# Apache Airflow

For Data Engineering
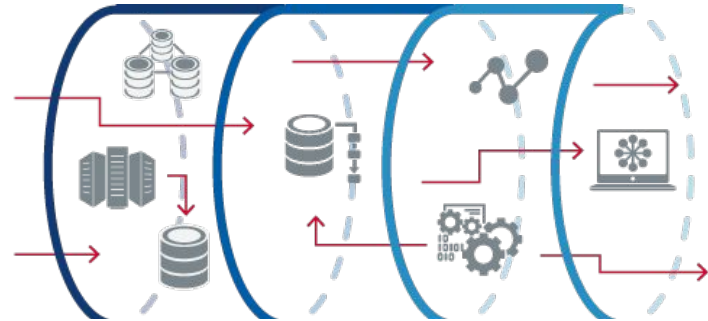
- Sampath Kumar, Principle Engineer

# Agenda

Goal of this session is to provide an overview of Airflow capabilities & how to use it for Data Engineering

- Introduction (DE, Airflow)
- Key Concepts (Key Words, Demo)
- Architecture (System Design)
- Features (Challenges, Loggins, Analytics)
- Challenges & Recommendations
- Q & A

PRAMATI

TECHTONIC
upshift your technology quotient

# DATA ENGINEERING

**Data Engineering** is the **aspect of data science** that focuses on **practical applications of data collection and analysis**. **Data Engineers** are tasked with

- **Designing, building, testing, integrating, managing, and optimizing data** from a variety of sources
- **Build the infrastructure** and architecture that enable data generation
- Primary focus is to build **free-flowing data pipelines** by combining a variety of big data technologies that enable real-time analytics
- Data engineers also write complex queries to ensure that data is easily accessible
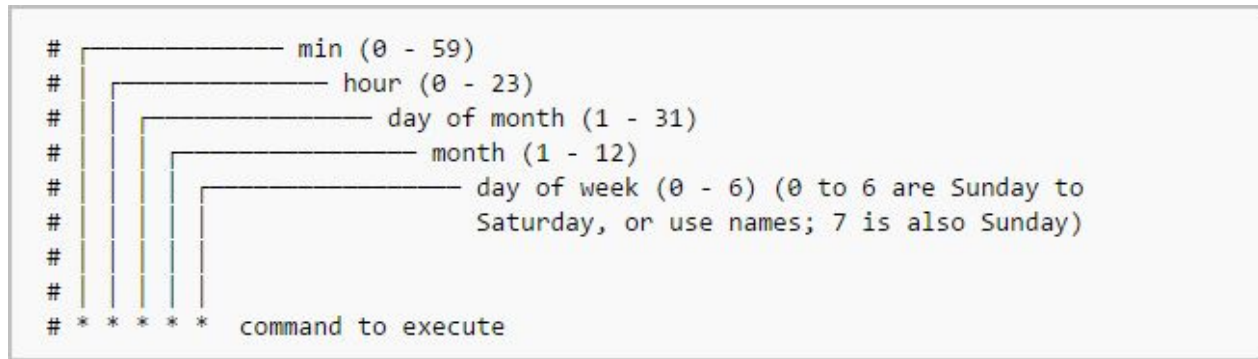
# AIRFLOW

- Open-source workflow automation and scheduling system that can be used to author and manage your data pipelines.
- It started at Airbnb in October 2014 as a solution to manage the company's increasing complex workflows.
- License: Apache License 2.0
- Written in: Python
- Operating system: Microsoft Windows, macOS, Linux
- Stable release: 1.10.5 / August 30, 2019; 3 months ago

