

# Converting Legacy Schedulers

## *& Announcing Astronomer Orbiter*

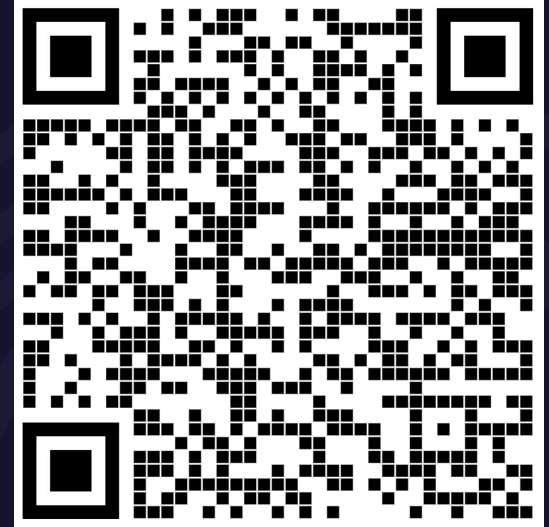
Fritz Davenport





**Fritz** @astronomer.io

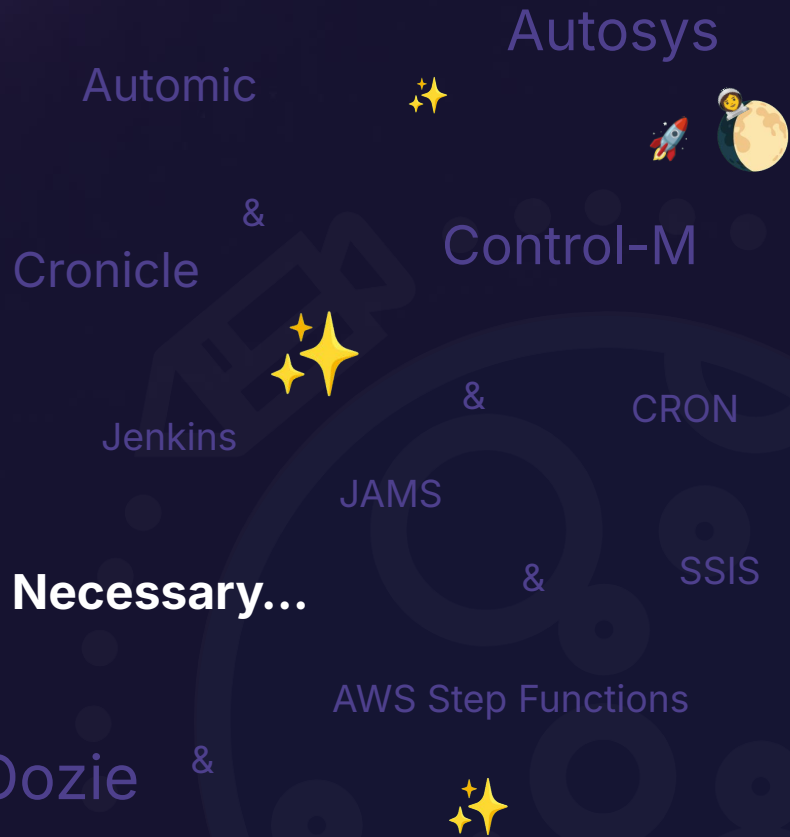
Principal Data Engineer





# Migrations

Difficult



**Necessary...**



# Migrations are...

Difficult

Risky

Boring

Complex

Slow

**Necessary...**



## Various Tools Exist

Control-M to Airflow

<https://github.com/GoogleCloudPlatform/dagify>

Oozie to Airflow

<https://github.com/GoogleCloudPlatform/oozie-to-airflow>

<https://www.astronomer.io/white-papers/migrating-from-oozie...>

Pentaho to Airflow

<https://medium.com/@swapnilspra/etl-code-migration-from-pentaho...>

& more



# Migration Process Overview

## 1. Setup

networking, infrastructure, authentication, etc

## 2. Identify, Group

common patterns, stakeholders, related workloads

## 3. Migrate

migrate migrate migrate migrate...



