# Pentaho Technical Overview

Max Felber
Solution Engineer
mfelber@pentaho.com
September 22, 2016

![Pentaho - A Hitachi Group Company logo]

# Industry Leader in Self-Service Big Data Preparation

**Gartner.**

# Market Guide for Self-Service Data Preparation

Published: 25 August 2016    ID: G00304870

Analyst(s): Rita L. Sallam, Paddy Forry, Ehtisham Zaidi, Shubhangi Vashisth

Learn how Gartner can help you succeed
Become a Client now ▸

## Summary

This report profiles 36 self-service data preparation products used by analysts and data scientists to accelerate data preparation for analysis, and increasingly by data engineers in data and analytics teams to create trusted, agile, curated data for a range of distributed analytics content authors.

## Overview

### Key Findings

▪ The trend toward ease of use and agility that has disrupted the BI and analytics and advanced analytics markets is also occurring for data integration for analytics.

▪ Most vendor offerings support broad data management capabilities, including interactive data preparation; data exploration, transformation, modeling and curation; and metadata support. Some also offer cataloging, enrichment and intelligent capabilities.

▪ The market is crowded with a range vendor choices, from stand-alone specialists to vendors that embed these tools into BI and analytics, data science and/or data integration platforms.

▪ Although accelerating the shift toward broadly deployed modern, agile BI and advanced analytics, these tools if unchecked can introduce multiple versions of the truth.

### Recommendations

Data and analytics leaders should:

- Gartner recently completed a study on 36 self-service data providers [Gartner Report]
- According to Gartner, a vendor should fulfill the following 4 pillars of self-service data preparation:
  1. Stand-Alone Self-Service Data Preparation
  2. Integrated With Existing Data Integration Platforms
  3. Integrated With Modern BI&A Platforms
  4. Integrated With Advanced Analytics/Data Science Platforms
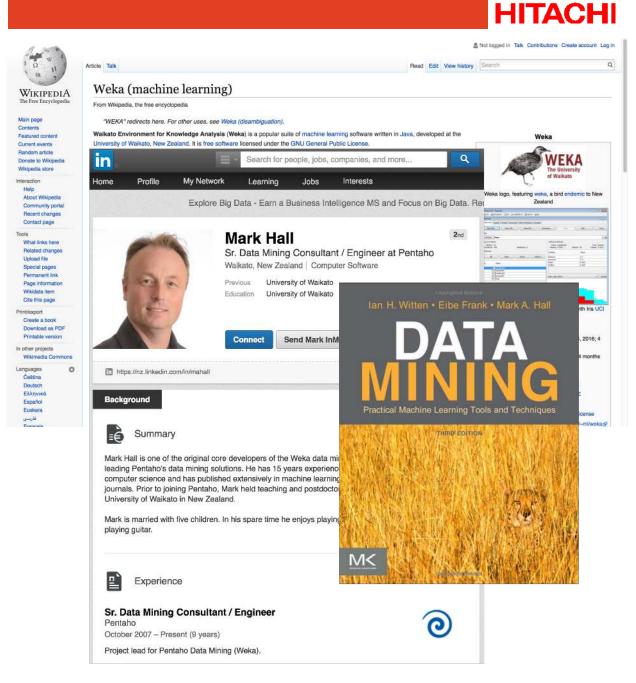- Only 3 vendors met each pillar: Oracle, IBM and Pentaho

- How often have you heard the term Data Lake?
- James Dixon, our founder, invented the term

- Maybe you've heard of Weka for machine learning?
- Mark Hall, the man who leads our data mining services not only developed it, he wrote the book on it

# So why do customers choose Pentaho?

1. Metadata Injection. By utilizing dynamic vs. static bindings, we can reduce the number of data transformations by 80-90 percent. YouTube Demo by Matt Casters, Chief Architect of Data Integration and Kettle Project Founder at Pentaho
2. Visual MapReduce. Graphically build Hadoop data transformations without coding. This enables you to reduce your Hadoop development time by over 85 percent. Additionally, Pentaho automatically manages the deployment and execution of Hadoop transformations with YARN. YouTube Demo by Doug Moran, Product Manager for Big Data Technologies and Co-Founder of Pentaho
3. Embedded Analytics. "White label" our reports, visualizations and dashboards directly into your web applications. Use Java and REST APIs to access Pentaho data transformations and reports. YouTube Demo by Anthony de Shazor, SVP of Customer Care and Principal Architect
4. Beyond ETL. Pentaho supports Enterprise Information Integration (EII), also known as data federation. No you can create ETL jobs that blend data from structured and big data sources and invoke it via JDBC with Teiid
5. Weka Scoring, Forecasting and R script execution with our Data Science Pack. This package helps Data Scientists reduce data preparation times by 60-80 percent
6. Deep Integration with your Cloudera, Hortonworks or MapR Hadoop ecosystem.  Pentaho offers real control over YARN jobs, Spark execution, Oozie, Sqoop and more
7. 2,000 Commercial Customers and 20,000 Production Deployments

# Best of Breed vs. Best Platform



VS.

