



Orchestrating Software-Defined Assets

Sandy Ryza (@s_ryz)

Lead Engineer, Dagster Project - Elementl

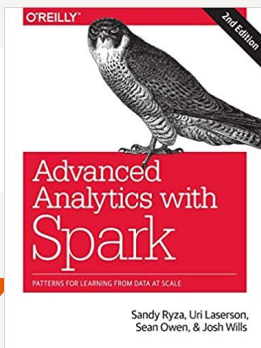
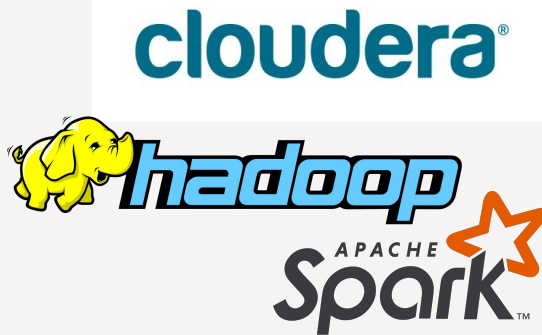


- Sandy

Building tools for
data people

Being a data
person

Building tools for
data people



**Clover
Health**



KEEP TRUCKIN

Frontend



Cluster
Orchestration



Dev Ops



Imperative → Declarative

Imperative

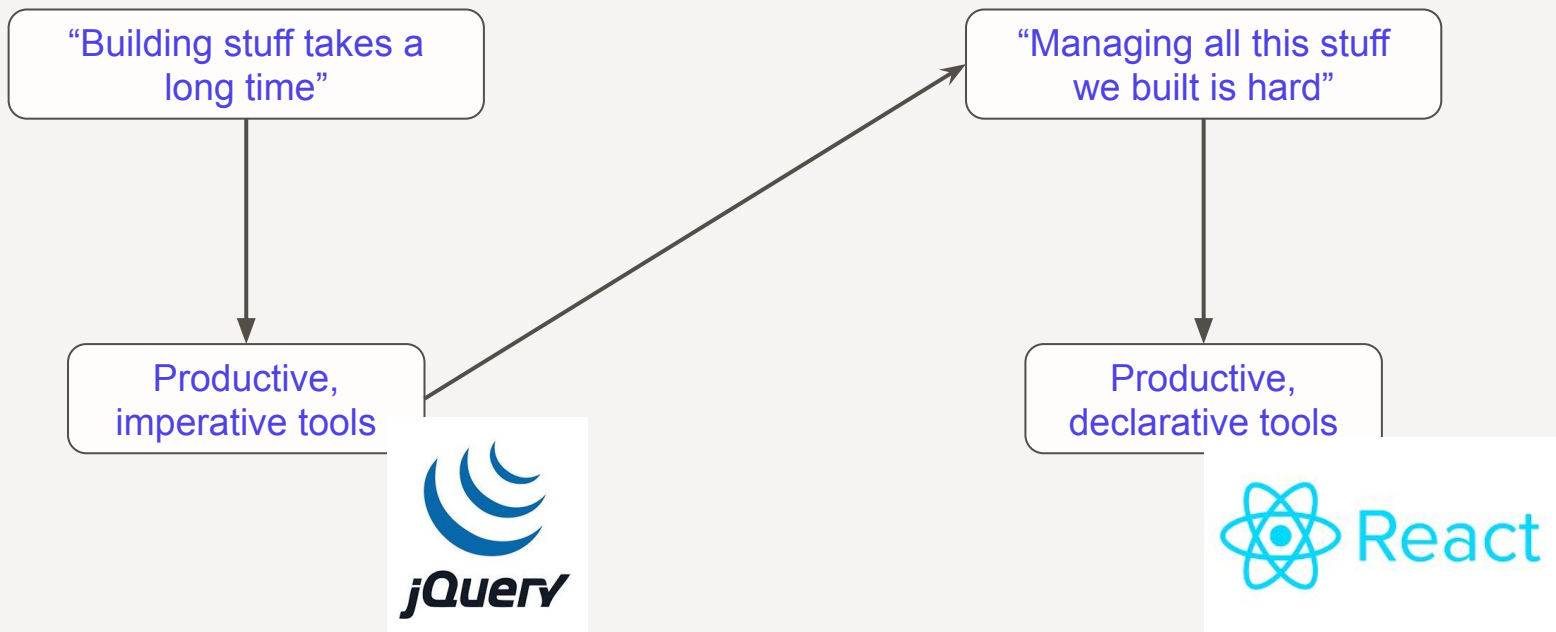
You give commands telling your system what to do.

Declarative

You describe the end state you want your system to be in.

What is going
on here?

The software domain trajectory



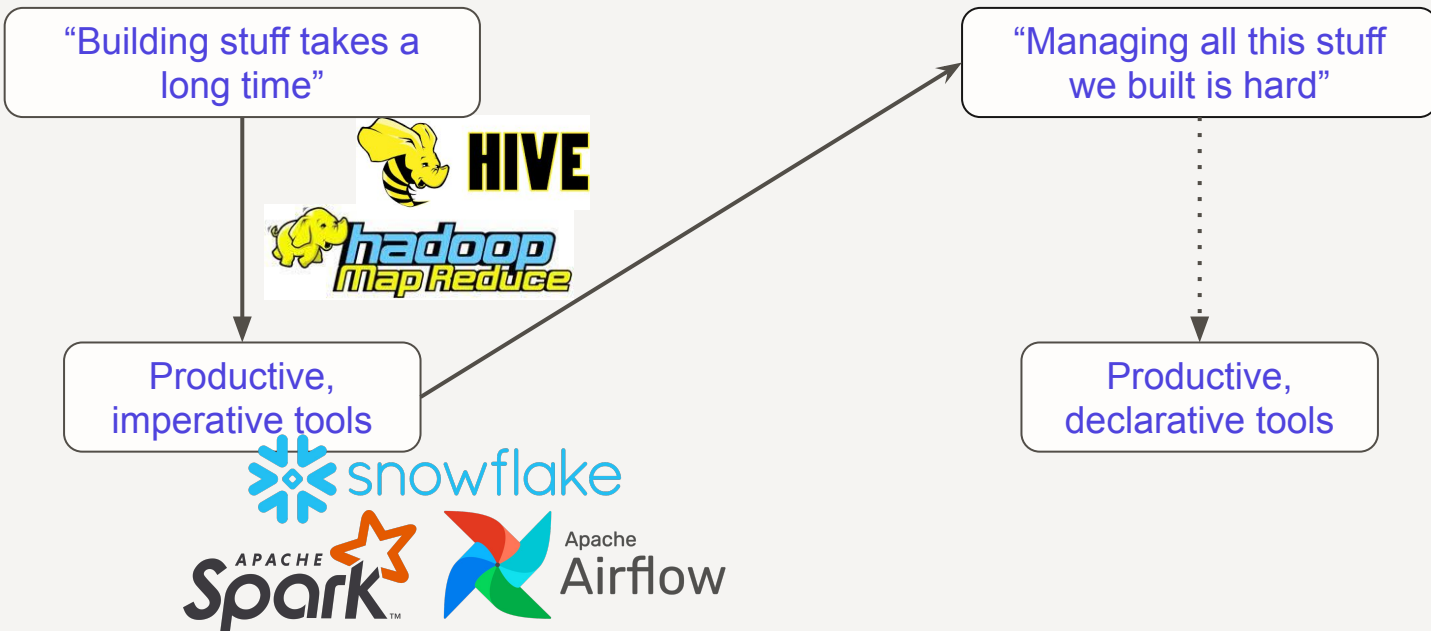
Declarative tools stop complexity from turning into chaos

