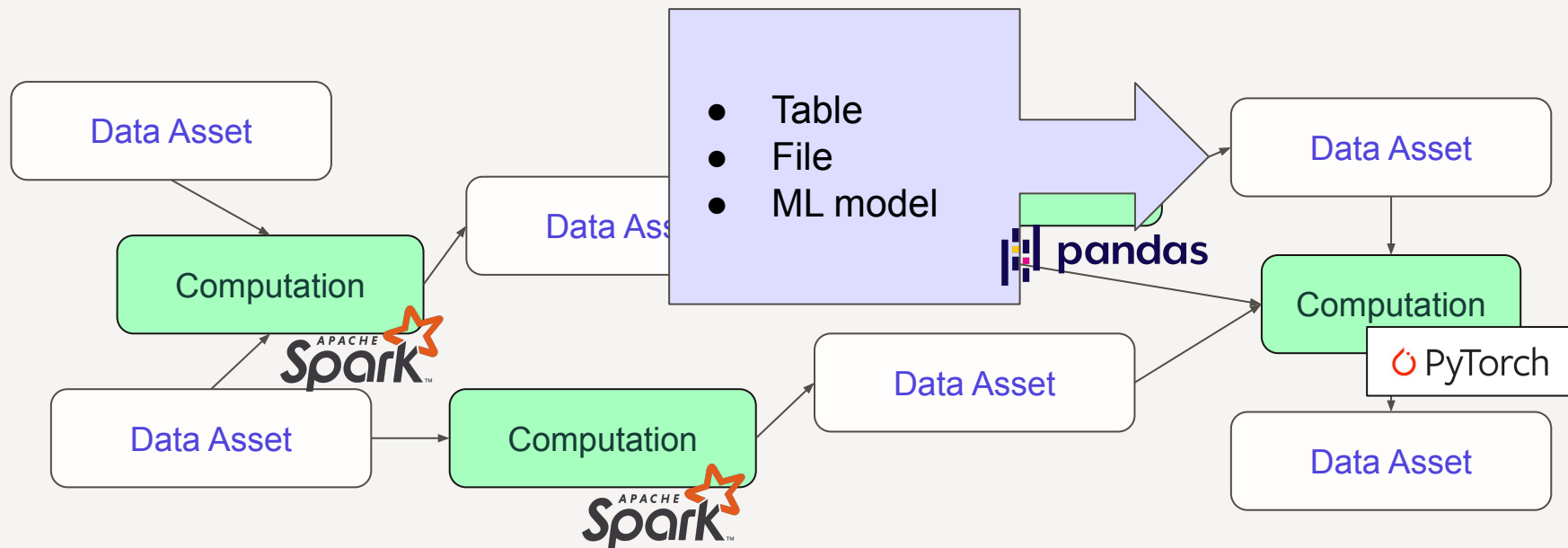# dagster

# Asset-Based Data Orchestration
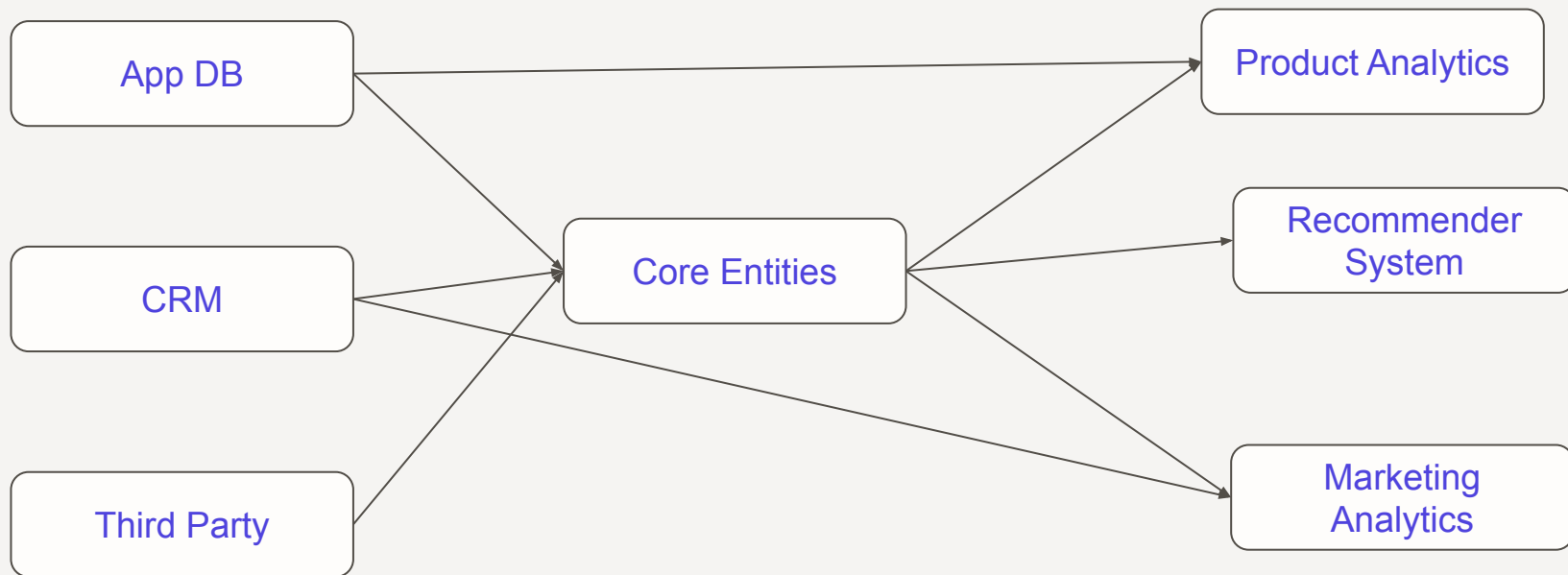
## Sandy Ryza (@s_ryz)
Lead Engineer, Dagster Project - Elementl

