

Multi-instance Asset Synchronization

-
Push or Pull ?

Sébastien Crocquevieille

3.0



Othmane EL Metioui
Head of Data & AI



Numberly

What we do

We generate value and growth for our clients by turning their digital marketing expenses into impactful and profitable data-driven investments.

CRM & Loyalty Management

- 1st Party Data Collection
- Customer Journey design
- Omnichannel campaign management
- **Clean Room management**
- Incrementality & CLV

Numberly Martech Platform

- Content builder
- Marketing Automation tool
- Messaging APIs
- Impactly

Digital Media

- Media Strategy
- Full funnel activation
- Data Collaboration & Retail Media
- Experimentation & measurement

Data Strategy & Insights

- Digital Consulting
- Insights & Analytics
- AI-driven solutions
- Numberly Academy

Numberly the Marketing Technologist

- **We are around since 2000, HQ in Paris France.**



Paris



Brussels



London



Amsterdam



Dubai



Tel Aviv



New-York



Montréal

- **95% of Our Infrastructure is Self-Hosted on Bare Metal Servers**
- **We are our own Internet Service Provider.**
- **We commit to Open Source & Developer Community**
 - Numberly Github, OSS financing, Europython Sponsor 11 years straight
- **Self hosted Airflow users since 2017 - version 1.x**

Numberly the Marketing Technologist

Airflow Key figures in Production:

- **1817** DAGs on 2 instances with **341** Datasets
- **16 500** daily DagRuns with **44 500** daily TaskInstances

WhoAml

Who Am I?

- Sébastien Crocquevieille
- Data Engineer @ Numberly
- French & Mexican
- Airflow user since 2018
- EuroPython 2023 Speaker
 - [Orchestrating Python Workflows in Apache Airflow](#)
- Pycon TW 2025 Speaker



You can call me **SEB**

Events vs Assets Scheduling

Disclaimer

THE CHARACTERS AND EVENTS DEPICTED IN THIS PHOTOPLAY
ARE FICTITIOUS. ANY SIMILARITY TO ACTUAL PERSONS,
LIVING OR DEAD, IS PURELY COINCIDENTAL.

External Events

It could be anything that happens outside Airflow

- Any message/post
- Any change in a value
- Any detectable action or signal

Event Driven Scheduling (Airflow 3.0)

Design :

- **Execution that strictly follows external events**
- **Individual action based on payload ?**
- **1 to 1 execution**

In practice :

- **A little complex to set up**
 - AssetWatcher: monitors external event source via triggers
- **Requires guarantees concerning:**
 - Ensuring deliverability
 - Avoiding duplicate events
 - Ordering ??

Asset-Aware Scheduling (Airflow 2.4)

Design :

- **Event that indicates a change in Data Source**
- **Data source changes are indistinguishable**

Constraints:

- **Similar to Event Driven but with lower constraints**
 - 1 to 1 execution **not** guaranteed
 - Order of events **shouldn't** be important
- **No triggers needed, just Airflow Assets**

Airflow (Data) Assets

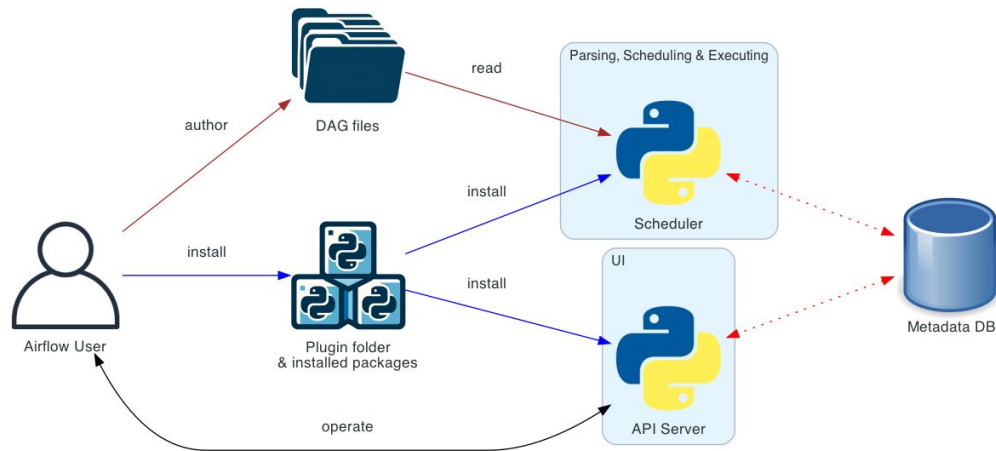
```
from airflow.sdk import Asset
```

```
example_asset = Asset("s3://asset-bucket/example.csv")
```

- It is just a string → representing a Data Source
- Stored in Airflow DB
- Maintaining link is YOUR responsibility
- Generated by DAG parsing

