

Josh Poduska, Sr. Business Analytics Consultant,  
Action Corporation



# Beyond Sampling: Fast, Whole-Dataset Analytics for Big Data on Hadoop

**October 2013 KNIME Day Boston**

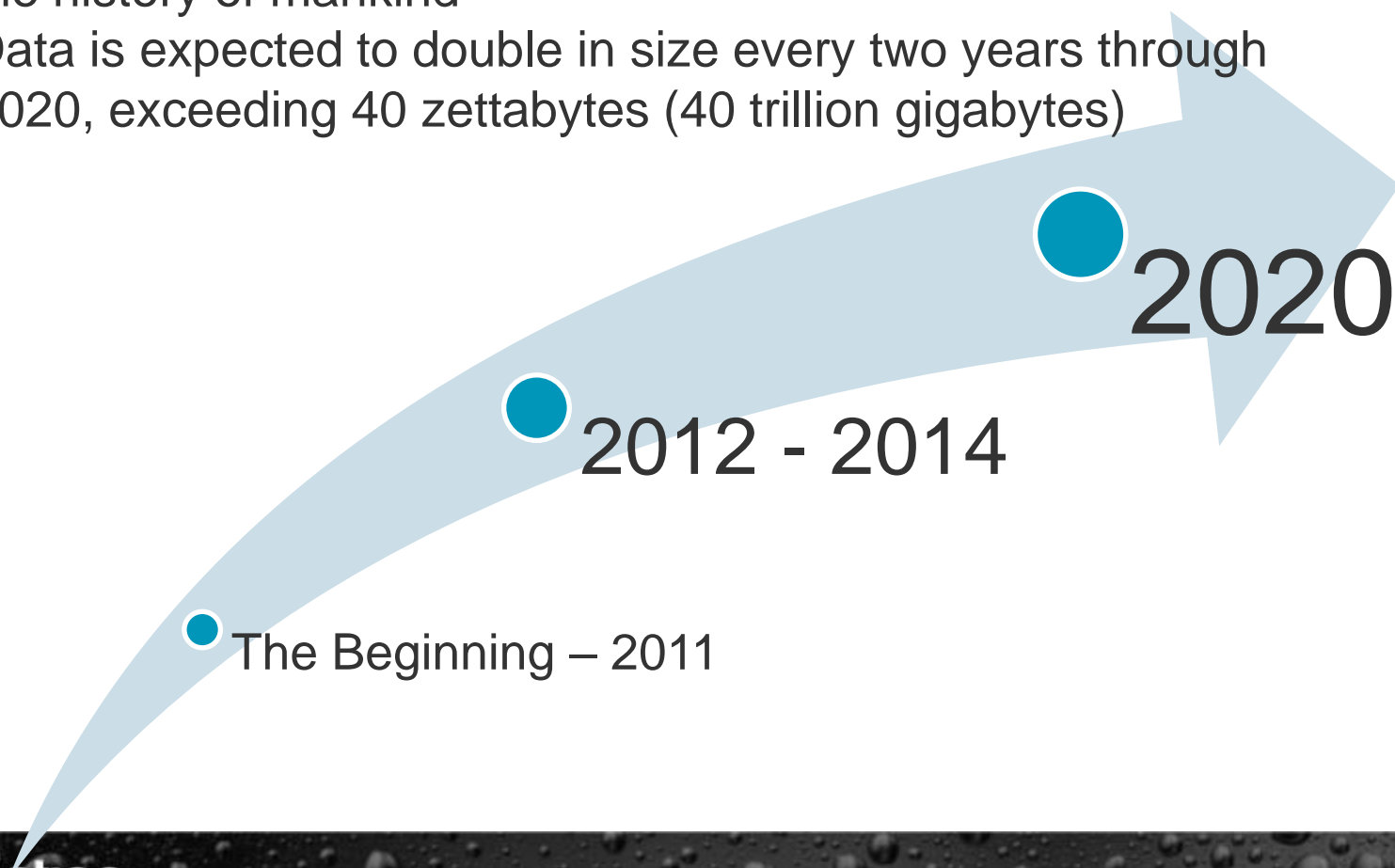


# > Agenda

- The Age of Data
- Scaling Gap
- Enter the Data Scientists
- POS Hadoop Analytics in KNIME
- Conclusions

# > The Age of Data

- In the last two years we have generated more data than in the history of mankind
- Data is expected to double in size every two years through 2020, exceeding 40 zettabytes (40 trillion gigabytes)



# > Entering the Age of Data

## ■ What's Changed?

- **Data is THE central business asset:**  
“Data are an organization's sole, non-depletable, non-degrading, durable asset. Engineered right, data's value increases over time because the added dimensions of time, geography, and precision.” (Peter Aitken)
- **Data generation has changed forever**
  - Instrumentation of ALL businesses, people, machines
- **Data is born digitally and flows constantly**
  - “All things are flowing..” (Heraclitus, 500 BC)

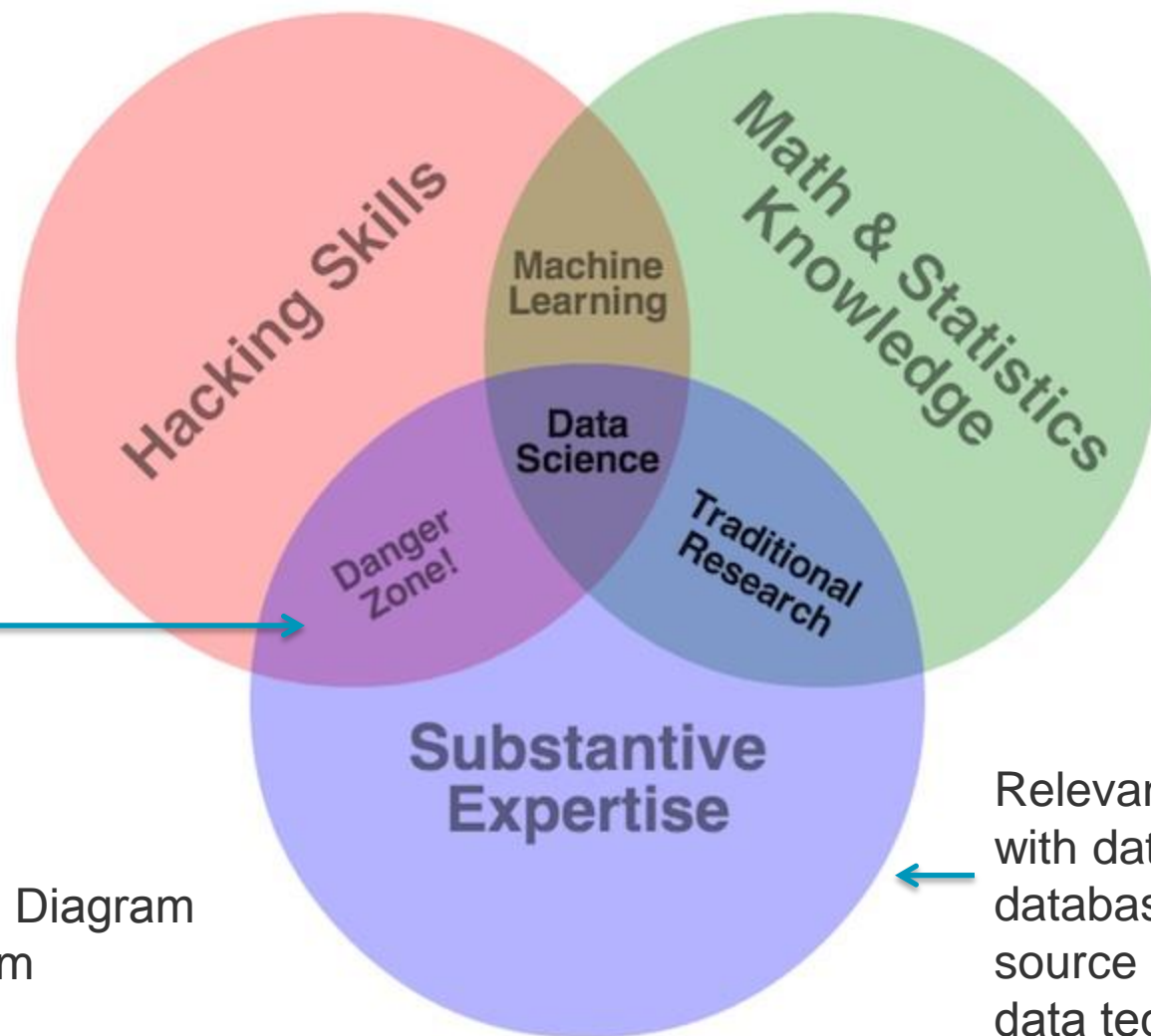


# > The Scaling Gap

- In the Age of Data, if you are not super-scaling you are failing. What does super-scaling entail?
- Software stacks to **consume, analyze** and **act** on event pipelines must be **frictionless** to set up
- Yet extremely **performant**: must **scale-up and scale-out (SUSO)** to fully exploit game-changing price/performance on modern commodity hardware
- And be **elastic**
- And still be **affordable**
- The **hard truth**:
  - Almost no legacy data/event processing stacks super-scale
  - And there is no path to reasonably (and economically) get there