# Data practitioners build and maintain data pipelines

# What's a data pipeline?

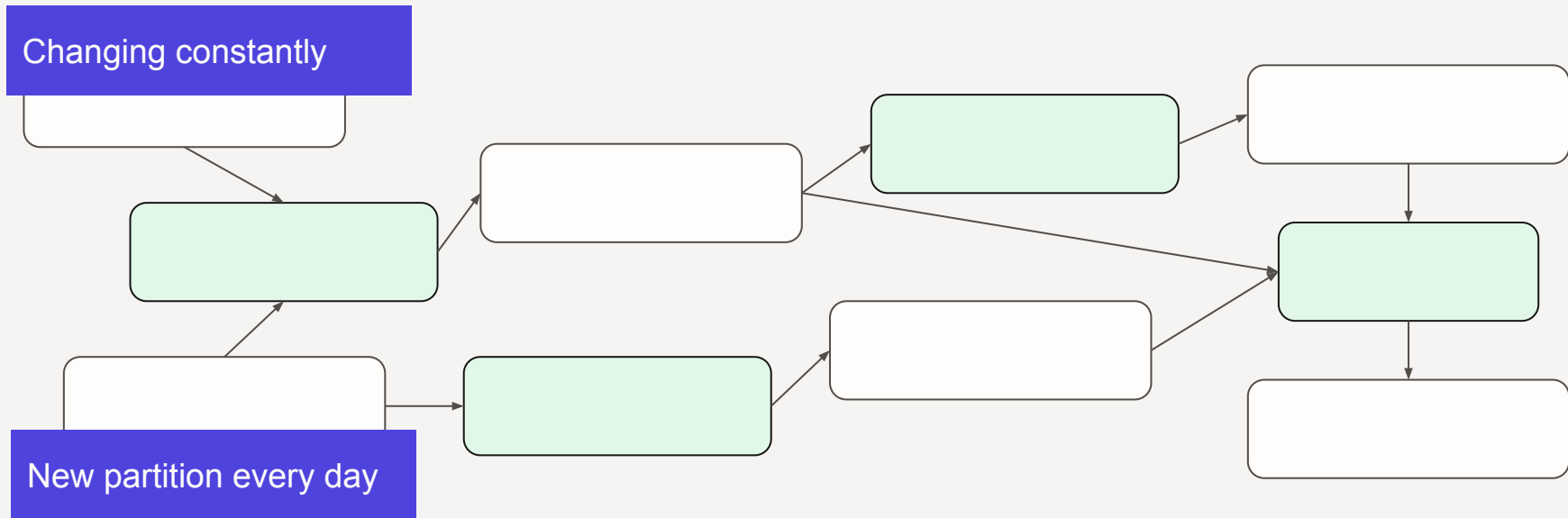# Data pipelines span entire organizations