Upscaling our Scheduling

Airflow Scheduler



Performance issues

- **Machine performance**

- High CPU and RAM usage
- Scheduler heartbeat failures
- dag_processing.total_parse_time is high

Performance issues

- **Machine performance**

- High CPU and RAM usage
- Scheduler heartbeat failures
- dag_processing.total_parse_time is high

- **Slow scheduling**

- Many tasks in “queued” or “scheduled” state
- Delayed DAG runs

Performance issues

- **Machine performance**

- High CPU and RAM usage
- Scheduler heartbeat failures
- dag_processing.total_parse_time is high

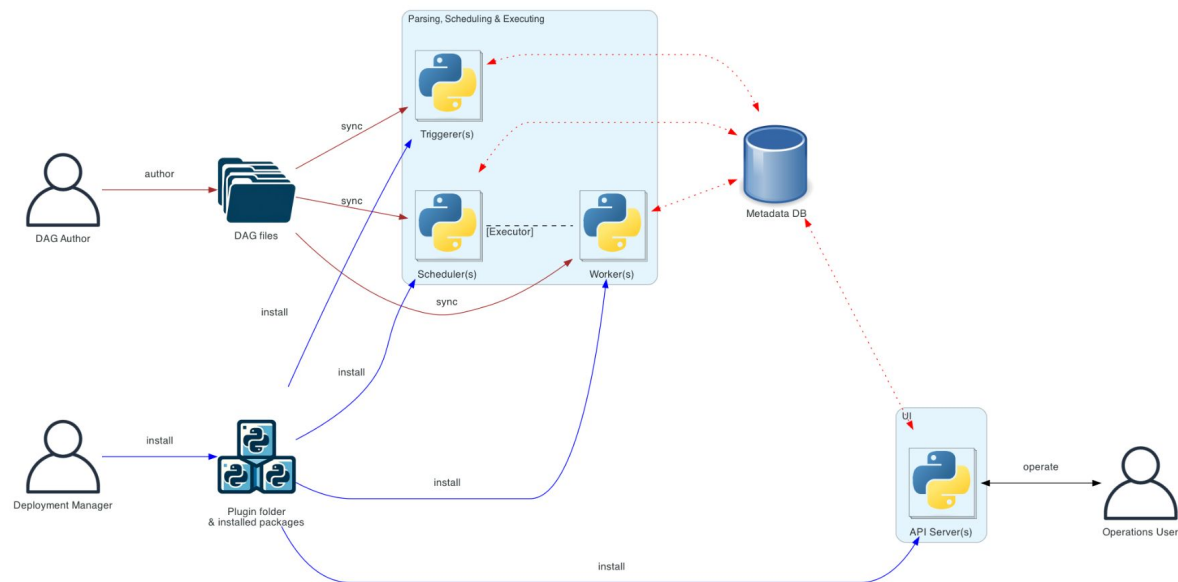
- **Slow scheduling**

- Many tasks in “queued” or “scheduled” state
- Delayed DAG runs

- **Database limitations**

- Too many connections
- Slow queries

Airflow Remote Executors



Airflow Remote Executors

Pros:

- Robust: decoupling workers from scheduler process
- Effective: Worker parallelization
- Available: Low latency workers always running

Airflow Remote Executors

Cons:

- Requires Infrastructure setup
- Can be Expensive: Cloud cost



Still having issues?

Scheduler Replication

Running More Than One Scheduler [\[link\]](#)

Airflow supports running more than one scheduler concurrently – both for performance reasons and for resiliency.

Scheduler Replication

Pros:

- Split scheduling load
- More resilient
- Not much work

Cons:

- More Database operations

Need more ?

Separate Instances

Just have 2 Airflows

Pros:

- Maximum scalability
- Separation of Concerns
- Better Access Control



Separate Instances

Just have 2 Airflows ?!

Cons:

- Even more maintenance
- Even more resources
- Redundant operations



Assets

+

Multiple instances

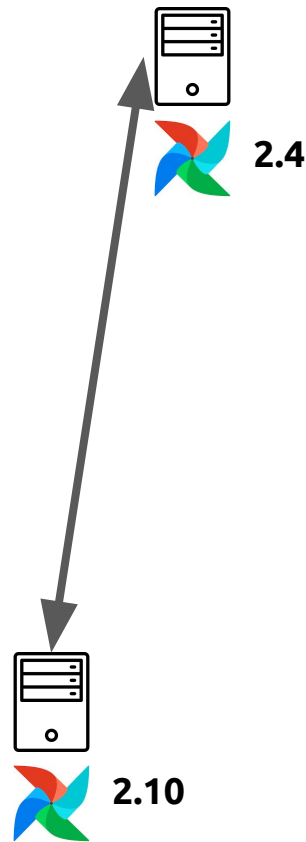
=

...