**Your legacy analytic software WILL fail in the Age of Data**

# > Venn Diagram of Data Scientists



Not statistically  
valid, lacking  
interpretation

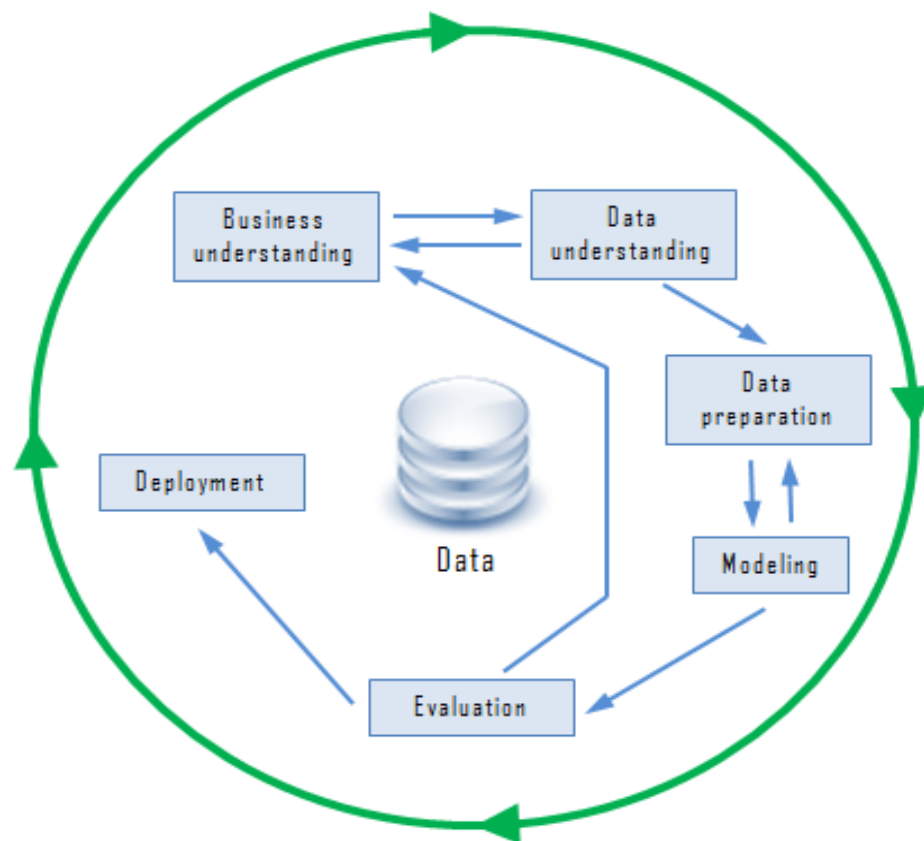
→

Drew Conway's Diagram  
[drewconway.com](http://drewconway.com)

← Relevant experience  
with data, tools,  
databases, open  
source software,  
data techniques, ...

# > Data Mining is a Process

- Successfully completing the process requires having a CONVERSATION with the data.



# > Tools Available (in big data context)

- Distributed Platforms

- Hadoop, Mesos, ...
- Dataflow, HPCC

- Analytic Databases

- ParAccel, Vectorwise, Vertica, Aster Data, GreenPlum, Netezza, ...
- Cassandra , HBase, MongoDB, ...

- Analytic Platforms

- Dataflow Analytics, SAS, SPSS, ...
- KNIME, R, RapidMiner, ...



# > Hadoop

- **Open source, distributed (scale out) platform for data processing on cheap hardware**
- Components
  - HDFS – Hadoop Distributed File System
  - MapReduce – computation framework
    - Broken into two phases
    - Map – takes input, produces a name/value mapping
    - Reduce – applies a final reduction of the mapping phase
  - HBase – name/value pair data store
  - Oozie, ZooKeeper, ...
  - Distributions: Apache, Cloudera, HortonWorks, MapR, Intel, ...

## Retail POS Application

# > Market Basket Analysis

## ■ The Data

- Retailers have Point of Sales (POS) data
- Items purchased in same basket are captured (line items)
- Summary of each basket (basket or order)
- Information about items - UPC, SKU, description, category hierarchy

## ■ Analysis

- Need to sell longer-held produce, labeled CLOSEOUT ITEMS
- Want to know what drives total receipt spend
- Want to know what items sell well together
- Which items drive purchase of other items

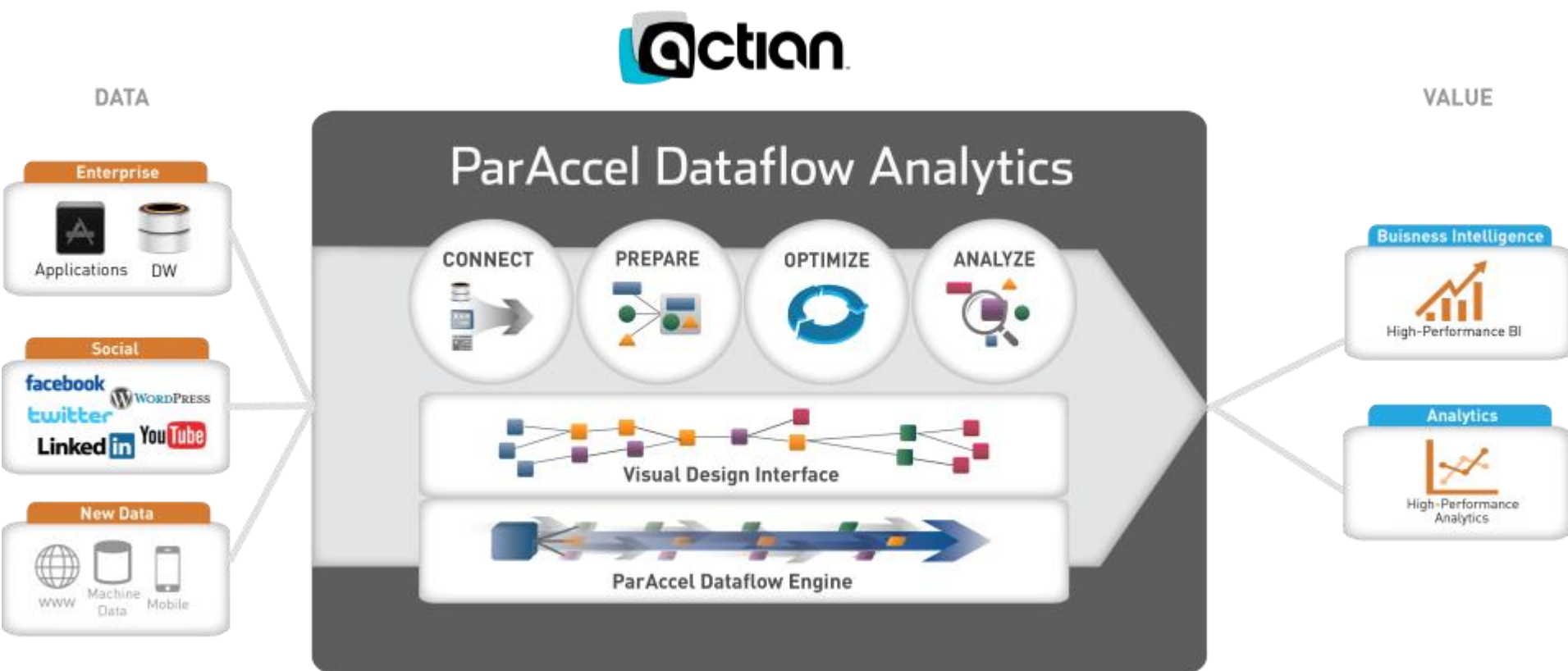
# > Configuration

- Hundreds of millions of rows of POS data
- Hadoop platform
  - 3 worker nodes, 1 head node
  - Running Cloudera CDH4
  - Distributed Dataflow 6.1
- Analysis Tools (run on desktop)
  - Actian Dataflow Analytics for Hadoop (KNIME installation library)
  - Gephi (open source graph visualization)

