

MapReduce

Distributed computing
can get **very complicated**

Managing resources and
memory across multiple nodes

MapReduce

Distributed computing
can get **very complicated**

What to do if a node
goes down?

MapReduce

MapReduce abstracts
the programmer from
all these complications

MapReduce

You just define 2
functions

map() reduce()

MapReduce

map() reduce()

The rest is taken care
of by Hadoop!

Let's take an example

We have a large text file

next up previous contents index
Next: Dynamic indexing Up: Index construction Previous: Single-pass in-memory indexing Contents Index

Distributed indexing

Collections are often so large that we cannot perform index construction efficiently on a single machine. This is particularly true of the World Wide Web for which we need large computer clusters [+] to construct any reasonably sized web index. Web search engines, therefore, use distributed indexing algorithms for index construction. The result of the construction process is a distributed index that is partitioned across several machines – either according to term or according to document. In this section, we describe distributed indexing for a term-partitioned index. Most large search engines prefer a document-partitioned index (which can be easily generated from a term-partitioned index). We discuss this topic further in Section 20.3 (page [+]).

The distributed index construction method we describe in this section is an application of MapReduce, a general architecture for distributed computing. MapReduce is designed for large computer clusters. The point of a cluster is to solve large computing problems on cheap commodity machines or nodes that are built from standard parts (processor, memory, disk) as opposed to on a supercomputer with specialized hardware. Although hundreds or thousands of machines are available in such clusters, individual machines can fail at any time. One requirement for robust distributed indexing is, therefore, that we divide the work up into chunks that we can easily assign and – in case of failure – reassign. A master node directs the process of assigning and reassigning tasks to individual worker nodes.

The map and reduce phases of MapReduce split up the computing job into chunks that standard machines can process in a short time. The various steps of MapReduce are shown in Figure 4.5 and an example on a collection consisting of two documents is shown in Figure 4.6. First, the input data, in our case a collection of web pages, are split into \$n\$ splits where the size of the split is chosen to ensure that the work can be distributed evenly (chunks should not be too large) and efficiently (the total number of chunks we need to manage should not be too large); 16 or 64 MB are good sizes in distributed indexing. Splits are not preassigned to machines, but are instead assigned by the master node on an ongoing basis: As a machine finishes processing one split, it is assigned the next one. If a machine dies or becomes a laggard due to hardware problems, the split it is working on is simply reassigned to another machine.

Figure 4.5: An example of distributed indexing with MapReduce. Adapted from Dean and Ghemawat. (2004).
\includegraphics[width=11.5cm]{art/mapreduce2.eps}
In general, MapReduce breaks a large computing problem into smaller parts by recasting it in terms of manipulation of key-value pairs. For indexing, a key-value pair has the form (termID, docID). In distributed indexing, the mapping from terms to termIDs is also distributed and therefore more complex than in single-machine indexing. A simple solution is to maintain a (perhaps precomputed) mapping for frequent terms that is copied to all nodes and to use terms directly (instead of termIDs) for infrequent terms. We do not address this problem here and assume that all nodes share a consistent term \$\rightarrow\$ termID mapping.

The map phase of MapReduce consists of mapping splits of the input data to key-value pairs. This is the same parsing task we also encountered in BSBI and SPIMI, and we therefore call the machines that execute the map phase parsers. Each parser writes its output to local intermediate files, the segment files (shown as \fbox{a-f\medstrut} \fbox{g-p\medstrut} \fbox{q-z\medstrut} in Figure 4.5).

For the reduce phase, we want all values for a given key to be stored close together, so that they can be read and processed quickly. This is achieved by partitioning the keys into \$j\$ term partitions and having the parsers write key-value pairs for each term partition into a separate segment file. In Figure 4.5, the term partitions are according to first letter: a-f, g-p, q-z, and \$j=3\$. (We chose these key ranges for ease of exposition. In general, key ranges need not correspond to contiguous terms or termIDs.) The term partitions are defined by the person who operates the indexing system (Exercise 4.6). The parsers then write corresponding segment files, one for each term partition. Each term partition thus corresponds to \$rs\$ segments files, where \$rs\$ is the number of parsers. For instance, Figure 4.5 shows three a-f segment files of the a-f partition, corresponding to the three parsers shown in the figure.

Collecting all values (here: docIDs) for a given key (here: termID) into one list is the task of the inverters in the reduce phase. The

Objective: Create a Frequency Distribution of words in the file

We have a large text file

next up previous contents index
Next: Dynamic indexing Up: Index construction Previous: Single-pass in-memory indexing Contents Index

Distributed indexing

Collections are often so large that we cannot perform index construction efficiently on a single machine. This is particularly true of the World Wide Web for which we need large computer clusters [+] to construct any reasonably sized web index. Web search engines, therefore, use distributed indexing algorithms for index construction. The result of the construction process is a distributed index that is partitioned across several machines – either according to term or according to document. In this section, we describe distributed indexing for a term-partitioned index. Most large search engines prefer a document-partitioned index (which can be easily generated from a term-partitioned index). We discuss this topic further in Section 20.3 (page [*]).