# Automatically updating data assets
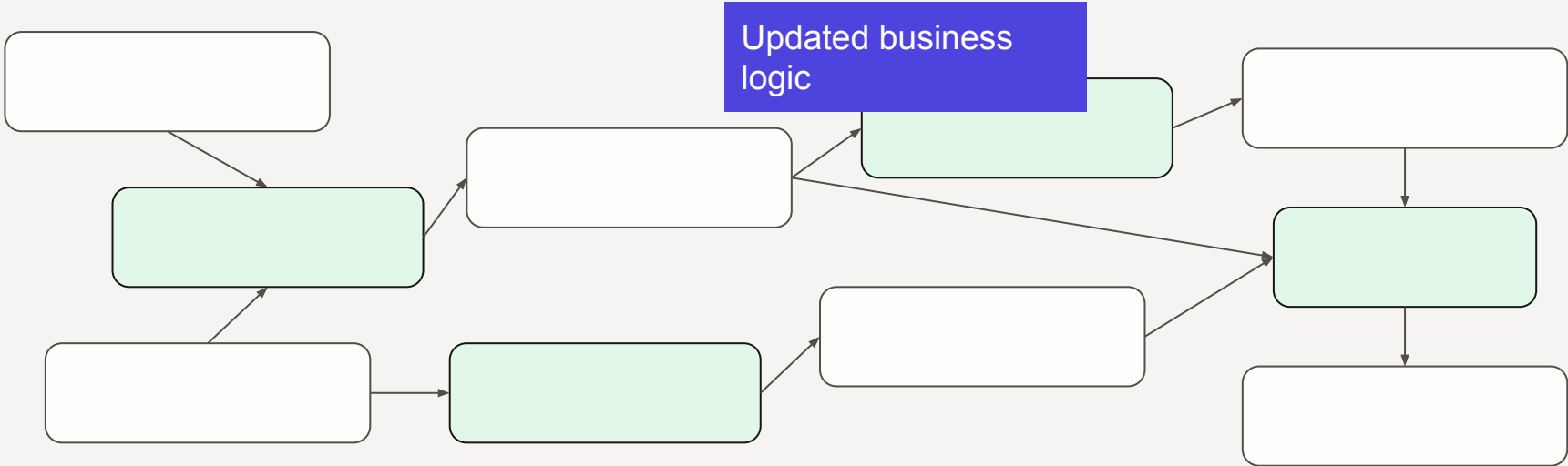
# Why update a data asset?

- Inputs have changed

# Changing inputs

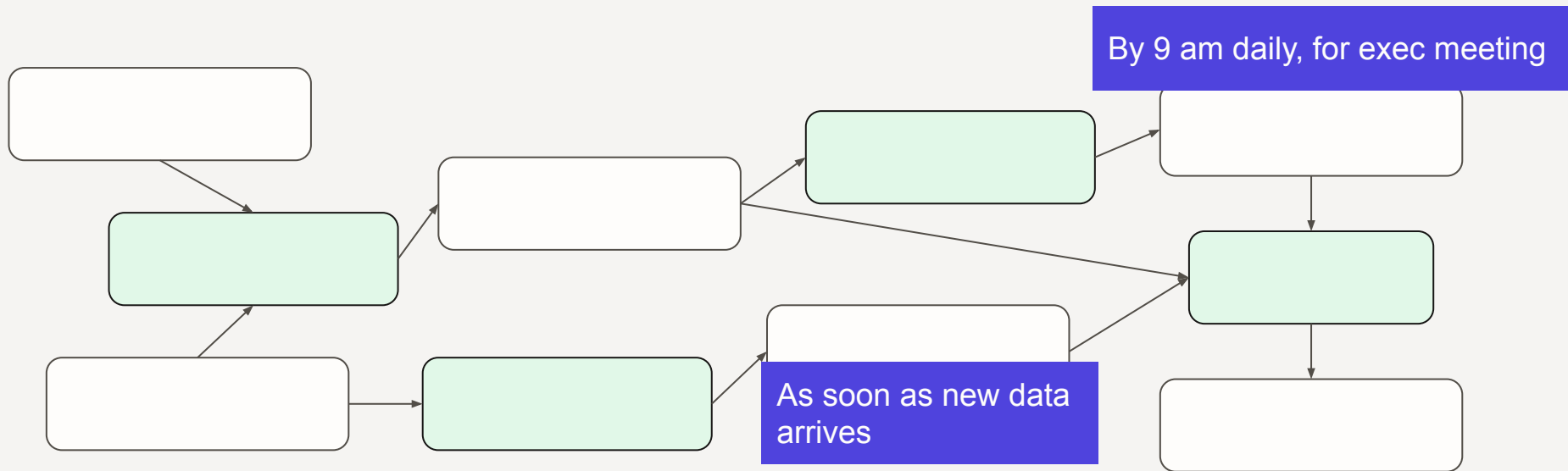# Why update a data asset?

- Inputs have changed
- Code has changed

# Code changes

Updated business logic

# Why update a data asset?

- Inputs have changed
- Code has changed
- Fresh data is needed

# Fresh data is needed

# Automatically updating data assets: how?

# The status quo: **workflow engines**

- DAG of tasks
- Run the DAG every hour/day/whatever

```python
from airflow.operators import BashOperator, DummyOperator
from airflow.models import DAG
from datetime import datetime

args = {
    'owner': 'airflow',
    'start_date': datetime(2015, 1, 1),
}

dag = DAG(dag_id='example1')

cmd = 'ls -l'
run_this_last = DummyOperator(
    task_id='run_this_last',
    default_args=args)
dag.add_task(run_this_last)

run_this = BashOperator(
    task_id='run_after_loop', bash_command='echo 1',
    default_args=args)
dag.add_task(run_this)
run_this.set_downstream(run_this_last)
for i in range(9):
    i = str(i)
    task = BashOperator(
```
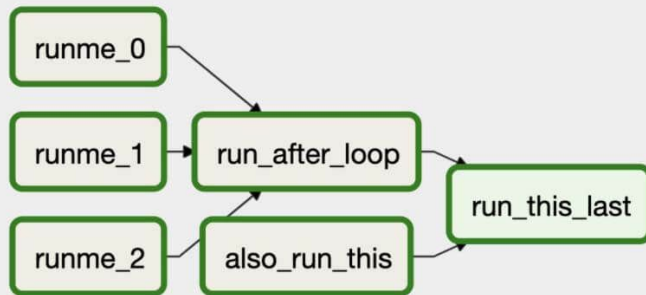
# Workflow engines: not actually the best way to schedule data pipelines?