Synchronizing Airflow Assets

Our use case

- **Airflow Migration 2.4 → 2.10 (latest version)**
- **Complete dependency overhaul (Python & Spark)**
- **Big dag library (~2000 DAGs)**
- **From **old** instance to **new****
- **Upgrading code as we migrate**



Our use case



- **Migrating DAGs by batches based on**
 - Complexity
 - Dependencies
 - Missing features
 - Urgency
- **Specific needs:**
 - Cross instance DAG Trigger
 - Cross instance Task sensor
- **What about Airflow Assets? ⇒ Synchronization !**

Push or Pull ?

Push

How do we start?

By getting the Asset Events!

GET

/datasets/events



Response samples

200

401

403

404

Content type
application/json

[Copy](#) [Expand all](#) [Collapse all](#)

```
{
  - "dataset_events": [
    - {
      "dataset_id": 0,
      "dataset_uri": "string",
      "extra": { },
      "source_dag_id": "string",
      "source_task_id": "string",
      "source_run_id": "string",
      "source_map_index": 0,
      + "created_dagruns": [ - ],
      "timestamp": "string"
    }
  ],
  "total_entries": 0
}
```

Design choices

- **Using Airflow API to get Asset Events**
 - Filter by -timestamp & store the “offset” in Airflow Variables

io Consumer

| | | | |
|---|---|---|------------|
| airflow-dataset-synchronizer.dataset-events | ✓ | Partition | Start From |
| | | <input checked="" type="checkbox"/> All <input checked="" type="checkbox"/> Earliest <input type="checkbox"/> Latest <input type="checkbox"/> Date (Paris time) <input type="checkbox"/> Timestamp or oid <input type="checkbox"/> Offset | |

Control

| | |
|--|---|
| <div> <div>▶</div> <div>⏸</div> <div>⏹</div> </div> <div>Msg limit</div> <div>Keyword filter</div> | <div>Filter-in</div> <div><input type="checkbox"/> Headers</div> <div><input type="checkbox"/> Key</div> <div><input checked="" type="checkbox"/> Value</div> |
|--|---|

