# Orchestrating & Optimizing a Batch Ingestion Data Platform for Americas #1 Sportsbook

## Gunnar Lykins

# AGENDA

FANDUEL GROUP

# Introduction to FanDuel Group

# Introduction to FanDuel Group

## Overview of FanDuel Group

- Innovative sports-tech entertainment company

- Diverse portfolio – gaming, sports betting, daily fantasy sports, advance deposit wager, TV/media

## Industry Leadership

- As of 2023…
  - #1 sports betting company in US
  - Fastest growing operator in iGaming
  - First US online operator to turn a profit for a full year

## Market Presence

- Operating in all 50 states

- Serving ~17 million customers

- Nearly 30 retail locations

## Workforce

- Over 4,000 employees

- 1,600+ in Technology

- ~100 in Data Engineering

# Batch Ingestion Data Platform Evolution

**2013 - Luigi**

Onset of development for batch pipelines

**2018 - Erie**

Libraries - code reusability

Alembic - data warehouse migrations are "self-served" Infrastructure-as-Code

Python 3

Airflow - improved scalability

**2019 - Automata**

Erie lacked standardization
Code reusability - enforced consistency

TOML files vs. Python - faster scalability of business

**2022 - Automata v2**

All Airflow deployments managed by Astronomer
Improved consistency in meeting internal SLAs for business-critical pipelines

**2023 - Automata v3**

Cost optimizations - separate deployments per business vertical
Increased data governance - serve regulatory compliance reporting & daily business reports

Present

# By the numbers…

Astronomer Workspaces

**3**

Astronomer Deployments

**17**

Monthly DAG Runs

**350K+**

Monthly Tasks Executed

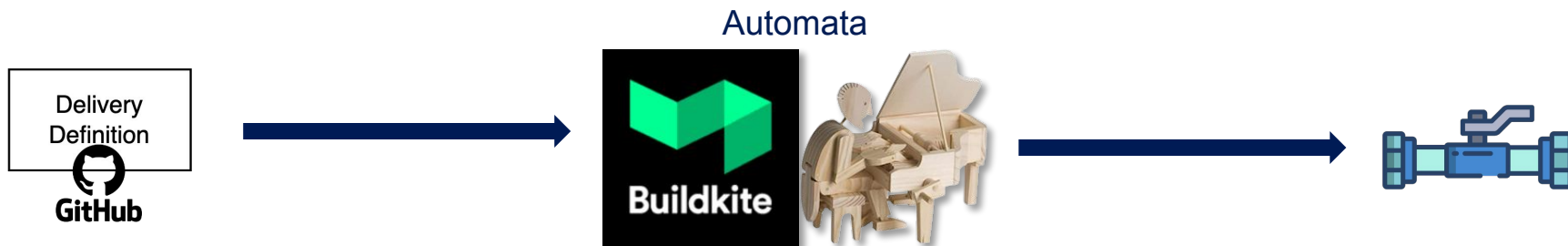**3.6 M**

# Creating & Managing DAGs

# Automata

- Greek for "self-acting, self-willed, self-moving"

- Internal name for FanDuel's batch ingestion data platform

- Orchestrator that serves both data and product needs – testament to the extensibility and generality of the batch ingestion platform
  - _Data_: ingest data from a source → perform data transformations → move data to storage
  - _Product_: perform custom-built, self-served steps as specified by a user

- Automated data processing managed by an internal team so stakeholders can place more efforts on business deliverables

# What has made it successful?

# Self-serving data platform with standardized way of pipeline development

Automata



"The **goal** of software architecture is to **minimize** the **human resources** required to build and maintain the required system."

# Delivery Definition Example

- Sources:
  - JDBC, file sources (.csv, .zip, .parquet, etc.), Kafka, Redshift, SFTP

- Destinations:
  - Delta Lake, FTP, Lake, Redshift, S3

- Cleansing/Light Transformations:
  - Time zone conversions, data type standardizations, permission settings (PII)



```
delivery_type = "file_to_lake_fullload"


owners = ["dummy-email@fanduel.com"]
description = "Sample .TOML for Airflow Summit 2024"

schedule_interval = "0 8 * * *"
source_timezone = "America/Los_Angeles"

[source]
name = "a_source_name"
bucket = "an-s3-bucket"
prefix = "data/DailyFiles"
fileregex = "fileregex-to-match"
file_type = "csv"
separator = ","

[[source.columns]]
name = "_c1"
type = "string"

[[source.columns]]
name = "_c2"
type = "timestamp"

[[steps]]
type = "transformation"
operations = "lower_case_string,empty_string_to_null"

[destination]
area = "formatted"
prefix = "finalData"
```