### a. Translate

Origin System → 🤔 → Airflow

### b. Test, Deploy



## 4. Adopt, Refactor



/input\_folder



# Structured Workflows

```
1 <coordinator-app name="hello-coord" frequency="${coord:days(1)}"  
2 ..... start="2009-01-02T08:00Z" end="2009-01-04T08:00Z"  
3 ..... xmlns="uri:oozie:coordinator:0.1">  
4 ..... <controls>  
5 ..... <timeout>10</timeout>  
6 ..... <concurrency>${concurrency_level}</concurrency>  
7 ..... <execution>${execution_order}</execution>  
8 ..... <throttle>${materialization_throttle}</throttle>  
9 ..... </controls>
```

```
1 <workflow-app xmlns="uri:oozie:workflow:1.0" name="demo-wf">  
2 ..... <start to="cleanup-node"/>  
3 .....  
4 ..... <action name="cleanup-node">  
5 ..... <fs>  
6 ..... <delete path="${nameNode}/user/${wf:user()}/${examplesRoot}/  
7 ..... </fs>  
8 ..... <ok to="hdfs-node"/>  
9 ..... <error to="fail"/>  
10 ..... </action>  
11 .....  
12 ..... <action name="hdfs-node">  
13 ..... <fs>  
14 ..... <move source="${nameNode}/user/${wf:user()}/${examplesRoot}/  
15 ..... <target="/user/${wf:user()}/${examplesRoot}/output-data  
16 ..... </fs>  
17 ..... <ok to="end"/>
```

Oozie  
(XML)



/input\_folder



# Structured Workflows

```
1 {
2   "Defaults" : {
3     "Application" : "SampleApp",
4     "SubApplication" : "SampleSubApp",
5     "RunAs" : "USERNAME",
6     "Host" : "HOST",
7     "Job": {
8       "When" : {
9         "Months": ["JAN", "OCT", "DEC"],
10        "MonthDays": ["22", "1", "11"],
11        "WeekDays": ["MON", "TUE", "WED", "THU", "FRI"],
12        "FromTime": "0300",
13        "ToTime": "2100"
14      }
15    }
16  },
17  "AutomationAPISampleFlow": {
18    "Type": "Folder",
19    "Comment" : "Code reviewed by John",
20    "CommandJob": [
21      {
22        "Type": "Job:Command",
23        "Command": "echo my 1st job"
24      },
25      {
26        "Type": "Job:Script",
27        "FilePath": "SCRIPT_PATH",
28        "FileName": "SCRIPT_NAME"
29      }
30    ],
31    "Flow": {
32      "Type": "Flow",
33      "Sequence": ["CommandJob", "ScriptJob"]
34    }
35  }
36 }
```

Client machine details

Schedule interval

Control-M  
(JSON, XML)

Job specifics

Dependency





```
8 with DAG(
9     dag_id="sample_workflow_with_external_event.demo_fin_sample_workflow",
10    schedule=None,
11    start_date=DateTime(year=1970, month=1, day=1),
12    catchup=False,
13    tags=["sample_workflow_with_external_event"],
14 ):
15     demo_fin_demo_service_task = EmptyOperator(
16         task_id="demo_fin_demo_service",
17         doc="Sample job flow with multiple branches and using site standards",
18     )
19     demo_fin_dynamic_event_task = EmptyOperator(
20         task_id="demo_fin_dynamic_event",
21     )
22     demo_fin_job000_task = BashOperator(
23         task_id="demo_fin_job000",
24         bash_command="sleep 300",
25     )
26     demo_fin_job001_task = BashOperator(
27         task_id="demo_fin_job001",
28         bash_command="sleep 300",
29         on_success_callback=send_smtp_notification(
30             to="alerts@astronomer.io",
31             smtp_conn_id="SMTP",
32             subject="Error occurred",
33             html_content="Job-{{ti.task_id}}-failed",
34         ),
35     )
```

/output\_folder

/dags

workflow\_a.py

...

