# Unified Batch & Stream Processing Platform

Himanshu Bari

Director Product Management

# Most Big Data Use Cases Are About

*Improving/Re-write EXISTING solutions To KNOWN problems…*

**DataTorrent**

# Current Solutions Were Built On

A. Imperfect information

B. Expensive s/w & h/w infrastructure

C. Relational data stores

**DATATORRENT**

# Inevitable Course of the Re-write

Specialized solutions

Near perfect information

Next gen data management platform

In Memory Processing
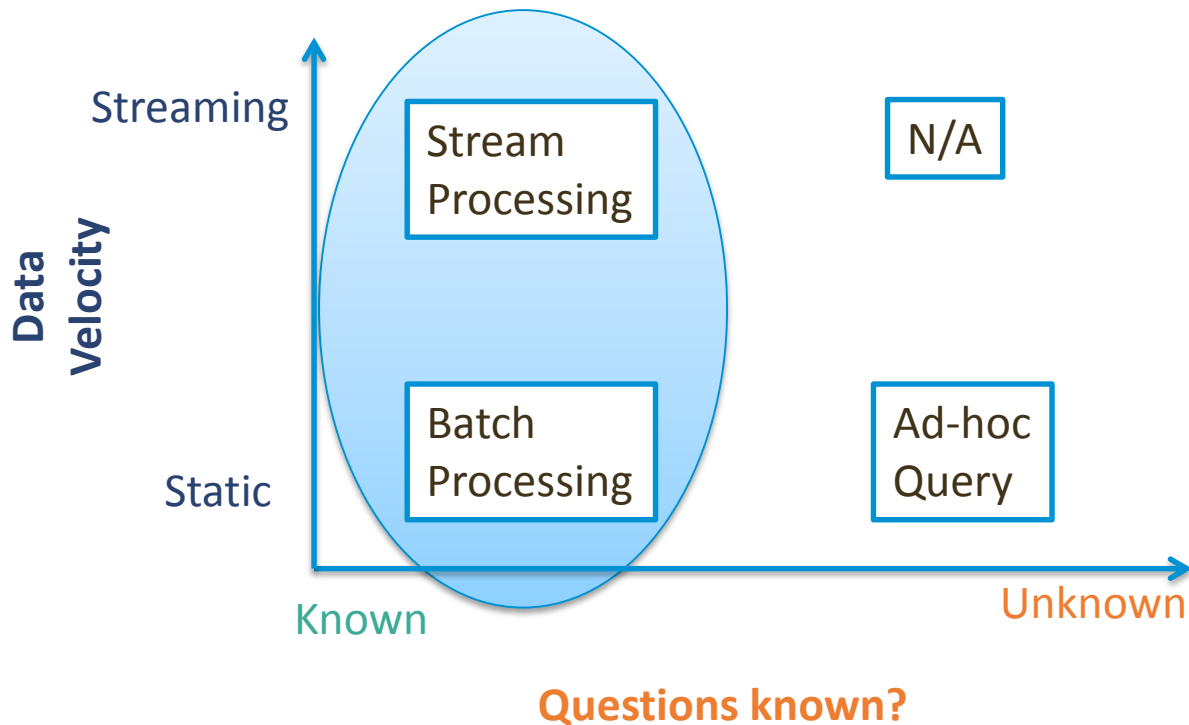
NoSQL    Graph    Hadoop    Search    EDW + RDBMS