- Make intentions explicit
- Offer a principled way of managing change

DATA



The data trajectory

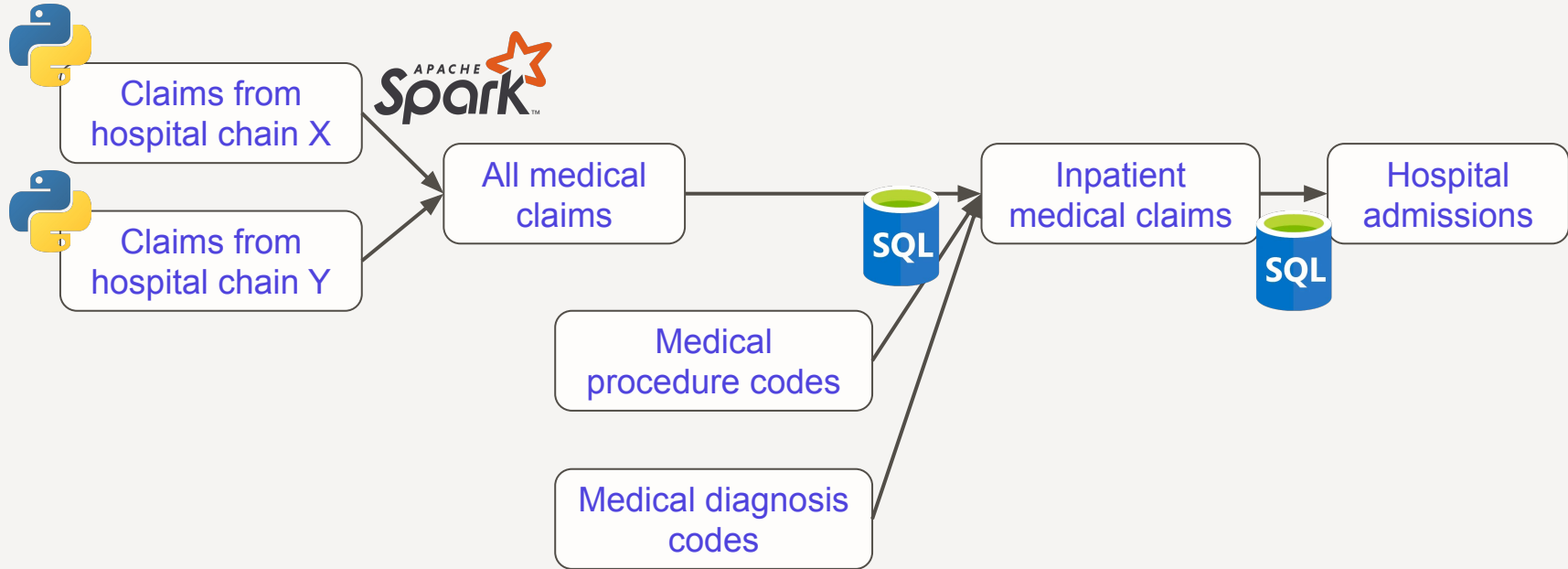


Big complexity

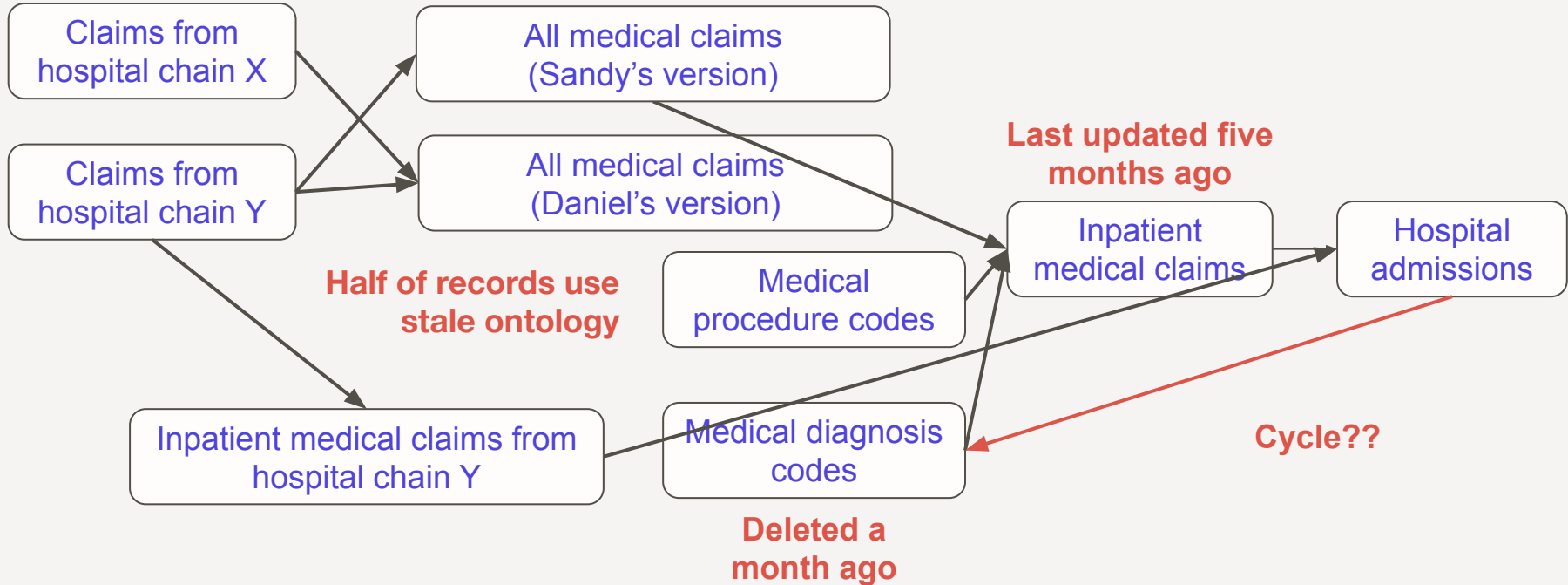
- 100s to 1,000s of tables, ML models, & datasets
- 10,000s of lines of code
- Multiple compute frameworks & storage systems



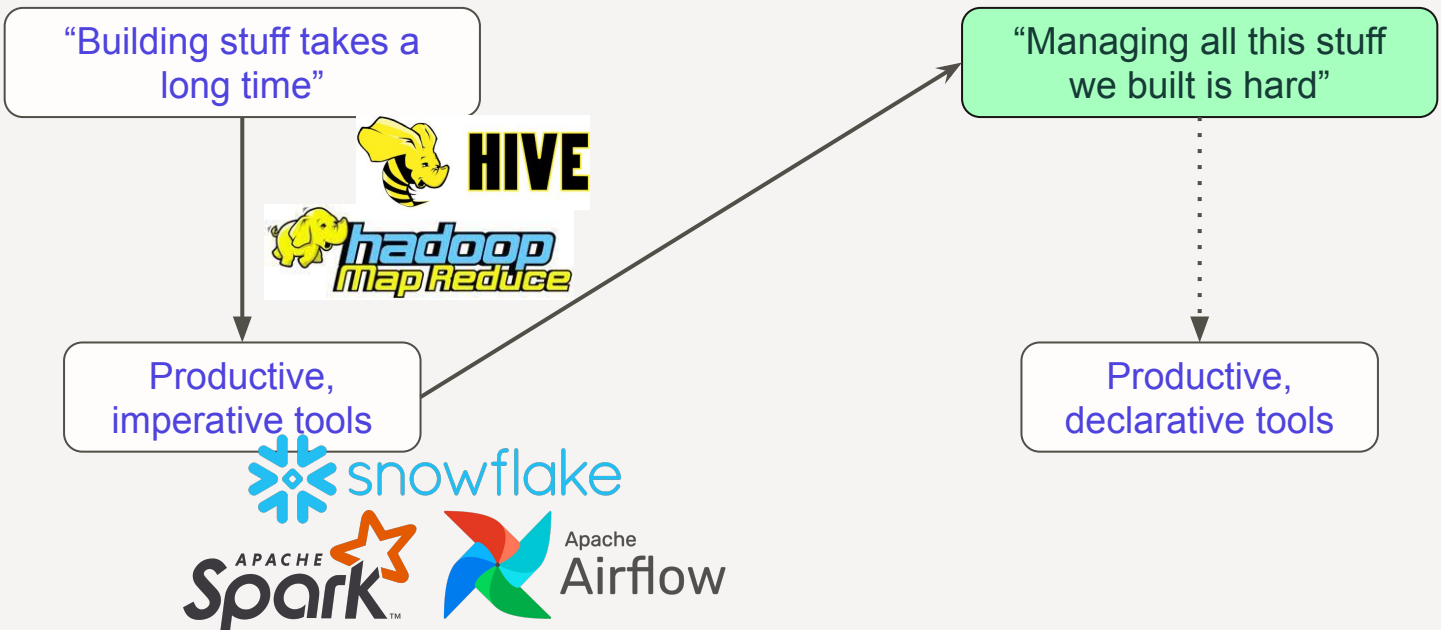
Example: how many hospital admissions?



