# PROCESSING BIG DATA

# Introduction to Processing Big Data

- **Moving Computation to Data**

- Parallel Operations

- Distributed Aggregations

- Accumulating and Broadcasting

- Working with Records of Pairs

- Multi-Dataset Operations

- From Concepts to Frameworks

# Moving Computation to Data

- Some typical properties of Big Data:
    - Data is split into multiple partitions residing in distant disks
    - The dataset is a huge collection of records
    - No orderings should be assumed a-priori
- Moving the data around is intolerably expensive

analyticscenter

**Input Split:**
R1, 1
R1, 2
...
R1, $n_1$

**Input Split:**
R2, 1
R2, 2
...
R2, $n_2$

**Input Split:**
R3, 1
R3, 2
...
R3, $n_3$

**Input Split:**
R4, 1
R4, 2
...
R4, $n_4$

## Distributed Collection
**The entire collection is split into partions, possibly residing on multiple machines**

analyticscenter

# Moving Computation to Data

- The key is moving computation to data, that is, performing local computations in parallel and reducing the amount of communication when data needs to be transferred across nodes

- The MapReduce programming model makes complex computations on distributed collections possible in a parallel and distributed fashion with a minimal amount of communication

  - Communication effectiveness (reducing the amount of data moving between nodes) will be key to implementing such operations on Big Data

analyticscenter

# Introduction to Processing Big Data

- Moving Computation to Data

- **Parallel Operations**

- Distributed Aggregations

- Accumulating and Broadcasting

- Working with Pairs of Records
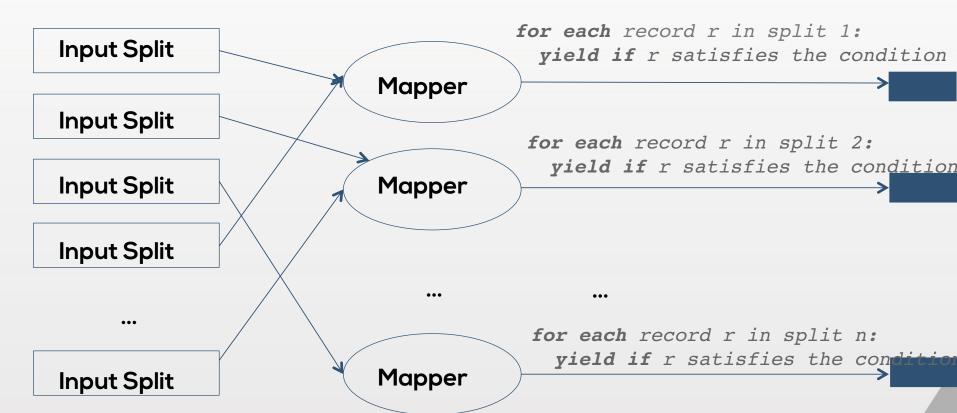
- Multi-Dataset Operations

- From Concepts to Frameworks

# Parallel Operations

- On a distributed collection of data, operations such as:

  - **Selection**

  - **Projection** can be implemented in parallel

```
//Pseudo-code for selection (filter out records that
//does not satisfy a boolean condition, i.e. where
//clause in SQL)

for each split i in parallel:
  for each record j in split i:
    if the record satisfies the condition
      yield the record
    end-if
  end-for
```

Input                    Read and process in parallel                           Output

| Input Split |

| Input Split |                    **Mapper**          *for each* record r in split 1:
                                                        *yield if* r satisfies the condition

| Input Split |

| Input Split |                    **Mapper**          *for each* record r in split 2:
                                                        *yield if* r satisfies the condition

...                                      ...              ...

| Input Split |                    **Mapper**          *for each* record r in split n:
                                                        *yield if* r satisfies the condition

**Select in parallel (WHERE predicate):**
**filter as a higher order function**

analyticscenter

Input    Read and process in parallel    Output

Input Split

Input Split

Input Split

Input Split

...