...

packages.txt

requirements.txt

# Apache Airflow Project



/input\_folder

/output\_folder



/dags

workflow\_a.py

...

...

packages.txt

requirements.txt

Structured  
Workflows

✨ Translate ✨

Apache  
Airflow Project



/input\_folder

/output\_folder



1

```
{  
  "workflowA": {  
    "some_property": ...  
    ...  
    "taskB": {  
      "other_property": ...  
      ...  
      ...  
    }  
    "taskC": { ... }  
    "dependency": {  
      "taskB": "taskC"  
      ...  
    }  
  }  
}
```

"Intermediate Representation"  
Python Dictionary



translate.py

/input\_folder

/output\_folder



```
{  
  "workflowA": {  
    "some_property": ...  
    ...  
    "taskB": {  
      "other_property": ...  
      ...  
    }  
    "taskC": { ... }  
    "dependency": {  
      "taskB": "taskC"  
      ...  
    }  
  }  
}
```

1

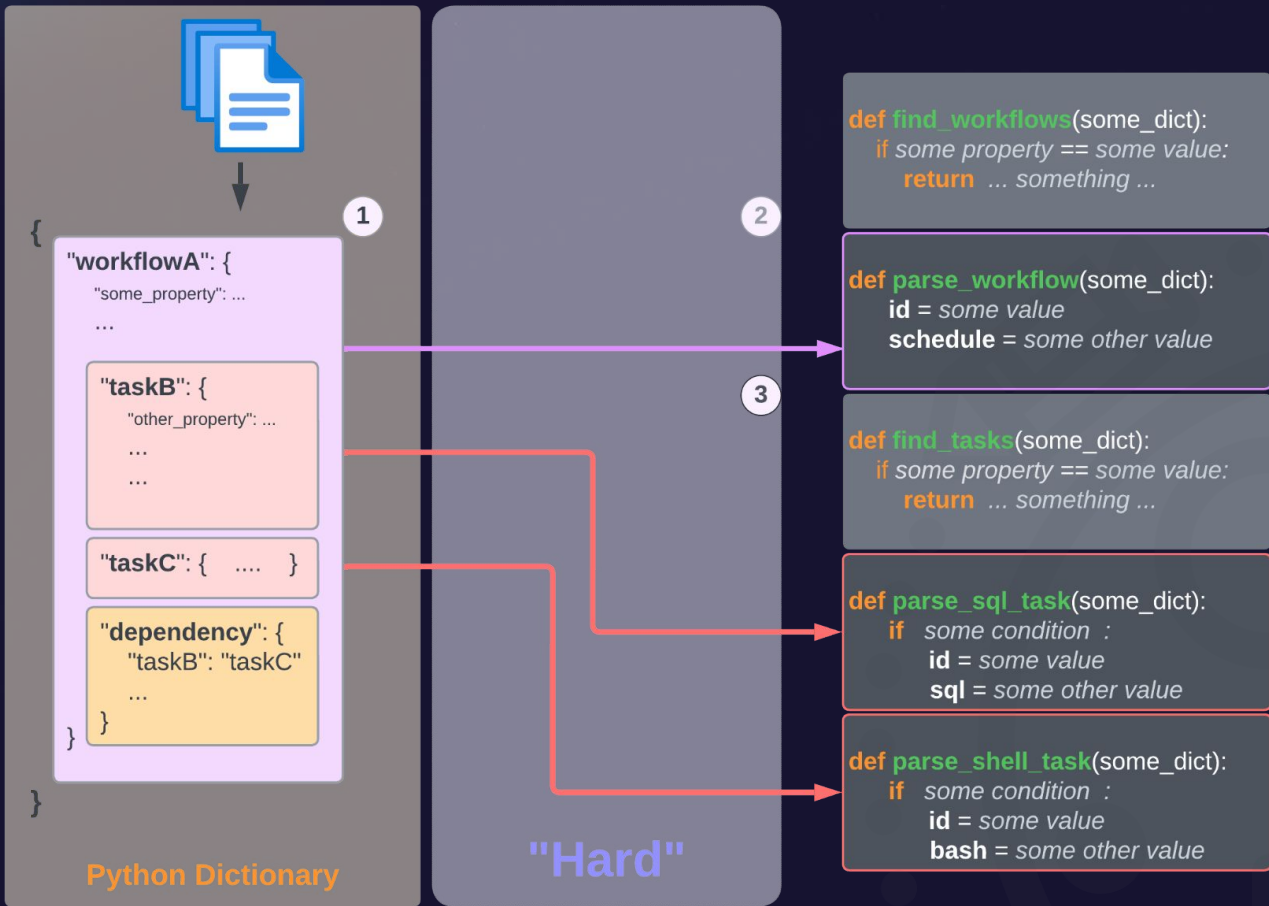
2

```
def find_workflows(some_dict):  
    if some_property == some_value:  
        return ... something ...
```

```
def find_tasks(some_dict):  
    if some_property == some_value:  
        return ... something ...
```