- Forces running in lockstep
  - Caught between doing redundant work and stale data
- Code management
  - What DAG should this new data asset be a part of?
  - Monolithic DAG objects
- Alerts when tasks fail vs. when data is late

# A different way: Asset-based orchestration

# Goals of asset-based orchestration

- Outcomes
  - Make data ready on time
  - Avoid redundant work
- Express scheduling in terms of the data assets
  - When does source data change?
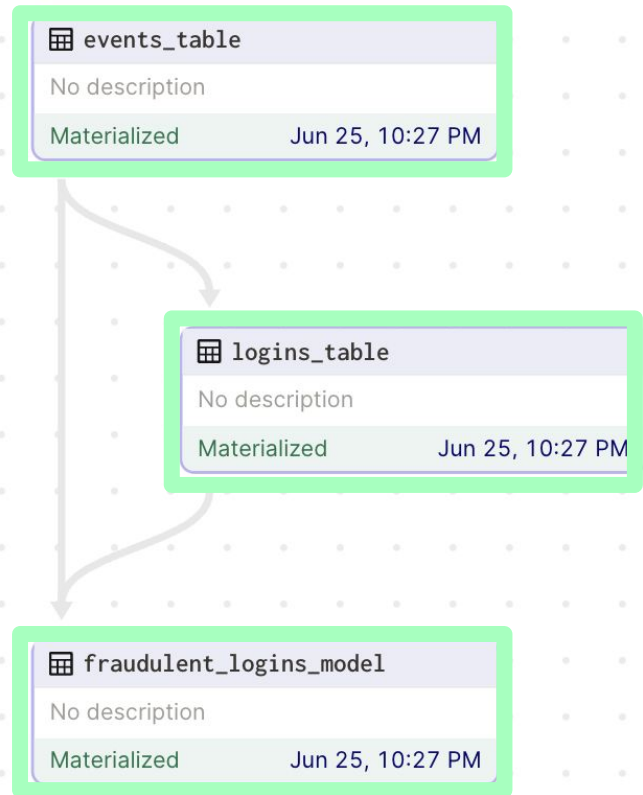  - How fresh do data assets need to be?
- Understand scheduling decisions

# Building a pipeline
aka defining some data assets

```python
from dagster import asset

@asset
def events_table():
    # code that writes the events table
    ...


@asset
def logins_table(events_table):
    # code that generates the logins table from the
    # events table
    ...


@asset
def fraudulent_logins_model(events_table, logins_table):
    # code that trains a bad login detection model
    ...
```

**default**   ⊞ Asset Group in basic-graph.py ↻                    ↻ Reload definitions

Lineage   List                                      ⊹ View global asset lineage

⊹ Type an asset subset... (ex: events_table+*)        0:04 ↻   ✦ Materialize all ▾

⊞ **events_table**
No description
Materialized          Jun 25, 10:27 PM

⊞ **logins_table**
No description
Materialized          Jun 25, 10:27 PM

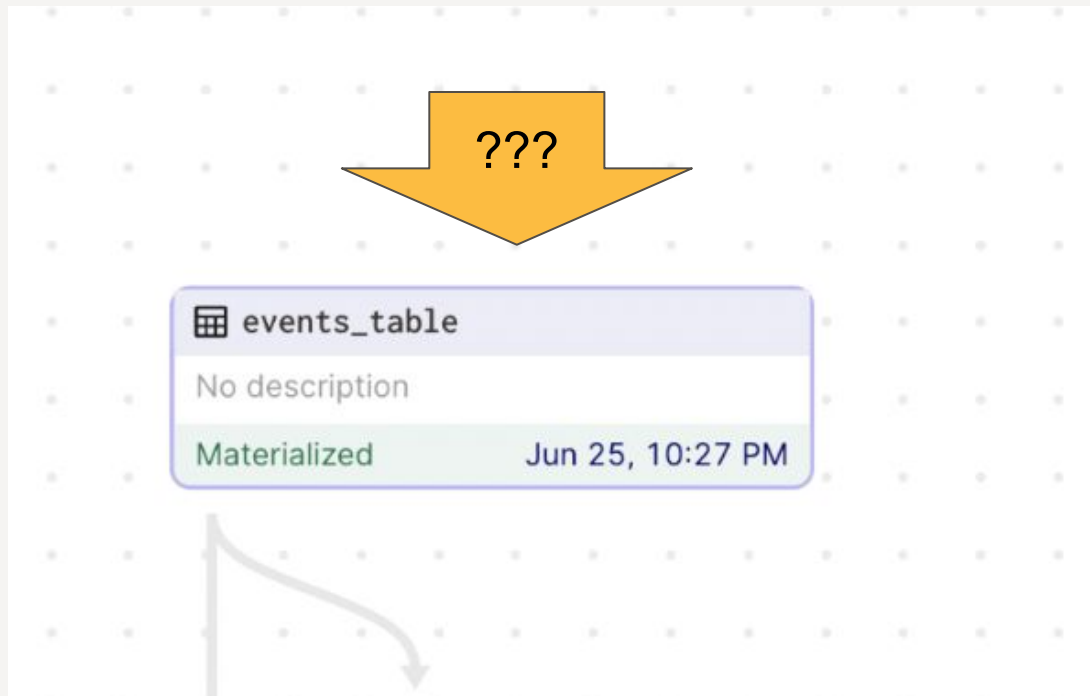⊞ **fraudulent_logins_model**
No description
Materialized          Jun 25, 10:27 PM

