# Globally Scalable Clickstream Analytics on Google Cloud

PyData Bangalore - Oct 2019
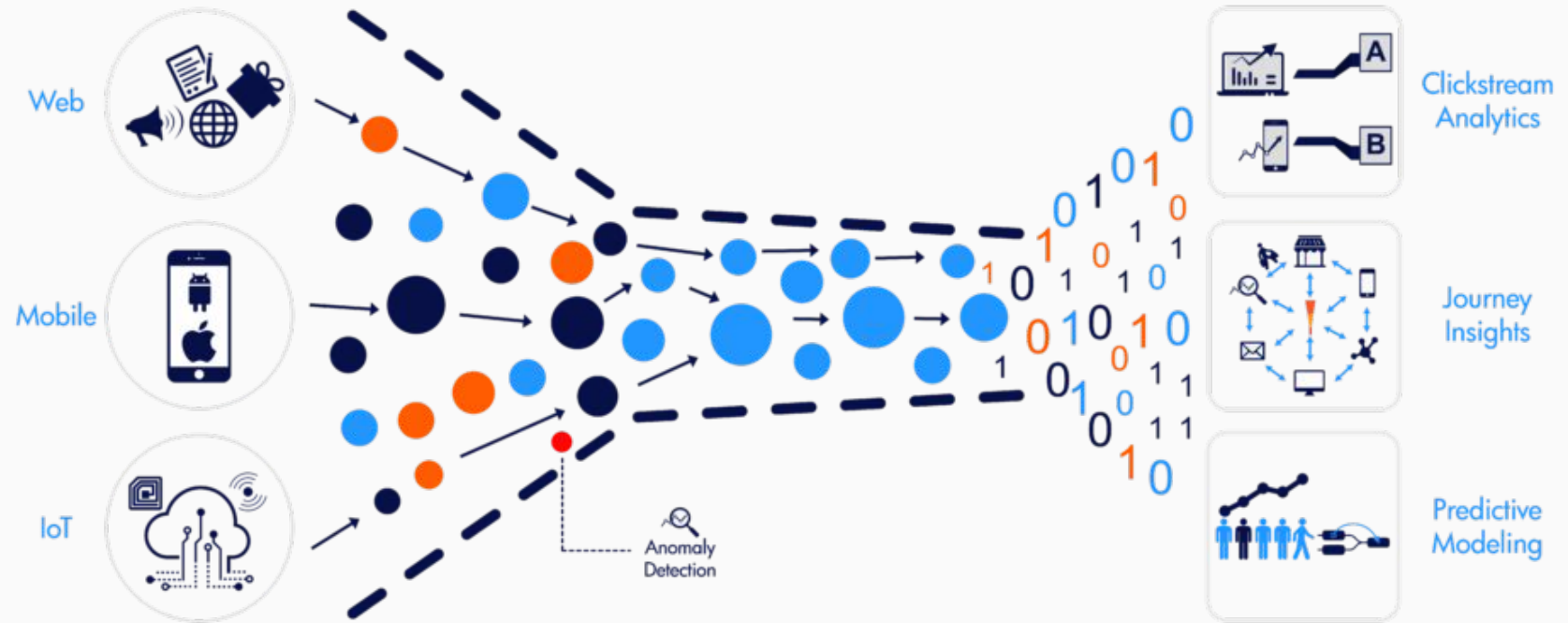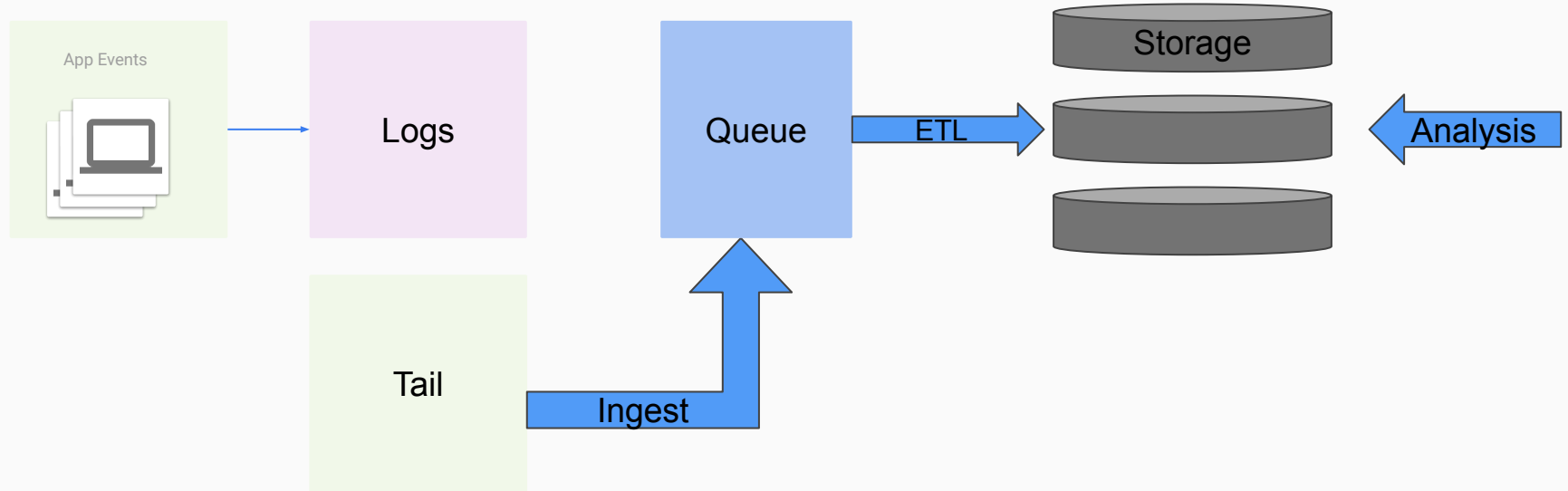
# Ramjee Ganti