Complexity » Chaos



The data trajectory

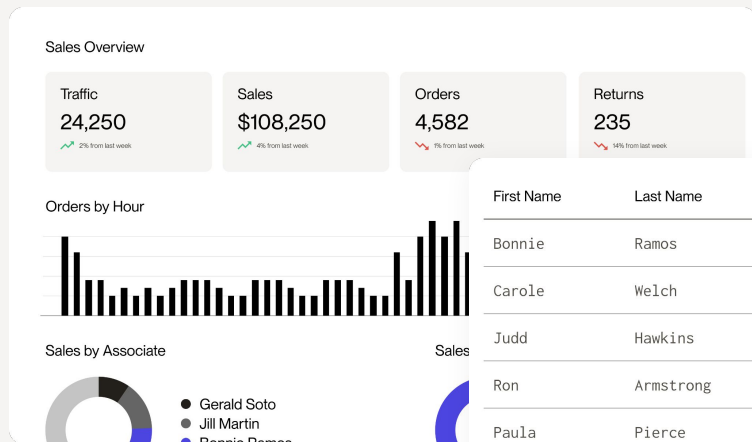


Declarative data management

What's the declarative entity?



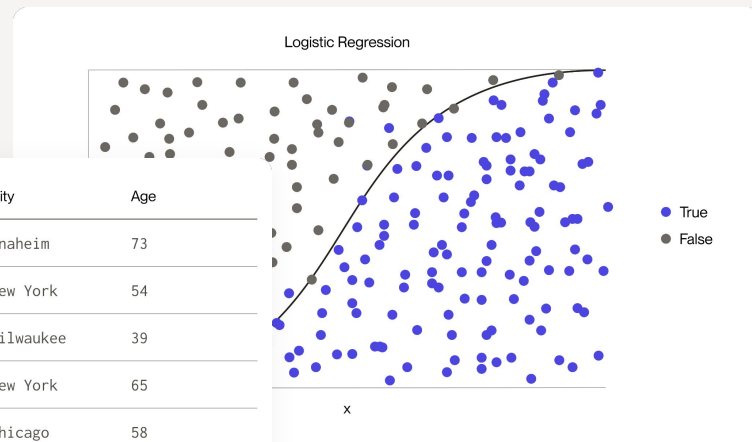
The asset!



Report

First Name	Last Name	Address	City	Age
Bonnie	Ramos	5530 E Pecan St	Anaheim	73
Carole	Welch	1931 McClellan Rd	New York	54
Judd	Hawkins	9532 Preston Rd	Milwaukee	39
Ron	Armstrong	6606 Cherry St	New York	65
Paula	Pierce	7674 Bruce St	Chicago	58
Bradley	Wallace	4896 W 6th St	Austin	61
Cat	Fowler	63 Pecan Acres Ln	New York	32
Courtney	Austin	874 Oak Lawn Ave	Chicago	28

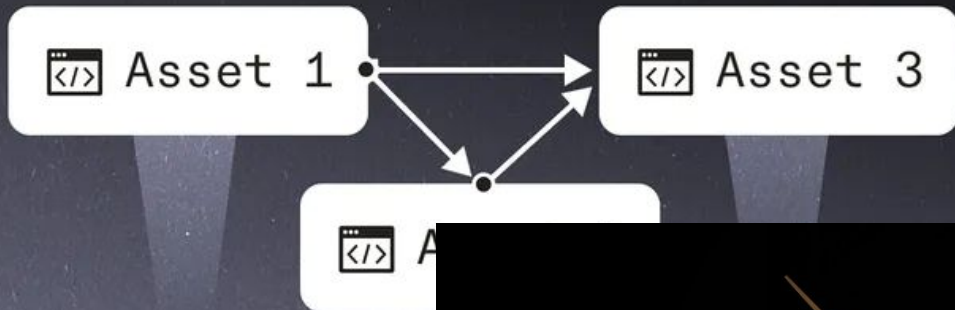
Table



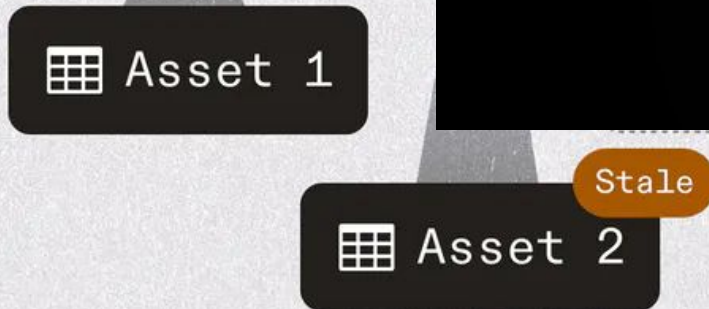
ML Model

Introducing the software-defined asset

Asset Definitions



Asset Materializations

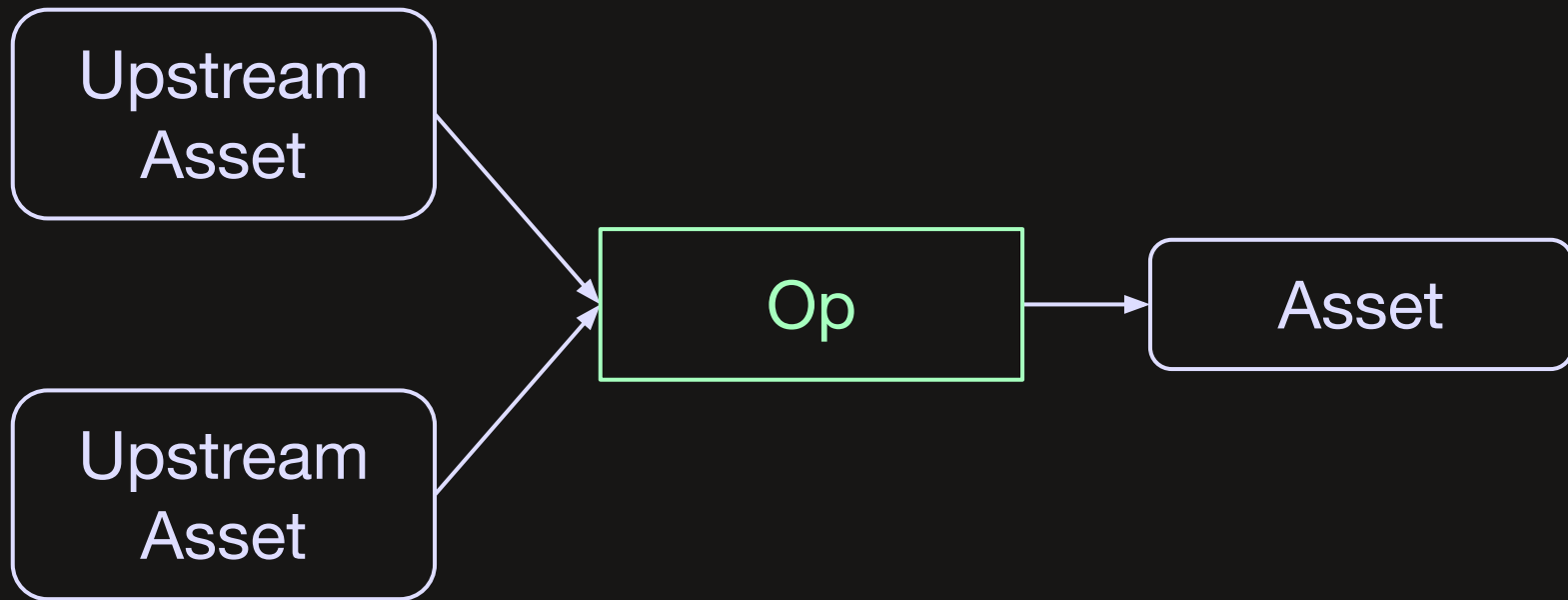




dagster

Software-defined asset

Specifies an asset that you intend to exist, and how to compute it.



Software-Defined Asset

Upstream
Asset Key

...

Upstream
Asset Key



Asset Key

Defining an asset in Python

asset key

upstream
asset key

```
from dagster import asset
```

```
@asset
```

```
def activity_forecast(activity_daily_stats: DataFrame) -> DataFrame:
```

```
    """Forecast of activity for the next 30 days"""
```

```
    start_date = activity_daily_stats.date.max()
```

```
    future_dates = date_range(start=start_date, end=start_date + DateOffset(days=30))
```

```
    predicted_data = 0.5 * np.exp(7 * (future_dates.astype(np.int64) / 10**18 - 1.6095))
```

```
    return DataFrame({"date": future_dates, "num_comments": predicted_data})
```

compute function

Python isn't the only way to
define an asset...