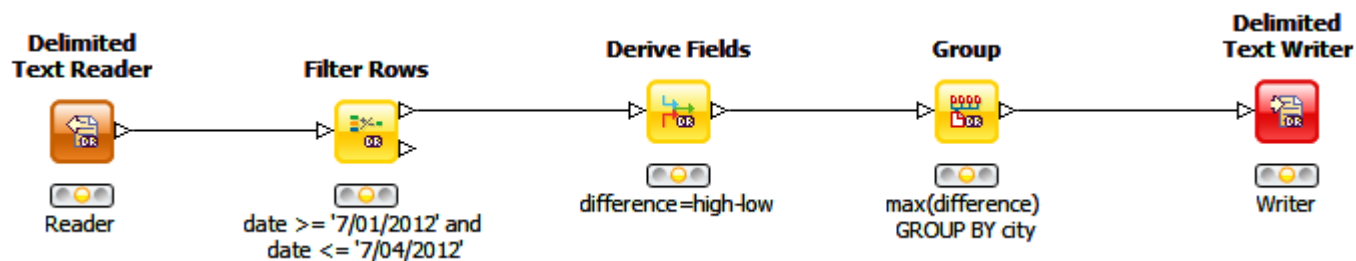
# > Dataflow Analytics for KNIME

- KNIME (knime.org)
  - Open source analytics workflow platform
  - Highly extensible
  - Active community of plugin contributors
  - Commercially available Server, Teamspace and Report products
- Dataflow Analytics
  - Actian developed extensions to KNIME
  - Includes scalable Dataflow technology
  - Large set of “nodes” based on Dataflow

# > Dataflow Analytics for KNIME



# > Dataflow Analytics for KNIME



Compiled to a set of physical graphs



## Phase 1



## Phase 2



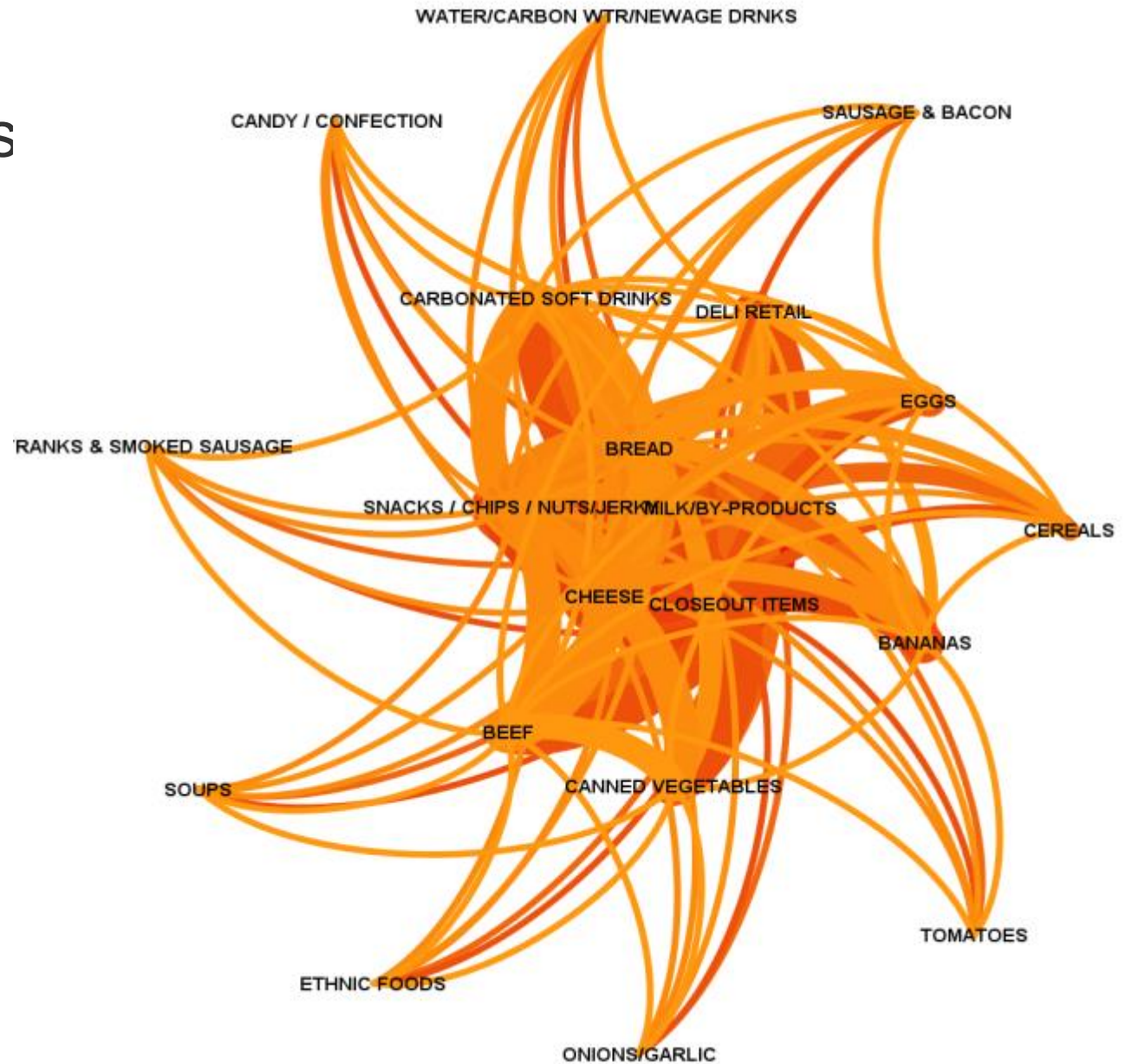


A high-speed photograph of a water splash, creating a horizontal wave of water across the middle of the frame. The water is clear and shows fine details of the splash.

## Demo



# Visualize Associations in Gephi



# > Big Data Conversation → Analytical Impact

- Retail Goal: Increase Spend on Closeout Produce
- High confidence antecedent:
  - Cheese + Bananas
  - Snacks + Bananas
  - Bread + Bananas ...
- High support antecedents:
  - Bananas
  - Tomatoes
  - Milk + Cheese ...