Input Split

Mapper

Mapper

...

Mapper

*for each* record r in split 1:
transform r and **yield**

*for each* record r in split 2:
transform r and **yield**

...

*for each* record r in split n:
transform r and **yield**

**Project in parallel (SELECT clause):**
**map as a higher order function**

analyticscenter

Input                    Read and process in parallel                                      Output

| Input Split |

*for each* record r in split 1:
*generate multiple records from r*
**yield each one of them**

( Mapper )

| Input Split |

*for each* record r in split 2:
*generate multiple records from r*
**yield each one of them**

( Mapper )

| Input Split |

| Input Split |

...

...

...

*for each* record r in split n:
*generate multiple records from r*
**yield each one of them**

| Input Split |

( Mapper )

**FlatMap in parallel**
**(in the functional programming sense)**

analyticscenter

# Parallel Operations

- Applying parallel operations is no hassle

    - Provided a parallel computation framework that runs higher order functions on records of a distributed collection

analyticscenter

# Introduction to Processing Big Data

- Moving Computation to Data

- Parallel Operations

- **Distributed Aggregations**

- Accumulating and Broadcasting

- Working with Pairs of Records

- Multi-Dataset Operations

- From Concepts to Frameworks

analyticscenter

# Distributed Aggregations

- Often, we need to compute aggregates over the collection

- We consider two strategies for aggregating distributed collections:

  - The aggregate operation is distributive

  - The aggregate operation is holistic

analyticscenter

# Distributed Aggregations

- If the aggregation is holistic (i.e. all data must be seen by the aggregate operation), the amount of communication is large

  - All record-level relevant data should be copied to the same processing unit that computes the aggregate

  - e.g. **median** is such a measure

analyticscenter

# Distributed Aggregations

```
//Local count function in parallel

for each split i in parallel:
  for each record j in split i:
    yield 1
  end-for
```

# Distributed Aggregations

```
//Global aggregate for count

input: [1, 1, 1, …, 1]
yield sum(input)
```

Input          Read

*for each* record r in split i:
*yield 1*

Input Split

Input Split

Input Split

Input Split

...

Input Split

[1,1, …, 1]

[1,1, …, 1]

∑

[1,1, …, 1]

...          ...

**Distributed Count**

analyticscenter

# Distributed Aggregations

- If the aggregation is distributive, we can perform aggregates with minimum amount of communication

- The trick is to

    - compute local aggregates

    - transfer all local aggregates to the same node

    - compute the global aggregate

# Distributed Aggregations

- We refer to

  - associative (same aggregate of local aggregates)

  - distributive (different aggregate of local aggregates)

  - algebraic (a function of multiple distributive aggregates)

operations as trivially distributive aggregations

# Distributed Aggregations

- Many aggregate operations can be distributed (communication) efficiently, such as:

    – **count** (sum of counts)

    – **max** (max of maxes)

    – **sum** (sum of sums)

    – **average** (sum/count)

# Distributed Aggregations

```
//Local count function in parallel

for each split i in parallel:
  for each record j in split i:
    yield 1
  end-for

  take a sum of all 1's
```

# Distributed Aggregations

```
//Global aggregate for count

input: [c1, c2, …, cm]
yield sum(input)
```

Input    Read    **for each** `record r in split i:`
                 **yield 1**

Input Split

Input Split

Input Split

Input Split

...

Input Split

Pre-aggregate

Σ    3

Σ    6

...    ...