Commodity hardware    Open source software

DATATORRENT

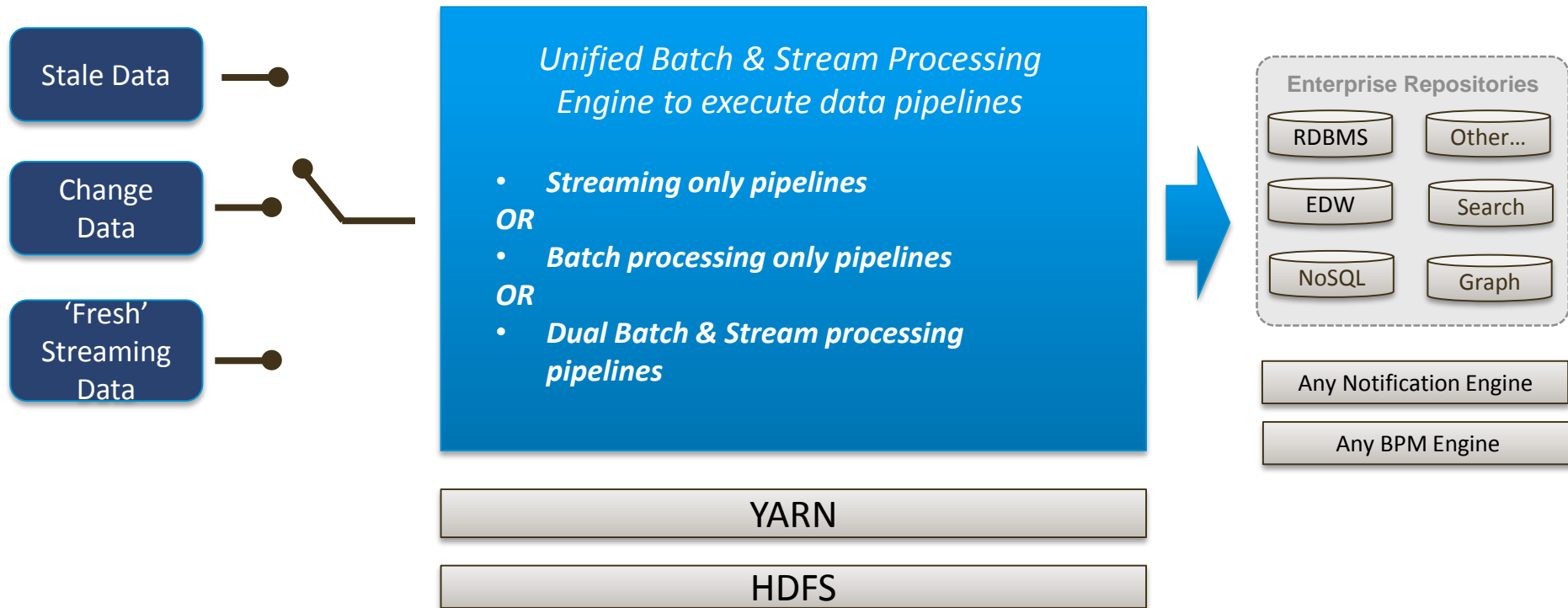# Data Processing Categories in Big Data Use Cases