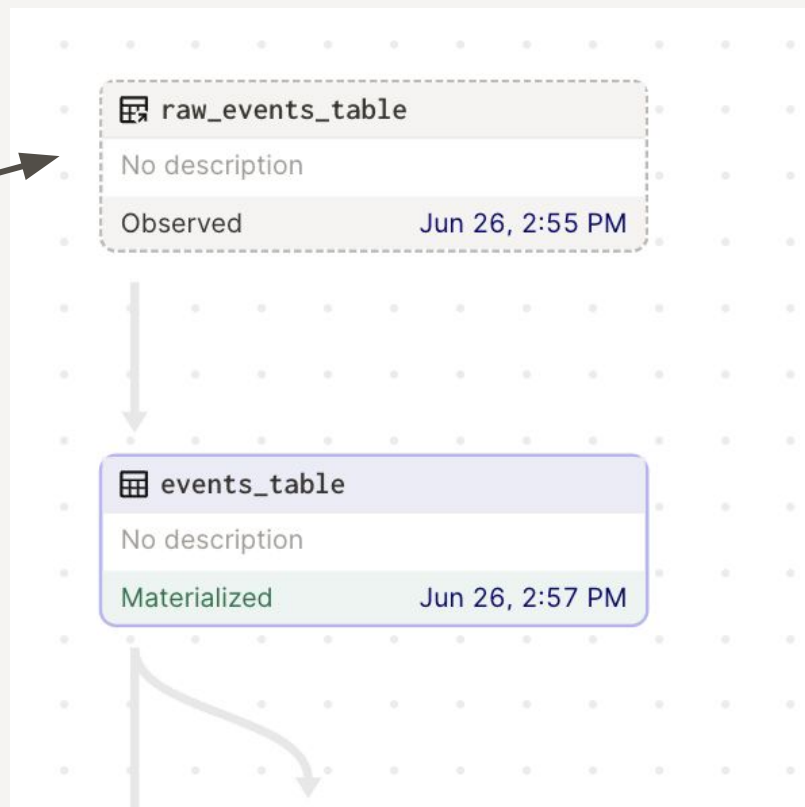# Asset-based orchestration in Dagster

# Auto-materialize policies

```python
@asset(auto_materialize_policy=AutoMaterializePolicy.eager())
def logins_table(events_table):
    ...
```
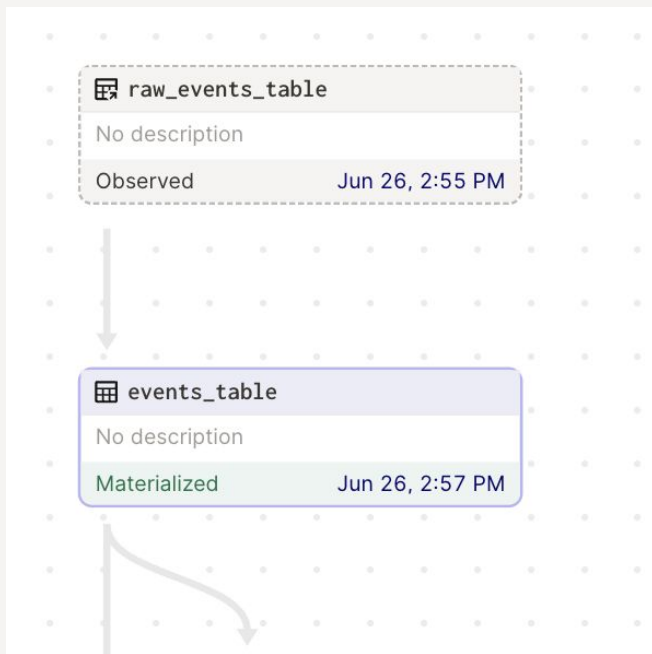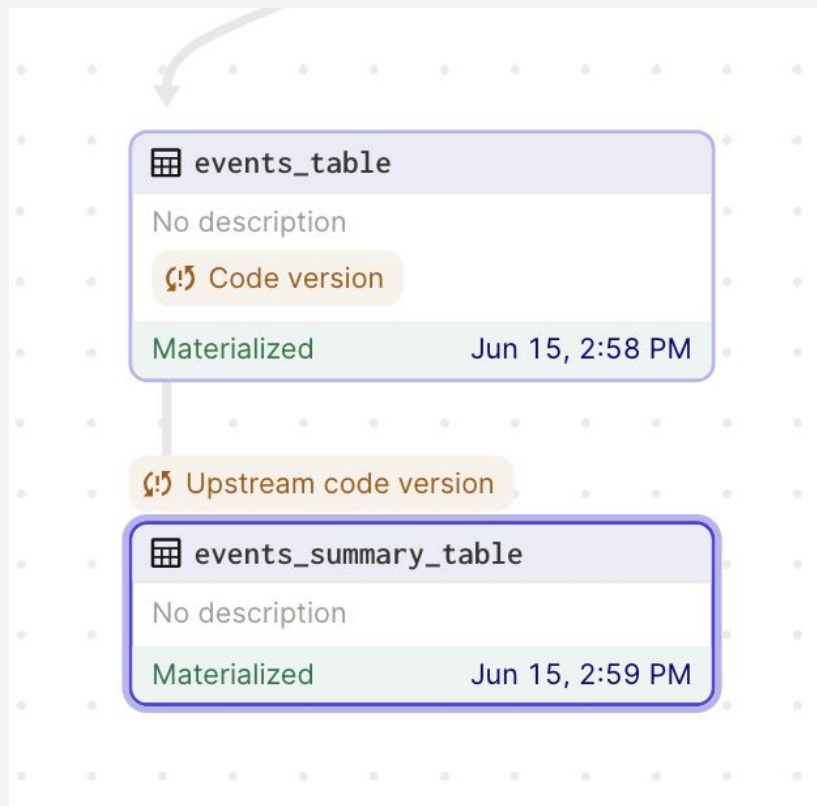
# The root of the graph