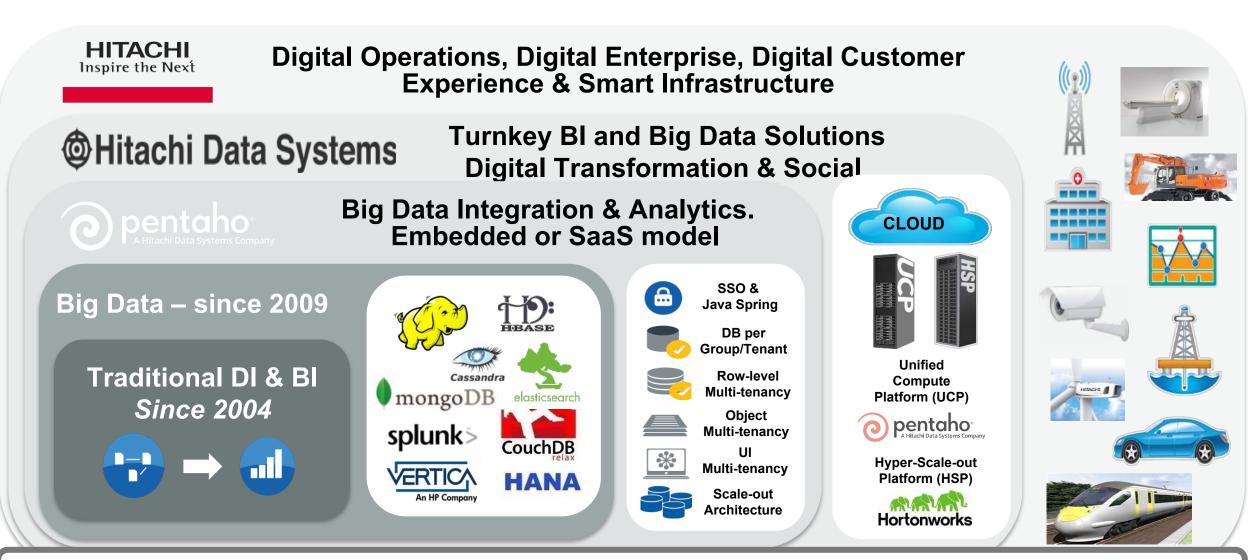# It's all about Balance. Focus on the Integration and Workflow.

# So why do customers buy Pentaho?
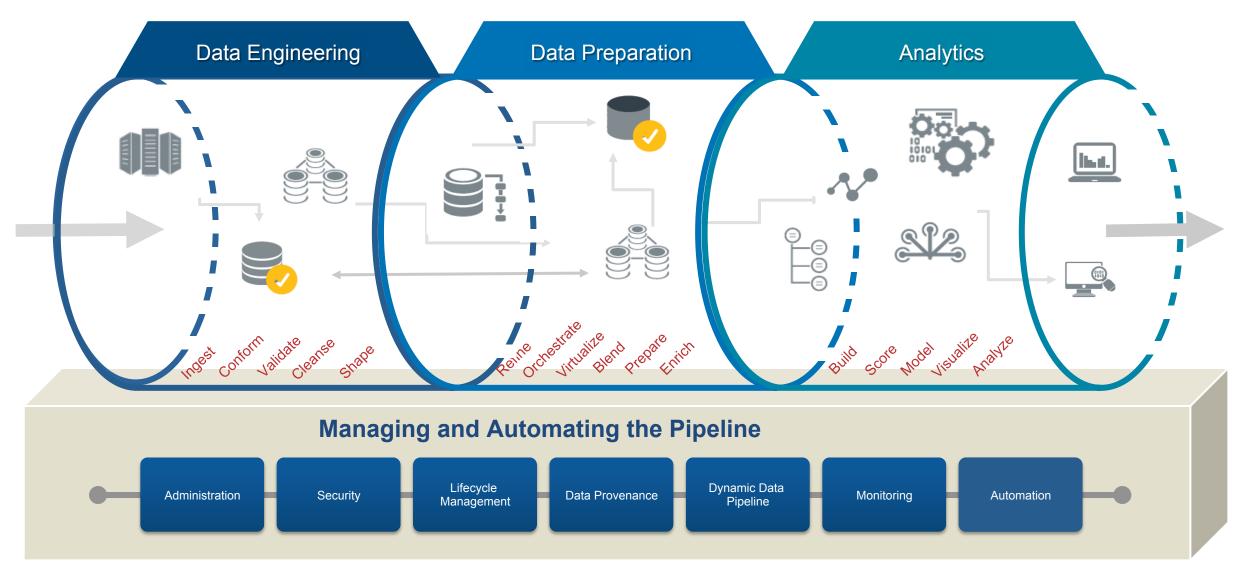
1. Metadata Injection. By utilizing dynamic vs. static bindings, we can reduce the number of data transformations by 80-90 percent. YouTube Demo by Matt Casters, Chief Architect of Data Integration and Kettle Project Founder at Pentaho
2. Visual MapReduce. Graphically build Hadoop data transformations without coding. This enables you to reduce your Hadoop development time by over 85 percent. Additionally, Pentaho automatically manages the deployment and execution of Hadoop transformations with YARN. YouTube Demo by Doug Moran, Product Manager for Big Data Technologies and Co-Founder of Pentaho
3. Embedded Analytics. "White label" our reports, visualizations and dashboards directly into your web applications. YouTube Demo by Anthony de Shazor, SVP of Customer Care and Principal Architect
4. Current Big Data Projects Struggling or Failing
5. Internet of Things (IoT) Initiatives
6. On-Boarding Initiatives (Consolidation, Data Warehousing, SaaS)
7. 360 Degree View (Blending Traditional & Big Data Sources)
8. Predictive Analytics ("R") & Machine Learning ("Weka")
9. Embedding (White-Labeling) Reporting & Analytics
10. ERP Migration
11. Cloud Deployments
12. Data Federation (a.k.a. Enterprise Information Integration)
13. Blended Application and Data Layer Integration (SOA, Web Services, etc.)

# Thank You