Python Dictionary

"Easy"





translate.py

/input\_folder

```
1 {
2     "Defaults" : {
3         "Application" : "SampleApp",
4         "SubApplication" : "SampleSubApp",
5         "RunAs" : "USERNAME",
6         "Host" : "HOST",
7         "Job": {
8             "When" : {
9                 "Months": ["JAN", "OCT", "DEC"],
10                "MonthDays": ["22", "1", "11"],
11                "WeekDays": ["MON", "TUE", "WED", "THU", "FRI"],
12                "FromTime": "0300",
13                "ToTime": "2100"
14            }
15        },
16    },
17    "AutomationAPISampleFlow": {
18        "Type": "Folder",
19        "Comment": "Code reviewed by John",
20        "CommandJob": {
21            "Type": "Job:Command",
22            "Command": "echo my 1st job"
23        },
24        "ScriptJob": {
25            "Type": "Job:Script",
26            "FilePath": "SCRIPT_PATH",
27            "FileName": "SCRIPT_NAME"
28        },
29        "Flow": {
30            "Type": "Flow",
31            "Sequence": ["CommandJob", "ScriptJob"]
32        }
33    }
```

```
def find_workflows(some_dict):
    if some property == some value:
        return ... something ...
```

```
def parse_workflow(some_dict):
    id = some value
    schedule = some other value
```

```
def find_tasks(some_dict):
    if some property == some value:
        return ... something ...
```

```
def parse_sql_task(some_dict):
    if some condition :
        id = some value
        sql = some other value
```

```
def parse_shell_task(some_dict):
    if some condition :
        id = some value
        bash = some other value
```



translate.py

/input\_folder



```
@dag_filter_rule  
def find_workflows(some_dict):
```

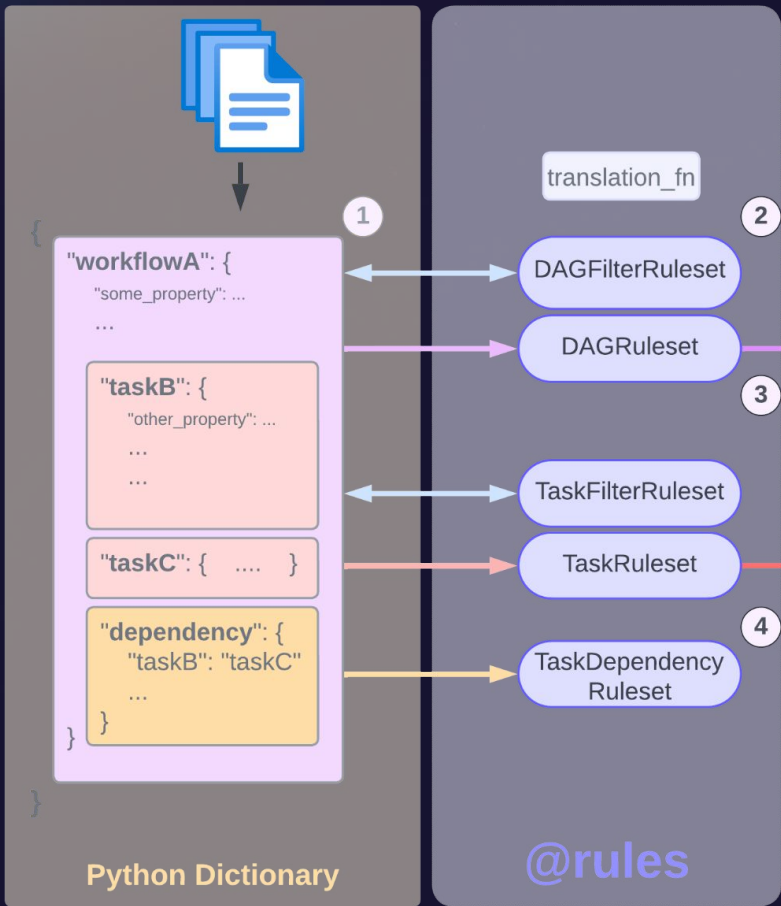
**Rule** = Filter or map input (e.g. this *specific input* converts to SQL Operator)

