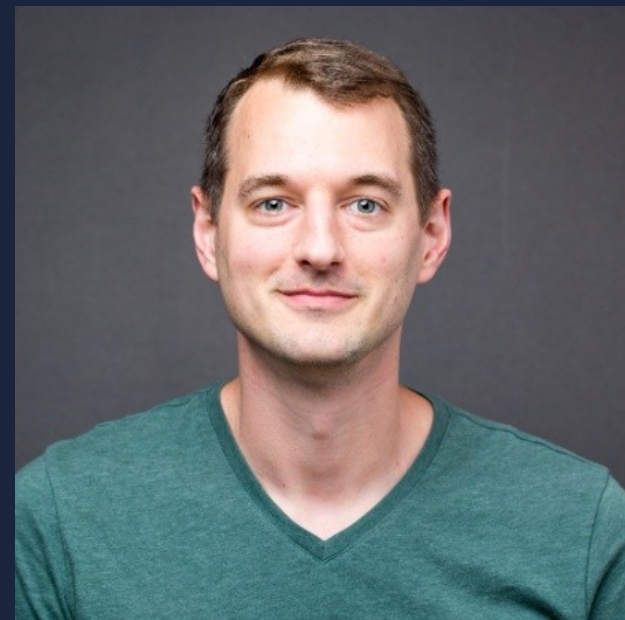
- Use your tools
- Preserve your code
- Deploy to your infrastructure
- Adopt incrementally

And there is a ton of work to do

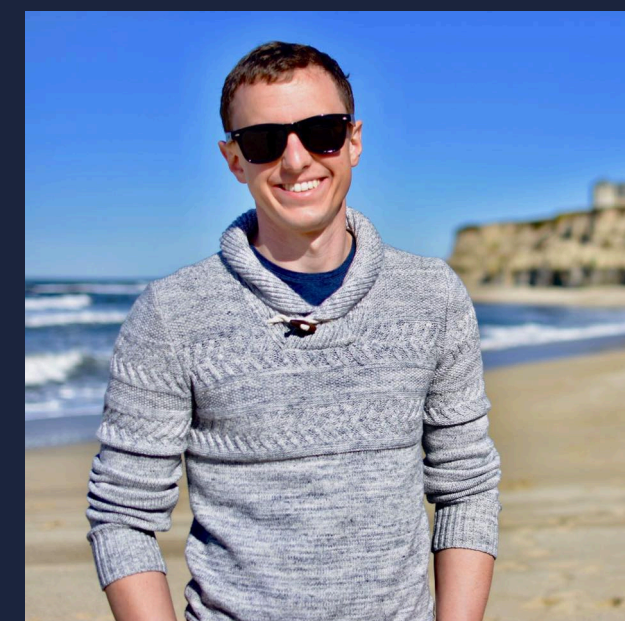
TEAM



Max Gasner



Nate Kupp



Alex Langenfeld

THANK YOU



Ben Gotow



Mikhail Novikov



Uma Roy

THANK YOU



Abe Gong

Superconductive Health



<https://github.com/dagster-io/dagster>

Join the team. Partner with us.

<https://elementl.com>

schrockn@elementl.com

@schrockn



DAGSTER