Defining an asset in SQL (with dbt)

asset key

upstream asset keys

activity_daily_stats.sql

```
select *  
from {{ ref('comment_daily_stats') }}  
full outer join {{ ref('story_daily_stats') }}  
using (date)
```

compute function

activity_daily_stats.sql

```
select *  
from {{ ref('comment_daily_stats') }}  
full outer join {{ ref('story_daily_stats') }}  
using (date)
```


@asset

```
def activity_forecast(activity_daily_stats: DataFrame) -> DataFrame:  
    """Forecast of activity for the next 30 days"""  
    start_date = activity_daily_stats.date.max()  
    future_dates = date_range(start=start_date, end=start_date + DateOffset(days=30))  
    predicted_data = 0.5 * np.exp(7 * (future_dates.astype(np.int64) / 10**18 - 1.6095))  
    return DataFrame({"date": future_dates, "num_comments": predicted_data})
```


```
from dagster_dbt.asset_defs import load_assets_from_dbt_manifest
```

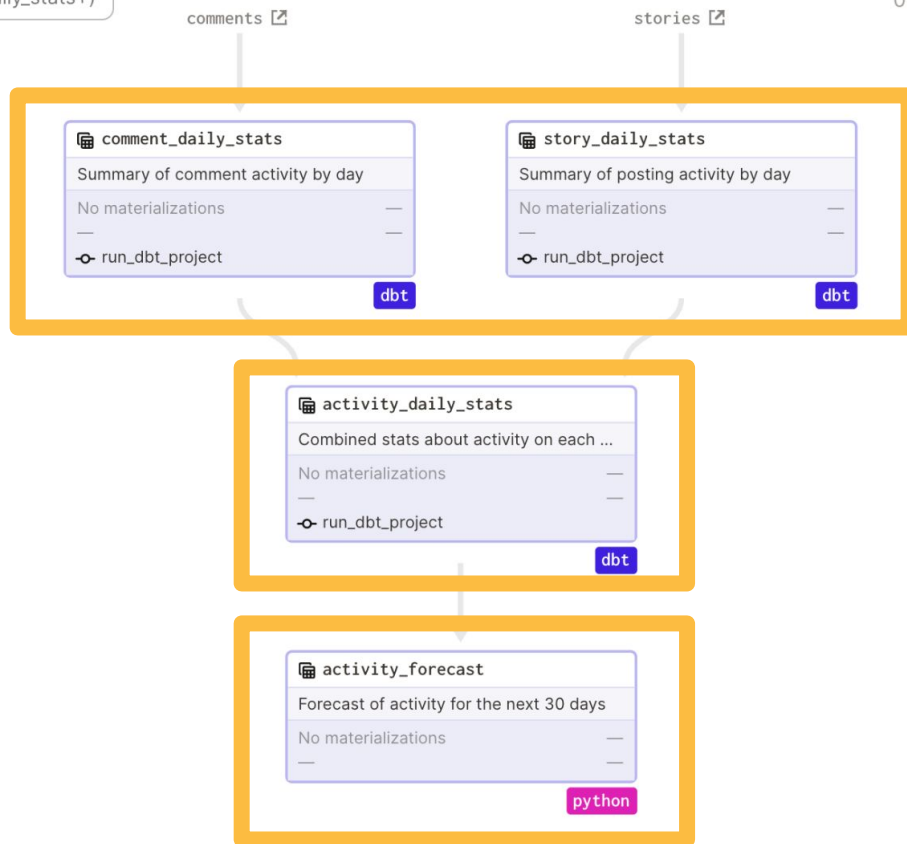
```
dbt_assets = load_assets_from_dbt_manifest(  
    json.load(open(os.path.join(DBT_PROJECT_DIR, "target", "manifest.json"))),  
    io_manager_key="warehouse_io_manager",  
)
```

```
@repository  
def activity_analytics():  
    return [  
        AssetGroup.from_package_name(  
            "hacker_news_assets.assets.activity_analytics", resource_defs=RESOURCES_LOCAL  
        )  
    ]
```

 Type an asset subset... (ex: activity_daily_stats+)


0:13 

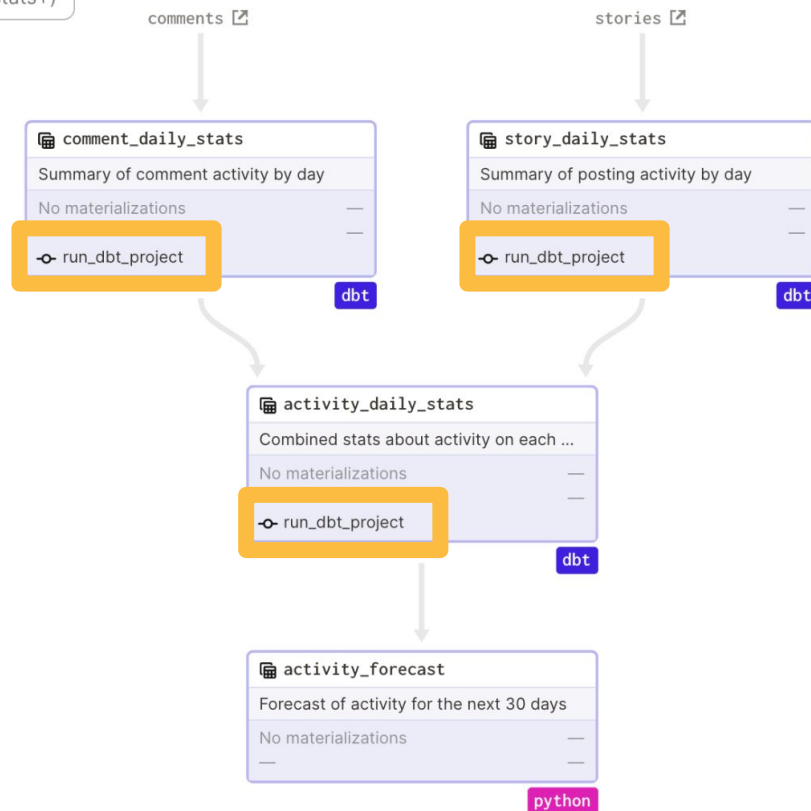
 Materialize All



 Type an asset subset... (ex: activity_daily_stats+)

0:13 

 Materialize All



Decentralized
dependencies →
more scalable graph

 Type an asset subset... (ex: activity_daily_stats+)

0:13 

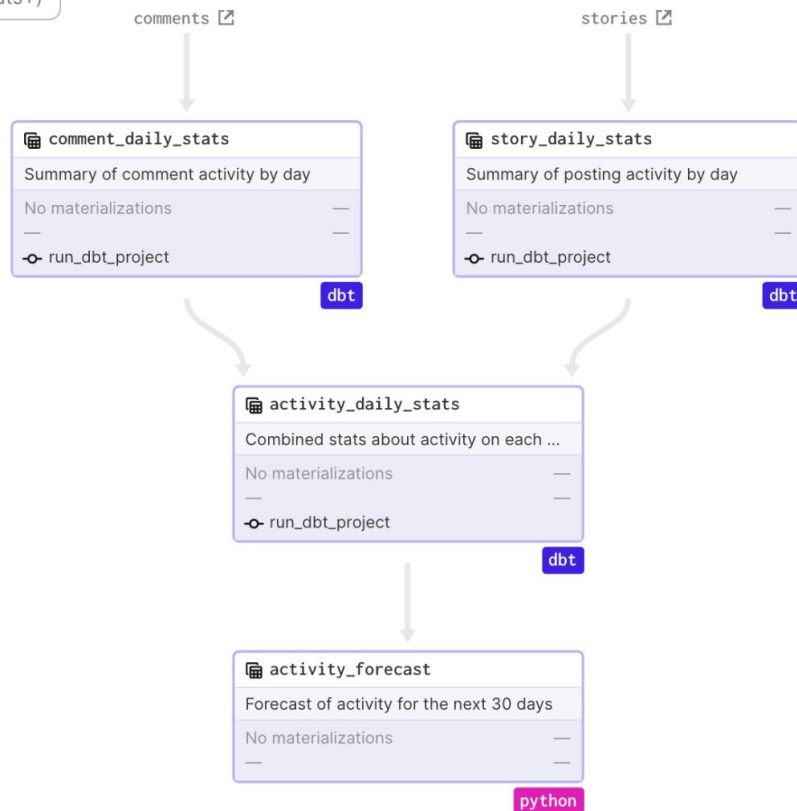
 Materialize All



 Type an asset subset... (ex: activity_daily_stats+)