**Ruleset** = Collection of rules (e.g. these rules filter to something DAG-like)

Python Dictionary

@rules

```
def parse_shell_task(some_dict):  
    if some_condition :  
        id = some_value  
        bash = some_other_value
```



```
@dag_filter_rule  
def find_workflows(some_dict):  
    if some_property == some_value:  
        return ... something ...
```

```
@dag_rule  
def parse_workflow(some_dict):  
    id = some_value  
    schedule = some_other_value
```

```
@task_filter_rule  
def find_tasks(some_dict):  
    if some_property == some_value:  
        return ... something ...
```

```
@task_rule  
def parse_sql_task(some_dict):  
    if some_condition :  
        id = some_value  
        sql = some_other_value
```

```
@task_rule  
def parse_shell_task(some_dict):  
    if some_condition :  
        id = some_value  
        bash = some_other_value
```





translate.py

/input\_folder

/output\_folder



```
"workflowA": {  
  "some_property": ...  
  ...  
  "taskB": {  
    "other_property": ...  
    ...  
  }  
  "taskC": { ... }  
  "dependency": {  
    "taskB": "taskC"  
    ...  
  }  
}
```

Python Dictionary

translation\_fn

DAGFilterRuleset

DAGRuleset

TaskFilterRuleset

TaskRuleset

TaskDependency Ruleset

@rules

class DAG

class Operator  
e.g. OrbiterBashOperator

class Task  
Dependency

Objects

```
with DAG(dag_id="workflowA"):  
  task_b = ExecuteSqlOperator(  
    task_id="task_b",  
    sql="select * from ..."  
  )  
  task_c = BashOperator(  
    task_id="task_c",  
    bash_command="sh script.sh"  
  )  
  task_b >> task_c
```

/dags

workflow\_a.py

...