- Engineer

- Scaled Tech at Startups

- Entrepreneur

- Anti Sugar Advocate

- @gantir

# What is Clickstream?

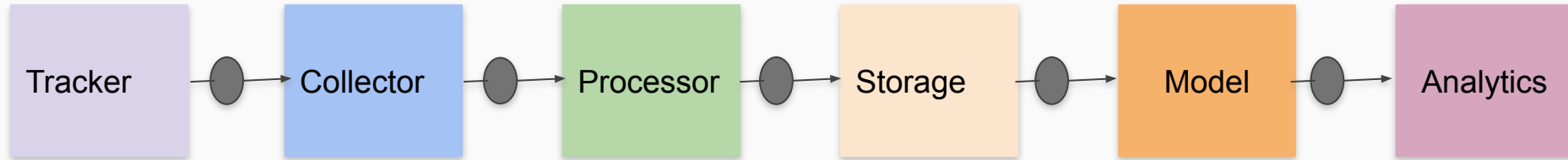# Logical Flow

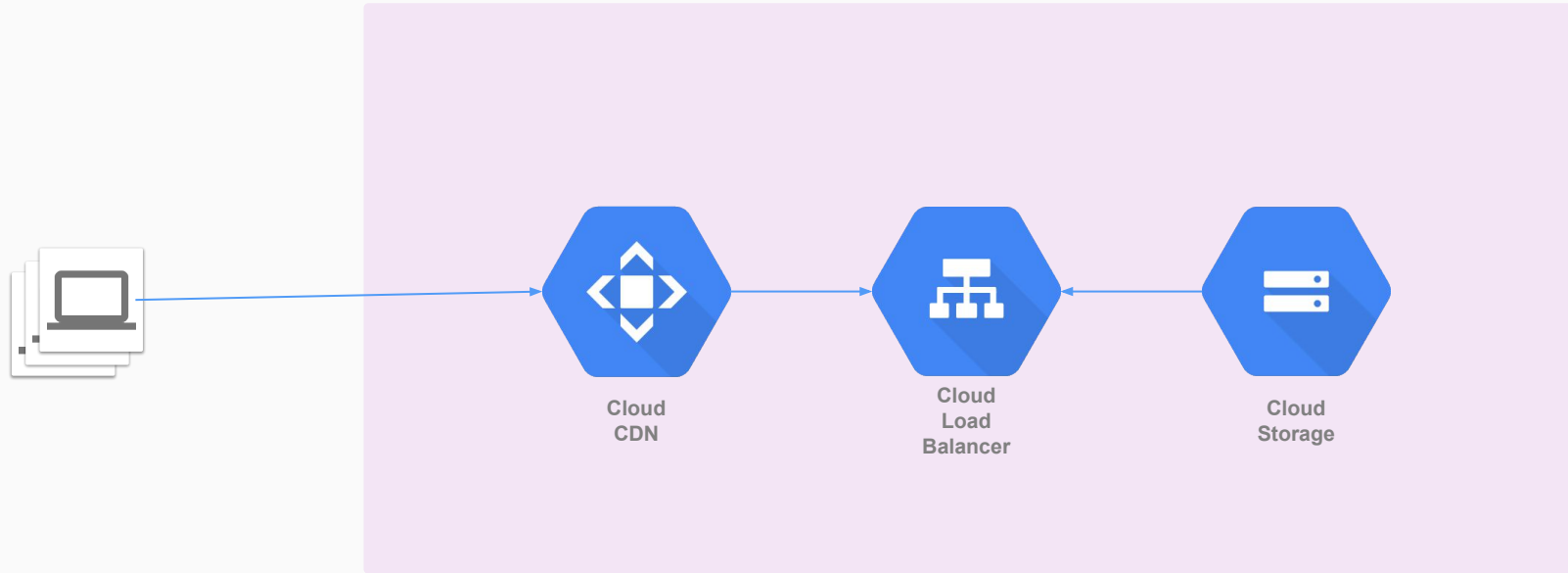# Concepts

- Tracking Protocol

- Versioned Schemas

# Components

Tracker → Collector → Processor → Storage → Model → Analytics