# Every Batch Process *Could* Have Been A Stream Process

- Every 'Static' data point was 'Streaming' at some point
- We choose to wait and collect a bunch of data points and then process them at once in 'Batch Mode'
- Move processing time closer to the data generation or 'Event' occurrence time
- Reduces time to insight and allows you to be proactive with timely actions
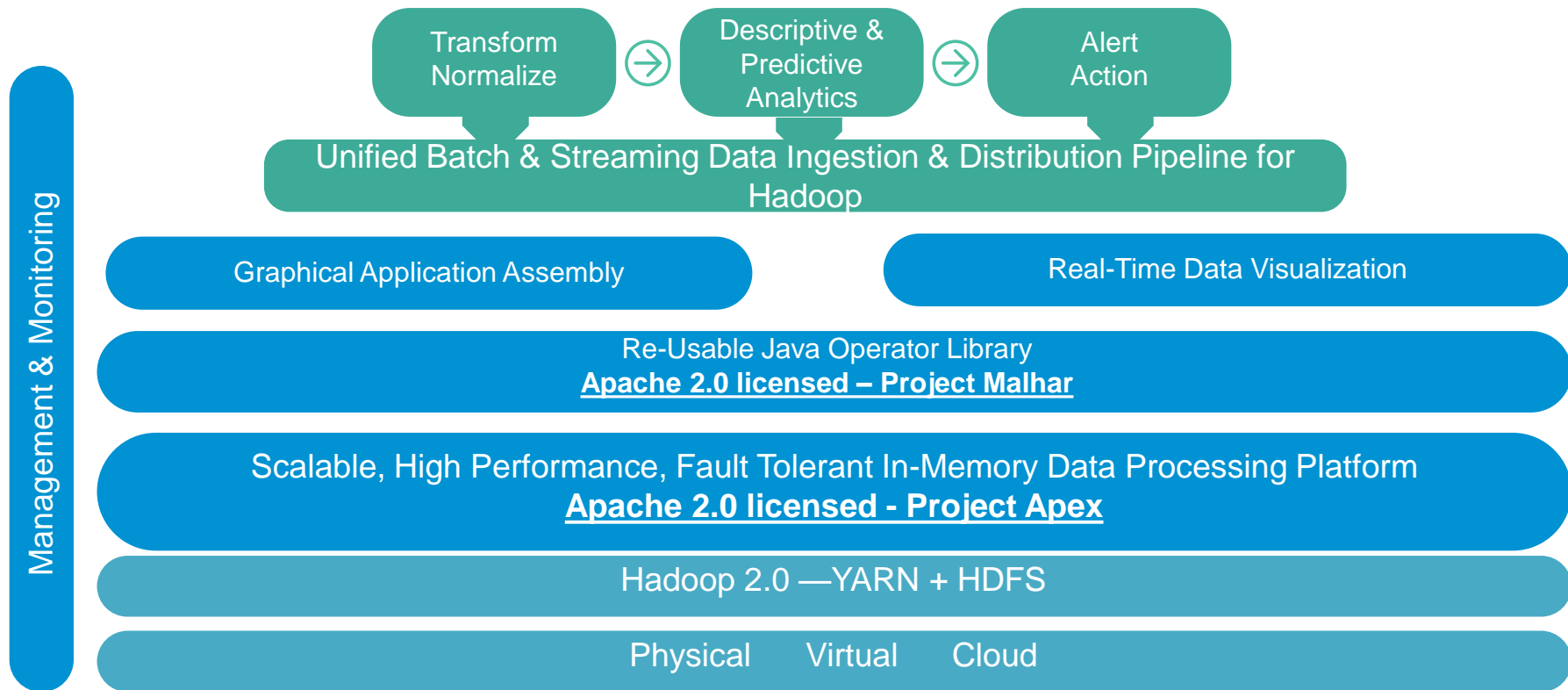
**DATATORRENT**

