# Scaling Data Science Platform as a Service

**Ameya Kanitkar**

Engineering Manager,  Data Science Platform

# What's common amongst these LinkedIn Products?

# This Talk

is about Metrics

- Linkedin has a data driven culture

- We measure and monitor everything

- Metrics and how we compute and serve them is at the heart of our data science platform as a service

# Once Upon a Time

there were great many metrics platforms

Experimentation

Marketing Metrics

Product Metrics

Executive Metrics

Running on my desktop-for-my-team metrics

# Which was bad

Apps　Apps　Apps　Apps

| Logic | Logic | Logic | Logic |

Tracking data, DB etc

- Inconsistent Data

- Lack of trust

- Closed systems with data and knowledge in silos

## Challenge 1: Data Consistency

# Challenge 2: Engineering Productivity & Operations



tracking data

segments

member info

Data Warehouse ~TB

OLAP cube

dashboards

- Identify data sources
- Build data pipelines
- Data Modeling
- Insights & Analysis
- Data Governance and Access management

# Challenges: Engineering Productivity (Contd..)



Build Insights

Data Scientists

Data Engineering

Infrastructure Engineering

Build and operate data pipelines

Build & scale data infrastructure

- Speed can slow down due to multiple teams are involved

- Data scientists productivity declines as they need to be dependent on other teams to build necessary insights

# Challenges: Operations (Contd..)



- Building and maintaining thousands of data pipelines is messy

- Consistently maintaining metadata, data lineage, data dependencies, SLA requires specialized systems and needs to be developed across all teams

# Solution

Our Approach | "Provide trusted repository of metrics, and build a self-serve platform for sustainable life cycle of metrics"

# Wish List

## What did we want?

- Same and consistent metrics/ insights in all data applications
- Same metadata everywhere
- Single definition of entity dimensions
- Single definition of event dimensions
- No duplication of metrics
- Allow data scientists to focus on their core skills/ job
- Increase engineering productivity, simplify and optimize operations

# Solution

## Metrics Platform as Service

"Platform that builds, manages and scales all metrics across all applications at LinkedIn"

# Metrics Platform as a Service



Data Scientists

Metrics Platform as a Service

Metrics

- Data scientists build metrics logic (config + code) into the central platform

- Platform automatically generates necessary data pipelines and centrally operate

- Platform computes and publishes metrics results into dashboard/ downstream apps

# Metrics Platform as a Service

```
metric-name:     daily-unique
input-dataset:   web-tracking, mobile-tracking
metric-formula:  DISTINCT_COUNT (members)
frequency:       daily
dimensions:      members, country, platform
logic:           daily-unique-users.sql
```

Data Scientists

Metrics Platform as a Service

auto-
generates

Metrics

mobile-tracking    web-tracking    dim-member    dim-country

daily-uniniq-users.sql

Platform-aggregation-logic

Hive-registration

metadata-registration

OLAP
cube

# Scaling Governance
# is not Easy

**metric-defs**

Single source/repository for all metrics

**ACL's**

ACL Management by Subject Matter Expertise

Tools for data lineage

Search & Discovery experience for your datasets and metrics

# Lessons along the way

Implementing default use cases should be simple, but more complex use cases possible

1. Added support for various compute engines: Pig, Hive, Spark, Spark SQL, Presto
2. Templatize common patterns such as cubing, multi datacenter availability, percentile support
3. Build downstream integrations
4. Explore plug-in architecture to avoid making platform team a bottleneck

# Lessons along the way

Leverage economies of
scale of being the central
platform

Simplified Compliance
- Platform managed compliance
- Data Retention and cleanup
- Data Quality Checks

Self-Serve & Operational Tooling
- Flow Management
- Backfills Portal
- Self-heal