# Source assets



raw_events_table

No description

Observed          Jun 26, 2:55 PM

source asset

events_table

No description

Materialized          Jun 26, 2:57 PM

```python
@observable_source_asset
def raw_events_table():
    path = "raw_data_bucket/raw_events.parquet"
    last_modified_datetime = get_last_modified_datetime(path)
    return DataVersion(str(last_modified_datetime))
```

raw_events_table

No description

Observed                    Jun 26, 2:55 PM


events_table

No description

Materialized                Jun 26, 2:57 PM

# What about code changes?

# Lazy auto-materialization





Upstream asset → Downstream asset

# Freshness policies

```
@asset(
    freshness_policy=FreshnessPolicy(
        cron_schedule="@daily",
        maximum_lag_minutes=24 * 60,
    )
)
def fraudulent_logins_model(events_table, logins_table):
    # code that trains a bad login detection model
    ...
```
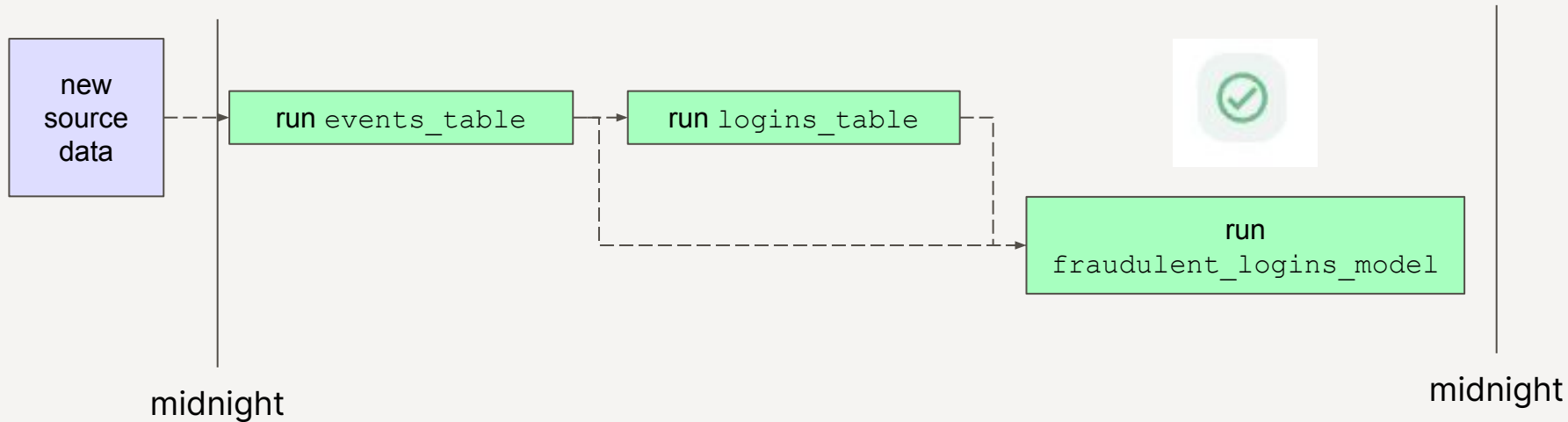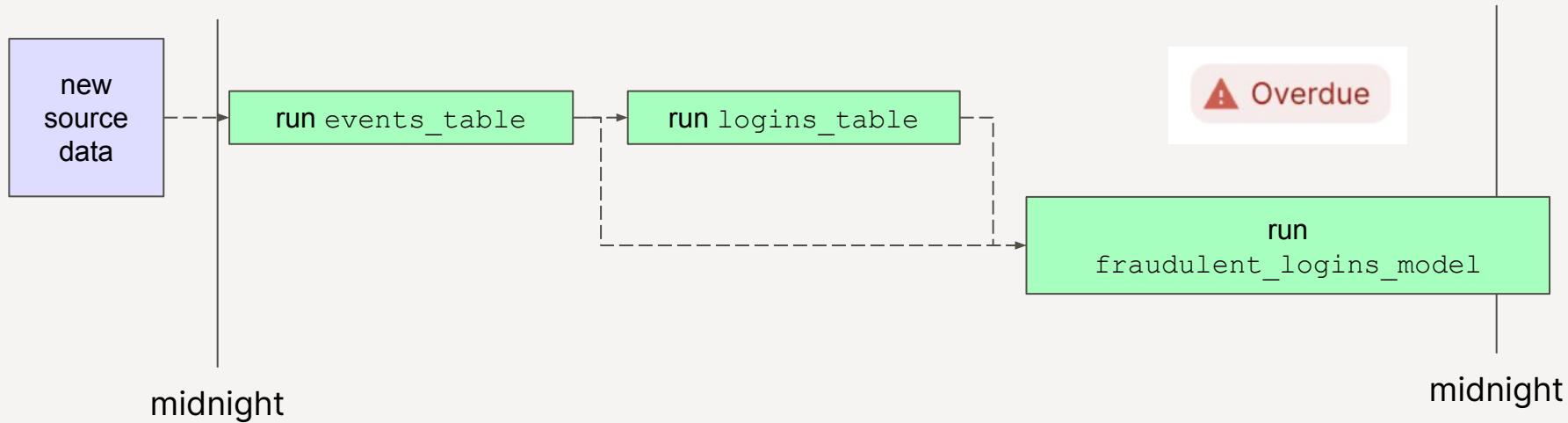
fraudulent_logins_model
View in Asset Catalog ↗