Σ    4

Σ

**Distributed Count**

analytics center
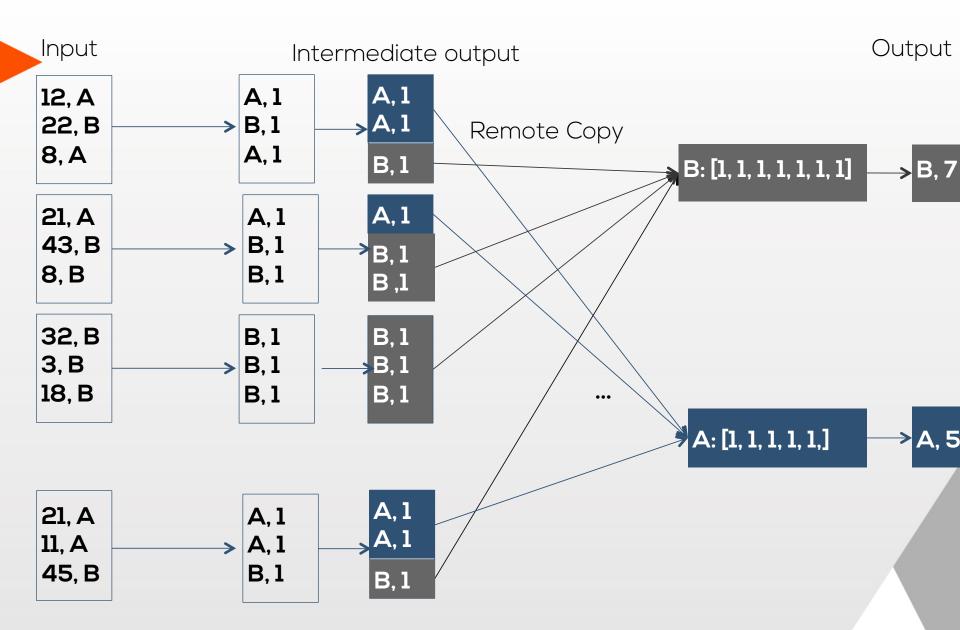
# Introduction to Processing Big Data

- Moving Computation to Data

- Parallel Operations

- Distributed Aggregations

- **Accumulating and Broadcasting**

- Working with Records of Pairs

- Multi-Dataset Operations

- From Concepts to Frameworks

analyticscenter

# Accumulating

- Sometimes a global counter, available to all parallel execution nodes, can be defined to be incremented by parallel operations, in an atomic way

- An example use case is counting the number of bad records in the collection

# Broadcasting

- A small data structure can be broadcast into all parallel execution node

- This might be a small data set or a configuration, for example

- Such variables are naturally read-only

# Introduction to Processing Big Data

- Moving Computation to Data

- Parallel Operations

- Distributed Aggregations

- Accumulating and Broadcasting

- **Working with Records of Pairs**

- Multi-Dataset Operations

- From Concepts to Frameworks

analytics center
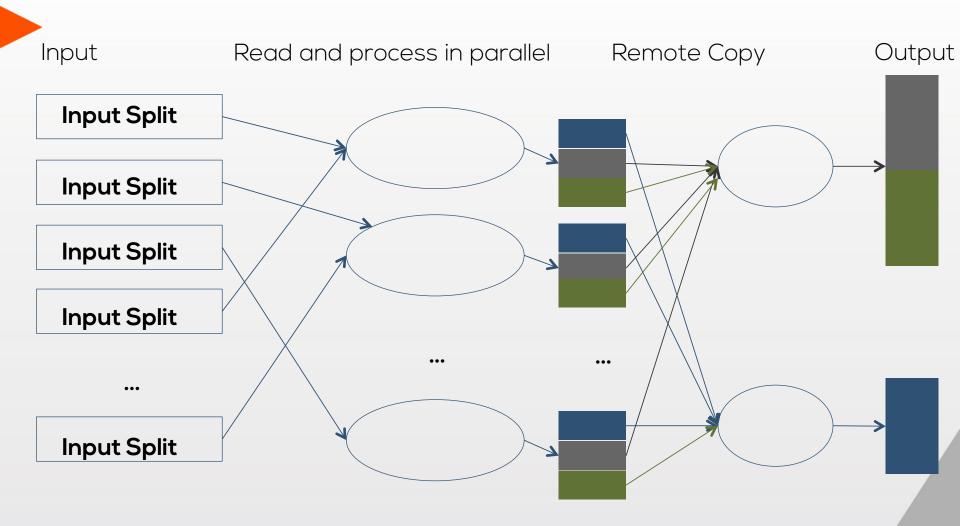
# Working with Records of Pairs