The distributed index construction method we describe in this section is an application of MapReduce, a general architecture for distributed computing. MapReduce is designed for large computer clusters. The point of a cluster is to solve large computing problems on cheap commodity machines or nodes that are built from standard parts (processor, memory, disk) as opposed to on a supercomputer with specialized hardware. Although hundreds or thousands of machines are available in such clusters, individual machines can fail at any time. One requirement for robust distributed indexing is, therefore, that we divide the work up into chunks that we can easily assign and – in case of failure – reassign. A master node directs the process of assigning and reassigning tasks to individual worker nodes.

The map and reduce phases of MapReduce split up the computing job into chunks that standard machines can process in a short time. The various steps of MapReduce are shown in Figure 4.5 and an example on a collection consisting of two documents is shown in Figure 4.6. First, the input data, in our case a collection of web pages, are split into \$n\$ splits where the size of the split is chosen to ensure that the work can be distributed evenly (chunks should not be too large) and efficiently (the total number of chunks we need to manage should not be too large); 16 or 64 MB are good sizes in distributed indexing. Splits are not preassigned to machines, but are instead assigned by the master node on an ongoing basis: As a machine finishes processing one split, it is assigned the next one. If a machine dies or becomes a laggard due to hardware problems, the split it is working on is simply reassigned to another machine.

Figure 4.5: An example of distributed indexing with MapReduce. Adapted from Dean and Ghemawat. (2004).
\includegraphics[width=11.5cm]{art/mapreduce2.eps}

In general, MapReduce breaks a large computing problem into smaller parts by recasting it in terms of manipulation of key-value pairs. For indexing, a key-value pair has the form (termID, docID). In distributed indexing, the mapping from terms to termIDs is also distributed and therefore more complex than in single-machine indexing. A simple solution is to maintain a (perhaps precomputed) mapping for frequent terms that is copied to all nodes and to use terms directly (instead of termIDs) for infrequent terms. We do not address this problem here and assume that all nodes share a consistent term \$\rightarrow\$ termID mapping.

The map phase of MapReduce consists of mapping splits of the input data to key-value pairs. This is the same parsing task we also encountered in BSBI and SPIMI, and we therefore call the machines that execute the map phase parsers. Each parser writes its output to local intermediate files, the segment files (shown as \fbox{a-f\medstrut} \fbox{g-p\medstrut} \fbox{q-z\medstrut} in Figure 4.5).

For the reduce phase, we want all values for a given key to be stored close together, so that they can be read and processed quickly. This is achieved by partitioning the keys into \$j\$ term partitions and having the parsers write key-value pairs for each term partition into a separate segment file. In Figure 4.5, the term partitions are according to first letter: a-f, g-p, q-z, and \$j=3\$. (We chose these key ranges for ease of exposition. In general, key ranges need not correspond to contiguous terms or termIDs.) The term partitions are defined by the person who operates the indexing system (Exercise 4.6). The parsers then write corresponding segment files, one for each term partition. Each term partition thus corresponds to \$rs\$ segments files, where \$rs\$ is the number of parsers. For instance, Figure 4.5 shows three a-f segment files of the a-f partition, corresponding to the three parsers shown in the figure.

Collecting all values (here: docIDs) for a given key (here: termID) into one list is the task of the inverters in the reduce phase. The

This is a pretty common task in Natural Language Processing

We have a large text file

next up previous contents index
Next: Dynamic indexing Up: Index construction Previous: Single-pass in-memory indexing Contents Index

Distributed indexing

Collections are often so large that we cannot perform index construction efficiently on a single machine. This is particularly true of the World Wide Web for which we need large computer clusters [+] to construct any reasonably sized web index. Web search engines, therefore, use distributed indexing algorithms for index construction. The result of the construction process is a distributed index that is partitioned across several machines – either according to term or according to document. In this section, we describe distributed indexing for a term-partitioned index. Most large search engines prefer a document-partitioned index (which can be easily generated from a term-partitioned index). We discuss this topic further in Section 20.3 (page [*]).

The distributed index construction method we describe in this section is an application of MapReduce, a general architecture for distributed computing. MapReduce is designed for large computer clusters. The point of a cluster is to solve large computing problems on cheap commodity machines or nodes that are built from standard parts (processor, memory, disk) as opposed to a supercomputer with specialized hardware. Although hundreds or thousands of machines are available in such clusters, individual machines can fail at any time. One requirement for robust distributed indexing is, therefore, that we divide the work up into chunks that we can easily assign and – in case of failure – reassign. A master node directs the process of assigning and reassigning tasks to individual worker nodes.

