**Jarek Potiuk**

Apache Airflow PMC member & committer

Member of Apache Software Foundation

Member of ASF Security Committee

𝕏 @jarekpotiuk

@potiuk
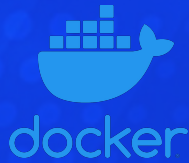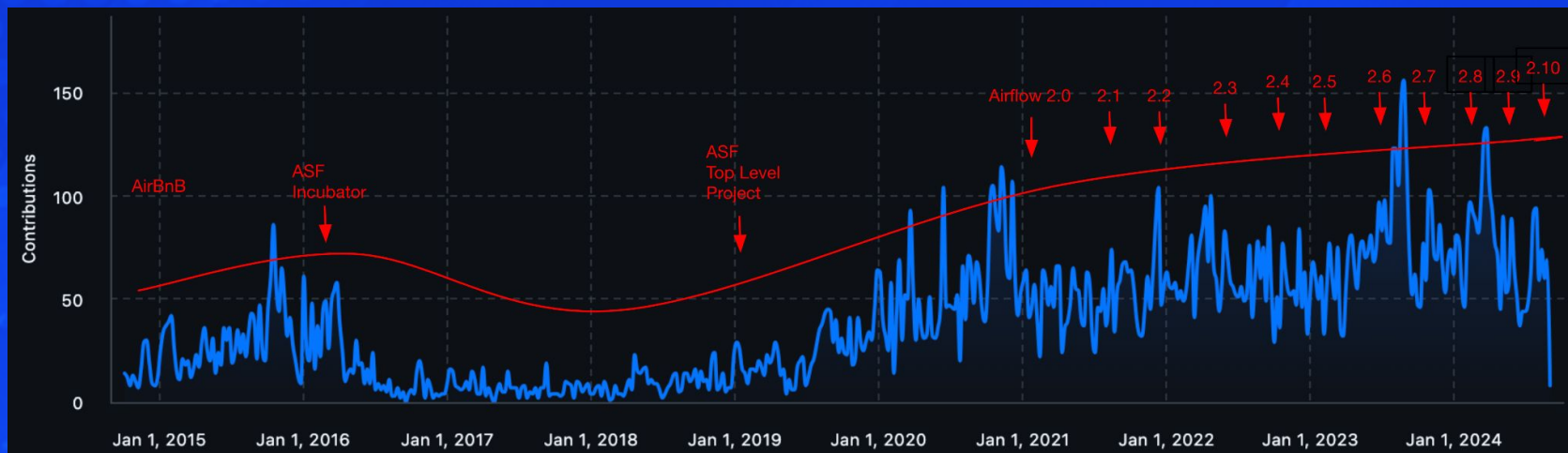
@jarekpotiuk

https://fosstodon.org/@jarekpotiuk

Airflow as Orchestrator

OSA CON | November 19-21, 2024

# Why Airflow 3

# Airflow timeline

# Airflow 2.0 strengths

- User base / Community
- Versatility
- Steady stream of new features
- Learning from the industry
  - implement what makes sense
  - when it makes sense

# Airflow 2 weaknesses

- Rooted in ETL
- Architecture UI not fully modernized (limits)
- Weak AI workflow support
- Weak DAG Change Management
- Deployment complexity (dependencies)
- Steep(ish) learning curve

# Why now?

- Focus on ML/AI workflows
- Users have new, different needs
- User skills are proliferating
- Python is winning - but is not the only in Data
- Technical debt piling
- Importance of streaming and inference

# Basic principles

# What to expect?

- Airflow 2 features continue to work (exceptions)
- Smooth migration for existing users is top priority
- AI/ML/Data workflows as first class citizens
- Enhance scalability, performance, enterprise-level security through architectural changes
- Make Airflow Asset and Partition native
- Make Airflow easier to deploy and learn

UI

# Modern Web Application/API

## (AIP-38, AIP-84, AIP-65)

OSA CON | November 19-21, 2024

OSA CON | November 19-21, 2024

# Architecture

# Task isolation

(AIP-72, AIP-82, AIP-83, AIP-69)