# But We Still Need Batch

- Historical data analysis

- What-if analysis

- Experimentation

- Data re-statement

- Transaction processing and re-conciliation

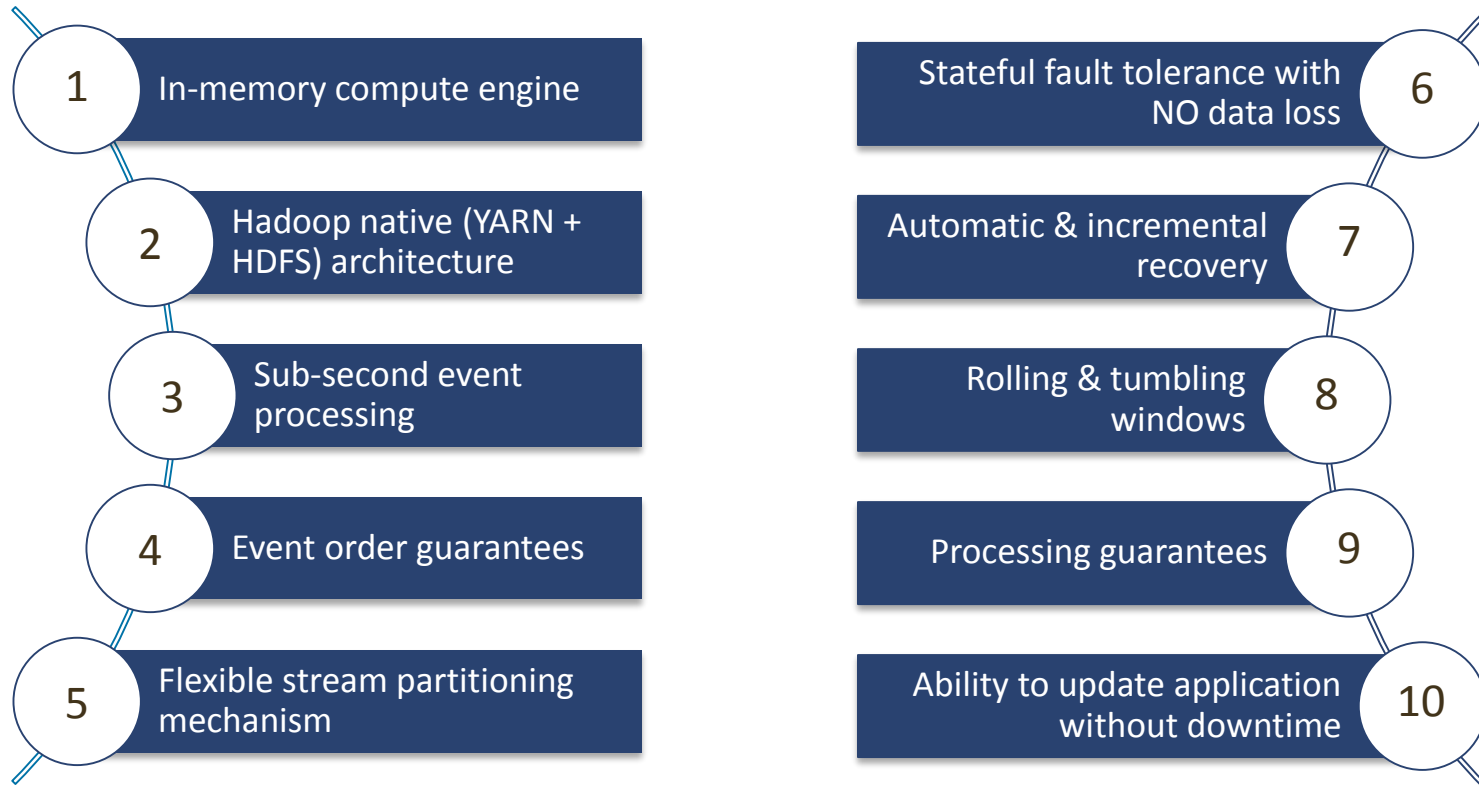- Audit

- Machine learning model training

- ……..

**DATATORRENT**

# Need A Unified Platform

Stale Data

Change Data

'Fresh' Streaming Data

**Unified Batch & Stream Processing Engine to execute data pipelines**

- **Streaming only pipelines**
- **OR**
- **Batch processing only pipelines**
- **OR**
- **Dual Batch & Stream processing pipelines**

YARN

HDFS

**Enterprise Repositories**

RDBMS

Other...

EDW

Search

NoSQL

Graph

Any Notification Engine

Any BPM Engine

**DATATORRENT**

# DataTorrent RTS Provides A Unified Batch & Streaming Platform

Transform Normalize →

Descriptive & Predictive Analytics →

Alert Action

**Unified Batch & Streaming Data Ingestion & Distribution Pipeline for Hadoop**

Management & Monitoring

Graphical Application Assembly

Real-Time Data Visualization

Re-Usable Java Operator Library
**Apache 2.0 licensed – Project Malhar**

Scalable, High Performance, Fault Tolerant In-Memory Data Processing Platform
**Apache 2.0 licensed - Project Apex**

Hadoop 2.0 —YARN + HDFS

Physical      Virtual      Cloud

**DataTorrent**

# DataTorrent Project Apex- Unified Batch & Stream Processing Engine

1. In-memory compute engine
2. Hadoop native (YARN + HDFS) architecture
3. Sub-second event processing
4. Event order guarantees
5. Flexible stream partitioning mechanism
6. Stateful fault tolerance with NO data loss
7. Automatic & incremental recovery
8. Rolling & tumbling windows
9. Processing guarantees
10. Ability to update application without downtime

