

CS 378 – Big Data Programming

Lecture 11

AVRO Formats

Review

- Assignment 5 – AVRO Objects
- Questions/Issues?
 - datum()
- Using the latest pom.xml

AVRO File Formats

- `TextOutputFormat`
 - How are various key and value types handled?
 - Recall that `TextOutputFormat` will cause `toString()` to be called
- `AvroKey<CharSequence>`
 - Acts like `Text`, so it just returns its string value
- `AvroValue<WordStatisticsData>`
 - Returns the value created by the `toString()` method

AVRO File Formats

- `TextOutputFormat`
- `AvroKey<Pair<CharSequence, WordStatisticsData>>`
 - Used as the key to the `write()` method, with value `NullWritable`
 - Generates a string representation in the `toString()` method of `Pair`
 - In this form: `{ "key": theKey, "value": theValue }`
 - *theKey* comes from `CharSequence`, so just a string
 - *theValue* comes from `WordStatisticsData`, so an AVRO text representation is generated (calls `toString()`)
 - `{ "document_count": }`

AVRO File Formats

- `AvroKeyValueOutputFormat`
 - Creates a generic Avro record with a “key” field and a “value” field
 - Like what we saw with `AvroKey<Pair< K, V >>`
 - Avro container file (binary)
 - Can be read in using: `AvroKeyValueInputFormat`

AVRO File Formats

- `AvroKeyOutputFormat<T>`
 - Extends
 - `AvroOutputFormatBase (AvroKey<T>, NullWritable>)`
 - Only the key is output, value is ignored
 - AVRO container file (binary format)
 - Can be read in using: `AvroKeyInputFormat`

AVRO File Formats

- `AvroSequenceFileOutputFormat`
 - Sequence file output format that can handle `AvroKey` and `AvroValue` in addition to `Writable`
 - Can be read with: `AvroSequenceFileInputFormat`

AVRO File Formats

- `AvroKeyValueInputFormat`
 - Reads generic Avro records with a “key” field and a “value” field
 - AVRO container file (binary)
 - Data should have been written with:
`AvroKeyValueOutputFormat`

AVRO File Formats

- `AvroKeyInputFormat`
 - Extends
 - `FileInputFormat(AvroKey<T>, NullWritable>)`
 - Only the key is read, value is ignored
 - Reads a AVRO container file (binary format)
 - Data should have been written with:
`AvroKeyOutputFormat`

AVRO File Formats

- `AvroSequenceFileInputFormat`
 - Input format that can read sequence files that support Avro types
 - Data should have been written with:
`AvroSequenceFileOutputFormat`

Design Pattern

- Structured to hierarchical design pattern
- Data sources linked by some foreign key
- Data is structured and row based
 - For example, from databases
- Data is semi-structured and event based
 - Web logs

Design Pattern

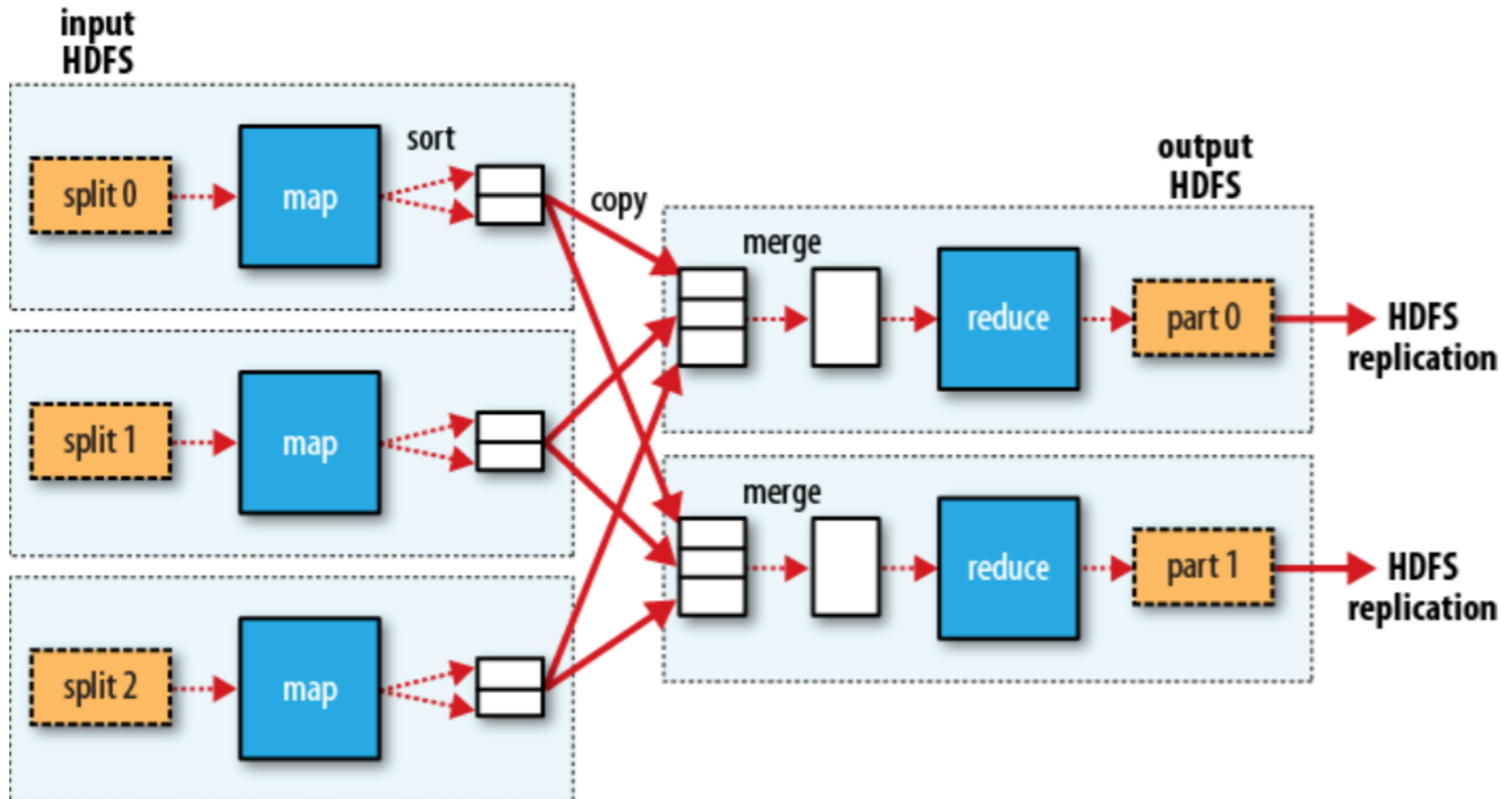
- Structured to hierarchical design pattern
- MultipleInputs
 - Able to accept data inputs from different formats
 - Mappers load and parse the input into a cohesive format
 - Prepared for work in the reducer
 - Map output key will be the unifying element of the hierarchical record
- Combiners don't help, as they don't "reduce" the data (make it smaller)

Design Pattern

- Structured to hierarchical design pattern
- Reducer takes all the data associated with a key
- Builds the structure to be output
- Example:
 - User session contains info about the user (IP, browser, ...)
 - An array of actions (page views, clicks, ...)

MapReduce in Hadoop

Figure 2.4, Hadoop - The Definitive Guide



Sessionizing Web Logs

- Create user sessions from individual web log entries
- Represents all the actions by a user
- Allows later analysis to “replay” the user actions
- Collect measures and metrics about user behavior
 - Pages viewed, time on page, clicks
 - Path through the site, entry to the site (from a search engine?)