# Which scenarios should the platform self-serve?

# 80/20 Principle

- Cover **most-common** data pipeline requirements

- Give the ability to users to **"hook up" custom code** in a <u>self-serving fashion</u>

```
delivery_type = "custom_to_lake_fullload"

owners = ["dummy-email@fanduel.com"]
description = "This is an example of a custom source

schedule_interval = "0 8 * * *"
source_timezone = "America/Los_Angeles"

[source]
source_type = "a_user_created_exceptional_source"
custom_argument_1 = "..."
custom_argument_2 = "..."
custom_argument_3 = "..."

[destination]
area = "formatted"
prefix = "my_destination_prefix"
```
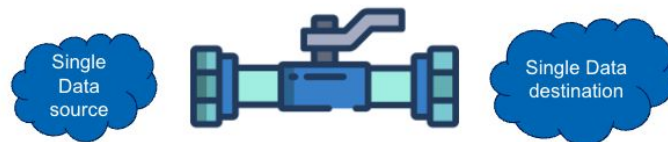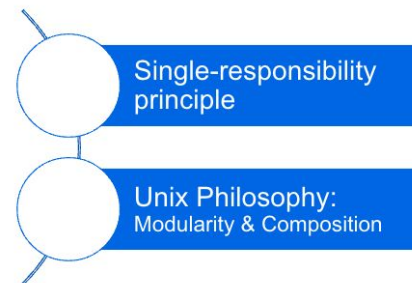
# Enforce consistency in the design of data pipelines

- "A data pipeline should have one, and only one, reason to change."
- "Make each data pipeline do one thing well."
- "Be able to compose more complex pipelines from simpler ones."

Single-responsibility principle

Unix Philosophy: Modularity & Composition

Single Data source

Single Data destination

Automata is a factory of data pipelines:

✔ Code is **re-used** inside the platform
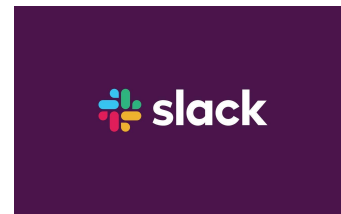
✔ Pipeline-generation code is **tested once and used many times**

# Monitoring and Troubleshooting

FANDUEL GROUP

# Monitoring, Alerting, & Observability Overview

An attractive feature to the batch ingestion data platform is structure monitoring, alerting, & observability

Technologies utilized:

- Terraform
- Slack & PagerDuty integrations
- Datadog
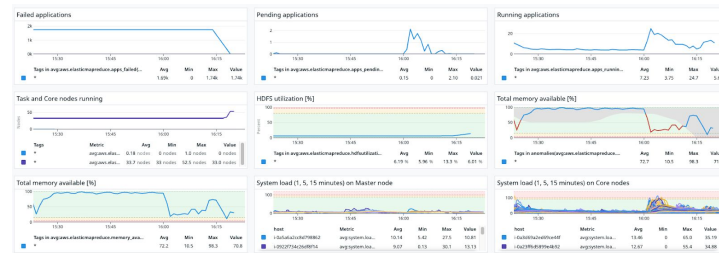- Databand

# Terraform: Infrastructure-as-Code

- Terraform allows the team to **manage infrastructure** in a safe, consistent, and repeatable way by defining resource configurations that can be **versioned and reused**

- Deploy configurations in different environments in a seamless manner

- Utilized to define:

  - Platform infrastructure for cloud compute resources

  - Datadog monitors and dashboards

# Datadog: Platform Compute Resource Monitoring

- Datadog is the primary tool utilized for data platform engineers to monitor various **performance and health metrics on compute resources** such as EMR, RDS, EKS, Airflow deployments, etc.