- Often, we work with datasets of records of pairs:

  – records are in `<k, v>` form

  – records are mapped into `<k, v>` form

  – the required aggregate is by `K`

- In this case, the aggregate operations can be performed on a per `K` basis

- To have multiple aggregation units running in parallel:

  – the results from the local operations should be somehow grouped by the `K`

  – all local results with the same `K` should be copied to the same aggregation unit

analytics center

Input

Intermediate output

Output

| 12, A |
| 22, B |
| 8, A |

A, 1
B, 1
A, 1

A, 1
A, 1
B, 1

Remote Copy

B: [1, 1, 1, 1, 1, 1, 1]

B, 7

| 21, A |
| 43, B |
| 8, B |

A, 1
B, 1
B, 1

A, 1
B, 1
B ,1

| 32, B |
| 3, B |
| 18, B |

B, 1
B, 1
B, 1

B, 1
B, 1
B, 1

...

A: [1, 1, 1, 1, 1,]

A, 5

| 21, A |
| 11, A |
| 45, B |

A, 1
A, 1
B, 1

A, 1
A, 1
B, 1

**Counting the numbers of A's and B's**
**GROUP-BY-COUNT**

29
analyticscenter

Input    Read and process in parallel    Remote Copy    Output

Input Split

Input Split

Input Split

Input Split

...

Input Split

...    ...

2 tasks, 3 groups

**Working with Records of Pairs (key-values)**

analytics center
30

# Working with Records of Pairs

- Again, we can avoid copying all per-record results to the aggregation units if the aggregation permits to do so

- We could have computed the local sums in the previous example, and reduce the amount of integers copied over network

analyticscenter

Input

Intermediate output

Output

| 12, A |
| 22, B |
| 8, A |

| A, 1 |
| B, 1 |
| A, 1 |

A, 2
B, 1

Remote Copy

B: [1, 2, 3, 1] → B, 7

| 21, A |
| 43, B |
| 8, B |

| A, 1 |
| B, 1 |
| B, 1 |

A, 1
B, 2

| 32, B |
| 3, B |
| 18, B |

| B, 1 |
| B, 1 |
| B, 1 |

B, 3

...

A: [2, 1, 2] → A, 5

| 21, A |
| 11, A |
| 45, B |

| A, 1 |
| A, 1 |
| B, 1 |

A, 2
B, 1

**Counting the numbers of A's and B's**
**(with pre-aggregation) -- GROUP-BY-COUNT**

analyticscenter

32

# Working with Records of Pairs

- Notice that when we work with records of pairs, the skew of the dataset by K's greatly affect:

    - the amount of data processed by

    - the amount of data copied into

    - the amount of work done by

an aggregation unit

# Introduction to Processing Big Data

- Moving Computation to Data

- Parallel Operations

- Distributed Aggregations

- Working with Records of Pairs
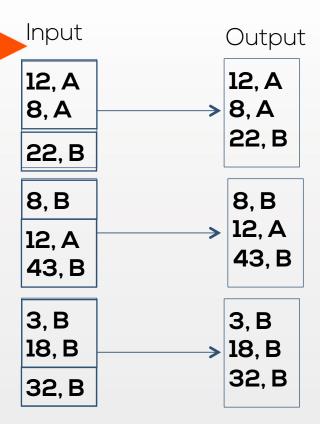
- **Multi-Dataset Operations**

- From Concepts to Frameworks

analyticscenter

# Multi-Dataset Operations

- Sometimes, we would like to perform operations yielding a single dataset from multiple datasets, such as:

  - Set operations like union and intersection

  - Cross-dataset operations like cartesian product and join

- Without loss of generality, we will assume 2 datasets, split into blocks on the same cluster

