# Astronomer

**The driving force behind Apache Airflow**

5 offices | 249 employees | 24×7 worldwide support

## 100%
Drives 100% of Airflow releases

## 55%
Of Airflow code contributed

## 18 of 25
18 of the top 25 committers on board, 8 PMC members

## 30K+
30K+ Airflow students in Academy ecosystem

ASTRONOMER

# Data Team

- Centralized Data Team
- Building critical operational and analytical pipelines with Airflow
- Product Influencers
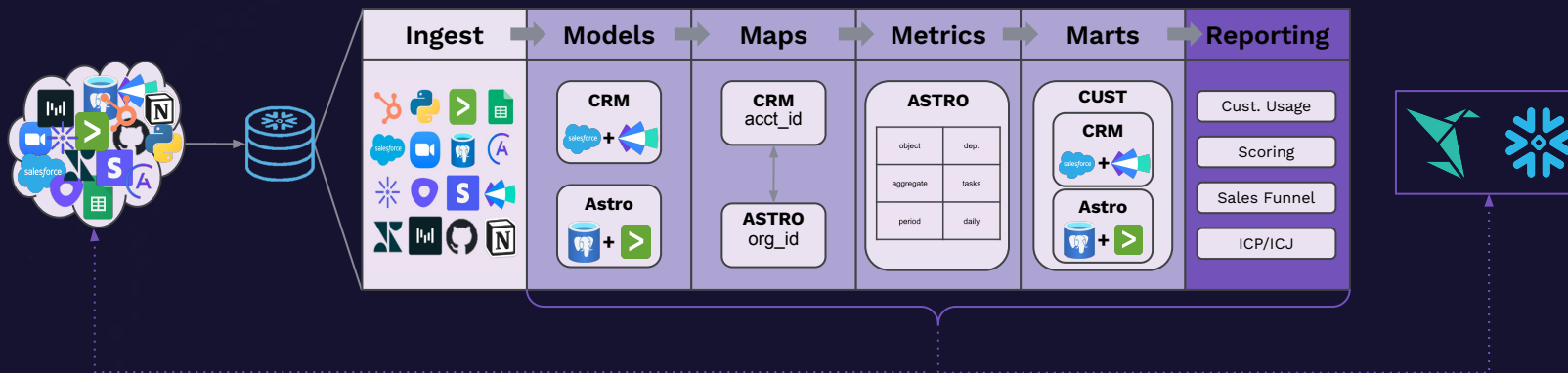  - Providing a Data-First Perspective

# Data Team Ecosystem

Our day-to-day work is standard ELT

**Ingest External Sources**

**Transform and Propagate**

**Deliver Data and Insights**



| Ingest | Models | Maps | Metrics | Marts | Reporting |
|---|---|---|---|---|---|

**CRM**

**CRM**
acct_id

**ASTRO**

| object | dep. |
|---|---|
| aggregate | tasks |
| period | daily |

**CUST**
**CRM**

Cust. Usage

Scoring

Sales Funnel

ICP/ICJ

**Astro**

**ASTRO**
org_id

**Astro**

Flow    Exposure

ASTRONOMER

| | | |
|---|---|---|
| **Application Integrations** | **Outbound Marketing** | **Account Scoring** |
| **Centralized Analytics** | **Billing** | **Embedded Dashboards** |
| **Sales Alerting** | **Warehouse Governance** | **rETL** |
| **Cost Controls** | **Ad Hoc Analysis** | **Product Testing** |

**To name a few...**

ASTRONOMER

# Our Initial Architecture

Focused on governance and onboarding.

### A DAG Factory

Quickly stand up pipelines.

- Abstracts Airflow
- Remain in a familiar context
  - Generate tasks from SQL, R, YAML, Notebooks, etc.

### Custom TaskGroups

Standardizes pipeline operations.

- Prioritizes code reusability
- Contract mechanism for production suitability

### Sensors

Manage cross DAG dependencies.