The map and reduce phases of MapReduce split up the computing job into chunks that standard machines can process in a short time. The various steps of MapReduce are shown in Figure 4.5 and an example on a collection consisting of two documents is shown in Figure 4.6. First, the input data, in our case a collection of web pages, are split into \$n\$ splits where the size of the split is chosen to ensure that the work can be distributed evenly (chunks should not be too large) and efficiently (the total number of chunks we need to manage should not be too large); 16 or 64 MB are good sizes in distributed indexing. Splits are not preassigned to machines, but are instead assigned by the master node on an ongoing basis: As a machine finishes processing one split, it is assigned the next one. If a machine dies or becomes a laggard due to hardware problems, the split it is working on is simply reassigned to another machine.

Figure 4.5: An example of distributed indexing with MapReduce. Adapted from Dean and Ghemawat. (2004).

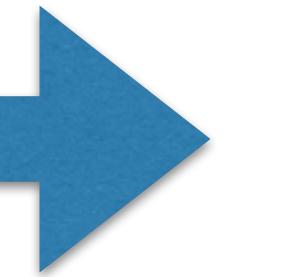
\includegraphics[width=11.5cm]{art/mapreduce2.eps}

In general, MapReduce breaks a large computing problem into smaller parts by recasting it in terms of manipulation of key-value pairs. For indexing, a key-value pair has the form (termID, docID). In distributed indexing, the mapping from terms to termIDs is also distributed and therefore more complex than in single-machine indexing. A simple solution is to maintain a (perhaps precomputed) mapping for frequent terms that is copied to all nodes and to use terms directly (instead of termIDs) for infrequent terms. We do not address this problem here and assume that all nodes share a consistent term \$\rightarrow\$ termID mapping.

The map phase of MapReduce consists of mapping splits of the input data to key-value pairs. This is the same parsing task we also encountered in BSBI and SPIMI, and we therefore call the machines that execute the map phase parsers. Each parser writes its output to local intermediate files, the segment files (shown as \fbox{a-f\medstrut} \fbox{g-p\medstrut} \fbox{q-z\medstrut} in Figure 4.5).