0:13 

 Materialize All



b46f63d2 🔔 Success 🔄 Run of activity_analytics_assets @ d1234f7b[🔗 View tags and config](#)[🐞 Debug file](#)☰ 🔍 🔄 🕒 ☐ Hide not started steps🔍 ☰ 🏠 Re-execute All (*) ▼

Not executed (0) ▼

No steps are waiting to execute

Executing (0) ▼

No steps are executing

Errored (0) ▼

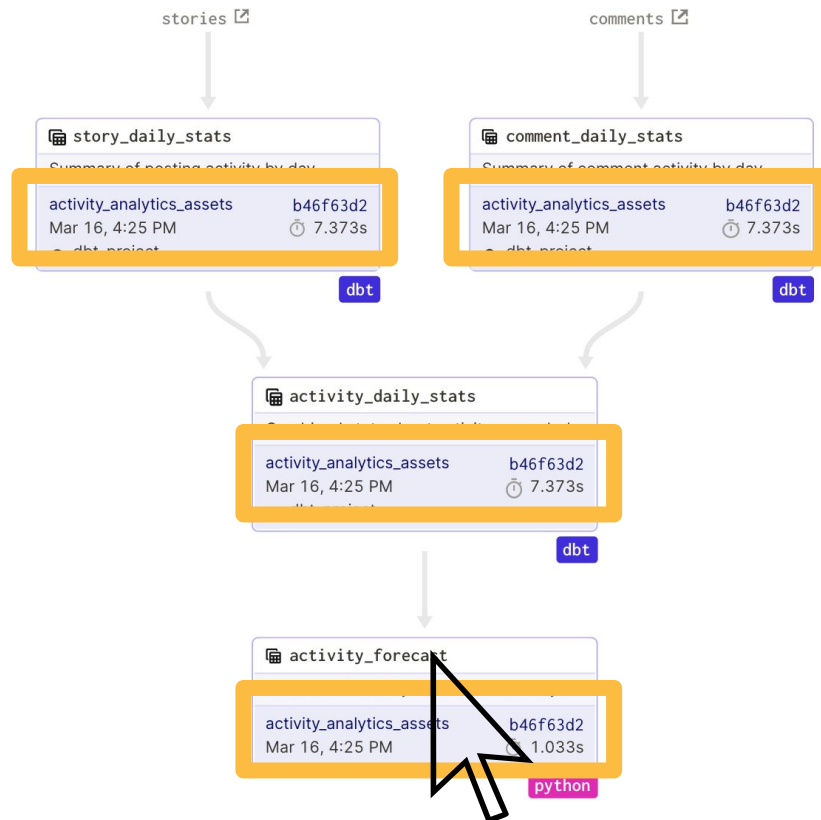
☰ 🔍 debug info warning error critical event[📄 Copy URL](#)

OP	EVENT TYPE	INFO	TIMESTAMP
activity_forecast	LOGS_CAPTURED	Started capturing logs for step: activity_forecast. <div><div>captured_logs</div><div>View stdout / stderr</div></div>	16:25:16.968
activity_forecast	STEP_START	Started execution of step "activity_forecast".	16:25:17.086
activity_forecast	LOADED_INPUT	Loaded input "activity_daily_stats" using input manager "warehouse_io_manager", from output "activity_daily_stats" of step "dbt_project"	16:25:17.208
activity_forecast	STEP_INPUT	Got input "activity_daily_stats" of type "DataFrame". (Type check passed).	16:25:17.261
activity_forecast	STEP_OUTPUT	Yielded output "result" of type "DataFrame". (Type check passed).	16:25:17.788
activity_forecast	ASSET_MATERIALIZAT...	Materialized value activity_forecast. <div><div>asset_key</div><div>activity_forecast</div><div>[View Asset]</div></div>	16:25:17.889
activity_forecast	HANDLED_OUTPUT	Handled output "result" using IO manager "io_manager"	16:25:18.067
activity_forecast	STEP_SUCCESS	Finished execution of step "activity_forecast" in 807ms.	16:25:18.119
-	ENGINE_EVENT	Multiprocess executor: parent process exiting after 38.84s (pid: 55765) <div><div>pid</div><div>55765</div></div>	16:25:18.833
-	RUN_SUCCESS	Finished execution of run for "activity_analytics_assets".	16:25:18.907
-	ENGINE_EVENT	Process for run exited (pid: 55765).	16:25:19.031

✱ Type an asset subset... (ex: story_daily_stats+)

0:01

✱ Rematerialize All



activity_forecast

Asset in activity_analytics@hacker_news_assets


0:11

Rematerialize

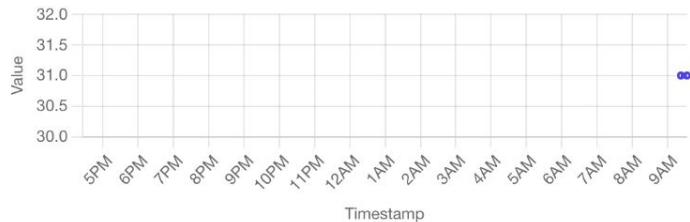
Activity Definition

Asset Events

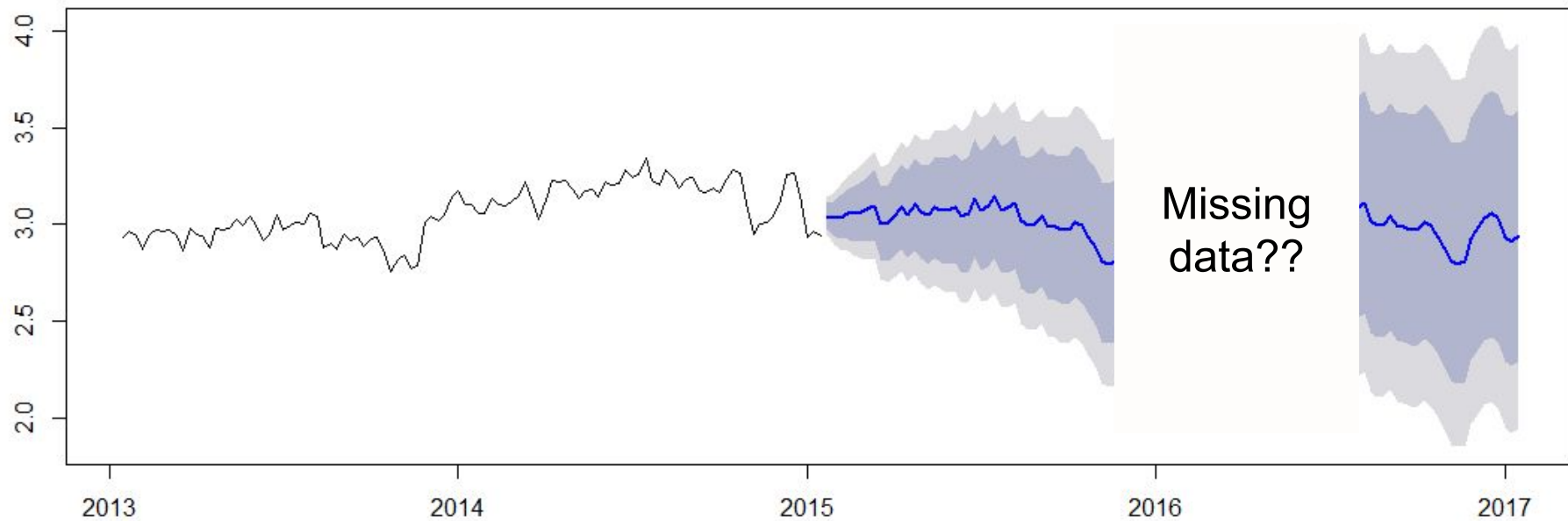
Timestamp	Job / Pipeline	Run
▶ Mar 17, 9:32 AM ⚡ Materialization		● 0e4f7d91
▼ Mar 17, 9:22 AM ⚡ Materialization		● 3f7a7370