- Asynchronous when possible

# Configuration as Code

```
/*
operator: include.task_groups.transform.CreateTable
description: >-
    Creates standard calendar, where each row
    represents a date with common date operations.
fields:
    date: Date (iso format yyyy-mm-dd).
    is_weekend: Date is Sat. or Sun.
    is_weekday: Date is weekday (MTWTF).
schema: !switch_value
    sandbox: env__sandbox_schema
    default: commons
primary_key:
    - date
tests:
    check_null:
        - is_weekend
validations:
    check_condition:
        - is_weekend != is_weekday
*/

SELECT
    gd.date::DATE AS date,
    IFF(gd.dow IN (0, 6), TRUE, FALSE) AS is_weekend,
    (NOT is_weekend) AS is_weekday
FROM generated_dates AS gd
ORDER BY date ASC
```

calendar.sql

**1 file = n tasks!**

Create a table; get
documentation and
testing come free.

calendar
- 7 tasks

validate
+ 2 tasks

drop_tmp
■ success
SnowflakeOperator

test_tmp
■ success
+ 1 tasks

swap
■ success
SnowflakeOperator

done
■ success
EmptyOperator

create_tmp
■ success
SnowflakeOperator

add_docs
■ success
+ 1 tasks

**And this was great!**

Successful Tasks by Week

26k Daily Tasks

# Until it wasn't...

# Until it wasn't...

**Marion Azoulai**
Pipeline Update
Oh no... Thankfully it looks like legit failures only, affecting both dev and prod.

**Marion Azoulai**
Daily Pipeline Update
Pipeline is a hot mess today!!
Dev:
- The GCS connexion broke following #1705 and as a result everything is red. Looks like there's an error in the constant

**Marion Azoulai**
Pipeline Update 🚰🚰 🐱
Oh n... !! Thankfully it looks like legit failures only, affecting both dev and prod.

**Marion Azoulai**
Daily Pipeline Update
Pipeline is a hot mess today!!
Dev:
- The GCS connexion broke following #1705 and as a result everything is red. Looks like...

**Marion Azoulai**
🐱🚰 Daily Pipeline Update 🚰🐱
💩💩💩💩💩💩💩💩
PROD: ⚠️ still behind (mart_cust + reporting)
- `snapshot_metrics_finance` failed due to wait sensor failing without logs (looks like our executor failure). As a result downstream failures prevented dailies from completing.

**Marion Azoulai**
🟢🟢🟩🐱🐱🐱🟩💩💩💩🎉🐱💩💩💥💩🔥
Pipeline Update 🟩 🐱
Oh n... !! Thankfully it looks like legit failures only, affecting both dev and prod.

**Marion Azoulai**
Daily Pipeline Update
Pipeline is a hot mess today!!
Dev:
• The GCS connexion broke followi...

**Marion Azoulai**
Daily Pipeline Update 🟩🐱🐶🐶🐶
PROD:
• ...
• Many task sensors timed out during daily run due to `metrics_finance` and `product_astro` dags taking too long to ...
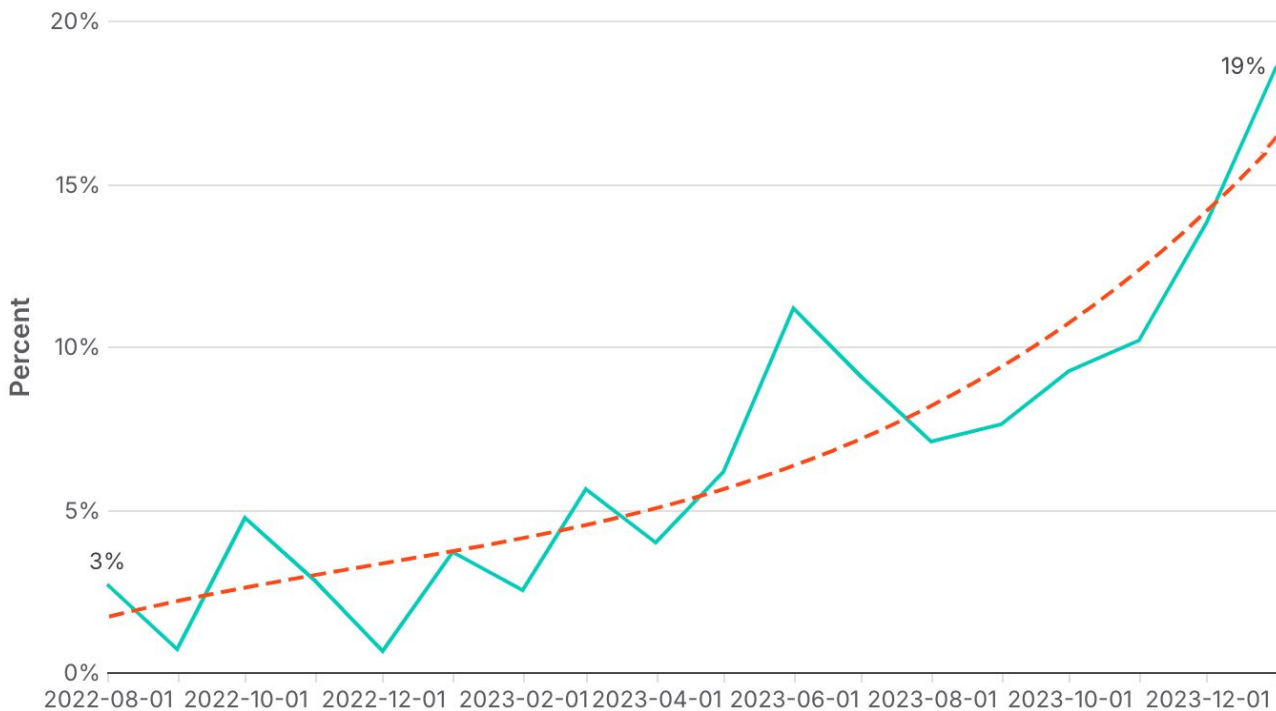...s red. Looks like ...