# > Big Data Conversation → Analytical Impact

- Retail Goal: Increase Spend on Closeout Produce
- What Closeout Produce associations exist with low lift?  
(↑ lift → bring consequent into basket)
  - Chicken or Cheese + Soft Drinks
  - **Action = Same visit coupon**

# > Big Data Conversation → Analytical Impact

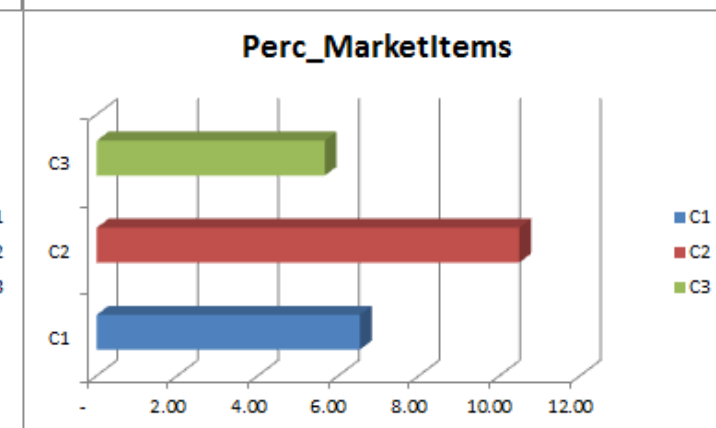
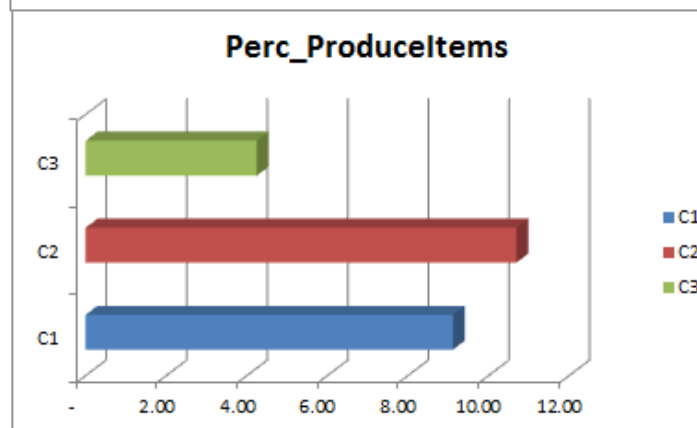
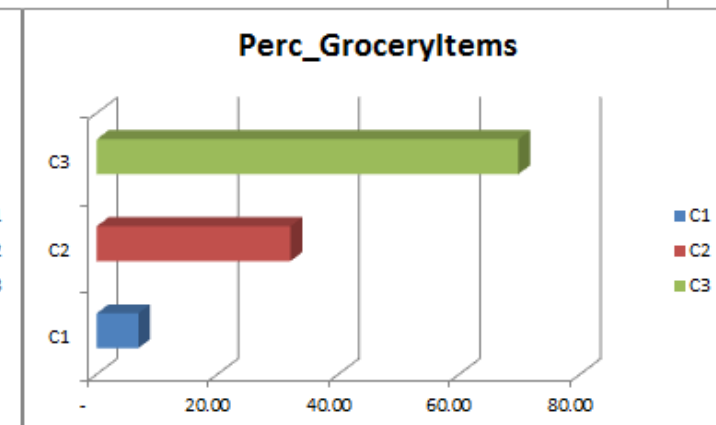
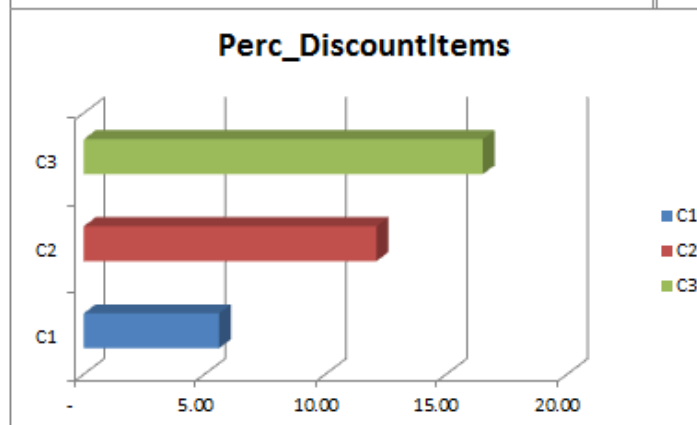
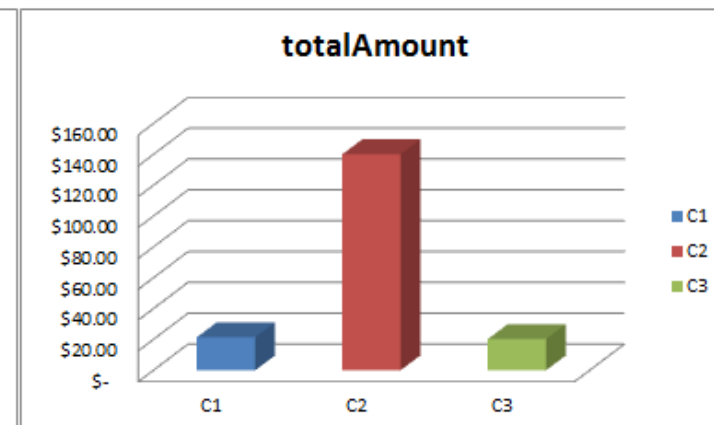
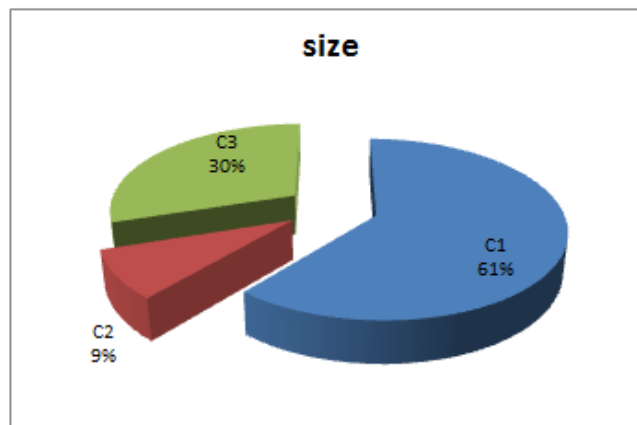
- Retail Goal: Increase Spend on Closeout Produce
- What has high support but no Closeout Produce association?
  - Candy, Canned Vegetables, Soup.
  - **Action = Same visit coupon, receipt coupon, delivered coupon, etc.**

# > Big Data Conversation → Analytical Impact

- Retail Goal: Increase Total Spend per Customer
- Clustering points to % Basket in Market Department being a key factor in total spend.



# Cluster Results



# > Big Data Conversation → Analytical Impact

- Retail Goal: Increase Total Spend per Customer
- Follow up Cluster with Linear Regression
  - % Basket in Market Department is indeed a key predictor even after factoring for:
    - DOW
    - Season
    - %Basket Market Department
    - %Basket Produce Department
    - %Basket Grocery Department
    - %Basket Discount Department



# Conclusion

- The Age of Data is here
  - Data is the central business asset
  - Data generation has changed forever
  - Shift of analysis focus to time-stamped events
  - Crisis of software that scales to meet demand
- Data Science is changing how data is:
  - Collected, discovered, analyzed, used, acted upon ...
- Big Data Conversations
  - Deep analysis is required to move beyond basic findings
  - Actionable results require very heavy lifting







# Questions

[www.action.com](http://www.action.com)