Row count	31																				
Path	/tmp/activity_forecast.pq 																				
Sample	<table><thead><tr><th></th><th>date</th><th>num_comments</th></tr></thead><tbody><tr><td>0</td><td>2022-03-16 00:00:00</td><td>0.65186</td></tr><tr><td>1</td><td>2022-03-17 00:00:00</td><td>0.652254</td></tr><tr><td>2</td><td>2022-03-18 00:00:00</td><td>0.652649</td></tr><tr><td>3</td><td>2022-03-19 00:00:00</td><td>0.653044</td></tr><tr><td>4</td><td>2022-03-20 00:00:00</td><td>0.653439</td></tr></tbody></table>			date	num_comments	0	2022-03-16 00:00:00	0.65186	1	2022-03-17 00:00:00	0.652254	2	2022-03-18 00:00:00	0.652649	3	2022-03-19 00:00:00	0.653044	4	2022-03-20 00:00:00	0.653439	
	date	num_comments																			
0	2022-03-16 00:00:00	0.65186																			
1	2022-03-17 00:00:00	0.652254																			
2	2022-03-18 00:00:00	0.652649																			
3	2022-03-19 00:00:00	0.653044																			
4	2022-03-20 00:00:00	0.653439																			
Columns	<table><tbody><tr><td>date</td><td>datetime64[ns]</td></tr><tr><td>num_comments</td><td>float64</td></tr></tbody></table>		date	datetime64[ns]	num_comments	float64															
date	datetime64[ns]																				
num_comments	float64																				

Row count



```
select * from activity_forecast
```



activity_forecast

Asset in activity_analytics@hacker_news_assets

0:02 Rematerialize

Activity Definition

Asset Events

Timestamp

Job / Pipeline

Mar 17, 9:32 AM

Materialization

Mar 17, 9:22 AM

Materialization

Row count

31

Path

/tmp/activity_forecast.pg

Sample	date	num_comments
0	2022-03-16 00:00:00	0.65186
1	2022-03-17 00:00:00	0.652254
2	2022-03-18 00:00:00	0.652649
3	2022-03-19 00:00:00	0.653044
4	2022-03-20 00:00:00	0.653439

Columns

date datetime64[ns]

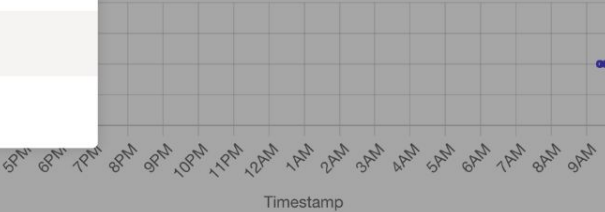
num_comments float64

activity_forecast

activity_analytics_assets
Job

activity_forecast
Asset

activity_daily_stats
Asset



Activity

Definition

Asset Events

Timestamp	Job / Pipeline	Run
▶ Mar 17, 9:32 AM ⚡ Materialization		● 0e4f7d91
▼ Mar 17, 9:22 AM ⚡ Materialization		● 3f7a7370

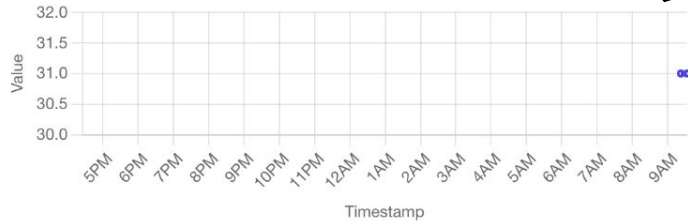
Row count	31
-----------	----

Path	/tmp/activity_forecast.pq
------	---------------------------

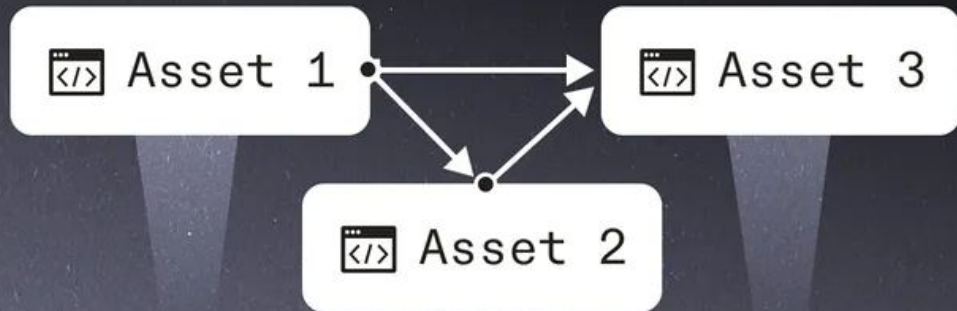
Sample	date	num_comments
0	2022-03-16 00:00:00	0.65186
1	2022-03-17 00:00:00	0.652254
2	2022-03-18 00:00:00	0.652649
3	2022-03-19 00:00:00	0.653044
4	2022-03-20 00:00:00	0.653439

Columns	date datetime64[ns]
	num_comments float64

Row count



Asset Definitions



Asset Materializations



Assets are not static

Orchestration



What does a traditional orchestrator do?



Invokes computations at the right time



Models the dependencies between computations



Tracks what computations ran

Orchestrators are experts on



When stuff happens



When stuff is going to happen



What it takes to make something happen

Critical questions about our assets

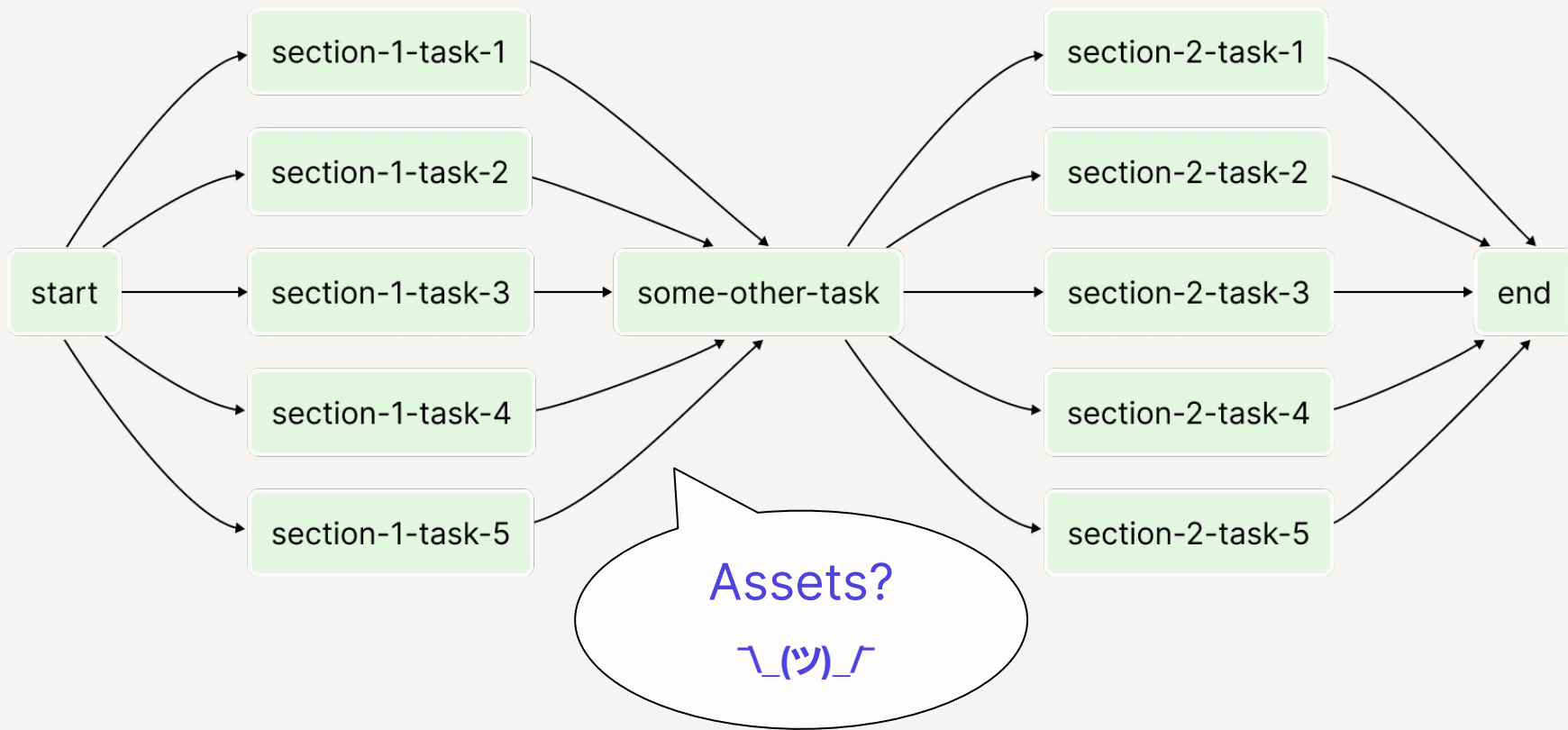
Is this asset up-to-date?