# Tracker (Web Tracker)

1.  Only 1st party cookie
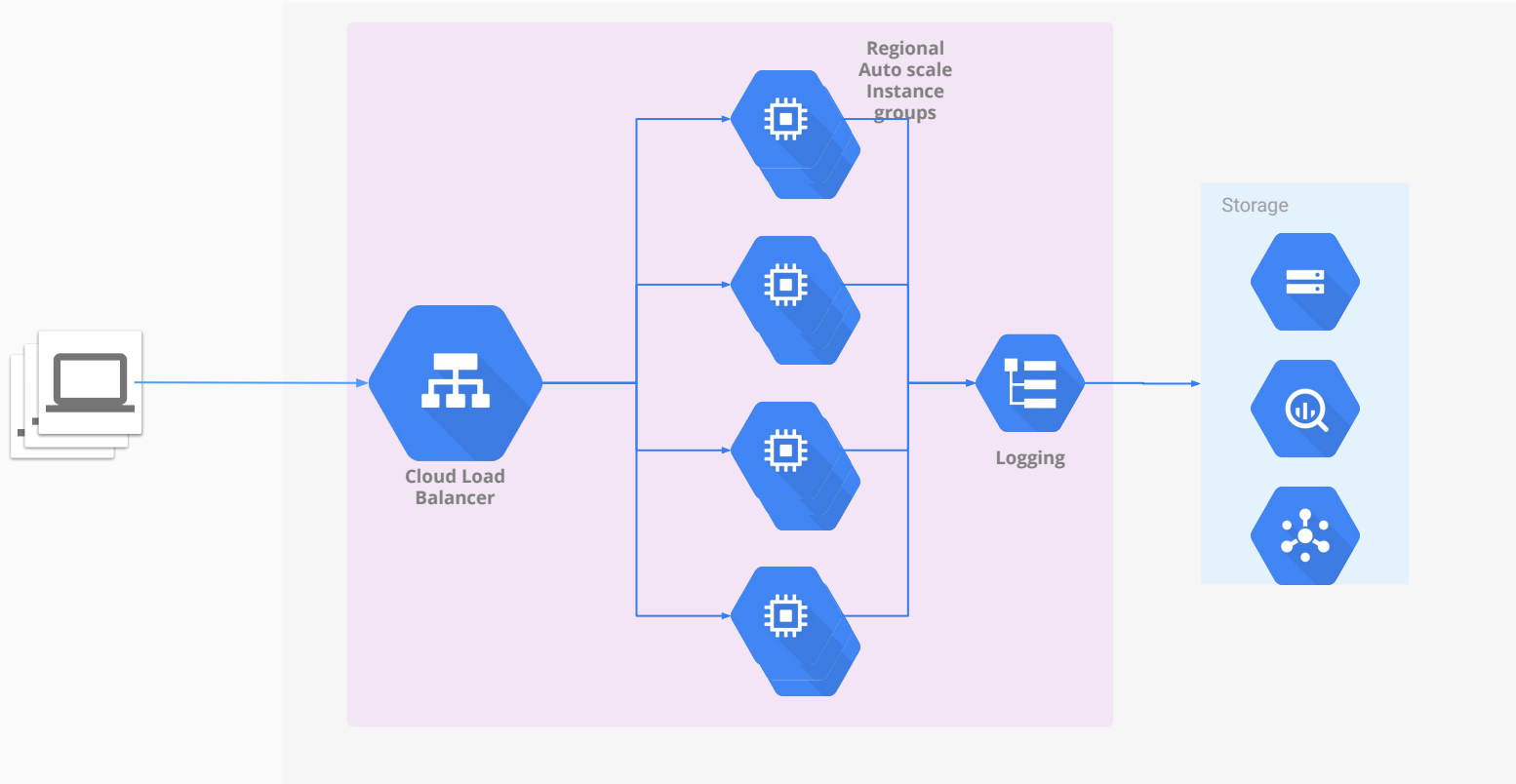
    a.   Cloud CDN

2.  3rd party cookie

    a.   GCLB + Instances

    b.   GCLB + Kubernetes

    c.   GAE

    d.   Cloud Functions

# Collector - Globally Scalable?

1. Only 1st party cookie

    a. Cloud CDN

2. 3rd party cookie

    a. GCLB + Instances