| Date | Timestamp | Partition | Offset | Headers | Key | Value |
|---------------------|---------------|-----------|--------|---------|--|---|
| 2025-09-01 22:35:42 | 1756758942053 | 0 | 26492 | | five.thalcloud.gr.br.clientsync.contactability | { "dataset_id": 86, "dataset_name": "five.thalcloud.gr.br.clientsync.contactability", "event_id": 23620 } |
| 2025-09-01 22:07:57 | 1756757277595 | 0 | 26491 | | five.thalcloud.gr.br.clientsync.campaign_targeting | { "dataset_id": 72, "dataset_name": "five.thalcloud.gr.br.clientsync.campaign_targeting", "event_id": 23619 } |
| 2025-09-01 22:07:57 | 1756757277595 | 0 | 26490 | | five.thalcloud.gr.br.clientsync.campaign_targeting | { "dataset_id": 78, "dataset_name": "five.thalcloud.gr.br.clientsync.campaign_targeting", "event_id": 23618 } |
| 2025-09-01 22:02:29 | 1756756949361 | 0 | 26489 | | five.thalcloud.gr.br.clientsync.contactability | { "dataset_id": 86, "dataset_name": "five.thalcloud.gr.br.clientsync.contactability", "event_id": 23617 } |
| 2025-09-01 21:51:55 | 1756756315091 | 0 | 26488 | | five.thalcloud.gr.br.clientsync.delete_users | { "dataset_id": 58, "dataset_name": "five.thalcloud.gr.br.clientsync.delete_users", "event_id": 23616 } |
| 2025-09-01 21:51:55 | 1756756315091 | 0 | 26487 | | five.thalcloud.gr.br.clientsync.delete_users | { "dataset_id": 77, "dataset_name": "five.thalcloud.gr.br.clientsync.delete_users", "event_id": 23615 } |
| 2025-09-01 21:41:20 | 1756755680563 | 0 | 26486 | | five.thalcloud.gr.br.clientsync.delete_users | { "dataset_id": 101, "dataset_name": "five.thalcloud.gr.br.clientsync.delete_users", "event_id": 23614 } |
| 2025-09-01 21:37:49 | 1756755469318 | 0 | 26485 | | five.thalcloud.gr.br.clientsync.campaign_targeting | { "dataset_id": 78, "dataset_name": "five.thalcloud.gr.br.clientsync.campaign_targeting", "event_id": 23613 } |
| 2025-09-01 21:37:49 | 1756755469318 | 0 | 26484 | | five.thalcloud.gr.br.clientsync.campaign_targeting | { "dataset_id": 72, "dataset_name": "five.thalcloud.gr.br.clientsync.campaign_targeting", "event_id": 23612 } |
| 2025-09-01 21:32:16 | 1756755136838 | 0 | 26483 | | five.thalcloud.gr.br.clientsync.contract | { "dataset_id": 80, "dataset_name": "five.thalcloud.gr.br.clientsync.contract", "event_id": 23611 } |
| 2025-09-01 21:31:47 | 1756755107354 | 0 | 26482 | | five.thalcloud.gr.br.clientsync.contract | { "dataset_id": 83, "dataset_name": "five.thalcloud.gr.br.clientsync.contract", "event_id": 23610 } |
| 2025-09-01 21:24:43 | 1756754683595 | 0 | 26481 | | five.thalcloud.gr.br.clientsync.aggregation_customer | { "dataset_id": 54, "dataset_name": "five.thalcloud.gr.br.clientsync.aggregation_customer", "event_id": 23609 } |
| 2025-09-01 21:18:41 | 1756754321064 | 0 | 26480 | | five.thalcloud.gr.br.clientsync.user | { "dataset_id": 15, "dataset_name": "five.thalcloud.gr.br.clientsync.user", "event_id": 23608 } |
| 2025-09-01 21:07:40 | 1756753660001 | 0 | 26479 | | five.thalcloud.gr.br.clientsync.campaign_targeting | { "dataset_id": 72, "dataset_name": "five.thalcloud.gr.br.clientsync.campaign_targeting", "event_id": 23607 } |
| 2025-09-01 21:07:40 | 1756753660001 | 0 | 26478 | | five.thalcloud.gr.br.clientsync.campaign_targeting | { "dataset_id": 78, "dataset_name": "five.thalcloud.gr.br.clientsync.campaign_targeting", "event_id": 23606 } |
| 2025-09-01 21:03:34 | 1756753414963 | 0 | 26477 | | five.thalcloud.gr.br.clientsync.contactability | { "dataset_id": 86, "dataset_name": "five.thalcloud.gr.br.clientsync.contactability", "event_id": 23605 } |
| 2025-09-01 20:42:26 | 1756752146414 | 0 | 26476 | | five.thalcloud.gr.br.clientsync.order | { "dataset_id": 230, "dataset_name": "five.thalcloud.gr.br.clientsync.order", "event_id": 23604 } |
| 2025-09-01 20:37:31 | 1756751851238 | 0 | 26475 | | five.thalcloud.gr.br.clientsync.campaign_targeting | { "dataset_id": 78, "dataset_name": "five.thalcloud.gr.br.clientsync.campaign_targeting", "event_id": 23603 } |
| 2025-09-01 20:37:31 | 1756751851238 | 0 | 26474 | | five.thalcloud.gr.br.clientsync.campaign_targeting | { "dataset_id": 72, "dataset_name": "five.thalcloud.gr.br.clientsync.campaign_targeting", "event_id": 23602 } |
| 2025-09-01 20:36:23 | 1756751783697 | 0 | 26473 | | five.thalcloud.gr.br.clientsync.corporate_client_id | { "dataset_id": 82, "dataset_name": "five.thalcloud.gr.br.clientsync.corporate_client_id", "event_id": 23601 } |
| 2025-09-01 20:34:52 | 1756751692813 | 0 | 26472 | | five.thalcloud.gr.br.clientsync.order_history | { "dataset_id": 230, "dataset_name": "five.thalcloud.gr.br.clientsync.order_history", "event_id": 23600 } |
| 2025-09-01 20:34:52 | 1756751692813 | 0 | 26471 | | five.thalcloud.gr.br.clientsync.user | { "dataset_id": 81, "dataset_name": "five.thalcloud.gr.br.clientsync.user", "event_id": 23599 } |

Design choices

- **Using Airflow API to get Asset Events**
 - Filter by -timestamp & store the “offset” in Airflow Variables
- **Push them to a Kafka topic ⇒ Single source of truth for all Asset Events**

Almost Done !

Some 403 issues

- **No POST endpoint for Dataset Events** - Solved in 2.9
 - Dirty DB editor DAG
- **Dataset Events POST doesn't create NEW Datasets** - "Solved" in 3.0 with Materialize EP
 - Make my own materialize DAG