**Marion Azoulai**
🟩🐱 Daily Pipeline Update 🟩🐱
💩💩💩💩💩💩💩💩💩
PROD: ⚠️ still behind (mart_cust + reporting)
• `snapshot_metrics_finance` failed due to wait sensor failing without logs (looks like our executor failure). As a result downstream failures prevented dailies from completing.

ASTRONOMER

Daily DAG failure Rate

# Pain Points

The challenges of a sensor-dependent setup

Over-reliance on Task and DAG Sensors

# Pain Points

The challenges of a sensor-dependent setup

Over-reliance on Task and DAG Sensors

Failure Cascades

# Pain Points

The challenges of a sensor-dependent setup

**Over-reliance on Task and DAG Sensors**

**Failure Cascades**

**Complex Error Diagnosis**

# Pain Points

The challenges of a sensor-dependent setup

**Over-reliance on Task and DAG Sensors**

**Failure Cascades**

**Complex Error Diagnosis**

**Recovery Delays**

# Pain Points

The challenges of a sensor-dependent setup

**Over-reliance on Task and DAG Sensors**

**Failure Cascades**

**Complex Error Diagnosis**

**Recovery Delays**

**Hitting System Limits**

# Pain Points

The challenges of a sensor-dependent setup

**Over-reliance on Task and DAG Sensors**

**Failure Cascades**

**Complex Error Diagnosis**

**Recovery Delays**

**Hitting System Limits**

**Excessive Compute**

# Pain Points

The challenges of a sensor-dependent setup

**Over-reliance on Task and DAG Sensors**

**Failure Cascades**

**Complex Error Diagnosis**

**Recovery Delays**

**Hitting System Limits**

**Excessive Compute**

**Scaling Challenges**

# Pain Points

The challenges of a sensor-dependent setup

Over-reliance on Task and DAG Sensors

Failure Cascades

Complex Error Diagnosis

Recovery Delays

Hitting System Limits

Excessive Compute

Scaling Challenges

Development Challenges

# Brainstorming the Solution

Focused on scalability and reliability

How do we solve our
sensor problem?

# Brainstorming the Solution

Focused on scalability and reliability

How do we solve our
sensor problem?


**Datasets**

# Brainstorming the Solution

Focused on scalability and reliability

How do we solve our sensor problem?

**Datasets**

How do we solve the visibility issue of datasets?

# Brainstorming the Solution

Focused on scalability and reliability

How do we solve our sensor problem?

**Datasets**

How do we solve the visibility issue of datasets?

**A Control DAG**

ASTRONOMER

# Brainstorming the Solution

Focused on scalability and reliability

How do we solve our sensor problem?

**Datasets**

How do we solve the visibility issue of datasets?

**A Control DAG**

How do we build a Control DAG that's scalable?