- Integrated with PagerDuty as well as dedicated alerting Slack channel

- Provides the ability to configure monitors with **adjustable thresholds** for iterative fine-tuning

- Aids in identifying issues before they escalate

- Dashboards provide **observability** benefits that have been particularly useful on **business-critical events** – 100% uptime throughout 2024 Super Bowl
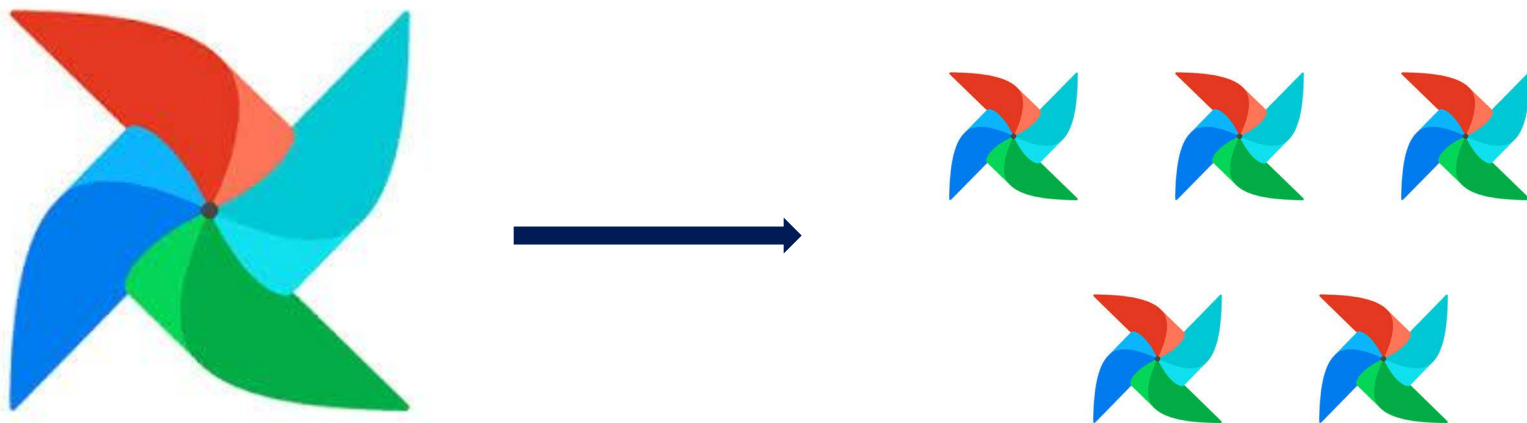
# Databand: Pipeline Specific Monitoring & Alerting

- Databand is the primary mechanism for pipeline owners to **define and customize alerts on a per pipeline and task basis**

- Integrated with PagerDuty as well as dedicated alerting channels within Slack

- Alerts are defined with respect to:

  - Pipeline run and state (running/success/failure)

  - Schema changes (column type change, column added, column removed, etc.)

  - Missing dataset operation (dependent operations to a pipeline aren't operational)

  - Custom task metrics (anomaly detection)

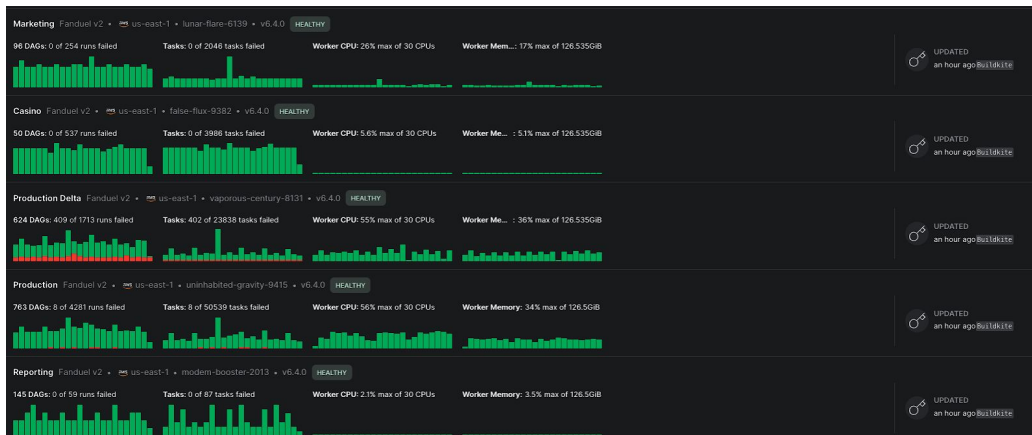# Scaling and Optimization

# Decoupling a Monolithic Deployment

- In 2022, saw over a 3x increase in the number of production pipelines

- Continuing to operate under one production monolithic deployment posed several risks to the **robustness of the platform**