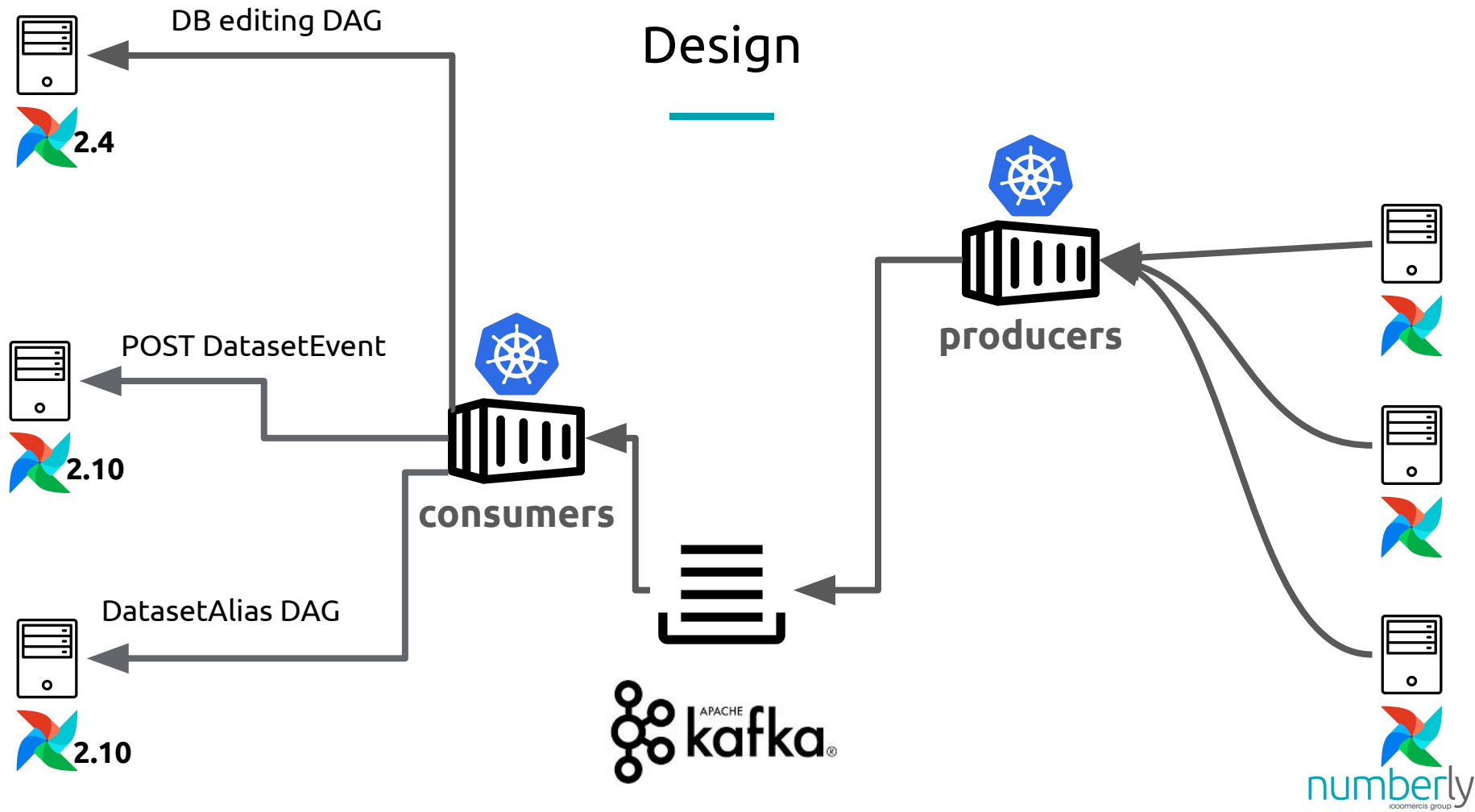
=====> **Conflict between parsed objects (Dataset) and created objects (Dataset Events)**

- **Datasets cannot be dynamically created** - Solved in 2.10 with DatasetAlias
 - Go to latest version (at the time) 2.10

Design choices

- **Using Airflow API to get Asset Events**
 - Filter by -timestamp & store the “offset” in Airflow Variables
- **Kafka topic ⇒ Single source of truth for all Asset Events**
- **Push back to Airflow through 3 different methods**

Design



```
1  {
2    "dataset_id": int,
3    "dataset_name": URI,
4    "event_id": int,
5    "event_timestamp": Datetime,
6    "origin": AIRFLOW INSTANCE URL,
7    "source_dag": DAG_ID,
8    "source_task": TASK_ID
9  }
```

It works !



Or Pull ?