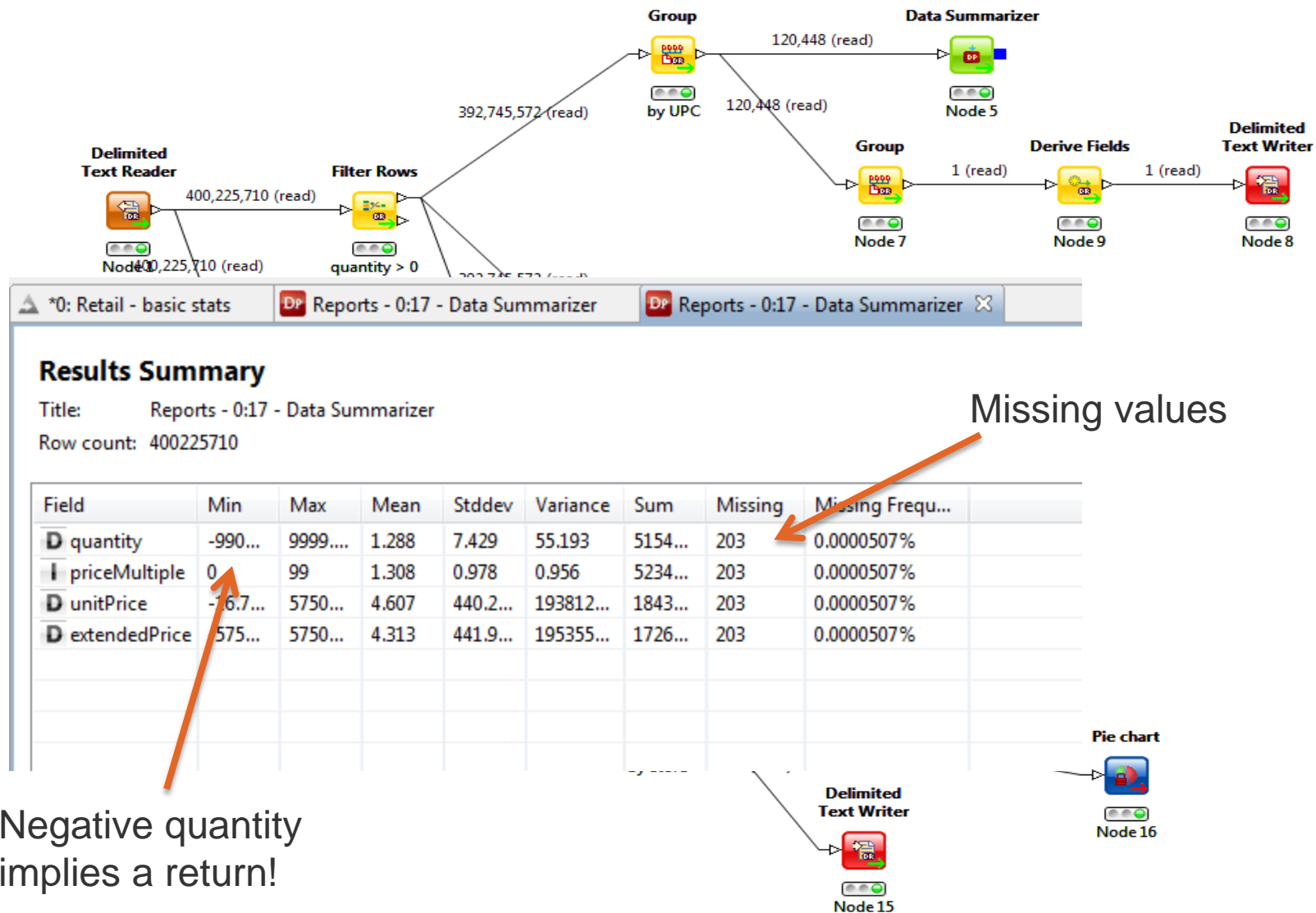
[facebook.com/actioncorp](https://facebook.com/actioncorp)

[@actioncorp](#)

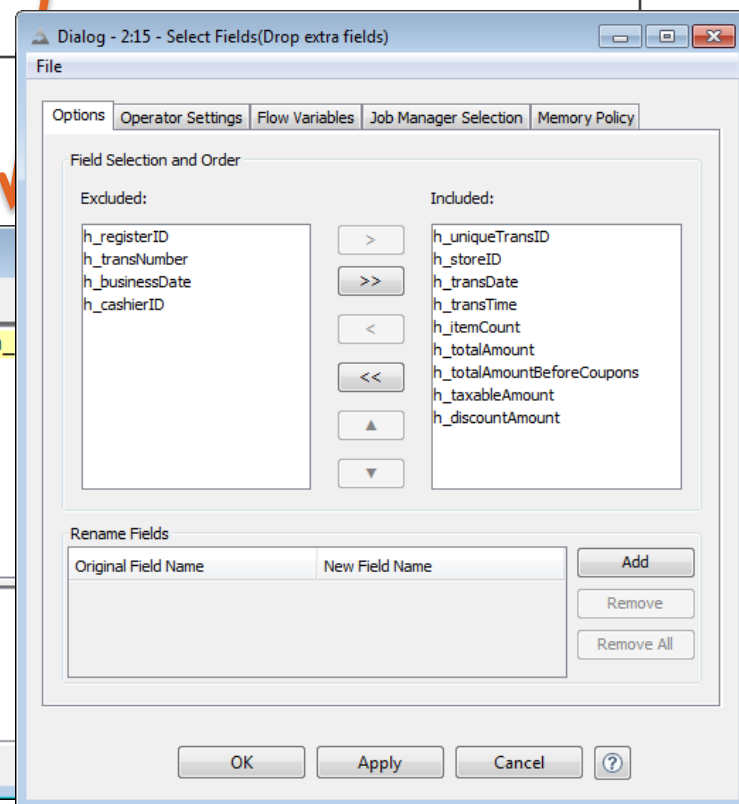
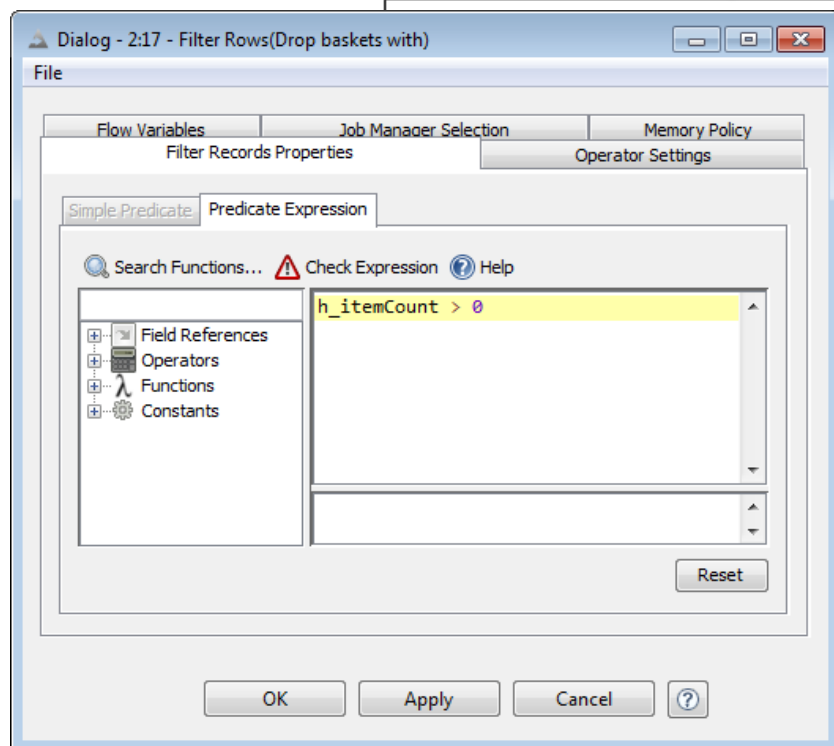
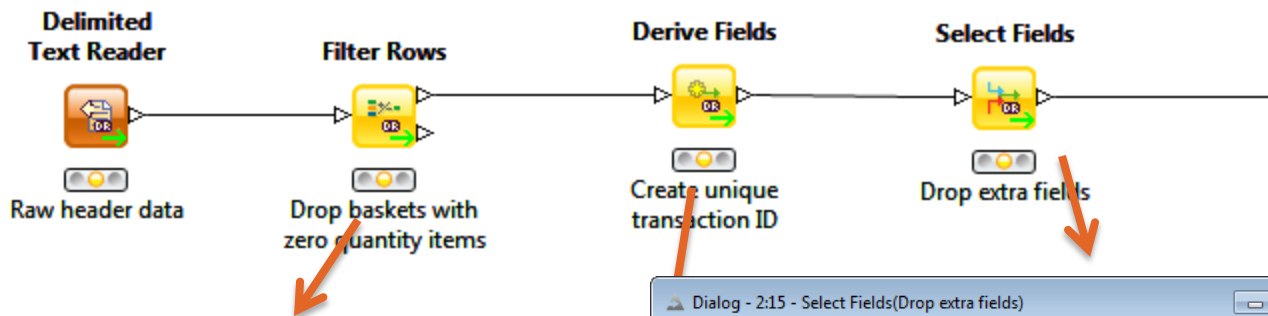
[joshua.poduska@action.com](mailto:joshua.poduska@action.com)

