ASCEND

Intelligent Orchestration: Data's missing link

Sean Knapp, Founder @ Ascend.io



Topics

- Quick Intro
- The State of Data Architectures
- Why Pipelines Suck
- Building a new Control Plane
- Making it Scale
- SaaS-ifying
- Future Topics



Quick Intro

About Ascend

- Founded 2015
- Team of 30
- We <3 Data Pipelines

About Me (Sean Knapp) 👋

ASCEND.10 COYALA Google

- 15 years building data platforms & teams
- Search Frontend TL @Google: first MapReduce in 2004
- Founder & CTO @Ooyala: 4B+ events/day
- Founder & CEO @Ascend: 1T+ events/day



Smarter Pipelines.





Despite advancements in every other part of the data lifecycle...

Building. Pipelines. SUCKS.

The Current State of Data Architectures

Why Pipelines Suck

Why Pipelines Suck

Databases & Warehouses

Where this...

Pipelines ...becomes 1,000s of lines of this...

Evolution of Pipeline Orchestration

	1.0	2.0	
Hosting Model	Roll-Your-Own	SaaS	
Code Generation	Manual	Templatized	
Interaction Model	Code	Code + GUI	
Control System	Scheduler	Scheduler	
Programming Model	Imperative	Imperative	
Examples		AWS Glue Azure Google	
		Data Factory Data Fusion	

We looked for ideas in adjacent spaces...

Who here has used a **Database**?

Who here has heard of **React**?

Who here uses **Kubernetes**?

What do they all have in common?

They're **Declarative**.

Declarative programming is a programming paradigm that expresses the logic of a computation without describing its control flow...

... [in an] attempt to **minimize or eliminate side** effects by describing what the program must accomplish...

... rather than describe how to accomplish it.

Declarative vs. Imperative

The **desired** outcome

What vs How

Declarative

Imperative

The **usual** outcome

Declarative

Imperative

Declarative vs. Imperative

		Declarative		Imperative
Pros		Less Code (like A LOT less)	\checkmark	Flexible
		Faster Dev Cycles	\checkmark	High Levels of Control
	Ø	Adaptive to Changes		
		Less Maintenance		
Cons	8	Domain specific Difficult to manually optimize Require annotations to override automated behaviors	× × × ×	State Management Stale assumptions in code Manual Optimizations Integrity Checks Failure Management

Imperative gives you the ability to do anything, and the responsibility to do gives you the responsibility to do giverything.

— Steven Parkes, CTO @ Ascend

Building a Control System for Declarative Pipelines

Our master plan...

What should a good control system do?

Then we spent $\frac{1}{2}$, 3 years building it!

So... how does it work?

Separation of Logic & Control

Logic Plane

User defined logic

Control Plane

Dynamic task generation to achieve desired state

Data Plane

Fully managed, portable cloud services

The Control System Answers

- **1** Is there anything I need to do?
- 2 What is the current state of my world?
- **3** What should it be?
- **4** What doesn't match?
- 5 How do I "fix" it?

1) Is there anything I need to do?

A Simple Analytics Pipeline

2) What is the current state of my world?

A Simple Analytics Pipeline

gs://ascend-io-dev-sean-dev-sean-record-fragments/21dd3c96_b039_4cc1_9566_2cfbb3ebe091/ASCEND_METADATA.json gs://ascend-io-dev-sean-dev-sean-record-fragments/21dd3c96_b039_4cc1_9566_2cfbb3ebe091/part-00000000.parquet ...

gs://ascend-io-dev-sean-dev-sean-record-fragments/21dd3c96_b039_4cc1_9566_2cfbb3ebe091/part-00000010.parquet

. . .

Data Plane

3) What should it be?

A Simple Analytics Pipeline

Data Plane ps://ascend-io-dev-sean-dev-sean-record-fragments/21dd3c96_b039_4cc1_9566_2cfbb3ebe091/part-00000000.parque ... ps://ascend-io-dev-sean-dev-sean-record-fragments/21dd3c96_b039_4cc1_9566_2cfbb3ebe091/part-00000010.parque

4) What doesn't match?

A Simple Analytics Pipeline

5) How do I "fix" it?

A Simple Analytics Pipeline

Scaling to **1B Partitions** and **1T+ records per day**

Scaling the Control Plane

- SHAs: lots and lots of SHAs (and SHAs of SHAs)
- Caching: lots and lots of caching
- Trees, not lists
 - Leverage time-series partitions
 - Aggregate metadata
 - SHAs at each node, not just leaves
- **Be Lazy**: Only do as much work as is useful right now

Scaling the Data Plane

• Do less work

- Intermediate Data Persistence
- Data & Task De-duplication
- Do the right kind of work
 - Small file aggregation
 - Small job optimizations (local mode)
 - Specialized compute pools for different tasks

• Do it efficiently

- Auto-Scaling Spark on Kubernetes w/ Spot/Preemptible Instances
- Single-zone clusters (reduce network costs)

SaaS-ifying the Control Plane

What we didn't discuss...

- Garbage Collection: background task
- Multi-cloud abstractions: k8s, MinIO
- **Data repair**: failure to retrieve data \rightarrow delete p-sha \rightarrow self-heal
- Part files: similarities & differences with other fragments
- **Resource Management**: capacity-aware Control Plane
- SLA Driven Scheduling & Job Prioritization: per-component priority + upstream inheritance
- Scaling Spark on Kubernetes
- Scaling to 100+ environments: terraform, templates, monitoring, & automation

Come ask us @ Office Hours!!! (4th floor in 476a) Declarative programming is a paradigm that expresses the logic of a computation without describing its control flow...

> ... [in an] attempt to **minimize or eliminate** side effects by describing what the program must accomplish...

... rather than describe how to accomplish it.

tl;dr

Declarative

The What Logic + Data \rightarrow Tasks

Imperative

The How State + Tasks \rightarrow Data

Smarter Pipelines.

See You in Office Hours

- Right after this, 4th floor, 476a
- Ask me anything
- Meet our CTO, Steven Parkes
- Visit our booth at the Partner
 Spotlight gallery for a live demo & free swag
- Visit us @ www.ascend.io