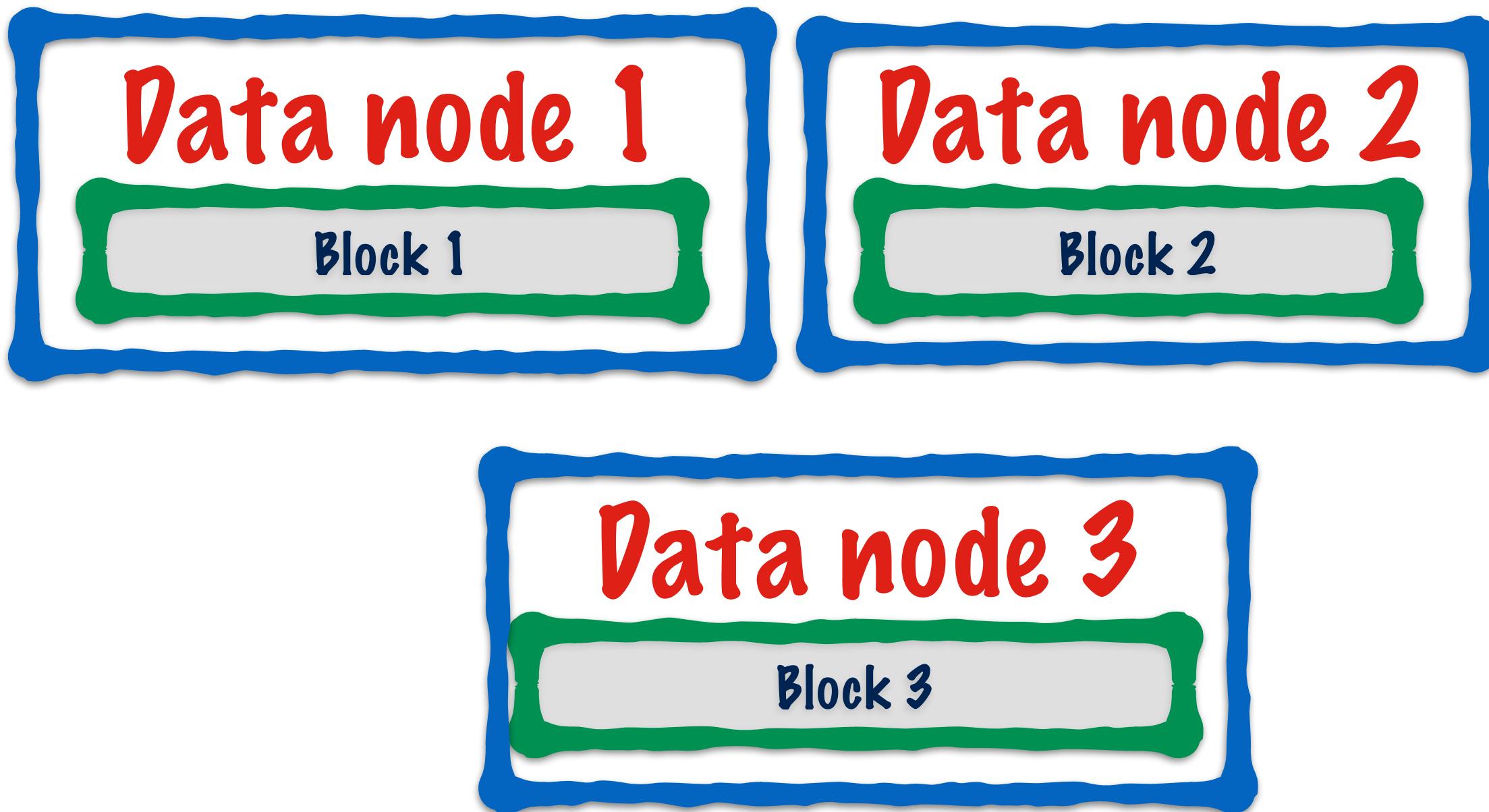
For the reduce phase, we want all values for a given key to be stored close together, so that they can be read and processed quickly. This is achieved by partitioning the keys into \$j\$ term partitions and having the parsers write key-value pairs for each term partition into a separate segment file. In Figure 4.5, the term partitions are according to first letter: a-f, g-p, q-z, and \$j=3\$. (We chose these key ranges for ease of exposition. In general, key ranges need not correspond to contiguous terms or termIDs.) The term partitions are defined by the person who operates the indexing system (Exercise 4.6). The parsers then write corresponding segment files, one for each term partition. Each term partition thus corresponds to \$sr\$ segments files, where \$sr\$ is the number of parsers. For instance, Figure 4.5 shows three a-f segment files of the a-f partition, corresponding to the three parsers shown in the figure.

Collecting all values (here: docIDs) for a given key (here: termID) into one list is the task of the inverters in the reduce phase. The



| Word | Count |
|---------|-------|
| because | 1 |
| each | 4 |
| figure | 9 |
| .. | .. |

The text file has been divided into blocks and stored in HDFS



Name node

The name node stores metadata

Here is how the data flows in a MapReduce job

Data node 1

Hey Diddle Diddle

Block 1

...

Data node 2

The cat and the fiddle

Block 2

....