## Backup slides

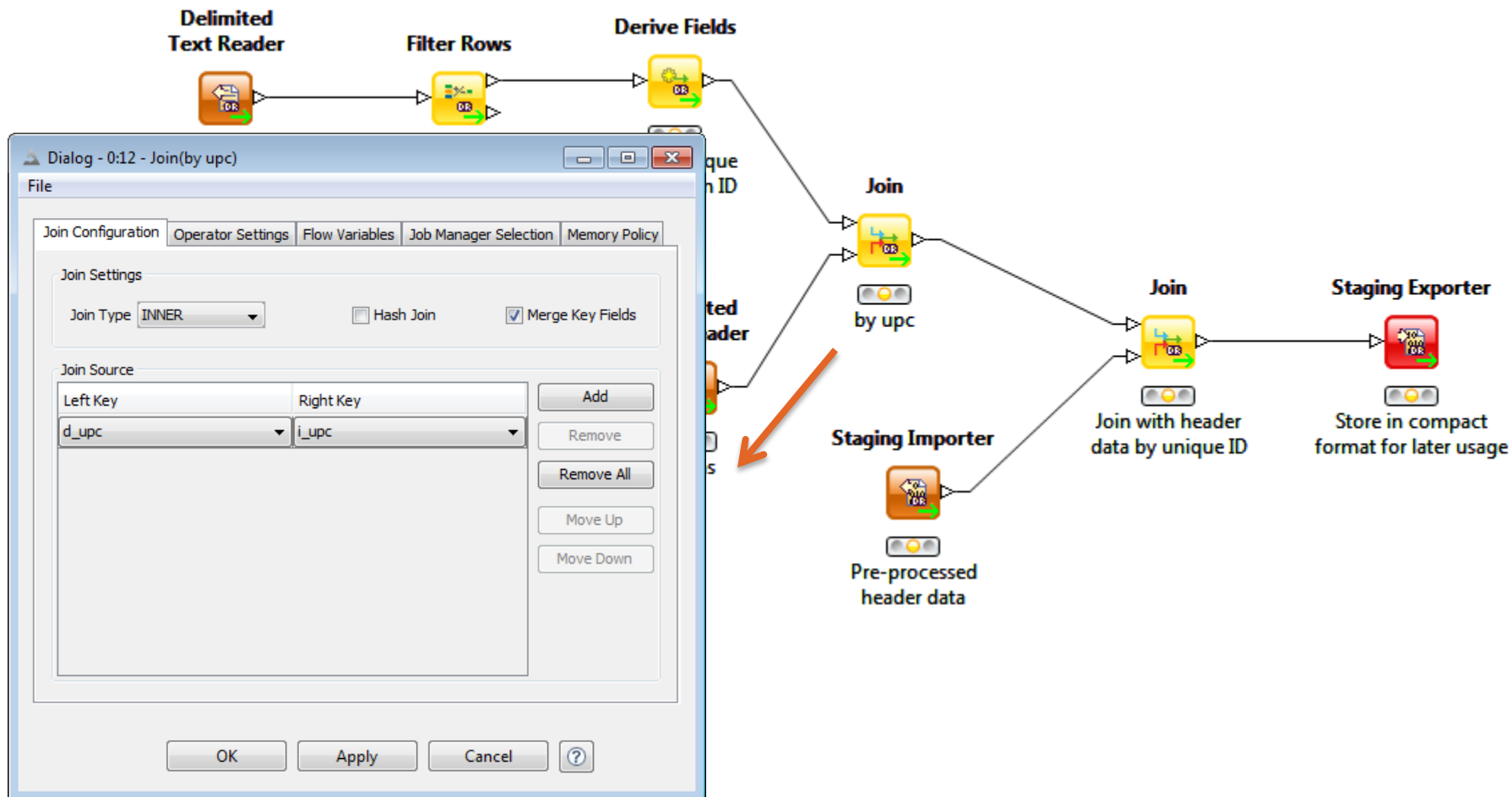
# > Start with data discovery



# > Cleanse, enrich & aggregate



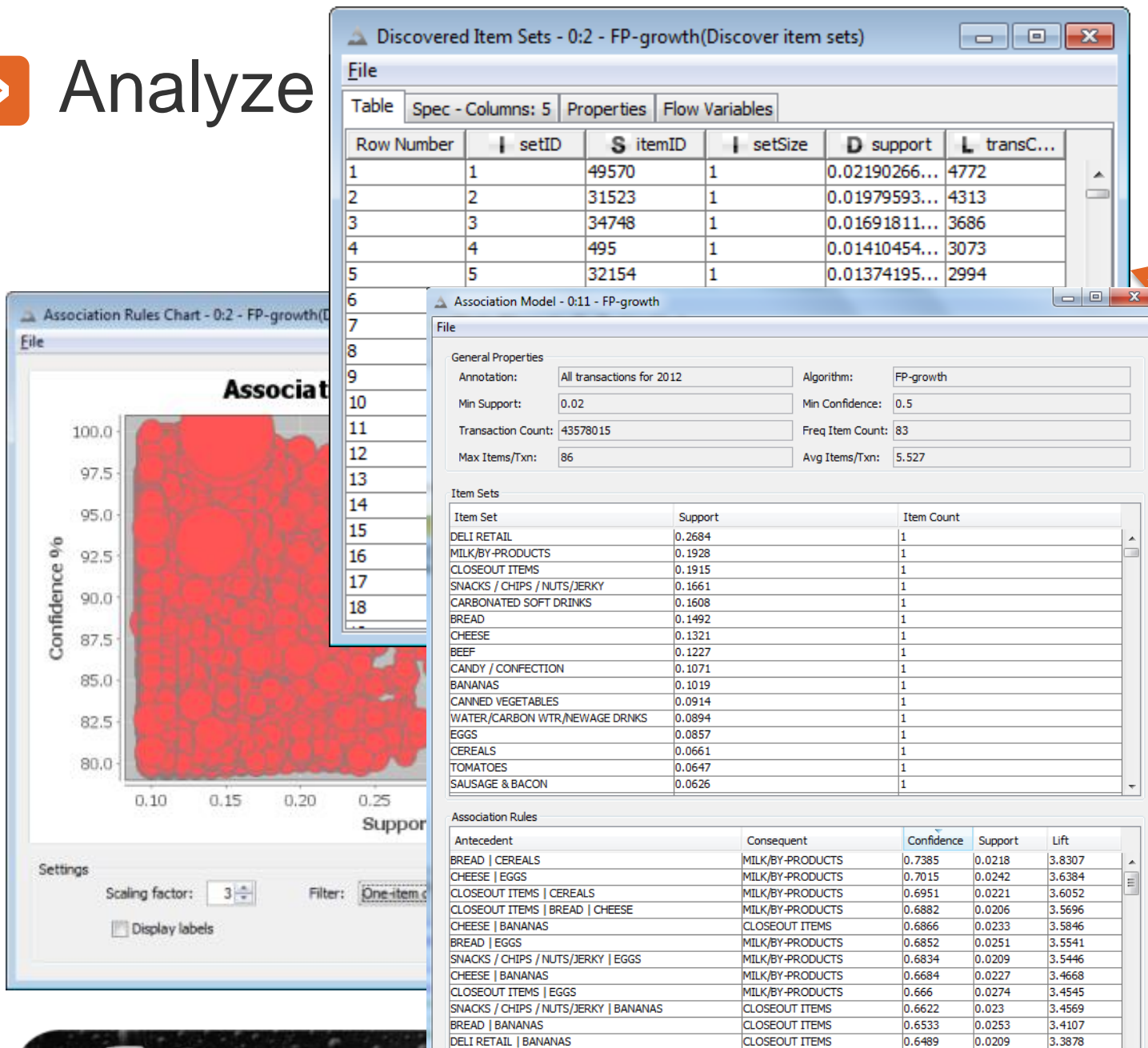
# > Cleanse and enrich







# Analyze



Delimited  
Text Writer



Node 14

PMML Writer



Node 12

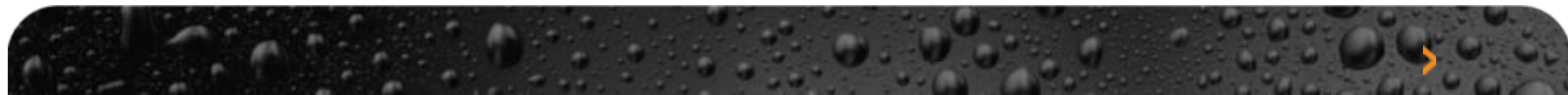
PM Model Converter



Node 15

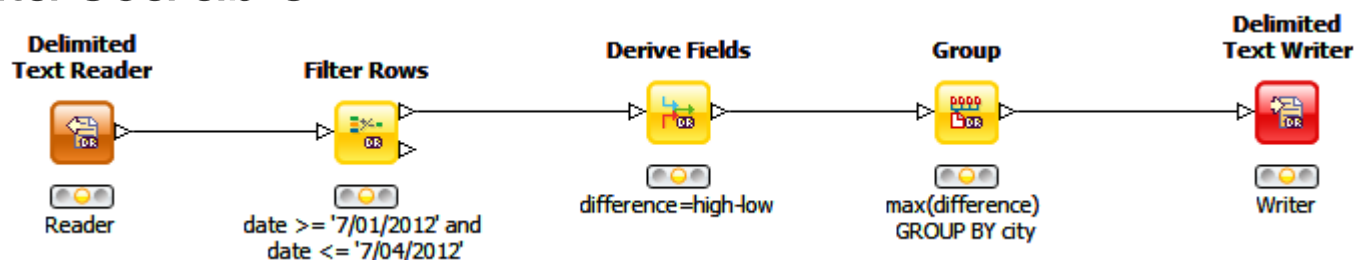
# > History of Dataflow

- Initially developed as next-gen data engine for integration
- Used to be DataRush
- Requirements
  - High data throughput
  - Scalable (data, multicore)
  - Based on dataflow concepts
  - Component based architecture
  - Easy to extend
  - Easily fits in visual development environment
- Embedded in Pervasive products (DataProfiler)
- Extended with SDK for more general use



# > Dataflow Concepts

- Operators (nodes) linked together in a directed graph
- Data flows along edges
- Shared nothing architecture
- Provides pipeline parallelism
- Supports data parallelism
- Data scalable





# Operator Library

- DR I/O
  - Read
    - ARFF Reader
    - Database Reader
    - Delimited Text Reader
    - HBase Reader
    - PMML Reader
    - Staging Importer
  - Write
    - ARFF Writer
    - BigQuery Writer
    - Database Writer
    - Delimited Text Writer
    - HBase Writer
    - PMML Writer
    - Staging Exporter
  - Force Staging
- DR Analytics
  - Association Rules
    - FP-growth
    - Frequent Items
  - Classifiers