Pull Based Scheduling

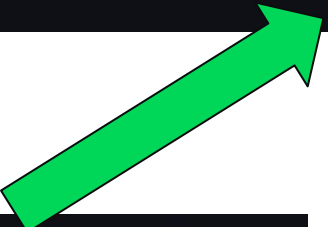
- **Airflow PULL from external instances directly**
- **Use Event based scheduling**
 - **Ignore** the fact that we need another migration
- **Poll Asset Events continuously to stay updated**

Asset Watcher ⇒ Asset Event ⇒ DAG execution

Let's code together

Using an Asset Watcher

```
from airflow.sdk import AssetWatcher  
  
my_watcher = AssetWatcher(name="my_watcher", trigger=???)
```



```
trigger: BaseEventTrigger | dict
```

Base Event Trigger

CLASS:

- Inherits from **BaseEventTrigger**

OBLIGATORY METHODS:

- **serialize**
 - Class path & arguments required for initialization must be serializable
- **run**
 - The actual trigger operation
- **__init__**

OPTIONAL METHODS:

- **cleanup**
- **_set_context**


```
7 from airflow.triggers.base import BaseEventTrigger, TriggerEvent
8
9 class SimpleEventTrigger(BaseEventTrigger):
10     def __init__(self, url: str, **kwargs):
11         self.url = url
12
13     def serialize(self) → tuple[str, dict[str, Any]]:
14         return (
15             "event_scheduling.SimpleEventTrigger",
16             {"url": self.url},
17         )
18
19     async def run(self) → AsyncIterator[TriggerEvent]:
20         while True:
21             event_str = await aiohttp.get(self.url).text()
22             if not event_str:
23                 continue
24             else:
25                 yield TriggerEvent(event_str)
26
```

```
27 from airflow.sdk import Asset, AssetWatcher
28
29 trigger = SimpleEventTrigger(url="http://magic/airflow_endpoint")
30 my_asset = Asset(
31     "simple_trigger_asset", watchers=[AssetWatcher(name="my_watcher", trigger=trigger)]
32 )
```

```
with DAG(dag_id="event_scheduled_job", schedule=[my_asset]) as dag:
```

Triggerer Process

[illegible]

Existing Event Triggers

Event Triggers:

- **KAFKA:** KafkaMessageQueueTrigger
- **GENERIC MQ:** MessageQueueTrigger



Base Triggers:

- **REDIS:** AwaitMessageTrigger
- **GOOGLE:** GCSTBlobTrigger
- **AWS:** EC2StateSensorTrigger, S3KeyTrigger, ...



