



Figure 2: Adaptive Execution With Spark and Visually Designed Hadoop MapReduce Jobs in Pentaho Data Integration

Broad Connectivity and Data Delivery

Pentaho Data Integration offers broad connectivity to a variety of diverse data, including all popular structured, unstructured and semi-structured data sources. Some examples include:

- Relational database management system (RDBMS): Oracle, IBM DB2, MySQL, Microsoft SQL Server, Postgres, IBM MQ
- Spark and Hadoop: Cloudera, Hortonworks, Amazon EMR, MapR (HPE Ezmeral Data Fabric), Microsoft Azure HDInsights, and Elastic Search
- NoSQL databases and object stores: MongoDB, Cassandra, HBase, Hitachi Content Platform, AWS S3, Google Cloud Storage, Microsoft Azure ADLS Gen 2
- Analytic databases: Redshift, Snowflake, Vertica, Greenplum, Teradata, SAP HANA, Amazon Redshift, Google Big Query
- Business applications: SAP, Salesforce, Google Analytics
- Files: XML, JSON, Microsoft Excel, CSV, txt, Avro, Parquet, ORC, EBCDIC (mainframe), unstructured files with metadata, including audio, video and visual files

To increase the performance of data extraction, loading and delivery processes, Pentaho offers the following capabilities:

- Native connectivity and bulk-loading to most common data sources, including Amazon Redshift and Snowflake.
- Data services to virtualize transformations without staging, making data sets immediately available to reports and applications.
- Automatic creation and publishing of metadata models to drive faster analytic results.
- Process streaming data in real-time.

Data Profiling and Data Quality

Pentaho technology provides data profiling capabilities, such as row counts, mathematical functions and identification of null values, as well as data quality operators, such as string manipulators, mapping functions, filtering and sorting. For name and address verification capabilities, Pentaho technology integrates with leading data quality vendors, such as Human Inference and Melissa Data. Pentaho data profiling and data quality capabilities help:

- Identify data that fails to comply with business rules and standards.
- Deduplicate and cleanse inconsistent and redundant data.
- Validate, standardize and correct name, address, email and telephone data.
- Replace file names and locations with simple business names by integrating with Pentaho Data Catalog.

Powerful Administration and Management

Pentaho Data Integration provides out-of-the box capabilities for managing operations for data integration projects. These capabilities include:

- Shared repository for collaboration among data analysts, developers and data stewards.
- Content management, versioning and locking to easily version jobs for roll-back to prior versions.
- Control over security privileges for users and roles and integration with third-party security systems; ability to set permissions for creating, reading or executing jobs and transformations.

[Discover Pentaho Data Integration](#) →

Pentaho Data Integration Technical Specifications and Compatible Components

[Learn More](#)