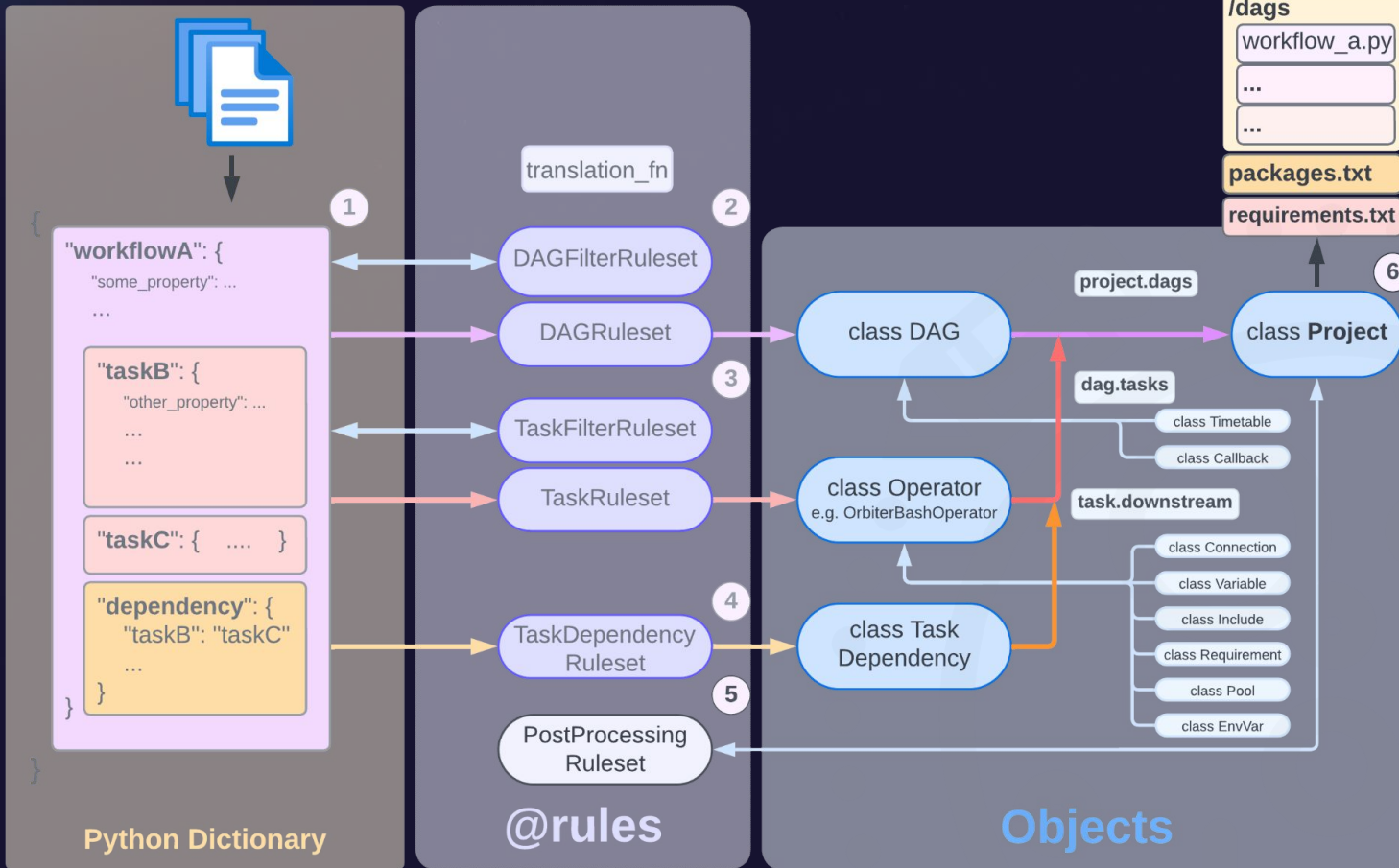
...