  - We may further assume a 'combined split' per node, including records from both splits of the underlying datasets
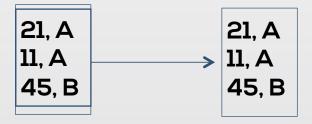
# Multi-Dataset Operations

- Union (duplicate records allowed) can trivially be parallelized, and there is no need of aggregation

    - Union can be implemented without a need of copying intermediate data over network
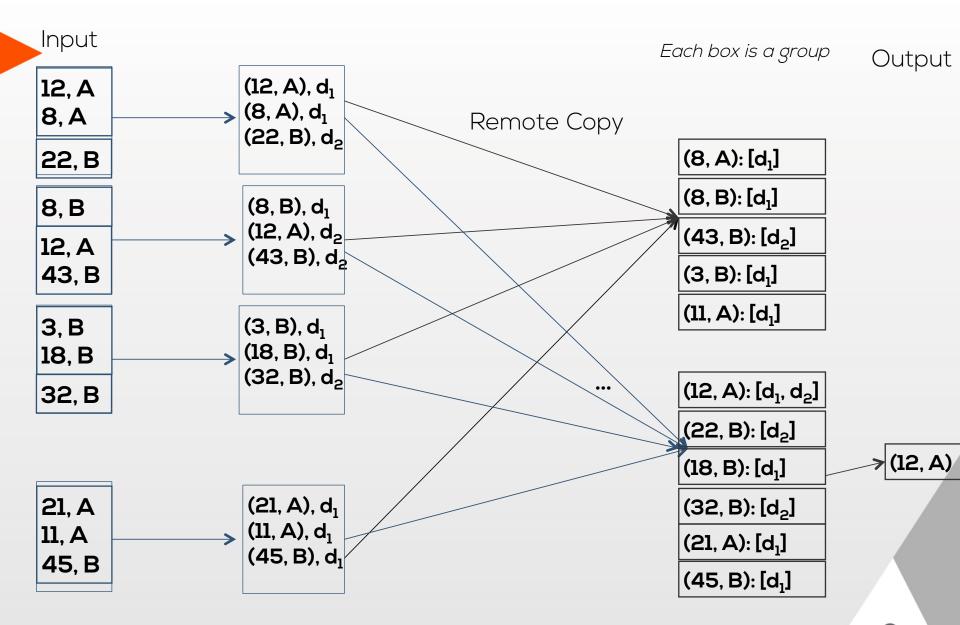
    - It is a **map-only** job

Input

Output

| 12, A |   | 12, A |
|---|---|---|
| 8, A | → | 8, A |
| 22, B |   | 22, B |

| 8, B |   | 8, B |
|---|---|---|
| 12, A | → | 12, A |
| 43, B |   | 43, B |

| 3, B |   | 3, B |
|---|---|---|
| 18, B | → | 18, B |
| 32, B |   | 32, B |

| 21, A |   | 21, A |
|---|---|---|
| 11, A | → | 11, A |
| 45, B |   | 45, B |

**UNION**

analytics center

# Multi-Dataset Operations: Union

- Things are not that easy with intersection

  - A record belonging to the resulting dataset can occur in different splits

  - thus we need to copy the intermediate records to the aggregation units, and only keep a record if it is contained in both datasets

    - This is a large amount of network

Input

Each box is a group

Output

12, A
8, A

22, B

8, B

12, A
43, B

3, B
18, B

32, B

21, A
11, A
45, B

(12, A), $d_1$
(8, A), $d_1$
(22, B), $d_2$

(8, B), $d_1$
(12, A), $d_2$
(43, B), $d_2$

(3, B), $d_1$
(18, B), $d_1$
(32, B), $d_2$

(21, A), $d_1$
(11, A), $d_1$
(45, B), $d_1$

Remote Copy

(8, A): [$d_1$]

(8, B): [$d_1$]

(43, B): [$d_2$]

(3, B): [$d_1$]

(11, A): [$d_1$]

...