**Freshness policy**                                                    ▼

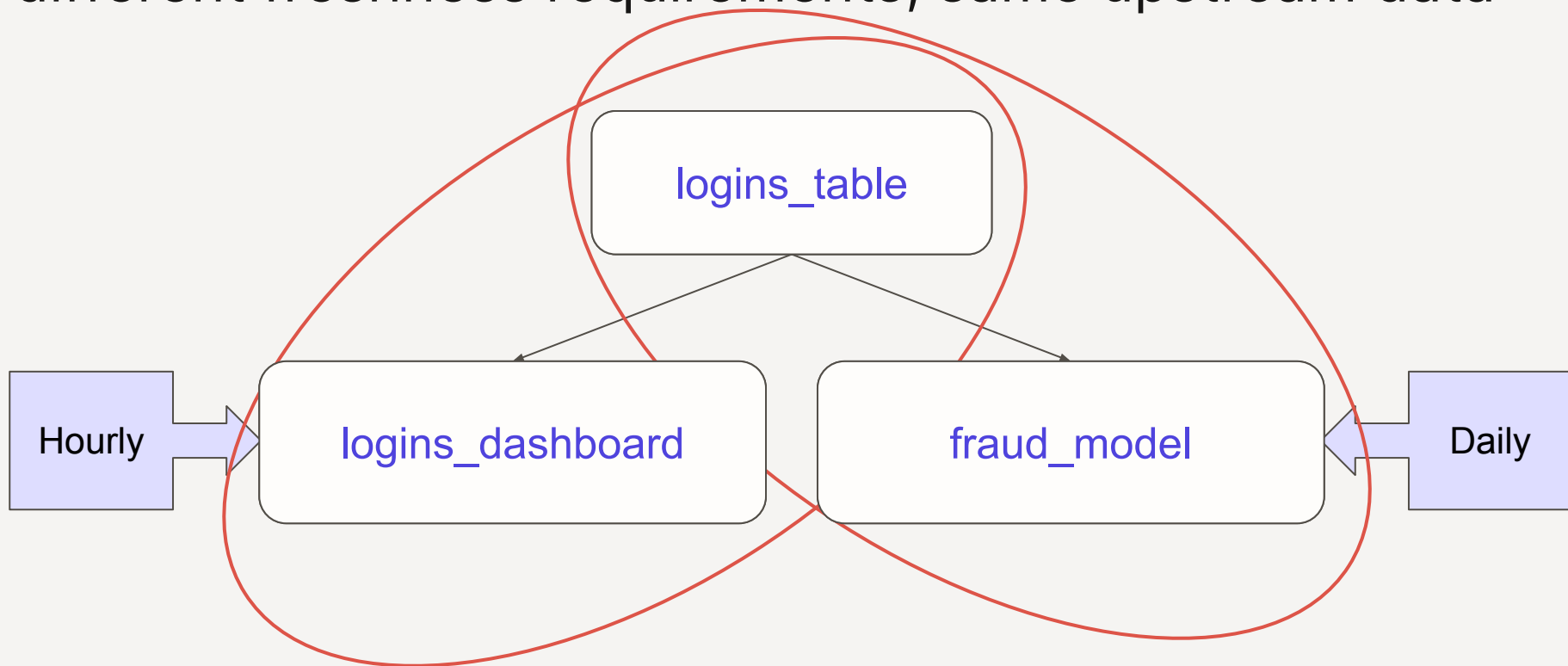By 12:00 AM UTC, this asset should incorporate    ⚠ Overdue
all data up to 24 hours before that time.

```python
@asset(auto_materialize_policy=AutoMaterializePolicy.lazy())
def logins_table(events_table):
    # code that generates the logins table from the
    # events table
    ...


@asset(
    freshness_policy=FreshnessPolicy(
        cron_schedule="@daily",
        maximum_lag_minutes=24 * 60,
    ),
    auto_materialize_policy=AutoMaterializePolicy.lazy(),
)
def fraudulent_logins_model(events_table, logins_table):
    # code that trains a bad login detection model
    ...
```
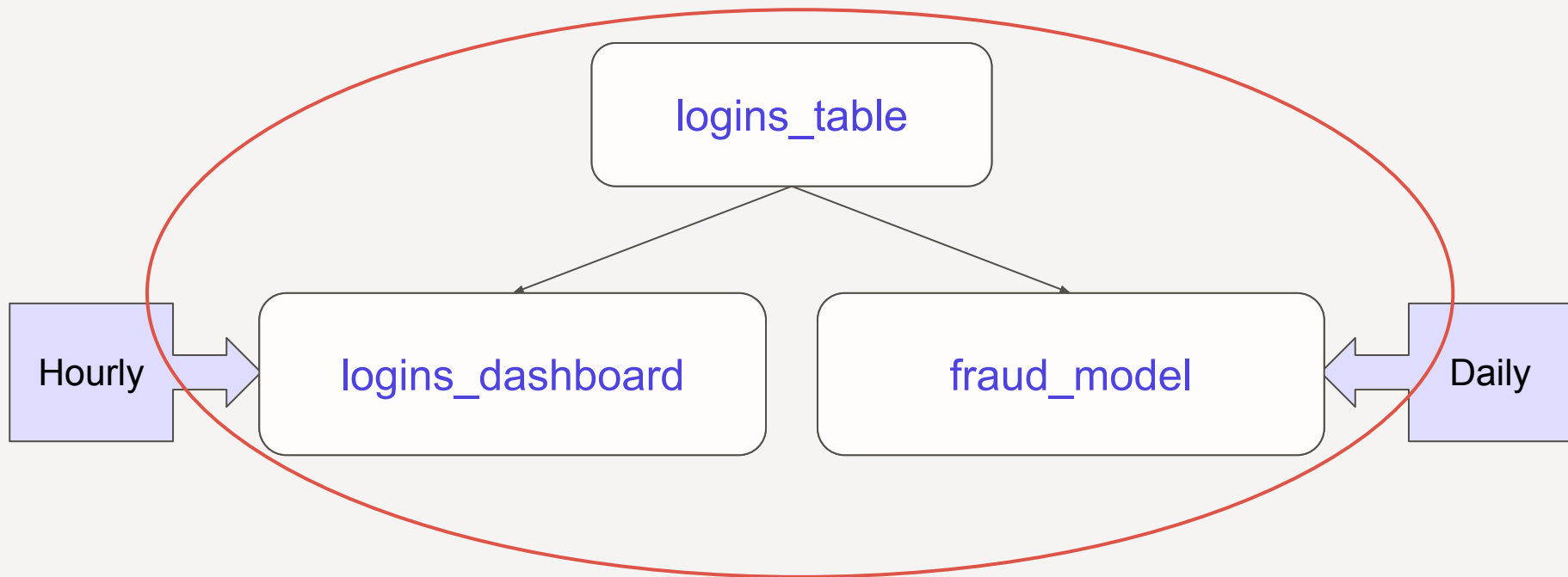
# Common scenario:
different freshness requirements, same upstream data