Data node 3

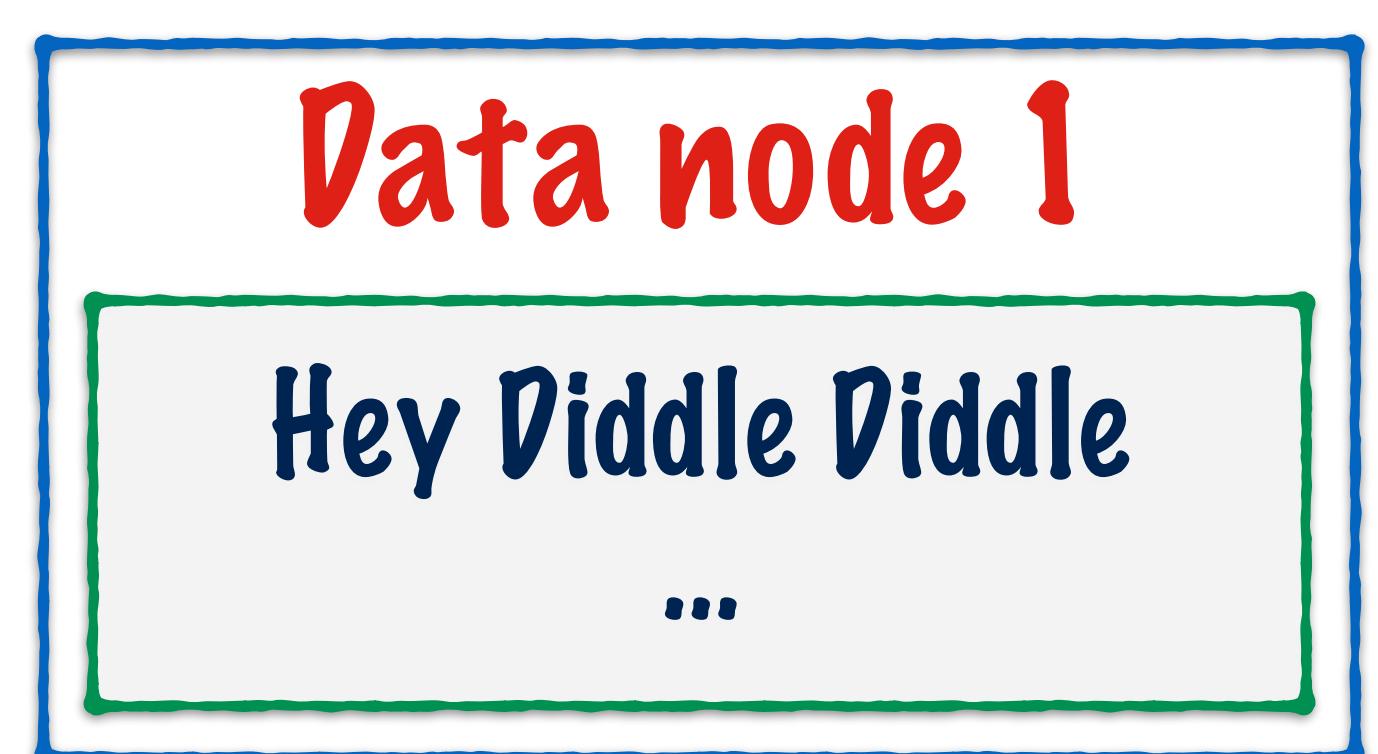
The cow jumped over the moon

Block 3

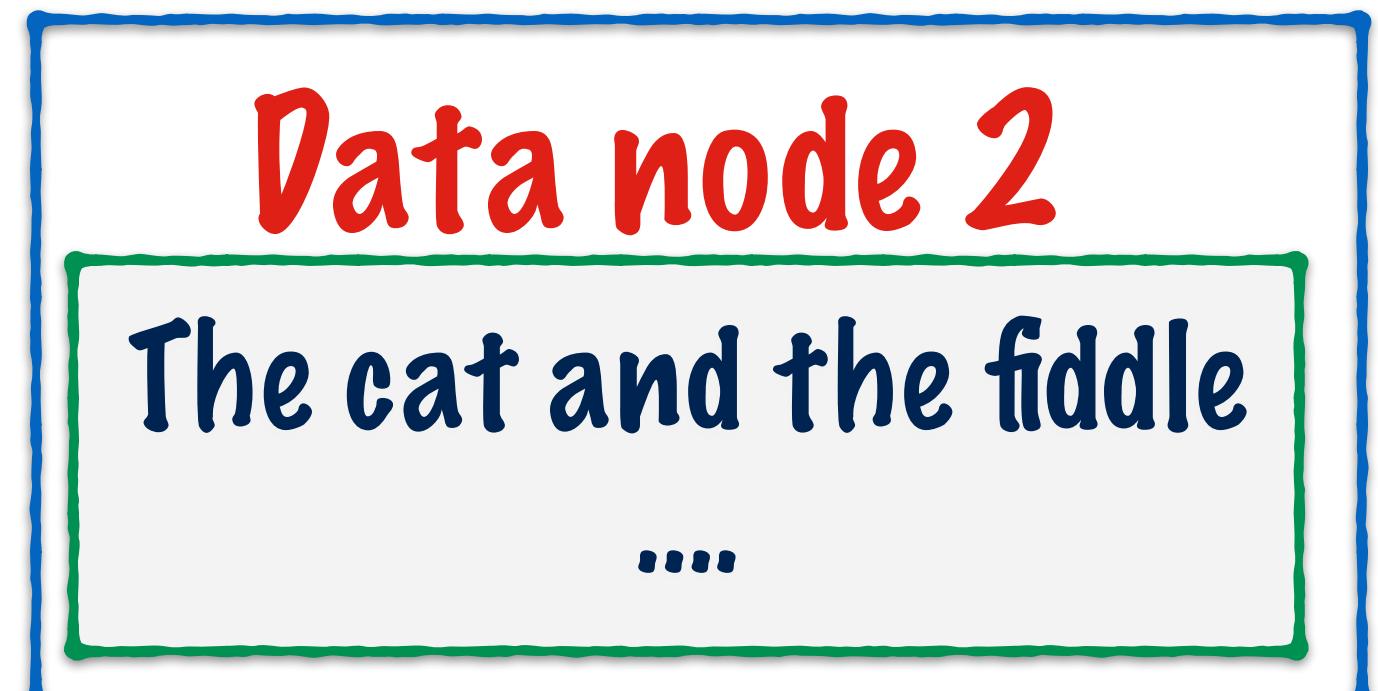
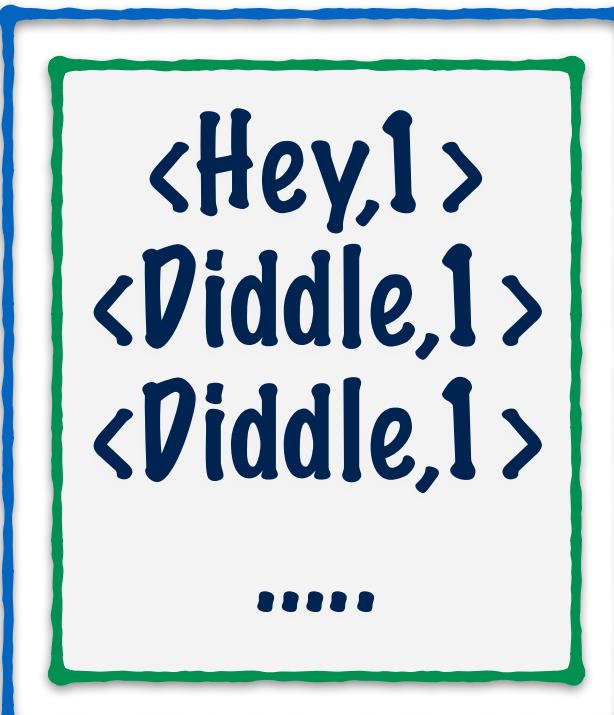
...