# CRON

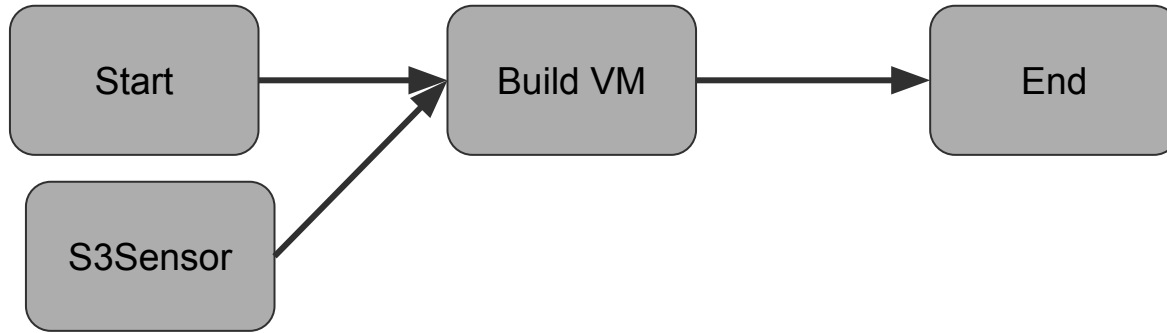CRON: Derived from work CRONOS(means time), is a software for unix-like systems to schedule jobs based on time.

```
#  ┌─────────────── min (0 - 59)
#  │ ┌───────────── hour (0 - 23)
#  │ │ ┌─────────── day of month (1 - 31)
#  │ │ │ ┌───────── month (1 - 12)
#  │ │ │ │ ┌─────── day of week (0 - 6) (0 to 6 are Sunday to
#  │ │ │ │ │         Saturday, or use names; 7 is also Sunday)
#  │ │ │ │ │
#  │ │ │ │ │
#  * * * * *  command to execute
```

Airflow is cron on steroids: it allows you to schedule tasks to run, run them in a particular order, and monitor / manage all of your tasks.

PRAMATI

TECHTONIC
upshift your technology quotient

# KEY CONCEPTS

- DAGS
  - Directed acyclic graphs that represent tasks workflow
- Task
  - Operators - Bash, Python, SSH, Http, MySql, SparkSubmit, Sensors(s3), Docker, Hive, Slack,..
- Hooks
  - To store credentials for services - AWS, GCP, DataBase, Email,..
- Vars & XCom
  - For sharing any global values or inter-task communication

PRAMATI

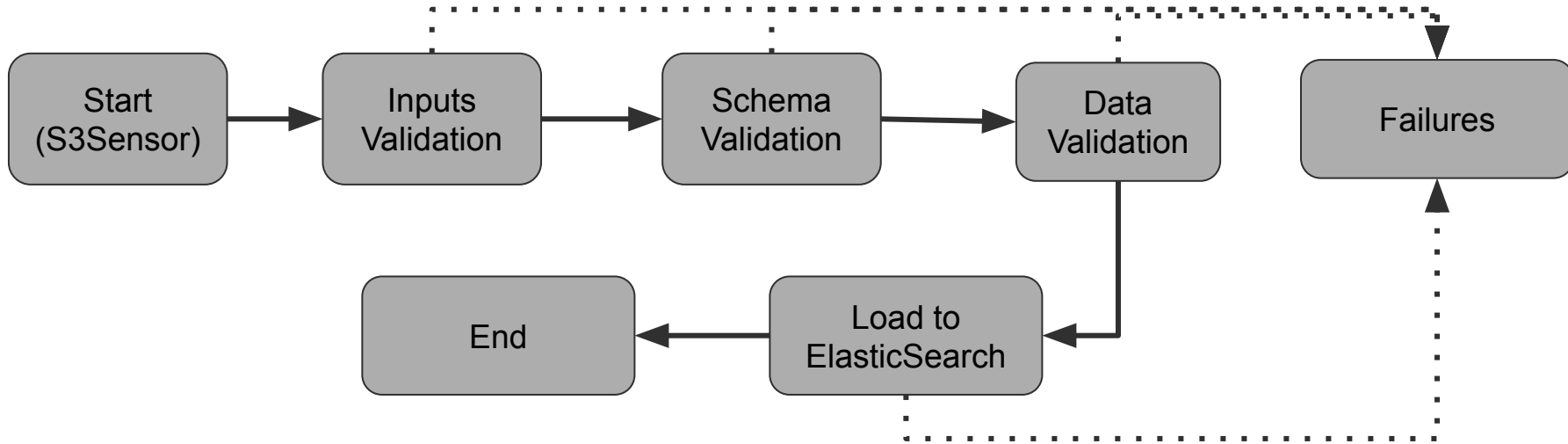TECHTONIC
upshift your technology quotient

