## Hadoop Open Source Projects

Hadoop is supplemented by an ecosystem of open source projects



PureData Ecosystem

IBM

## How to Analyze Large Data Sets in Hadoop

- Although the Hadoop framework is implemented in Java, MapReduce applications do not need to be written in Java
- To abstract complexities of Hadoop programming model, a few application development languages have emerged that build on top of Hadoop:
  - Pig
  - Hive
  - Jaql







# Pig, Hive, Jaql – Similarities

- Reduced program size over Java
- Applications are translated to map and reduce jobs behind scenes
- Extension points for extending existing functionality
- Interoperability with other languages
- Not designed for random reads/writes or low-latency queries



© 2013 IBM Corporation

TEN

PureData Ecosystem

27

# Pig, Hive, Jaql – Differences

Characteristic	Pig	Hive	Jaql
Developed by	Yahoo!	Facebook	IBM
Language	Pig Latin	HiveQL	Jaql
Type of language	Data flow	Declarative (SQL dialect)	Data flow
Data structures supported	Complex	Better suited for structured data	JSON, semi structured
Schema	Optional	Not optional	Optional

# Pig

- The Pig platform is able to handle many kinds of data, hence the name
- 9 Canactering and a second s

PureData Ecosystem

## Pig

30

- Three steps in a typical Pig program:
  - LOAD
    - Load data from HDFS
  - TRANSFORM
    - · Translated to a set of map and reduce tasks
    - Relational operators: FILTER, FOREACH, GROUP, UNION, etc.
  - DUMP or STORE
    - · Display result on to the screen or store it in a file

#### • Pig data types:

- Simple types:
  - int, long, float, double, chararray, bytearray, boolean
- Complex types:
  - tuple: ordered set of fields

(John,18)`

bag: collection of tuples

{(John,18), (Mary, 29)}

[name#John, phone#1234567]

• map: set of key/value pairs



## Pig

#### Example: wordcount.pig

input = LOAD `./all\_web\_pages' AS (line:chararray); -- Extract words from each line and put them into a pig bag -- datatype, then flatten the bag to get one word on each row words = FOREACH input GENERATE FLATTEN(TOKENIZE(line)) AS word; -- create a group for each word word\_groups = GROUP words BY word; -- count the entries in each group word\_count = FOREACH word\_groups GENERATE COUNT(words) AS count, group; -- order the records by count ordered\_word\_count = ORDER word\_count BY count DESC; STORE ordered\_word\_count INTO `./word\_count\_result';

- How to run wordcount.pig?
  - Local mode:

/bin/pig -x local wordcount.pig

- Distributed mode (MapReduce):

hadoop dfs -copyFromLocal all\_web\_pages input/all\_web\_pages bin/pig -x mapreduce wordcount.pig

31

32

PureData Ecosystem

### Hive

- What is Hive? (1)
  - Data warehouse infrastructure built on top of Hadoop
  - Provides an SQL-like language called HiveQL
  - Allows SQL developers and business analysts to leverage existing SQL skills
  - Offers built-in UDFs and indexing
- What Hive is not?
  - Not designed for low-latency queries, unlike RDBMS such as DB2 and Netezza
  - Not schema on write
  - Not for OLTP
  - Not fully SQL compliant, only understand limited commands



© 2013 IBM Corporation

## Hive

- Components:
  - Shell
  - Driver
  - Compiler
  - Engine
  - Metastore
    - Holds table definition, physical layout

### Data models:

- Tables
  - Analogous to tables in RDBMS, composed of columns
- Partitions
  - · For optimizing data access, e.g. range partition tables by date
- Buckets
  - Data in each partition may in turn be divided into Buckets based on the hash of a column in the table

© 2013 IBM Corporation

PureData Ecosystem

## Hive

33

### Example: movie ratings analysis



TEM

### Jaql

Designed for easy manipulation and analytics of semi-structured data, support formats such as JSON, XML, CSV, flat files, etc. Developed by TEM. Flexibility with optional schema

Easy extensibility

35 © 2013 IBM Corporation

PureData Ecosystem

## Jaql – How does a Jaql query work?

- A Jaql query can be thought of as a pipeline
- Data manipulation through operators

   FILTER, TRANSFORM, GROUP, JOIN, EXPAND, SORT, TOP



## Jaql

In addition to core operators, Jaql also provides built-in functions

#### Data models:

- Atomic types: boolean, string, long, etc.
- Complex types: array, record

### Where to run Jaql queries?

- Shell: jaqlshell
  - Cluster mode
  - Local mode: for testing, prototyping purposes
- Eclipse
- Embedded in Java



© 2013 IBM Corporation

IRM

PureData Ecosystem

## Jaql

37

Example: find employees who are manager or have income > 50000



### **Query**

read(hdfs "employees"));

### <u>Data</u>

## Jaql

### Example: find employees who are manager or have income > 50000



PureData Ecosystem

### Jaql

Example: find employees who are manager or have income > 50000



## Jaql

### Example: find employees who are manager or have income > 50000



PureData Ecosystem

IBM

# Other Hadoop Related Projects

#### Data serialization:

Avro

· uses JSON schemas for defining data types, serializes data in compact format

#### • Monitoring:

- Chukwa
  - Built on top of HDFS and MapReduce
  - · Structured as a pipeline of collection and processing stages

#### Distributed coordination:

- ZooKeeper
  - Distributed configuration, synchronization, naming registry

#### Jobs management:

- Oozie
  - · Simplifies workflow and coordination of MapReduce jobs

#### Structured data storage:

- HBase
  - · Column-oriented non-relational distributed database built on top of HDFS

### Summary

- Apache Hadoop is a software framework for distributed processing of large data sets across clusters of computers
- Two major components of Hadoop:
  - MapReduce programming model
  - Hadoop Distributed File System
- Hadoop is supplemented by an ecosystem of projects



Questions?

E-mail: Subject: impe.biginsights@ca.ibm.com Big data bootcamp