What do I need to run to refresh this asset?

When will this asset be updated next?

What code and data were used to generate this asset?

After pushing a change, what assets need to be updated?



Orchestration + Assets = ?

Asset orchestrator:
manages change in assets

Triggered asset jobs

Re-materialize your assets on a schedule, or via a sensor

Reconciliation

Update your assets when they don't match their definitions

```
ScheduleDefinition(  
    job=assets.build_job(  
        "daily_stats_job", selection=["*activity_daily_stats"]  
    ),  
    cron_schedule="0 0 * * *",  
),
```

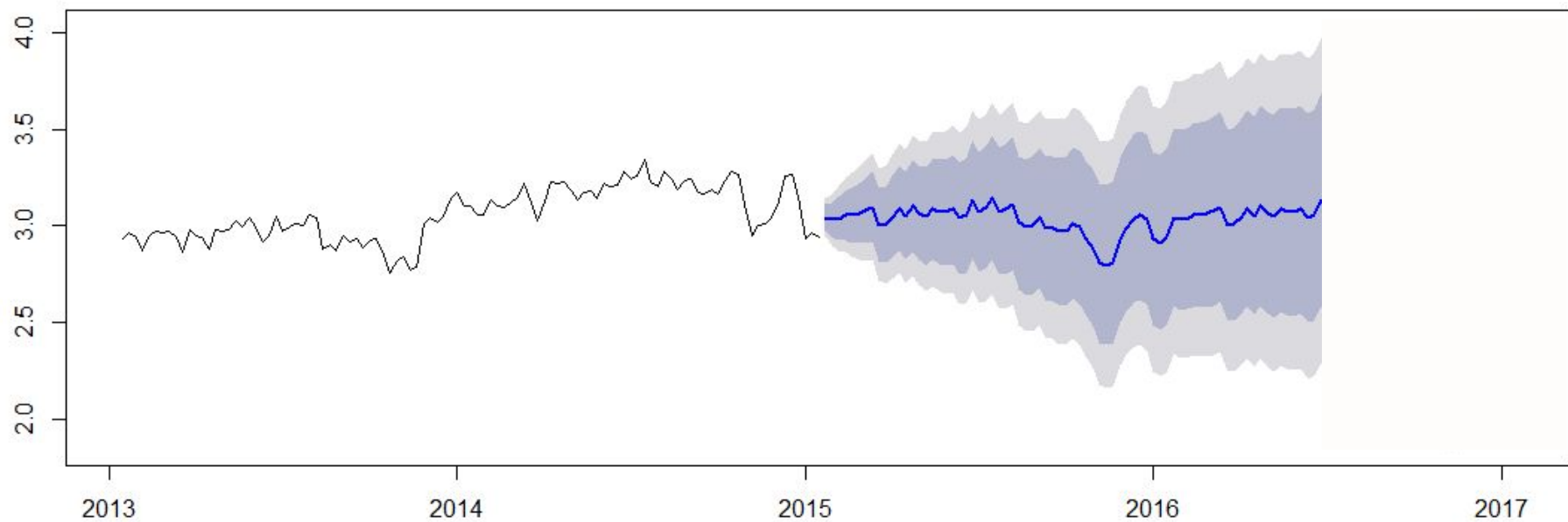
✖ Type an asset subset... (ex: story_daily_stats+*)

0:08

Materialize All



```
select * from activity_forecast
```



activity_daily_stats

Asset in activity_analytics@hacker_news_assets

Schedule: At 12:00 AM

Activity Definition

Description

daily_stats_job run_dbt_project

Combined stats about activity on each day

Raw SQL:

```
select *
from {{ ref('comment_daily_stats') }}
full outer join {{ ref('story_daily_stats') }}
using (date)
```

Upstream Assets (2)

[View upstream graph](#)

story_daily_stats

No materializations

run_dbt_project

comment_daily_stats

No materializations

run_dbt_project

Downstream Assets (1)

[View downstream graph](#)

activity_forecast

No materializations

Critical questions about our assets

Is this asset up-to-date?

What do I need to run to refresh this asset?

When will this asset be updated next?

What code and data were used to generate this asset?

After pushing a change, what assets need to be updated?