Distributed Airflow Architecture

# Task SDK

- No direct DB access from task

- Small set of dependencies

- Non-Python native support

- Remote task execution

# Ad-hoc / Fast tasks

- No schedule involved

- Low overhead for  tasks

- Inference ready

- "Almost" streaming

- External-event driven scheduling

# Edge execution

- Agent execution

- Very low overhead

- Enterprise-security

- Multiple Platform support

- Multi-language support

# Dag Authoring

# Data Assets and Partitions
## (AIP-74, AIP-75, AIP-76)

# Dataset changes to Asset

```python
# The "uri" argument is optional but strongly encouraged.
# The asset's name defaults to the function name if none is explicitly given.
@asset(uri="s3://aws_conn_id@bucket/raw_bus_trips.parquet")
def raw_bus_trips():
    # Write bus trips data to asset...
```
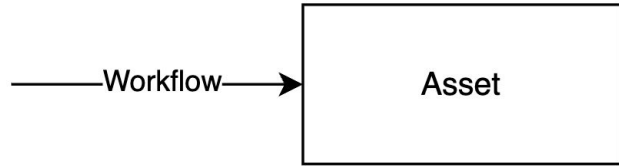
```
@asset(..., schedule="@yearly")
def asset1():  # Write happens when a new calendar year is entered.
    ...


@asset(..., schedule=asset1)
def asset2():  # Write happens whenever asset1's write completes.
    ...

@asset(..., schedule=asset1)
def asset3():  # Runs parallel to asset 2.
    ...


@asset(..., schedule=asset2 | asset3)
def asset4():  # Write happens whenever EITHER asset2 or asset3 finishes.
    ...
```
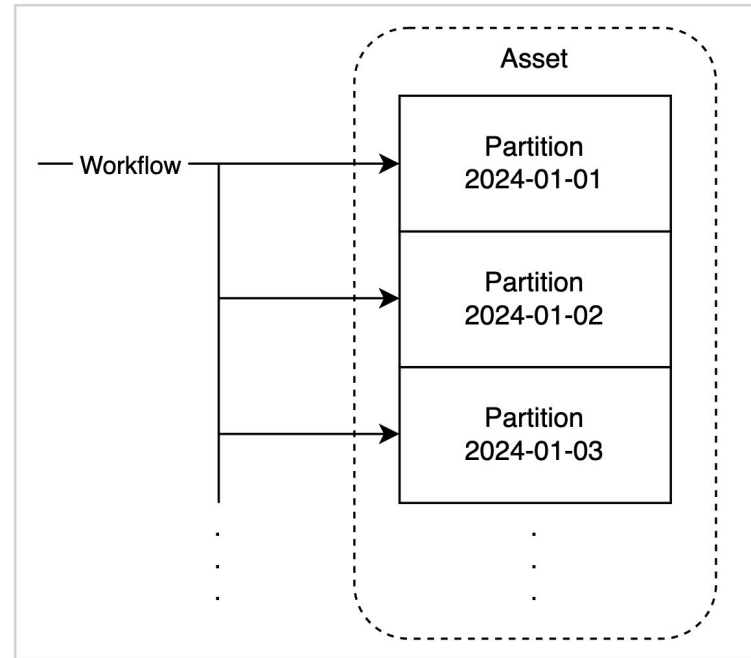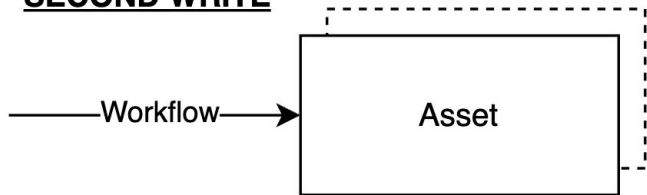
Easy Deployment / Running

# Easy deployment / running

- Standalone Airflow

- Improve messaging and debugging (Survey !!!)

- Separate dependencies for Core/Providers/SDK/UI

- Decouple Authentication and Authorization

- Enterprise-grade deployment options

# Security

# SLSA Threat model

# More / Future

- React UI plugins
- UI-controlled backfills
- Explicit Template Fields
- Lineage auto-hook instrumentation
- Enhanced CLI Security
- Multi-team
- Asset auto read/write

# Airflow is more active than ever



November 7, 2024 – November 14, 2024 | Period: 1 week

**Overview**

187 Active pull requests | 86 Active issues

| 138 Merged pull requests | 49 Open pull requests | 37 Closed issues | 49 New issues |

Excluding merges, **60 authors** have pushed **132 commits** to main and **182 commits** to all branches. On main, **685 files** have changed and there have been **22,045 additions** and **8,647 deletions**.

Used by 12.9k
+ 12,918

Contributors 3,135

+ 3,121 contributors

How users are affected?

When ?

Early
Q2 2025