    - Decision Tree Learner
    - Decision Tree Predictor
    - Decision Tree Pruner
    - K-Nearest Neighbors Classifier
    - Naive Bayes Learner
    - Naive Bayes Predictor
    - SVM Learner
    - SVM Predictor
  - Clustering
    - k-Means
  - Regression
    - Linear Regression (Learner)
    - Logistic Regression (Learner)
    - Logistic Regression (Predictor)
    - Regression (Predictor)
- DR Transformations
  - Aggregate
    - Group
    - Join
    - Semi/Anti-Join
    - Union All
  - Filter
    - Filter Rows
    - Limit Rows
    - Random Sample
    - Select Fields
  - Manipulation
    - Assert Sorted
    - Date Value Extraction
    - Missing Value
    - Normalize Values
    - Partition Data
    - Rank Fields
    - Regular Expression
    - Run Script
    - Sort
    - Substring
    - Time Difference
    - Trim Whitespace
- DR Data Explorer
  - Data Quality Analyzer
  - Data Summarizer
  - Data Summarizer Viewer
  - Distinct Values
  - DataMatcher
    - Cluster Duplicates
    - Cluster Links
    - Discover Duplicates
    - Discover Links
    - Encode

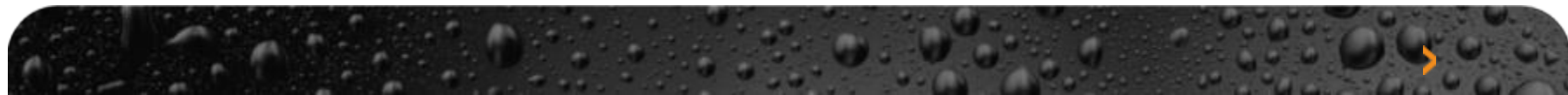
# > Integration with Hadoop

## ■ Data Level

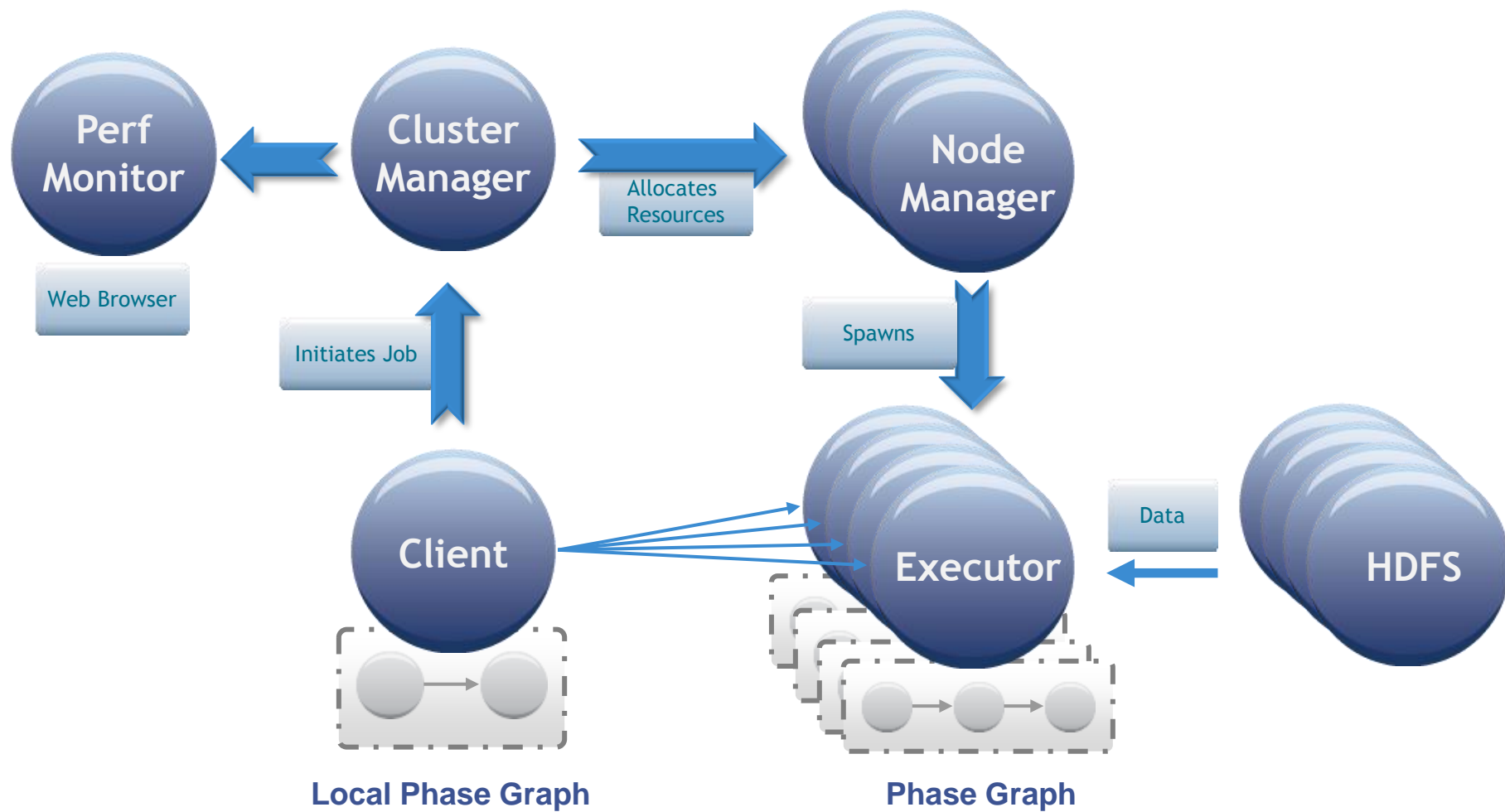
- HDFS access
  - File system abstraction – works with all I/O operators
  - Distributed execution – uses splits much like MR
- HBase
  - Temporal key-value data store based on column families
  - Fast loading using HFile integration
  - Fast temporal queries