# Brainstorming the Solution

Focused on scalability and reliability
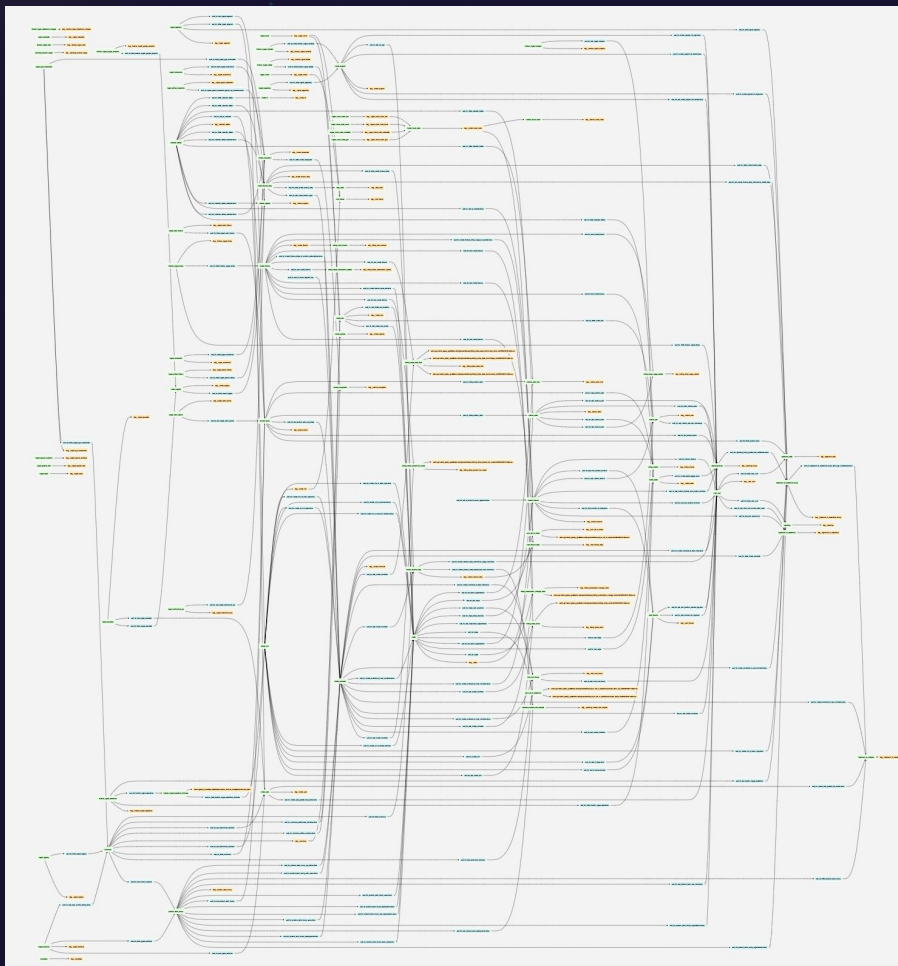
How do we solve our sensor problem?

**Datasets**

How do we solve the visibility issue of datasets?

**A Control DAG**

How do we build a Control DAG that's scalable?

**Airflow**

# But... make it functional

# Pain Points

The challenges of a sensor-dependent setup

| | | | |
|---|---|---|---|
| **Over-reliance on Task and DAG Sensors** | **Failure Cascades** | **Complex Error Diagnosis** | **Recovery Delays** |
| **Hitting System Limits** | **Excessive Compute** | **Scaling Challenges** | **Development Challenges** |

# Pain Points

The challenges of a sensor-dependent setup

| | | | |
|---|---|---|---|
| Over-reliance on Task and DAG Sensors | Failure Cascades | Complex Error Diagnosis | Recovery Delays |
| Hitting System Limits | Excessive Compute | Scaling Challenges | Development Challenges |

ASTRONOMER

# Our Re-Architecture

Focused on scalability and reliability

Dataset Scheduling

=

Increased Reliability

ASTRONOMER

# Pain Points

The challenges of a sensor-dependent setup

| | | | |
|---|---|---|---|
| Over-reliance on Task and DAG Sensors | Failure Cascades | Complex Error Diagnosis | Recovery Delays |
| Hitting System Limits | Excessive Compute | Scaling Challenges | Development Challenges |