Each block here would represent a part of the text file

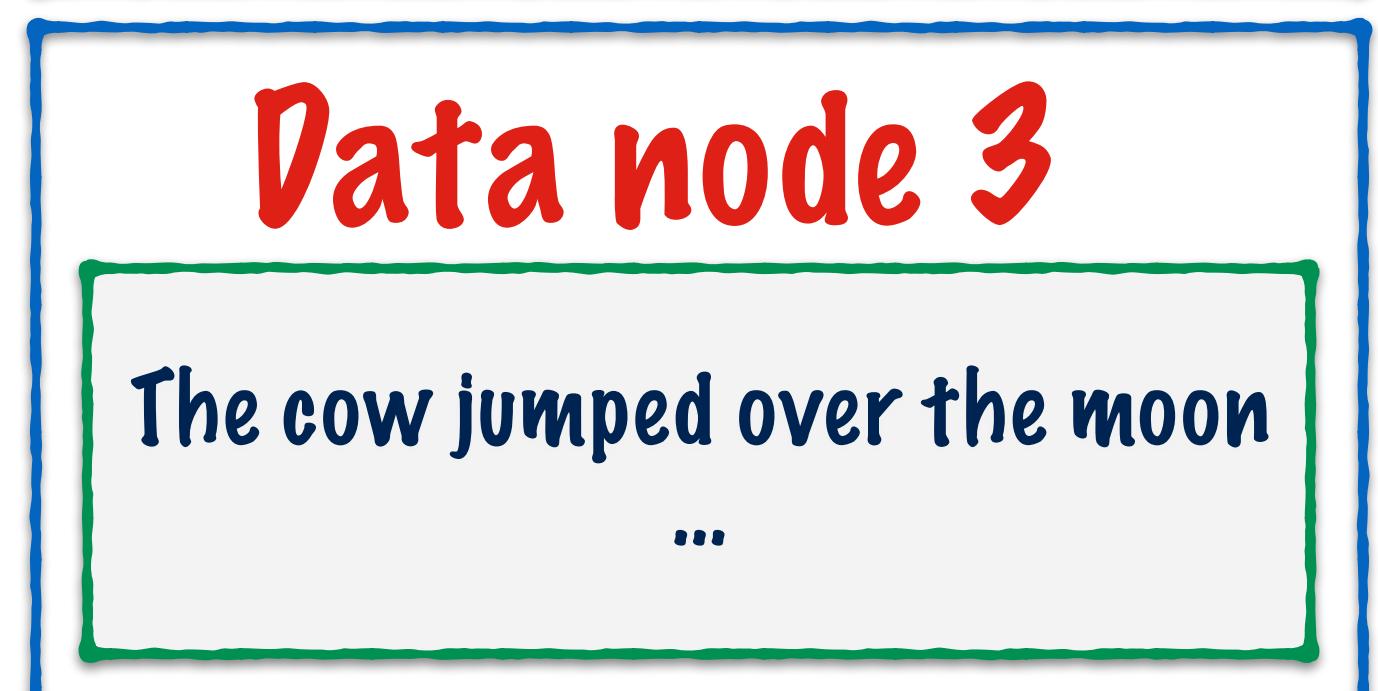
Here is a snapshot of some text from each block