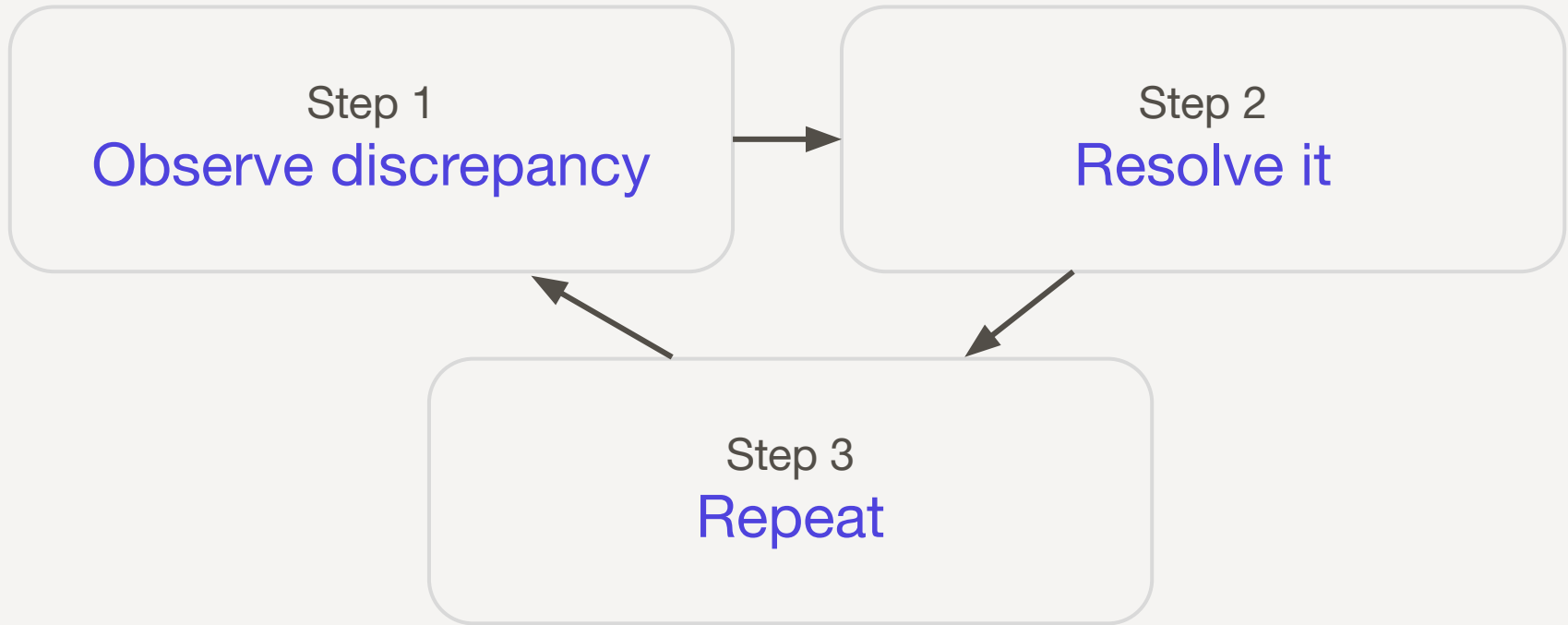
Scheduled asset jobs

Re-materialize your assets on a schedule

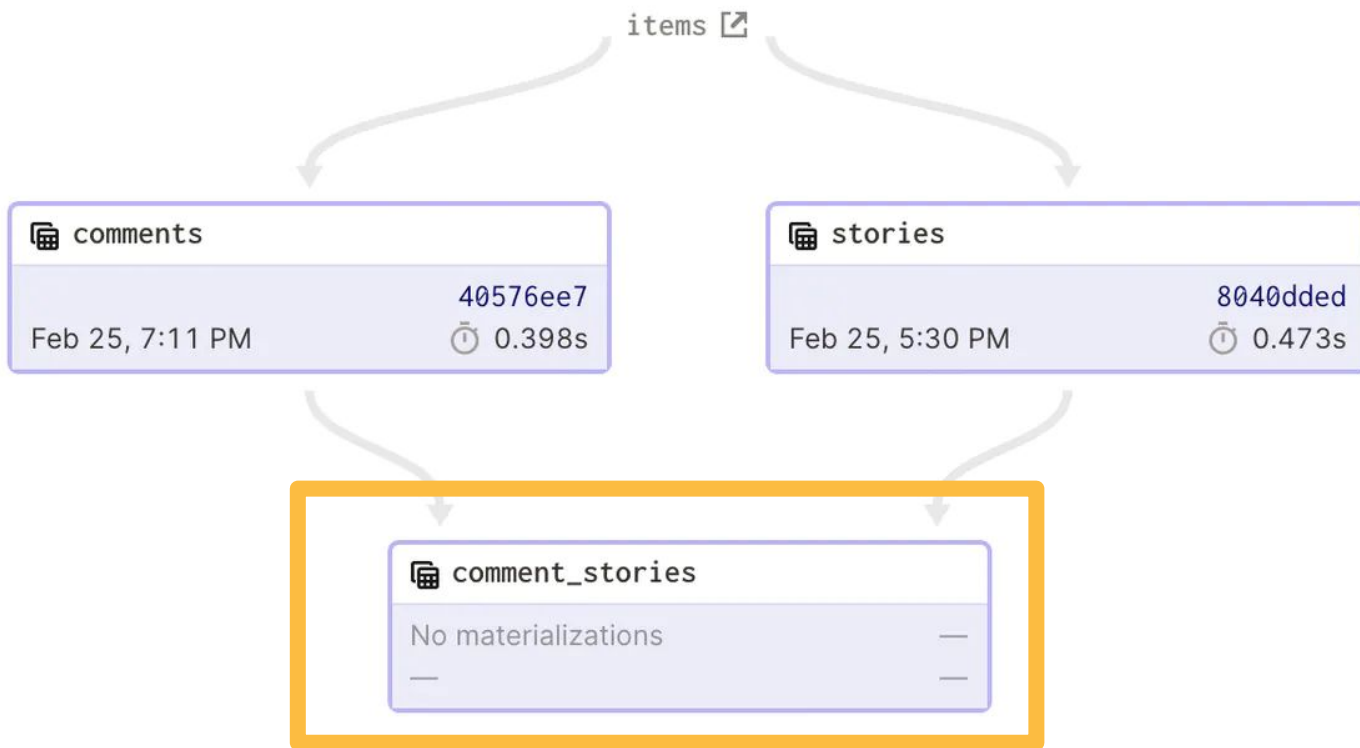
Reconciliation

Update your assets when they don't match their definitions

Reconciliation: how it works



Discrepancy: definition w/o materialization



Discrepancy: partition definition w/o materialization

items  Asset in core@hacker_news_assets 

Activity **Definition**

Description

Items from the Hacker News API: each is a story or a comment on a story.

Partitions

Hourly, starting 2022-03-18-12:00 UTC.

3/10 Partitions

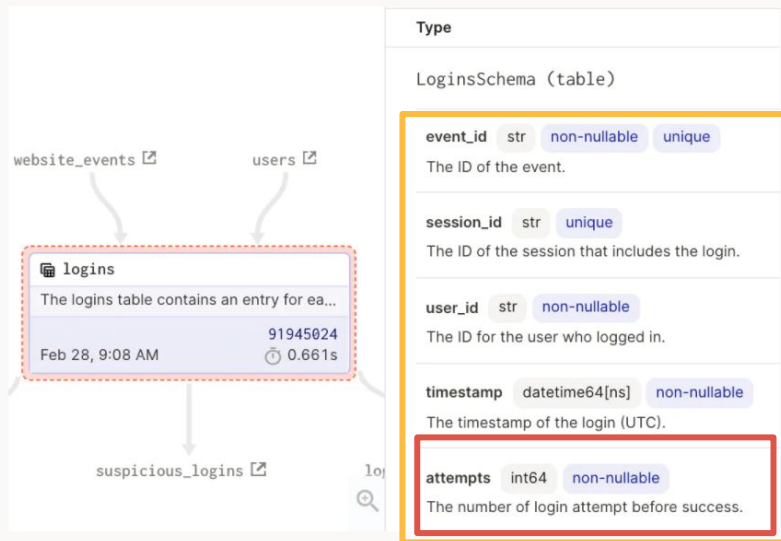


Discrepancy: upstream changed



Discrepancy: columns

Definition



Latest Materialization

Asset Events	
Timestamp	
Feb 28, 9:08 AM	
Materialization	
schema	
event_id	str non-nullable unique
The ID of the event.	
session_id	str unique
The ID of the session that includes the login	
user_id	str non-nullable
The ID for the user who logged in.	
timestamp	datetime64[ns] non-nullable
The timestamp of the login (UTC).	

The fruits of reconciliation

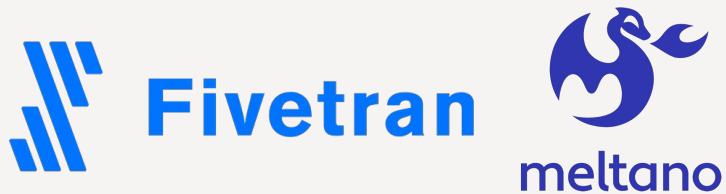
- Move bookkeeping about what's stale out of your head
- Avoid unnecessary computation
- Understand why computations occurred

Software-defined assets & the modern data stack

A set of tools and practices
that have dramatically
simplified common patterns
for working with data.



Define derived
tables



Define source
tables

Pre-modern world

Execution

Ingest

SQL
transforms

Python
transforms

ML

Orchestration



Modern world

Execution



Orchestration

Sync external tables @ midnight

Rebuild derived tables @ 1 am



We're forced to choose...

Unified but

(Pre-modern)

but siloed

(Data Stack)

Unified and Declarative



dagster



```
airbyte_assets = build_airbyte_assets(  
    connection_id="my-airbyte-connection-id",  
    destination_tables=["table1", "table2"],  
)
```



```
dbt_assets = load_assets_from_dbt_project(  
    project_dir="my/dbt/project/dir/",  
)
```



```
@asset  
def predicted_usage(daily_usage_stats: DataFrame) -> DataFrame:  
    ...
```


runs
analytics_assets_job
Feb 16, 10:21 AM
093c029a
0:00:21

organizations
analytics_assets_job
Feb 16, 10:21 AM
093c029a
0:00:21

airbyte

airbyte

event_logs
analytics_assets_job
Feb 16, 10:21 AM
093c029a
0:00:21

runs_annotated
analytics_assets_job
Feb 16, 10:21 AM
093c029a
5.493s

airbyte

dbt

steps_annotated
analytics_assets_job
Feb 16, 10:21 AM
093c029a
5.493s

dbt

step_events_aggregated
analytics_assets_job
Feb 16, 10:21 AM
093c029a
5.493s

dbt

step_stats
analytics_assets_job
Feb 16, 10:21 AM
093c029a
5.493s

dbt

run_stats
analytics_assets_job
Feb 16, 10:21 AM
093c029a
5.493s

dbt

daily_step_stats
analytics_assets_job
Feb 16, 10:21 AM
093c029a
5.493s

dbt

daily_run_stats
analytics_assets_job
Feb 16, 10:21 AM
093c029a
5.493s

dbt

usage_stats
analytics_assets_job
Feb 16, 10:21 AM
093c029a
5.493s

dbt

daily_usage_stats
analytics_assets_job
Feb 16, 10:21 AM
093c029a
5.493s

dbt

predicted_usage
analytics_assets_job
Feb 16, 10:21 AM
093c029a
5.183s

python

Boiling it all down

Missing abstraction: the software-defined asset

Declarative approach → less chaos, more trust

Asset orchestrators manage change in assets

Python: first-class citizen of the Modern Data Stack

Unified control plane for the Modern Data Stack

Thank you

Sandy Ryza
@s_ryz





dagster

Asset Definitions

Asset Materializations

Software-Defined Asset

Upstream Asset Key



Asset Key

Software-Defined Asset

Upstream Asset Key



Asset Key

The casualties of chaos

Trust

Productivity