ASTRONOMER

# Our Re-Architecture

Focused on scalability and reliability

**Dataset Scheduling**

**=**

**Increased Reliability**
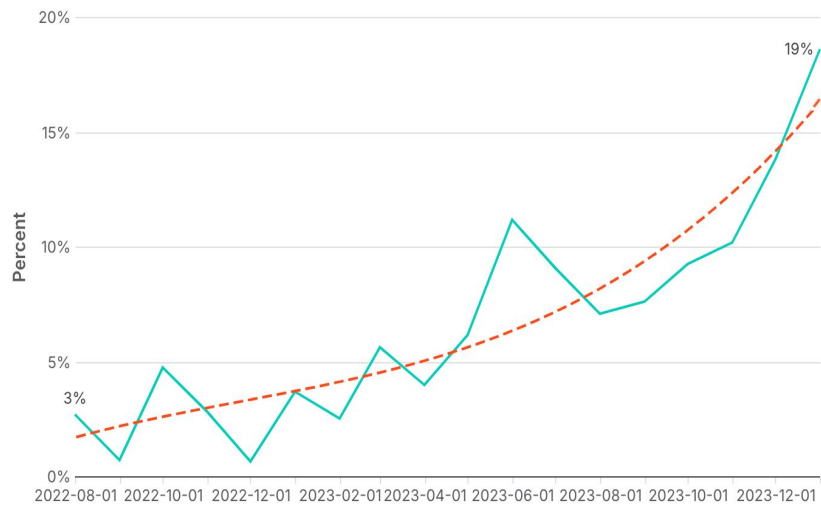
**End to End Visibility**

**=**

**Increased Confidence**

# Pain Points

The challenges of a sensor-dependent setup

| | | | |
|---|---|---|---|
| Over-reliance on Task and DAG Sensors | Failure Cascades | Complex Error Diagnosis | Recovery Delays |
| Hitting System Limits | Excessive Compute | Scaling Challenges | Development Challenges |

# Our Re-Architecture

Focused on scalability and reliability

| Dataset Scheduling | End to End Visibility | Micro-Pipelines |
|:---:|:---:|:---:|
| = | = | = |
| Increased Reliability | Increased Confidence | Failure Minimization |

ASTRONOMER

Daily DAG failure Rate

Daily DAG failure Rate

Daily DAG failure Rate

↓37.9%

Avg. Hourly Run Duration

8.4 min improvement

ASTRONOMER

# Pop Quiz

# What is the most commonly used operator?

**A:** BashOperator

**B:** SnowflakeOperator

**C:** PythonOperator

**D:** LazyAutomationOperator

# Most Commonly Used Operators

Last 30 days

| | |
|---|---|
| PythonOperator | 21% |
| _PythonDecoratedOperator | 12% |
| BigQueryInsertJobOperator | 11% |
| BranchPythonOperator | 3% |
| SnowflakeOperator | 3% |

# What is the most failure prone operator?

**A:** BigQueryCheckOperator

**B:** EmptyOperator

**C:** SSHOperator

**D:** PythonSensor

# Most Commonly Failing Operators

Last 30 days

| Operator | Successful | Failed |
|----------|-----------|--------|
| BigQueryCheckOperator | 80% | 20% |
| DbtCloudRunJobOperator | 83% | 17% |
| SSHOperator | 91% | 9% |
| PythonSensor | 93% | 7% |
| SqlSensor | 94% | 6% |
| DatabricksRunNowOperator | 94% | 6% |
| DatabricksSubmitRunOperator | 96% | 4% |
| _PythonDecoratedOperator | 96% | 4% |
| BashOperator | 97% | 3% |
| MsSqlOperator | | 2% |

■ Failed Tasks  ■ Successful Tasks

ASTRONOMER

## Longest Running Operators

Median task duration in minutes in the last 30 days

| Operator | Duration |
|---|---|
| DatabricksSubmitRunOperato… | 8.15 |
| DatabricksSubmitRunDeferra… | 8.12 |
| _PythonVirtualenvDecorated… | 7.43 |
| DatabricksRunNowOperator | 6.55 |
| DatabricksSubmitRunOperator | 5.53 |
| DbtCloudRunJobOperatorAsy… | 5.18 |
| BatchOperator | 3.78 |
| DataprocSubmitPySparkJobO… | 2.78 |
| EcsRunTaskOperator | 2.63 |
| DbtCloudRunJobOperator | 2.45 |