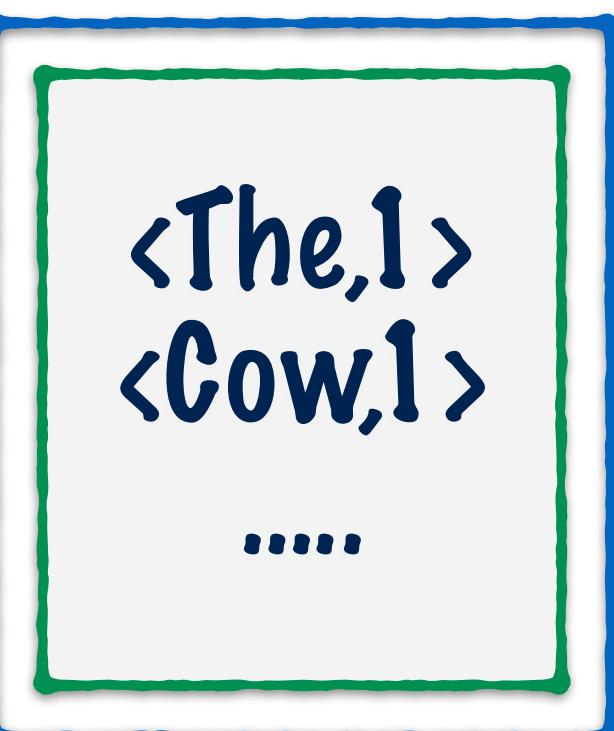
MAP →



MAP →