- Set out to create separate deployments segmented by business vertical, which provided several benefits
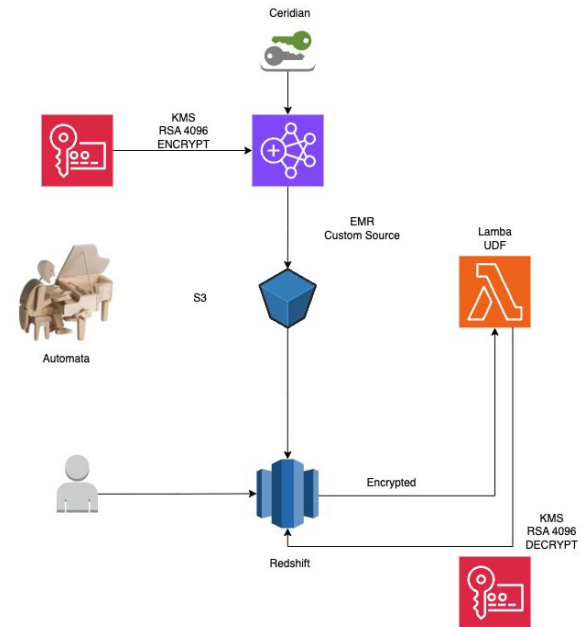
# Decoupling a Monolithic Deployment

✔ Security – restricting access to individuals on a **"need-to-have" basis**

✔ Stability – as the volume of data ingestion increases, the **blast radius is reduced** on production issue

✔ Scalability – **enhanced governance** on Astronomer deployment configuration parameters on a per deployment basis



~$10k/mo savings in EC2 instances

# Additional Optimizations:
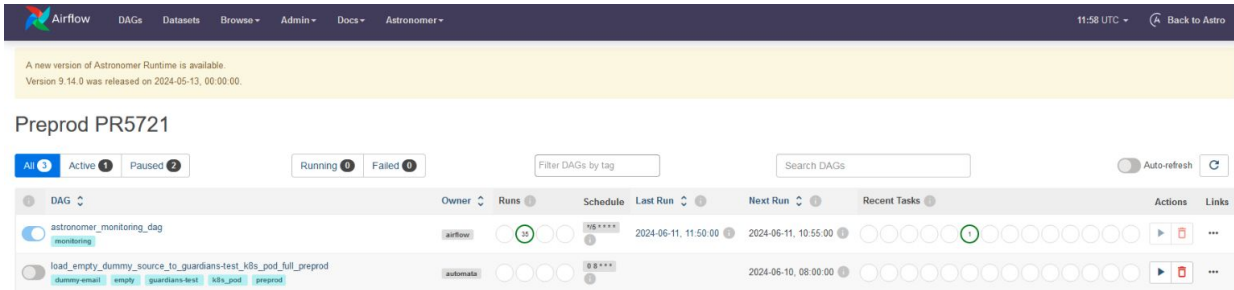# Sensitive PII Data with Encryption Requirements

- Restricted access to only stakeholders

- Separate IAM permissions to connections and secrets per workspace

- Users without access to the deployment will not be able to turn on/off the pipelines



Data remains encrypted throughout the whole process, even when it's copied to Redshift

# Additional Optimizations:
# Auto-Scaling of Development Environments

- Developers can provision dedicated Airflow deployments for testing changes in Pull Requests prior to merging to production

- Integrates with CI/CD steps (Buildkite) – Pull Request and deployment IDs documented in Postgres database

- Maintenance DAG runs to delete deployments and database records that are 8+ hours old

# Future Trends and Considerations

# How else can Automata be leveraged outside data engineering?

Batch ingestion data platform provides value to the organization for being the primary tool for **orchestrating scheduled jobs**

- Scheduling logging & reporting via Buildkite for business analytics & **extrapolating insights**

- Automating process for sending scheduled emails to customers on performed transactions for tax purposes

- Migrating Casino iOS game files and automating the process of delivering them to the iOS App Store

# Conclusion

# Key Principles of Automata

✔ **Self-Service**: allows users to interact with and create data pipelines through writing easily digestible .TOML files

✔ **Standardization**: enhances the reliability and maintainability of the platform as well as improvements to consistently meeting SLAs

✔ **Scalability**: yields substantial performance improvements and cost reductions with multiple workspaces and deployments

✔ **Observability**: enables monitoring for a centralized platform at a granular level with circa 100 engineers actively contributing to it

✔ **Orchestration Diversity**: provides the framework to the creation of data pipelines in a streamlined manner and enlightens other business cases outside of data

# Questions?