(12, A): [$d_1$, $d_2$]

(22, B): [$d_2$]

(18, B): [$d_1$]

(32, B): [$d_2$]

(21, A): [$d_1$]

(45, B): [$d_1$]

(12, A)

**INTERSECTION**

analyticscenter

# Multi-Dataset Operations on Records of Pairs: Joins

- Join-like operations can greatly be optimized depending on:
  - The size of the datasets
  - Whether or not the datasets are sorted, or how they are initially partitioned into the nodes
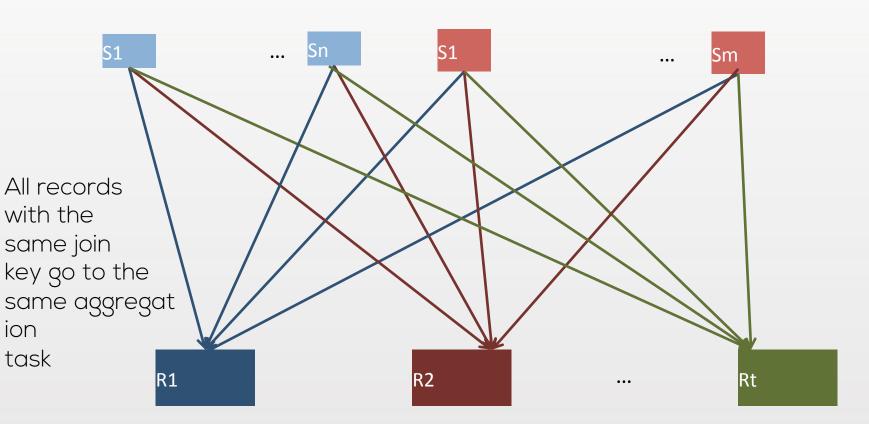
# Multi-Dataset Operations on Records of Pairs: Shuffle Joins

- When we have no idea of how datasets are partitioned, and the sizes of both are big, we resort to shuffle joins
  - The parallel step would yield
    - the join key
    - the appropriate projection of the row
    - the dataset id
  - The records with the same join key would be copied into the same aggregation units, and these units would yield the cartesian product of all the records after repartitioning the group of projected rows with the same keys per dataset

analyticscenter

# Multi-Dataset Operations on Records of Pairs: Shuffle Joins



All records with the same join key go to the same aggregation task
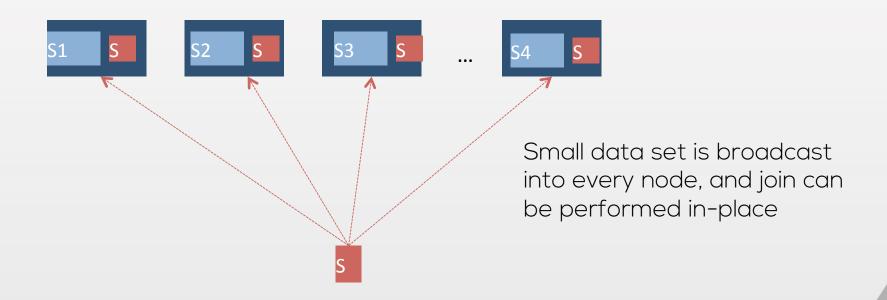
# Multi-Dataset Operations on Records of Pairs: Broadcast Joins