## ■ Execution

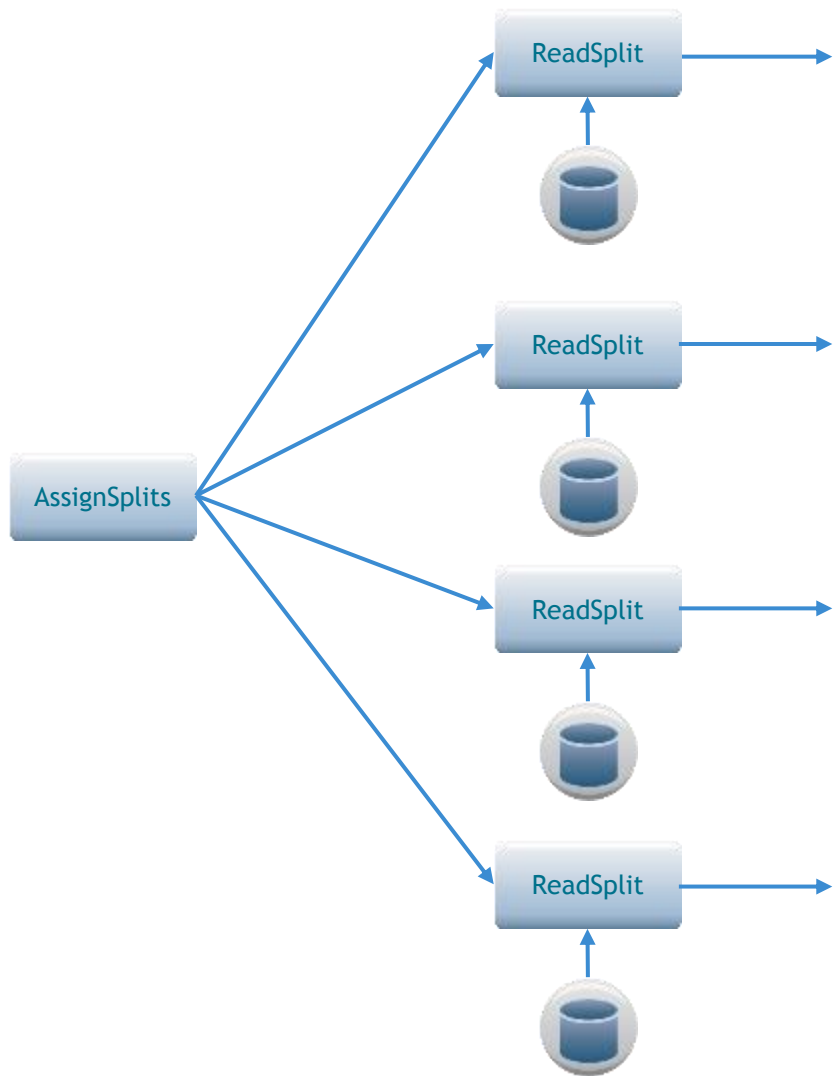
- Distributed execution uses distribute DataRush engines (not MapReduce)
- Integrating with YARN for resource sharing



# > Distributed Execution



## > Distributed I/O



- Allows downstream operators to be parallelized
- Parallelization concepts are the same whether the graph is run locally or distributed

# > Performance Test

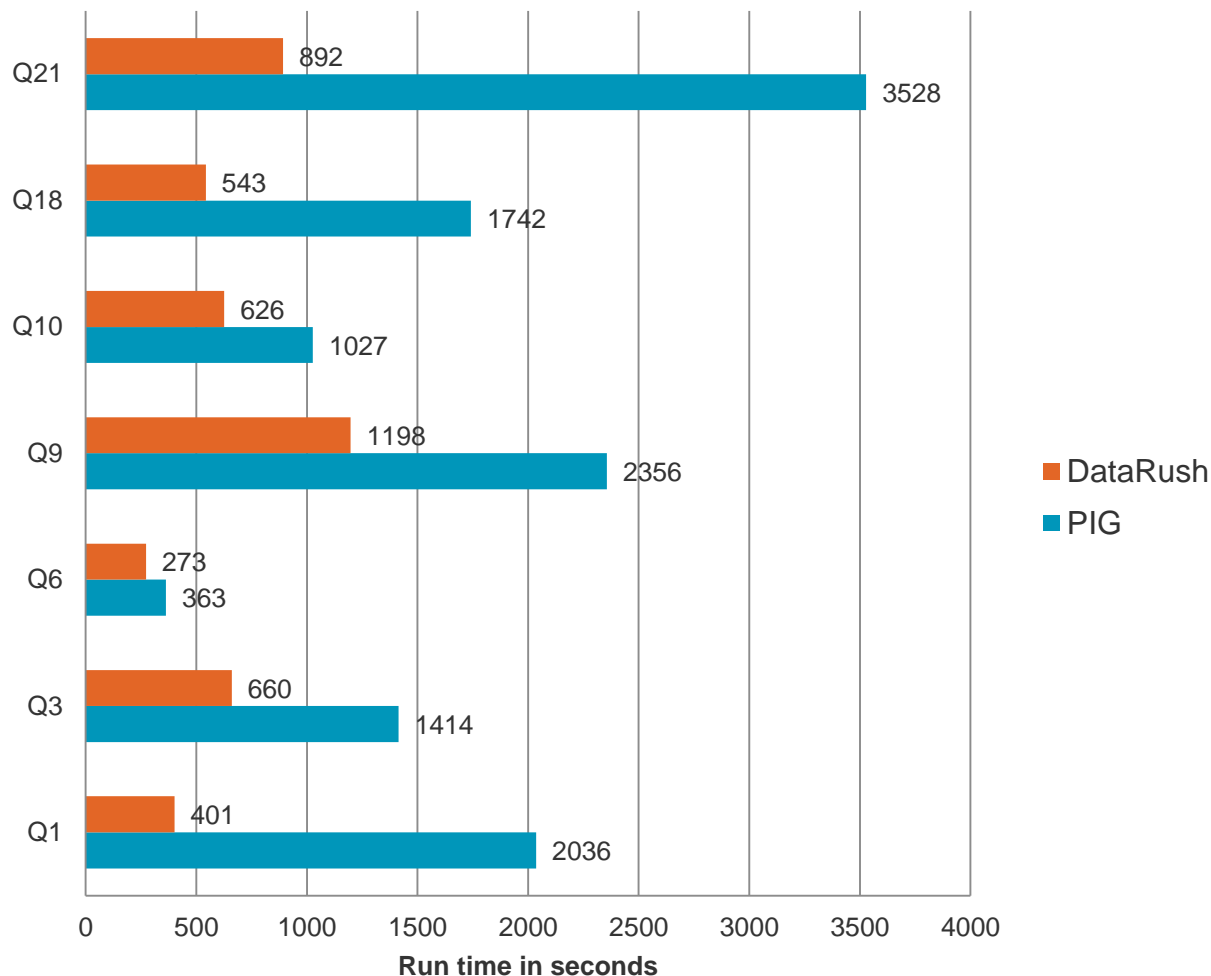
## ■ Dataflow versus PIG

- Used TPC-H data
- Generated 1TB data set in HDFS
- Ran several “queries” coded in Dataflow and PIG
- Run times in seconds (smaller is better)

### Cluster Configuration:

- 5 worker nodes
- 2 X Intel E5-2650 (8 core)
- 64GB RAM
- 24 X 1TB SATA 7200 rpm

## TPC-H : 1 Terabyte Test : Run times



# > Dataflow Analytics Solutions

- Opera Solutions
  - Data science solutions provider
  - Embedding DataRush in engineered solutions
- Healthcare
  - Claims cleansing & processing
- Retail
  - Market basket analysis
  - Product category resolution (MDM)
- Telecom
  - CDR processing & analysis

***“[Dataflow’s] efficiency and ability to automatically scale, whether on a single server or a Hadoop cluster, supports our vision for consistent, reusable, scalable Big Data analytics.”***

– Armando Escalante, Chief Operating Officer, Opera Solutions