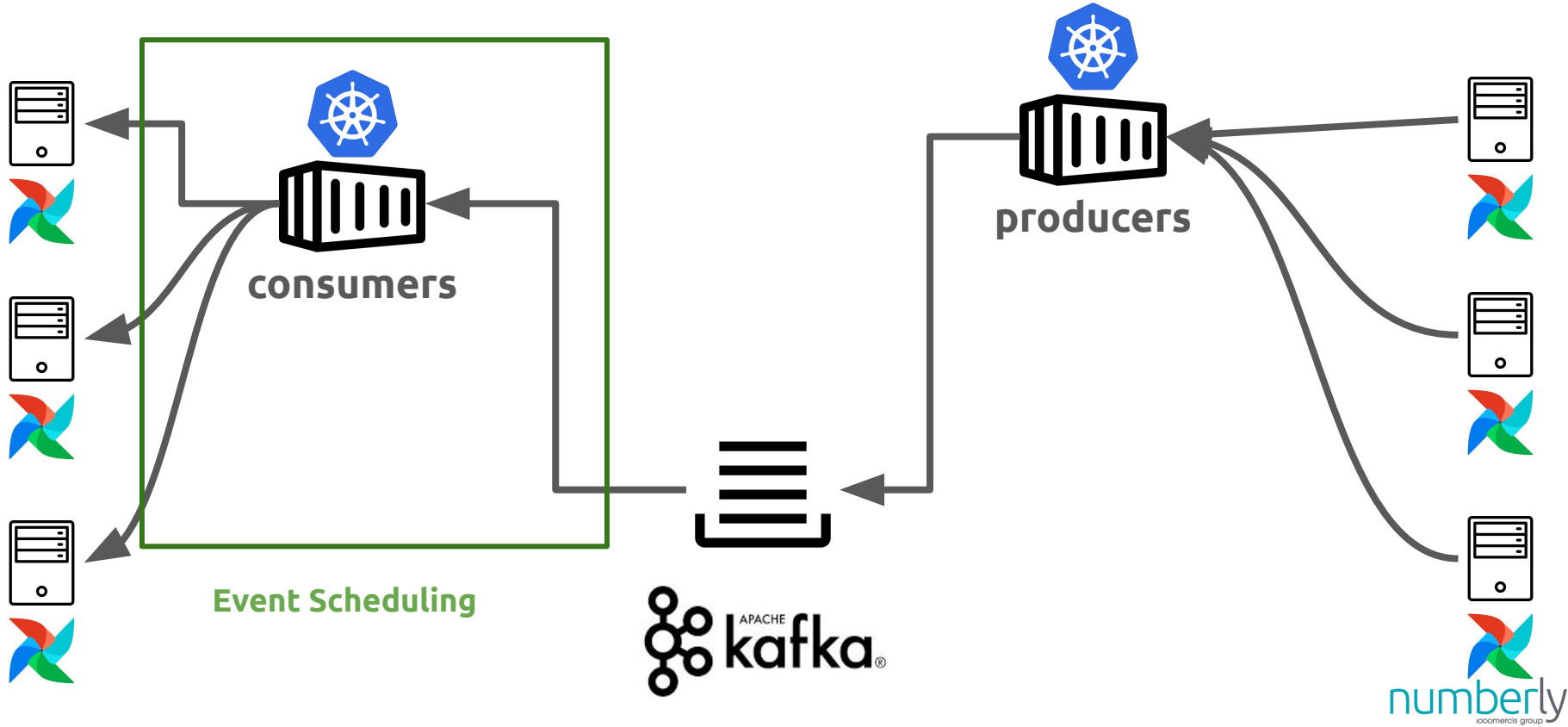
Pull based scheduling

Caveats:

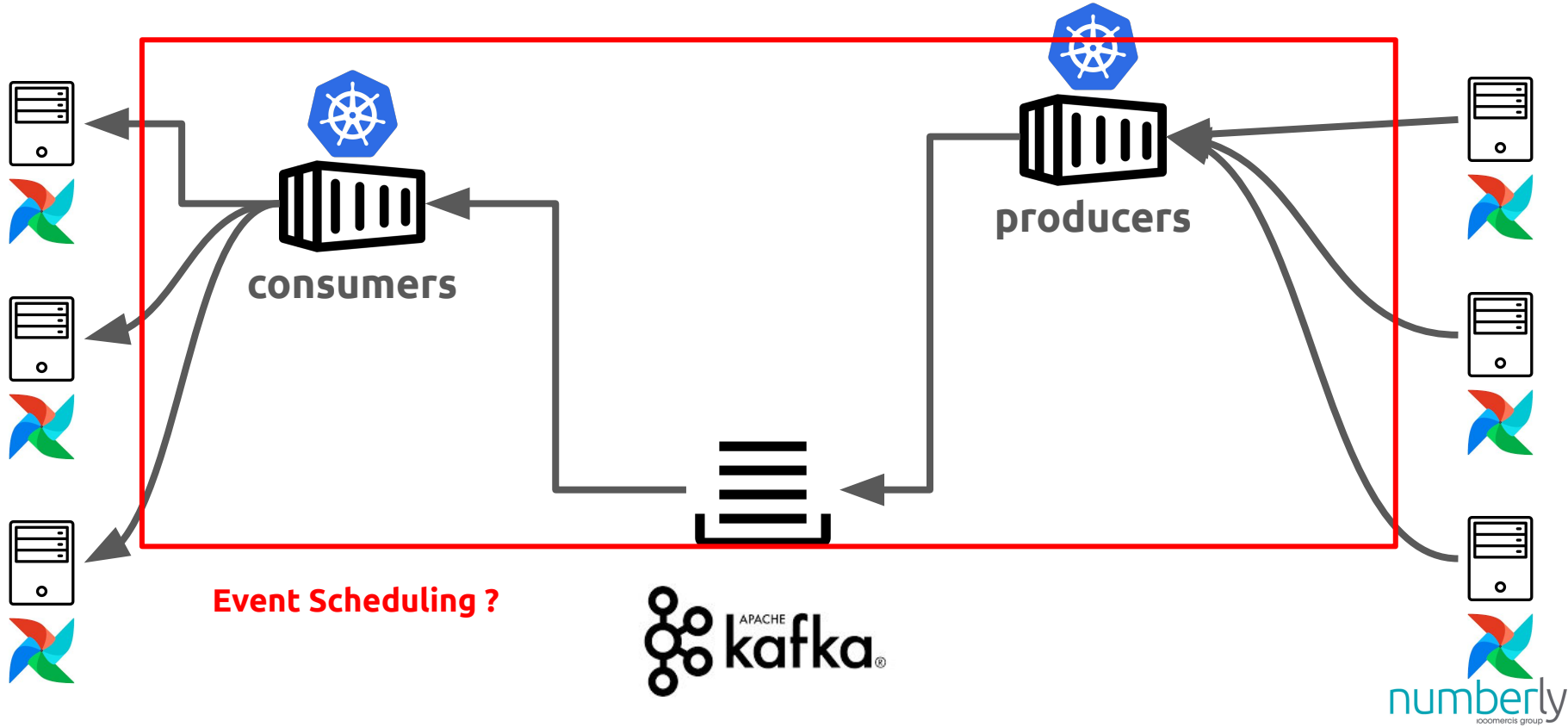
- **“run” method is king**
 - Error handling
 - Pagination
 - Authentication
 - No blocking
 - Not heavy
- **Know where your Triggerers are running**
 - Do you have one ?
 - Is it sharing scheduler resources ?
- **Use cases outside message queues are “exploratory” if possible**

So... can we pull ?