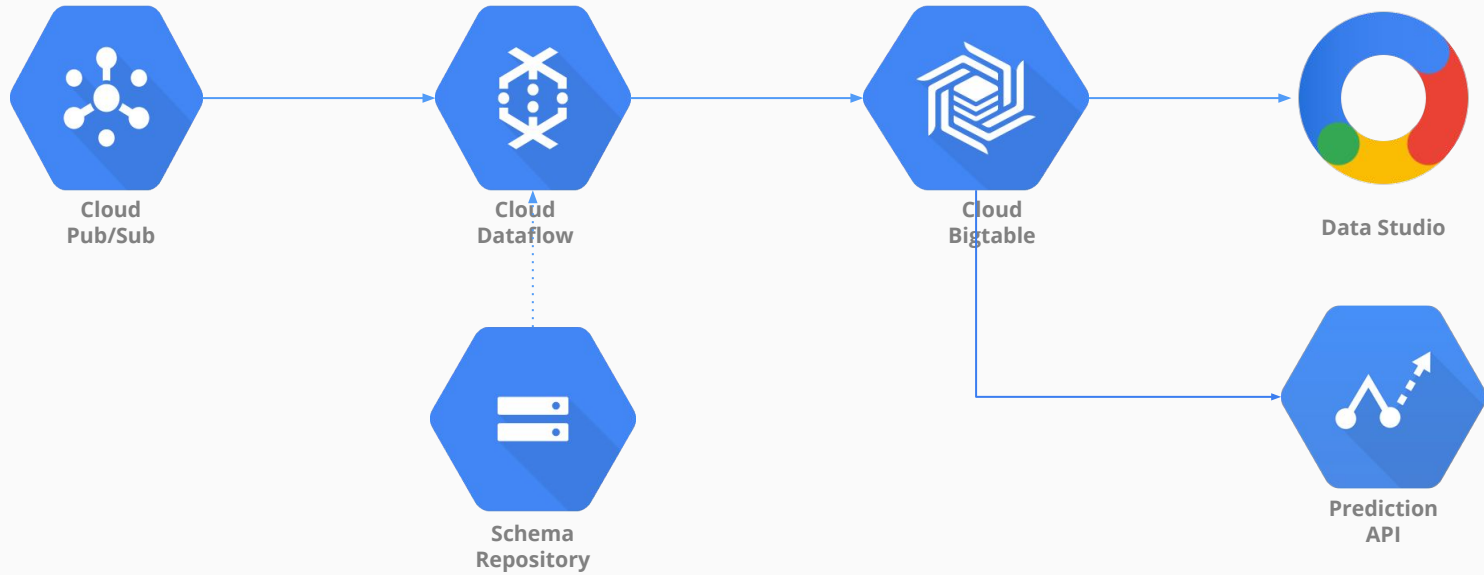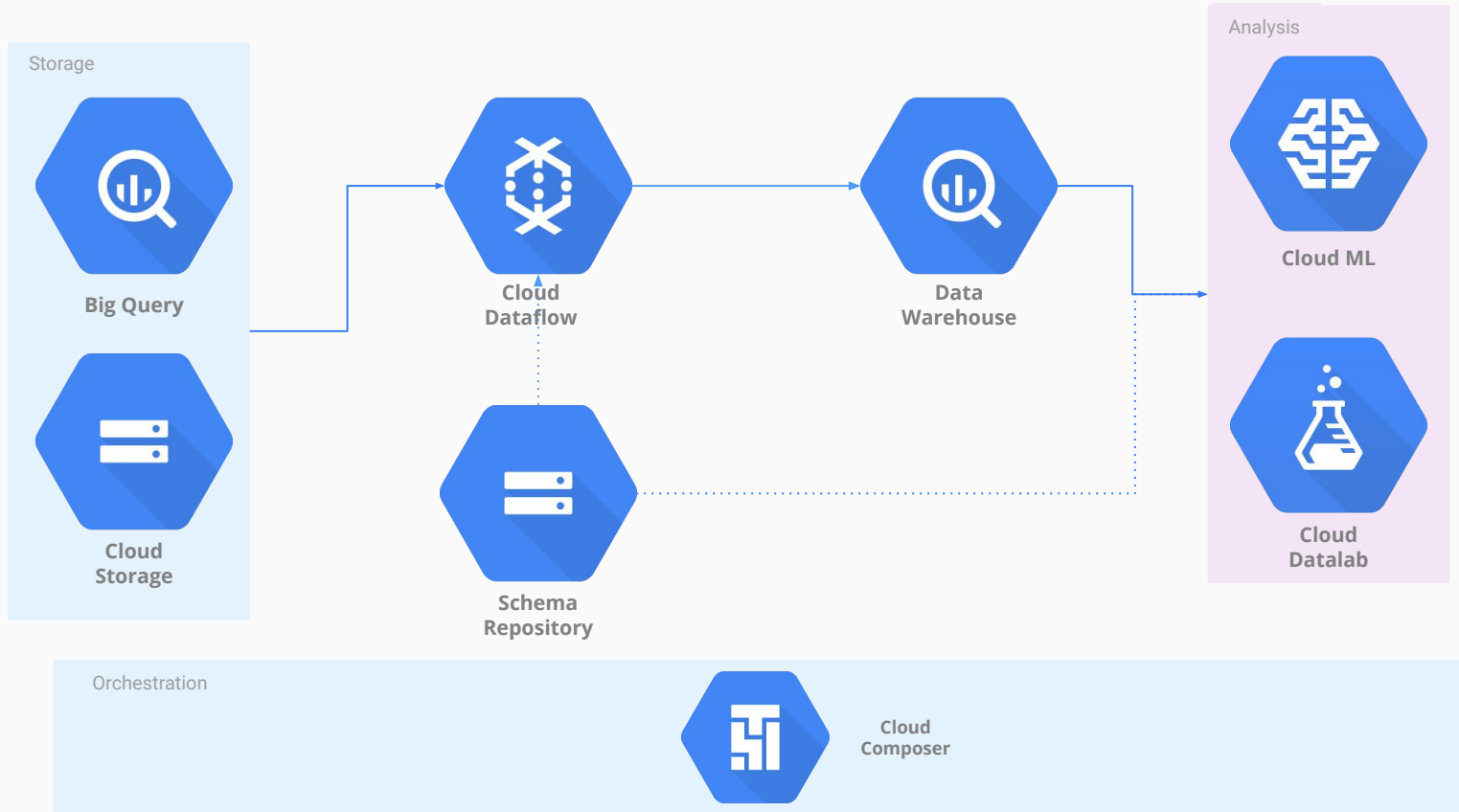    b. GCLB + Kubernetes

# Collector

# Apache Beam (Dataflow)

- Unified programming model for processing both batch and streaming data

- Cloud Dataflow is Google's managed service on GCP.

- Infinitely scalable

- Only focus on the data and not the underlying infrastructure.

- Support for different runners (Spark, Flink, Hadoop MapReduce etc)

# Real time

# Batch



Storage
- Big Query
- Cloud Storage

Cloud Dataflow

Data Warehouse

Schema Repository

Analysis
- Cloud ML
- Cloud Datalab

Orchestration
- Cloud Composer

# References

- [Snowplow Analytics](#)

- [Google Cloud Solutions](#)

- [Streaming 101](#)

# Questions

Thanks

# Annexure