How many tasks in a DAG?
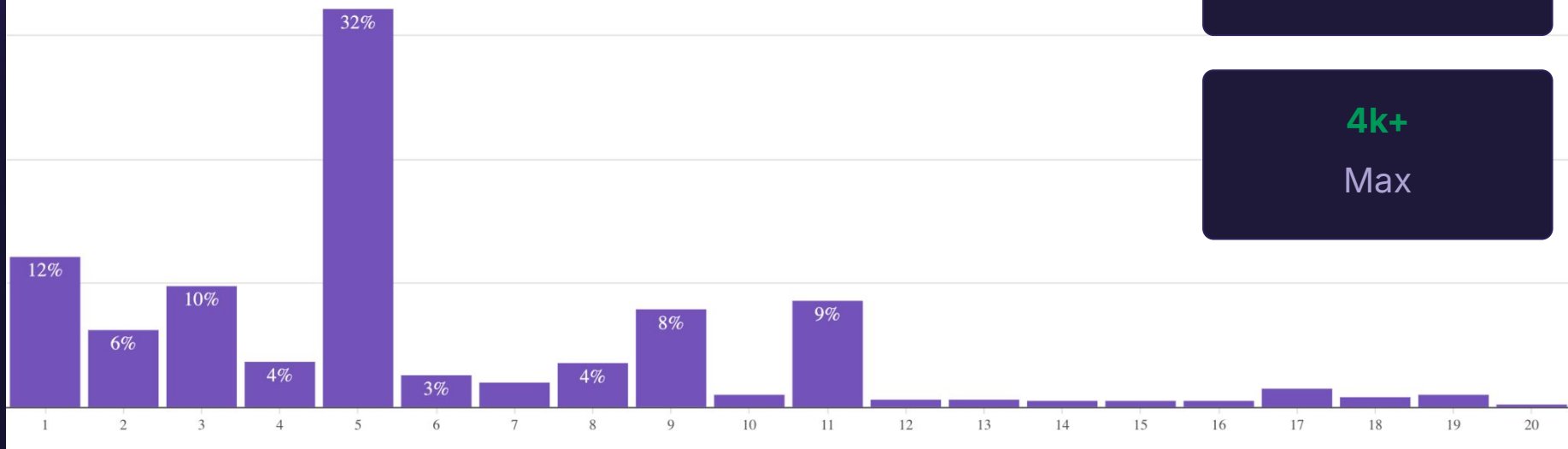
How many tasks does a DAG typically have?

A: 1

B: 5

C: 10

D: 20

DAGs by Unique Task Count

90+%
<= 20 tasks

4k+
Max

Which hour (UTC) are daily DAGs most often scheduled?

A: 0

B: 8

C: 12

D: 20

# DAGs by Scheduled Hour

Daily scheduled dags in the past 30 days

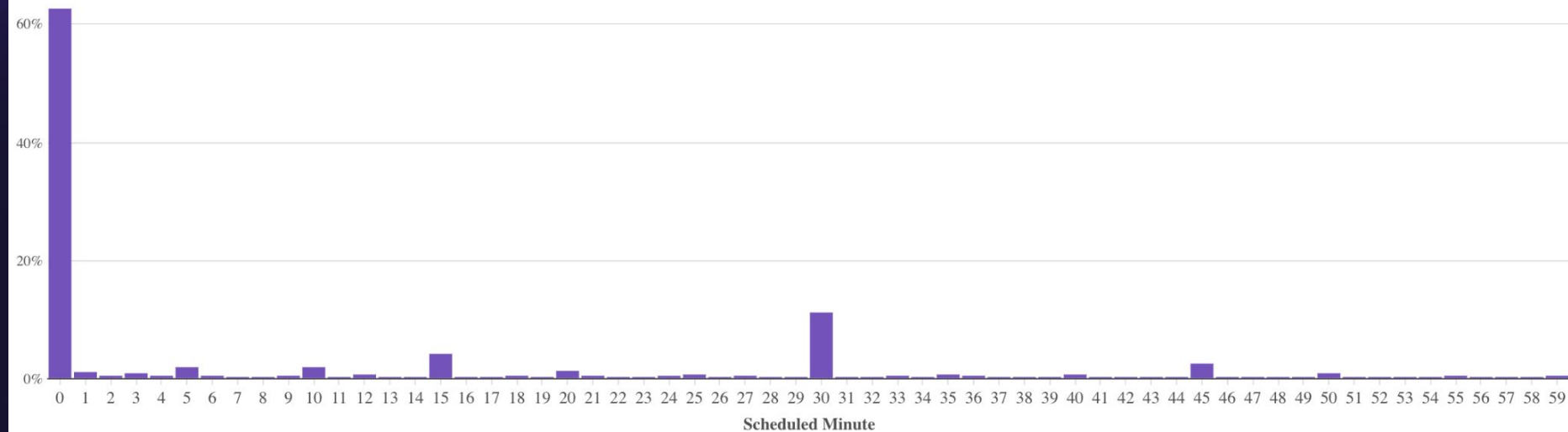**DAGs by Scheduled Minute**

Daily scheduled dags in the past 30 days

# Thank you!
# Any questions?

ASTRONOMER