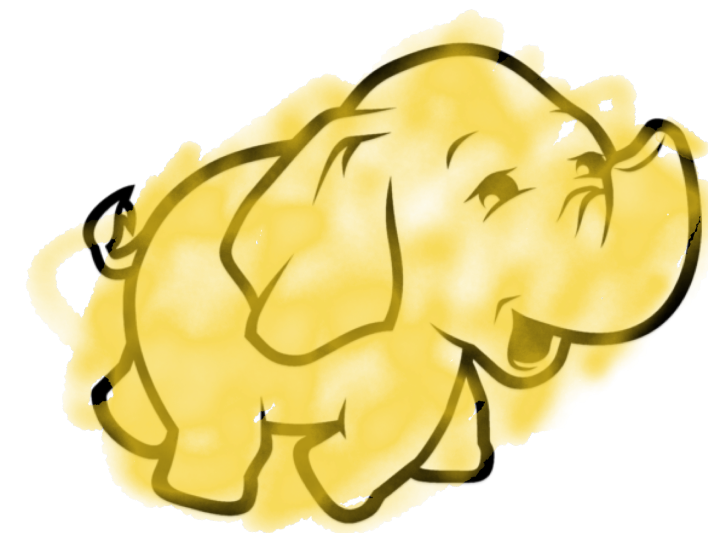# Lessons along the way

Innovations are a lot more impactful with platforms

Batch -> Realtime Convergence

- Automatic Hadoop MR -> Kafka/ Samza convergence via apache calcite

Hadoop
Offline Metrics
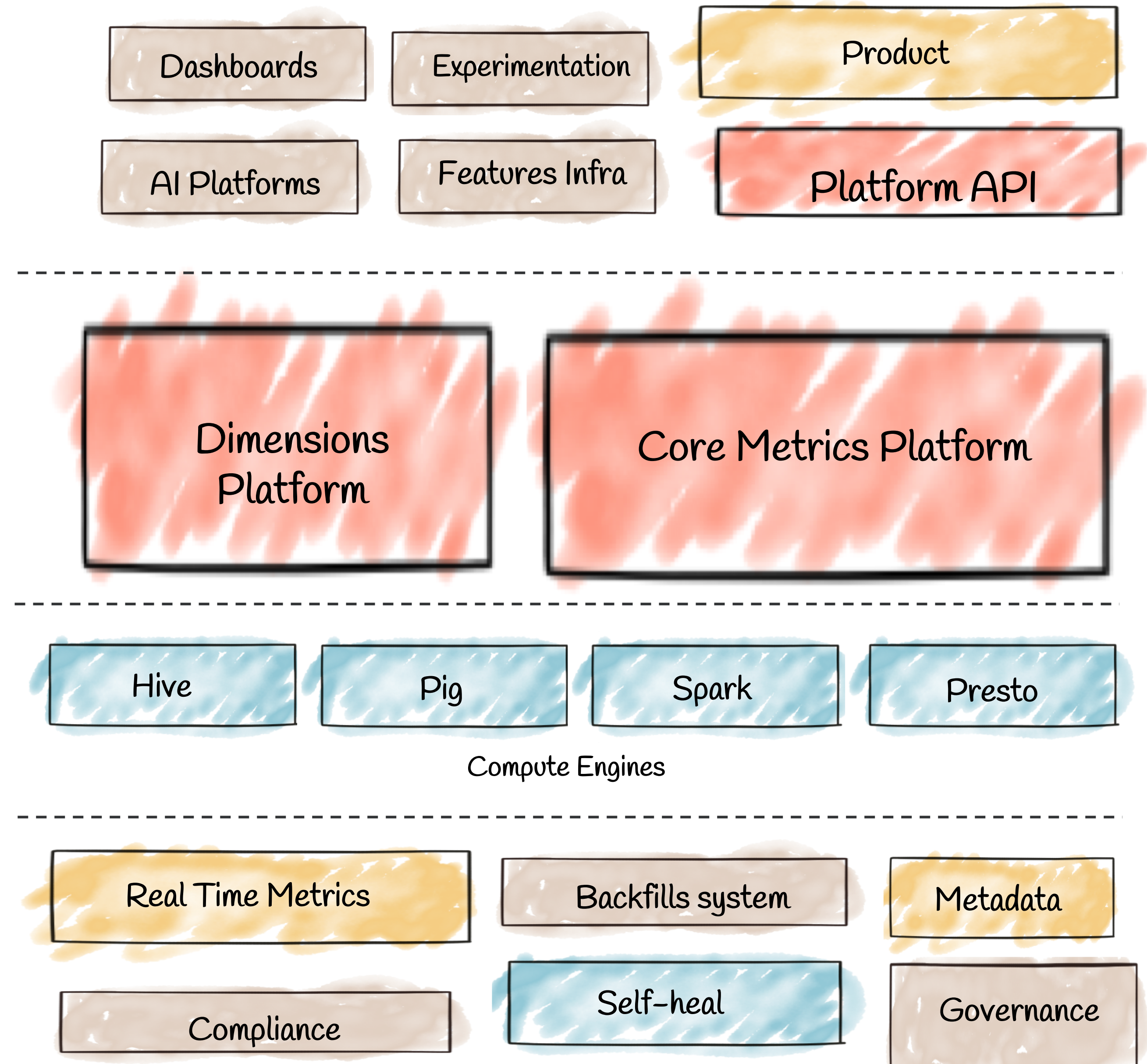
kafka/ samza
Real TimeMetrics

Invest in scaling governance

Build for customizations

Leverage economies of scale

Evolve and Innovate

## Metrics Platform as Service Ecosystem

Dashboards

Experimentation

Product

AI Platforms

Features Infra

Platform API

Dimensions Platform

Core Metrics Platform

Hive

Pig

Spark

Presto

Compute Engines

Real Time Metrics

Backfills system

Metadata

Compliance

Self-heal

Governance

# Thank you

akanitka@linkedin.com