**DataTorrent**

# Unified Data Ingestion & Distribution Pipeline

**Input & Output Variety**
- FTP, S3 etc.
- Kafka & JMS
- Change data capture

**Tackle data size & speed fluctuations**
- Overcome HDFS small file problem
- Skew management

**Hadoop Native**
- Runs within the Hadoop cluster over YARN & HDFS

**Easily customizable**
- Easily extend and insert operations for data preparation

**Run-time updates**
- Parameters like filtering criteria, bandwidth utilization & polling interval should be updateable at runtime

**Simple to build & Operate**
- Graphical UI & API
- End to end metrics

**DATATORRENT**

# Hadoop Data Ingestion & Distribution Application

# Data Prep & Analytics Layer Requirements

- Truly scale horizontally across the Hadoop cluster

- Pre-built operators
  - Re-ordering
  - Normalization
  - Transpose
  - De-duplication
  - Tagging
  - Filtering
  - Enrichment

- Operators work seamlessly in both streaming & batch mode
  - Data local HDFS read & process
  - Ability to do computations on time window as well as file boundaries

- Ability to re-use existing business logic

- Simple workflow & scheduling capabilities
  - Built-in or integrations with Oozie or other schedulers

**DATATORRENT**

# Development API Requirements

- Consistent between Streaming & Batch pipelines
- No mapreduce
- No exposure of low level processing engine concepts
- Easily extendible

**DataTorrent**

# Malhar Operator Library Overview

## Malhar Operator Library

### Input / Output Connectors

| File Systems | RDBMS | NoSQL | Messaging |
|---|---|---|---|
| Notification | In Memory Databases | Social Media | Protocol read/write |

### Compute Operators

| Pattern Matching | Stats & Math | Machine Learning & Algorithms |
|---|---|---|

| Parsers | Stream Manipulators | Query & Scripting |
|---|---|---|

**DATATORRENT**

# Visual Application Assembly

- **Easy to Use**
  - Web based drag-n-drop development environment
  - Automatic port compatibility validation
  - Simple schema management
  - Generic property configurator

- **Easy to Operate**
  - No external component dependency - Runs natively in Hadoop
  - Integrated with DataTorrent management platform

- **Simple to extend**
  - Simple API to enable any existing DataTorrent operator
  - Ability to plug any business logic using a custom operator

# Streaming or Batch Data Processing Visualization

- Intuitive
  user interface
  - Auto-generate or custom create
  - One dashboard for multiple apps
  - Supports bar, line, pie, area
    charts & data tables

- Easy to Operate
  - No external component
    dependency - Runs natively
    in Hadoop
  - Integrated with DataTorrent
    management platform

- Simple to extend
  - Any DAG operator can be
    made a real-time datasource

# Summary



**Project Apex**
https://www.datatorrent.com/product/project-apex/

**DataTorrent RTS Sandbox**
https://www.datatorrent.com/download/

- World is moving from 'Batch' to 'Streaming' BUT both are required
- Need a Hadoop native in memory compute engine that is scalable & fault tolerant in BOTH batch & streaming modes
- With out -of-the box data prep & analytics operators
- Using a consistent & functional development API
- Operationalized through a common management & monitoring layer

DataTorrent

Big Data. Big Actions. Now.

# Some Verticals & Use Cases

| Ad-Tech | Telco & Cable |
|---|---|
| • Real-time customer facing dashboards on key performance indicators<br>• Click fraud detection<br>• Billing optimization | • Call detail record (CDR) & extended data record (XDR) analysis for<br>  • Service quality improvement<br>  • Capacity planning/optimization<br>• Understanding customer behavior AND context<br>• Packaging and selling anonymous customer data |
| **Financial Services** | **IoT** |
| • Fraud & risk monitoring<br>• Sentiment based trading strategies<br>• Usage based insurance<br>• Improved credit risk assessment<br>• Improving turn around time of trade settlement processes | • Process optimization<br>• Proactive maintenance prediction<br>• Remote monitoring & diagnostics |

**DATATORRENT**