Using Pull ?



Using Pull ?



Cost of pulling

With the queue ?

- **Constant triggerer polling**
 - Less accurate than push
 - Or noisy
- **Still need to push to queue**
 - Keep dedicated service?
 - Use custom "high frequency" DAGs ?

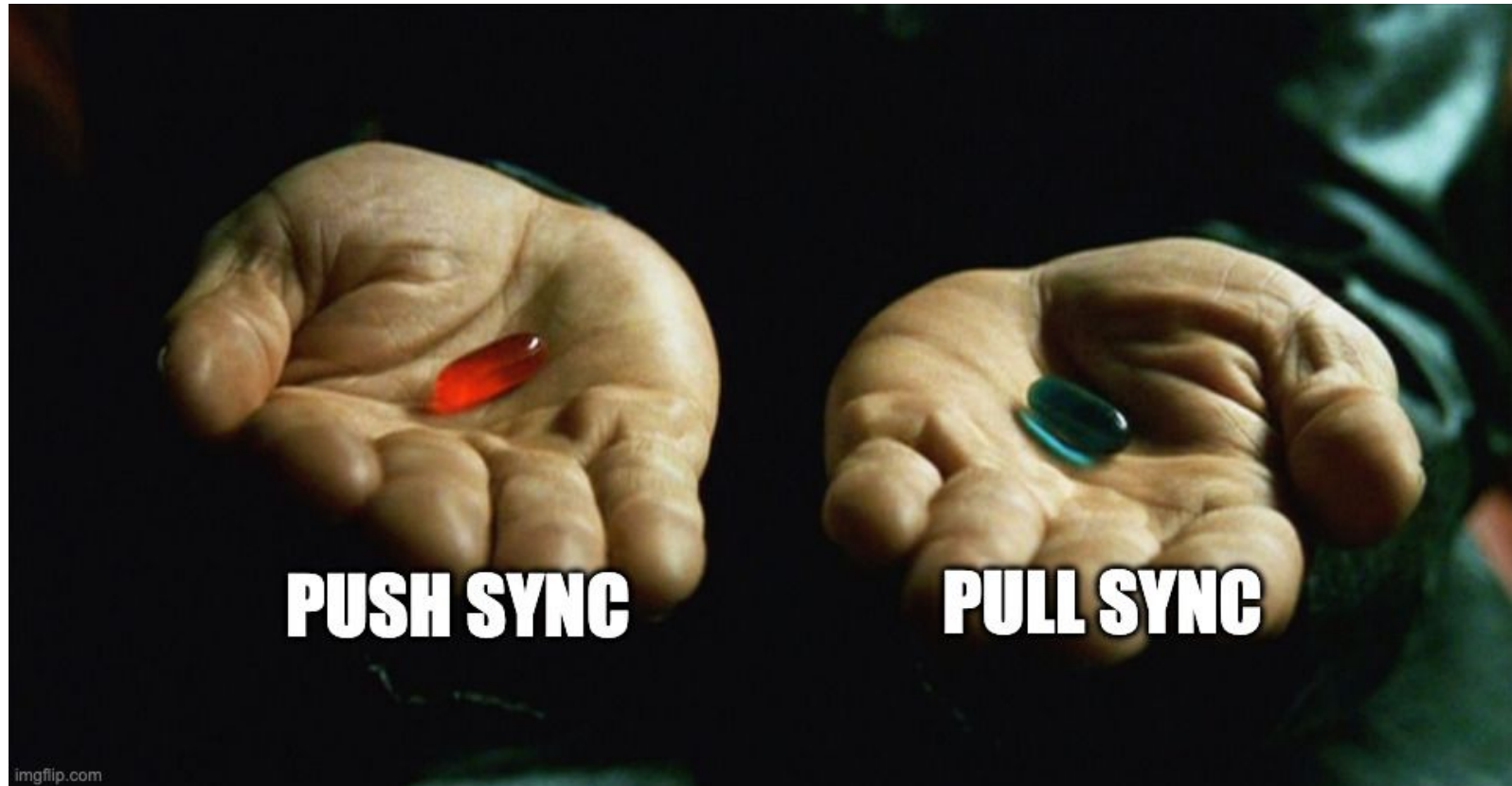
And without the queue?

- **API request + Airflow Variable**
- **$O(n^2)$ network calls**
- **Even More Airflow Activity**



So the choice is

I'm still on Airflow 2 !



Question time !

Bonus

ASTRONOMER

THE DATA FLOWCAST



Scaling On-Prem Airflow
With 2,000 DAGs
at Numberly
with Sébastien Crocquevieille