# Sample - DAG

```
┌─────────────┐          ┌─────────────┐          ┌─────────────┐
│             │          │             │          │             │
│    Start    │───────▶  │   Build VM  │───────▶  │     End     │
│             │      ▲   │             │          │             │
└─────────────┘      │   └─────────────┘          └─────────────┘
                     │
┌─────────────┐      │
│             │      │
│  S3Sensor   │──────┘
│             │
└─────────────┘
```

**Goal:** As and when user files come into AWS S3 bucket, start a high-spec docker VM and process the job. Prior to start and post process, job status to be notified.

Solutions required to be cost efficient and need to auto-scale if required.

PRAMATI

TECHTONIC
upshift your technology quotient

# DEMO - UI

# Sample - DAG

```
Start          Inputs          Schema          Data
(S3Sensor)  →  Validation  →  Validation  →  Validation          Failures

                                              ↓                      ↑
              End      ←    Load to      ←
                            ElasticSearch
```

A data engineering pipeline, where an S3 sensor is used to identify the arrival of input file, following by several validation checks and then load into ElasticSearch which will be used for serving clients.

PRAMATI

TECHTONIC
upshift your technology quotient

# ARCHITECTURE

- MetaDB
- Message Broker
- Airflow Webserver
- Airflow Scheduler
- Airflow Workers

# FEATURES

- Airflow Workers are Horizontally scalable
- Airflow Messaging Broker - Celery
- Airflow Integrations - GCP, Azure, AWS, Qubole & Databricks
- Hooks, Connections & Pools - Environment(dev/test/prod) friendly
- DAG - Dynamic sub dags & Branching

PRAMATI

TECHTONIC
upshift your technology quotient

# ANALYTICS

- Part of being productive with data is having the right weapons to profile the data you are working with.
- Airflow provides a simple query interface to write SQL and get results quickly, and a charting application letting you visualize data.

# LOGGING

- Logging is visible from UI

# CHALLENGES

- Airflow is not Apache NIFI or Apache Spark
  - is not - data routing or data transformation system.
- Airflow Workers
  - requires identical access(network, authentication & authorisation)
  - requires similar hardware capabilities
- Airflow Webserver & Scheduler
  - not scalable (Work-around - supervisor, docker health checks)

PRAMATI

TECHTONIC
upshift your technology quotient

# RECOMMENDATIONS

Airflow is a best suited where

- Agility is important
- Portability
- Segregation b/w compute and workflow mgmt
- Scalability on demand
- Pool/Connection management
- Job Analytics
- Logging visibility

PRAMATI

TECHTONIC
upshift your technology quotient

# Q&A

Open for discussions,..

PRAMATI

TECH TONIC
upshift your technology quotient

# REFERENCES

- Roles of Data Engineer - https://inlovewithcode.wordpress.com/2019/05/15/roles-of-data-engineer-required-skill/
- System Design - https://inlovewithcode.wordpress.com/system-design/
- Airflow Tutorials - https://airflow-tutorial.readthedocs.io/en/latest/airflow-intro.html
- Airflow Documentation - https://airflow.apache.org/docs/stable/
- Airflow Dockers - https://towardsdatascience.com/getting-started-with-airflow-using-docker-cd8b44dbff98

PRAMATI

TECHTONIC
upshift your technology quotient

Thank you

PRAMATI

TECHTONIC
upshift your technology quotient