- If we knew that one of the datasets is small

  - We can broadcast the small dataset into all nodes

  - The join can be performed in parallel, without need of any copying

analyticscenter

# Multi-Dataset Operations on Records of Pairs: Broadcast Joins

| S1 | S | | S2 | S | | S3 | S | ... | S4 | S |

Small data set is broadcast into every node, and join can be performed in-place
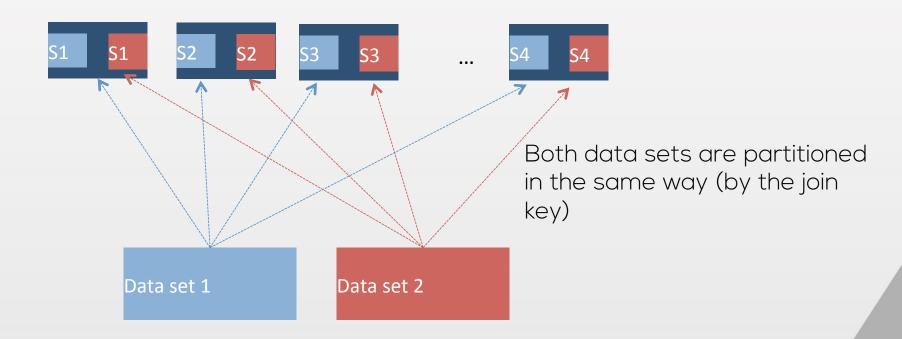
S

analyticscenter

# Multi-Dataset Operations on Records of Pairs: Merge Joins

- If the datasets are partitioned (and sorted) in the same way, we ensure that the rows that should be joined are always included in the splits on the same node

- Then we can perform local joins and report the resulting dataset

analyticscenter

# Multi-Dataset Operations on Records of Pairs: Merge Joins

S1  S1    S2  S2    S3  S3    ...    S4  S4

Both data sets are partitioned in the same way (by the join key)

Data set 1

Data set 2

# Introduction to Processing Big Data

- Moving Computation to Data

- Parallel Operations

- Distributed Aggregations

- Accumulating and Broadcasting

- Working with Records of Pairs

- Multi-Dataset Operations

- **From Concepts to Frameworks**

analyticscenter

# From Concepts to Frameworks

- The described concepts are implemented by many frameworks, including

  – Apache Hadoop MapReduce

  – Apache Spark

- We only need to implement the local operations and the rest are performed automatically by the framework:

  – Parallel execution

  – Copying intermediate data

  – Ensuring all records with the same key are collected together

  – Execution of the aggregate operations

analyticscenter

# MapReduce

- Hadoop MapReduce allows us writing distributed programs, which process vast amounts of data

  - in a parallel and distributed manner

  - on large (of thousands of nodes) clusters

  - reliably

# MapReduce

- MapReduce defines an InputFormat for splitting data into chunks, and iterating over records of pairs of these chunks

- MapReduce framework then;

  - Runs **user defined map functions** on the input records in a completely parallel manner

  - **Shuffles** (ensuring map outputs with the same key are kept together) the intermediate map outputs and copies them to the **Reducers**

  - **Sorts** and passes the map outputs with the same key **to the user defined reduce** functions

  - Reports the **Reducer** outputs

- The framework takes care of scheduling and monitoring tasks, and re-executing the failed ones

# MapReduce

- Typically, the input(s) and output formats are files stored in the HDFS

  - Although this is the typical case, as long as the correct I/O formats defined, other sources and sinks are available

  - Typical MapReduce inputs are in the form of:

    - Plain text files (possibly compressed)

    - SequenceFiles

    - Avro Files

    - Parquet Files

    - ORC Files

    - HBase tables

  each of which is associated with different `InputFormat`s

# MapReduce

- Typically, the storage nodes (i.e., the DataNodes in the HDFS cluster) and compute nodes of MapReduce are the same

  - The MapReduce framework and the HDFS are running on the same cluster of nodes

  - That allows Map tasks to run on local data –**data locality**

analyticscenter

# MapReduce

- Applications specify

  - the I/O Formats,

  - input/output locations

  - *map* and *reduce* functions by implementing the appropriate interfaces (`Mapper#map`, `Reducer#reduce`)

    - To transform/filter records, `map` functions are defined to run on each input record

    - `reduce` functions are defined such that they run on groups of values attached to a key

analyticscenter

# MapReduce

- In MapReduce, broadcasting can be performed using the **DistributedCache**

- Accumulators can be implemented using MapReduce **Counters**

# Processing Big Data    End of Chapter

analyticscenter