# Common scenario:
different freshness requirements, same upstream data

# Common scenario:
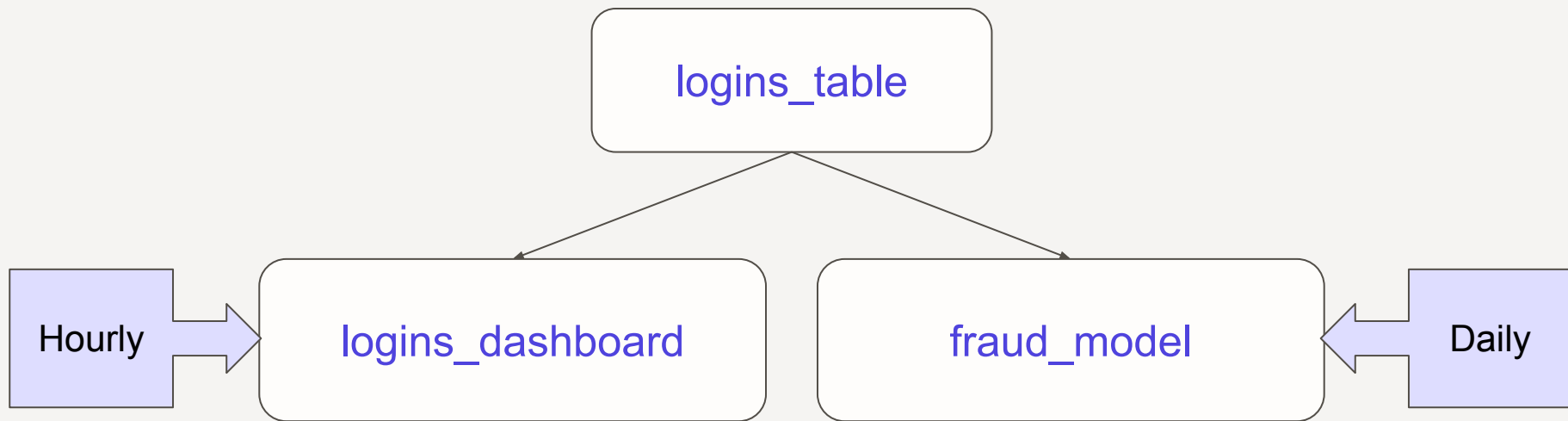different freshness requirements, same upstream data

**Common scenario:**
different freshness requirements, same upstream data

# Common scenario:
different freshness requirements, same upstream data

# Asset-based orchestration: observability

# To sum up…

- Data pipeline = graph of data assets connected by computations
- Workflows are not an adequate scheduling abstraction
- Asset-based orchestration
  - Express intentions more clearly
  - Avoid redundant computations
  - Debug scheduling decisions

# Thank you

Sandy Ryza
@s_ryz

elementl

dagster