translate.py

/input\_folder

/output\_folder



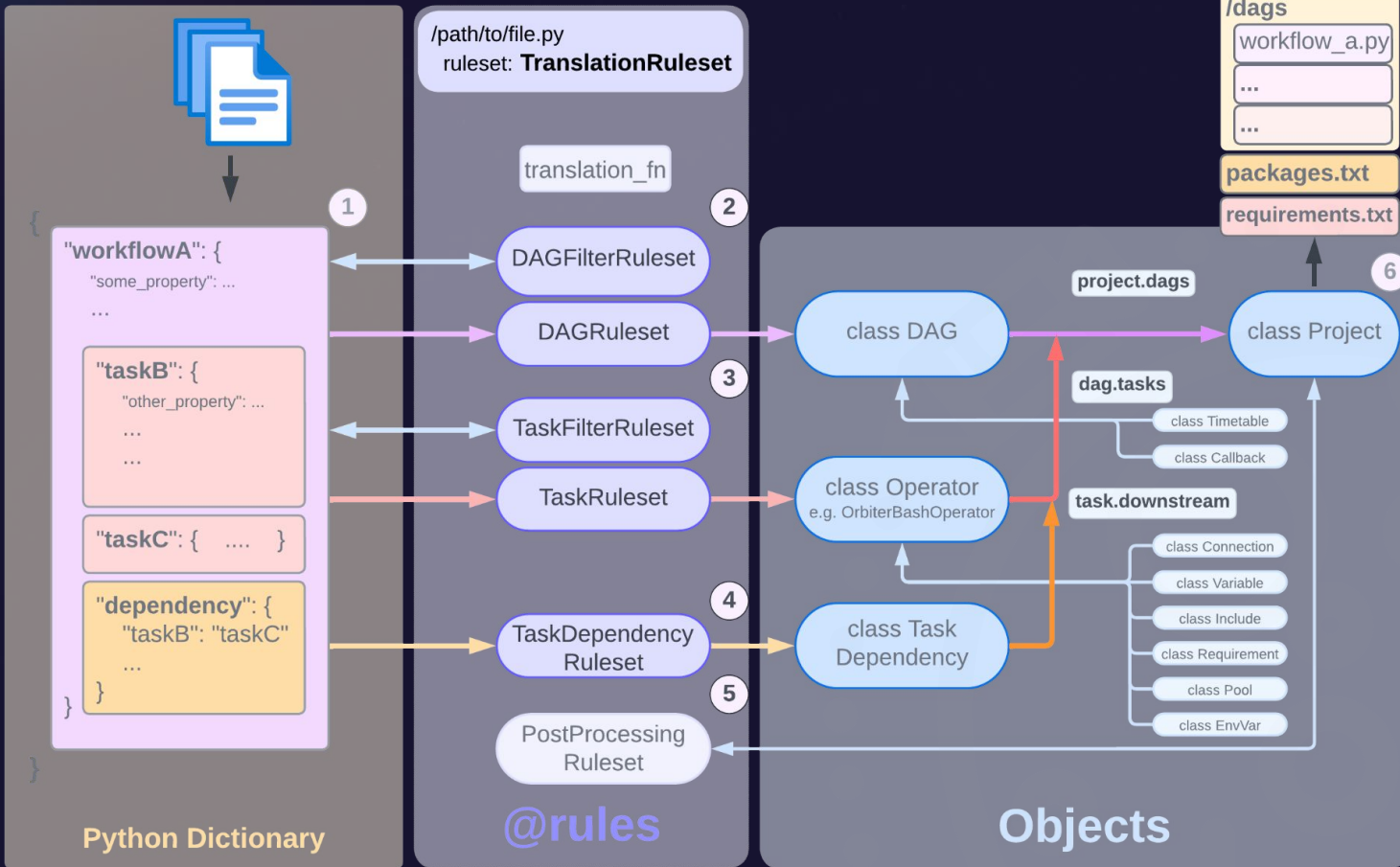


`translate.py`

`/input_folder`

`--ruleset path.to.file.ruleset`

`/output_folder`



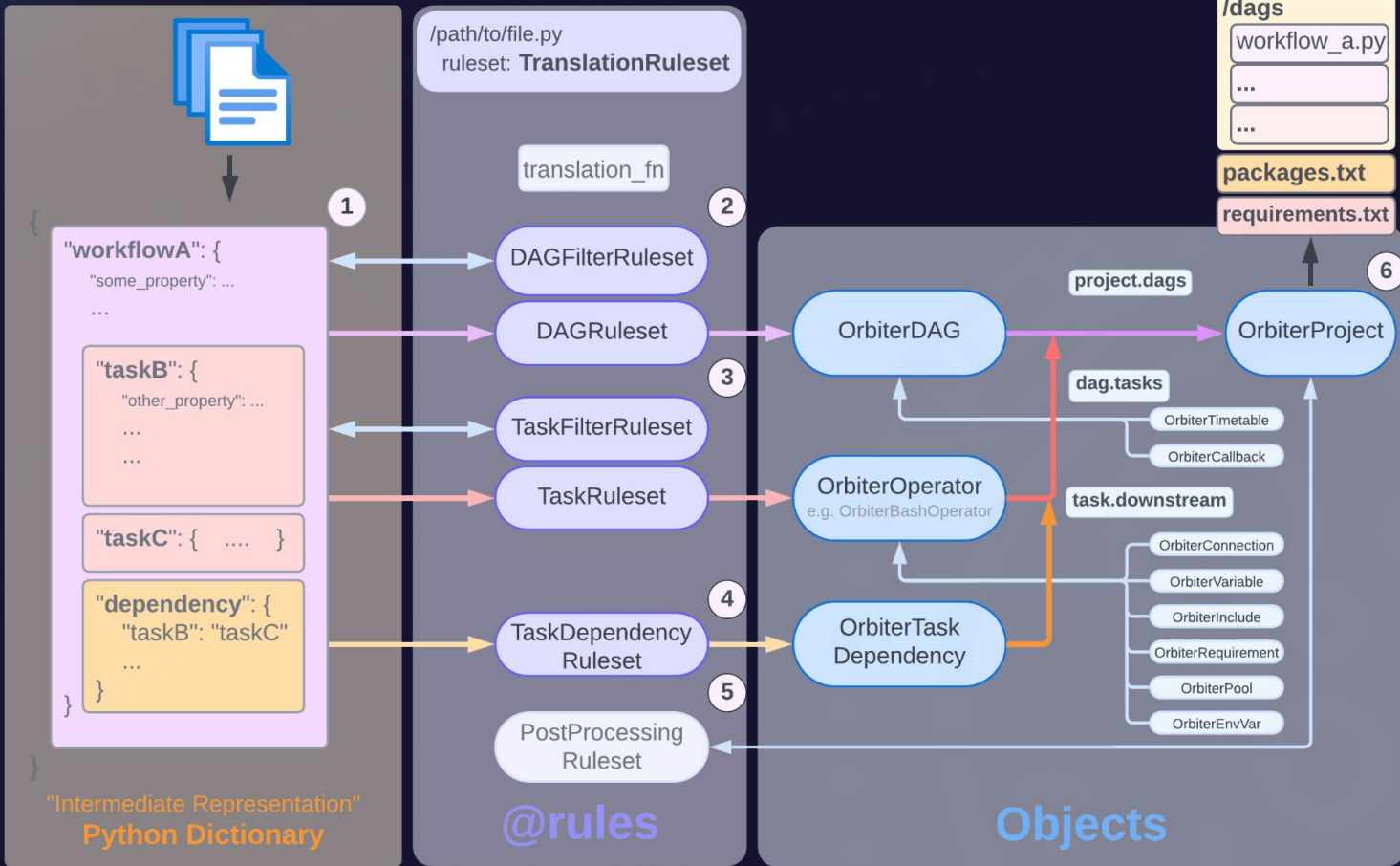


**orbiter translate**

/input\_folder

--ruleset path.to.file.ruleset

/output\_folder



"Intermediate Representation"  
Python Dictionary

@rules

Objects

# Orbiter

Land legacy workloads safely down in a new home on Apache Airflow!



Any system\*

Framework

Extendable & Open Source



Batteries included

```
1. pip install astronomer-orbiter
```

```
2. orbiter list-rulesets
```

```
3. orbiter install ...
```

```
4. orbiter translate
```

```
    <input dir> --ruleset <ruleset> <output dir>
```

CLI







## Orbiter

*Migrate <sup>(m)</sup>any legacy workloads to Airflow!*

<https://astronomer.github.io/orbiter/>

# Questions?

## Community Translations

*Contribute yours!*

<https://github.com/astronomer/orbiter-community-translations>



## Astronomer Translations & Assistance

*A Team of Airflow Experts to Guide Your Migration!*

<https://www.astronomer.io/professional-services/>







## Acknowledgements

["Choosing Apache Airflow over other Proprietary Tools..."](#)

[Parnab Basak, Airflow Summit 2022](#)

["Observations from migrating away from Control-M to Airflow"](#)

[Arjun Anandkumar](#)

[Astronomer Migration White Papers](#)

["DAGify - Enterprise Scheduler Migration Accelerator for Airflow"](#)

[Konrad Schieban](#)

& more