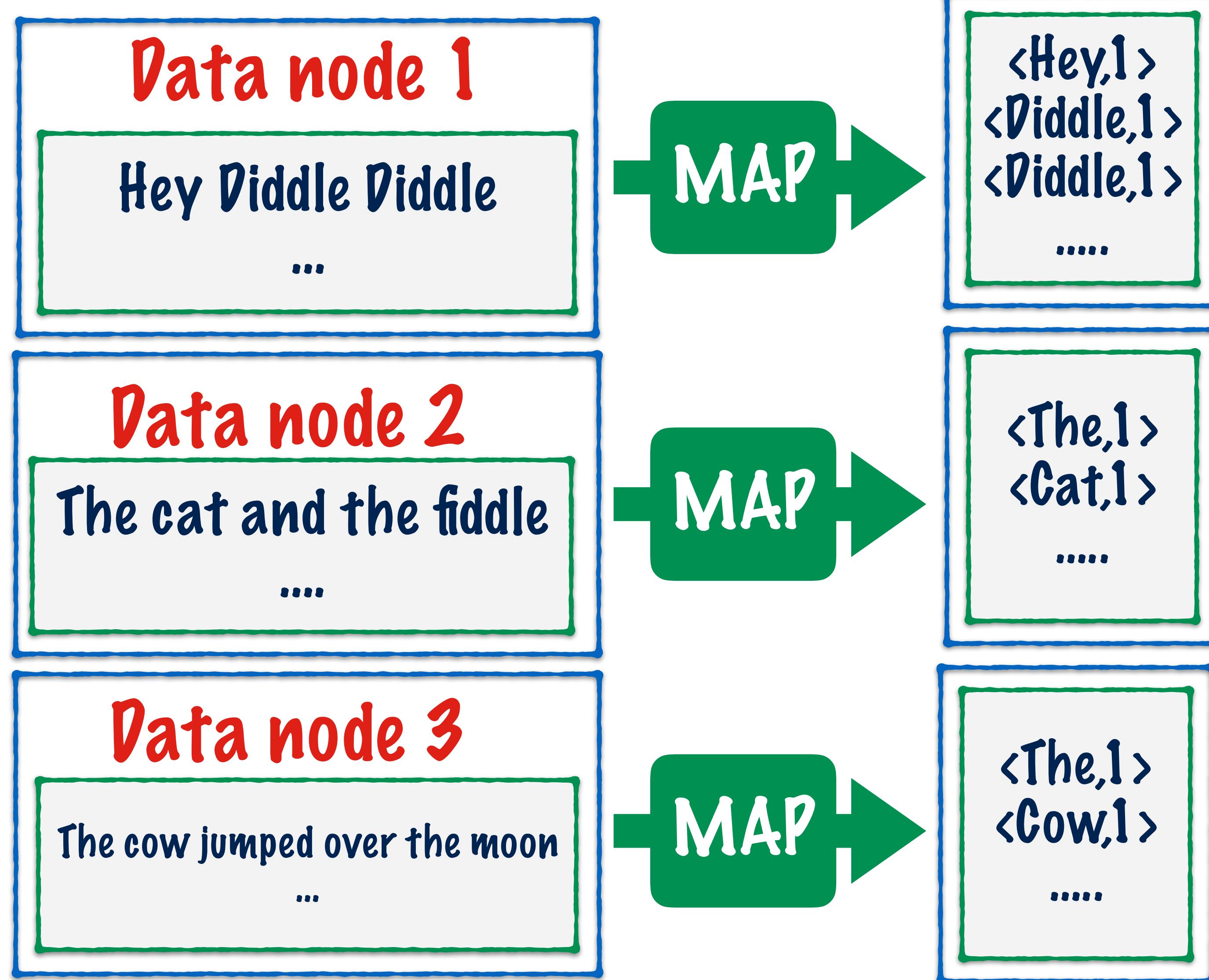


MAP →



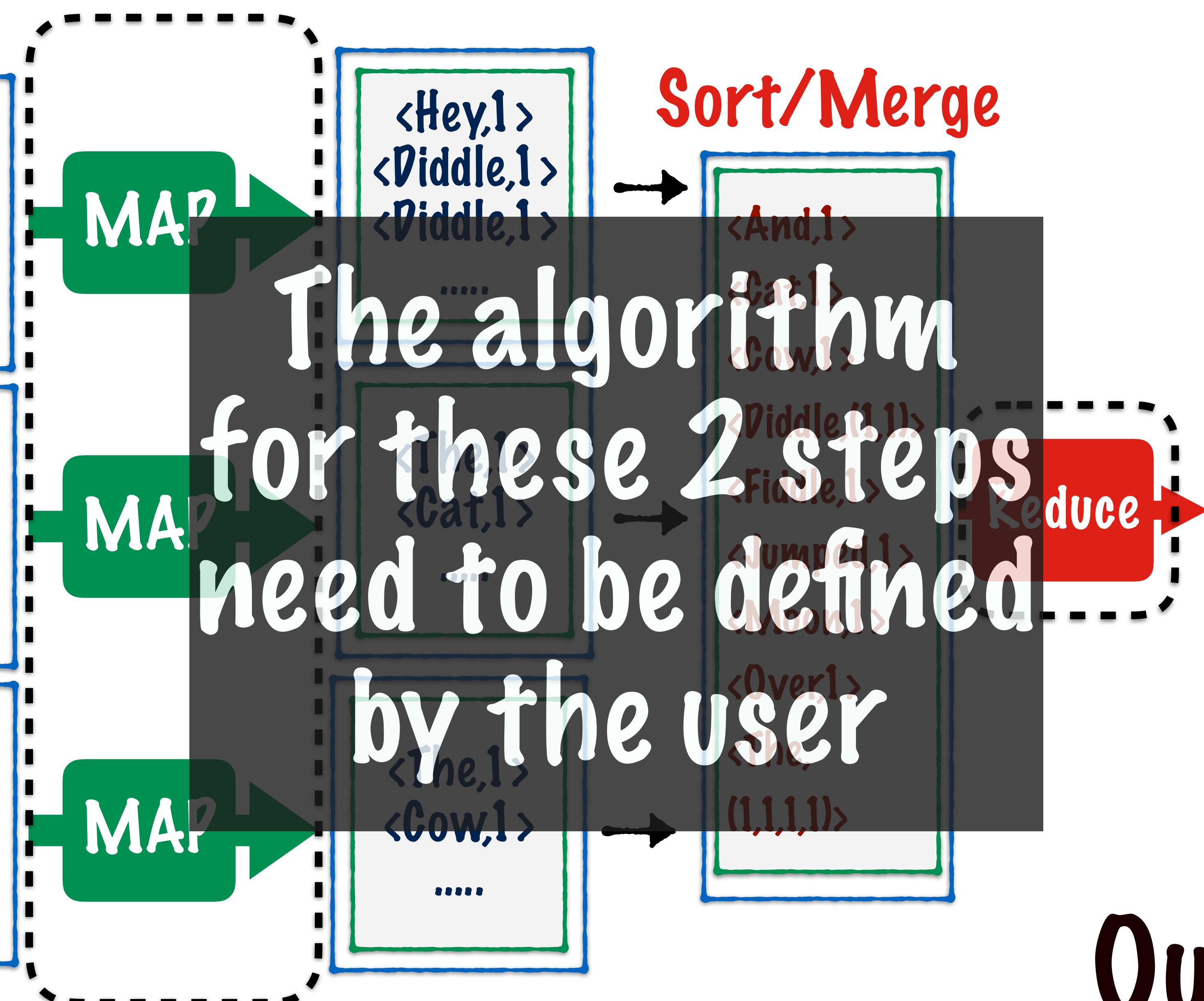
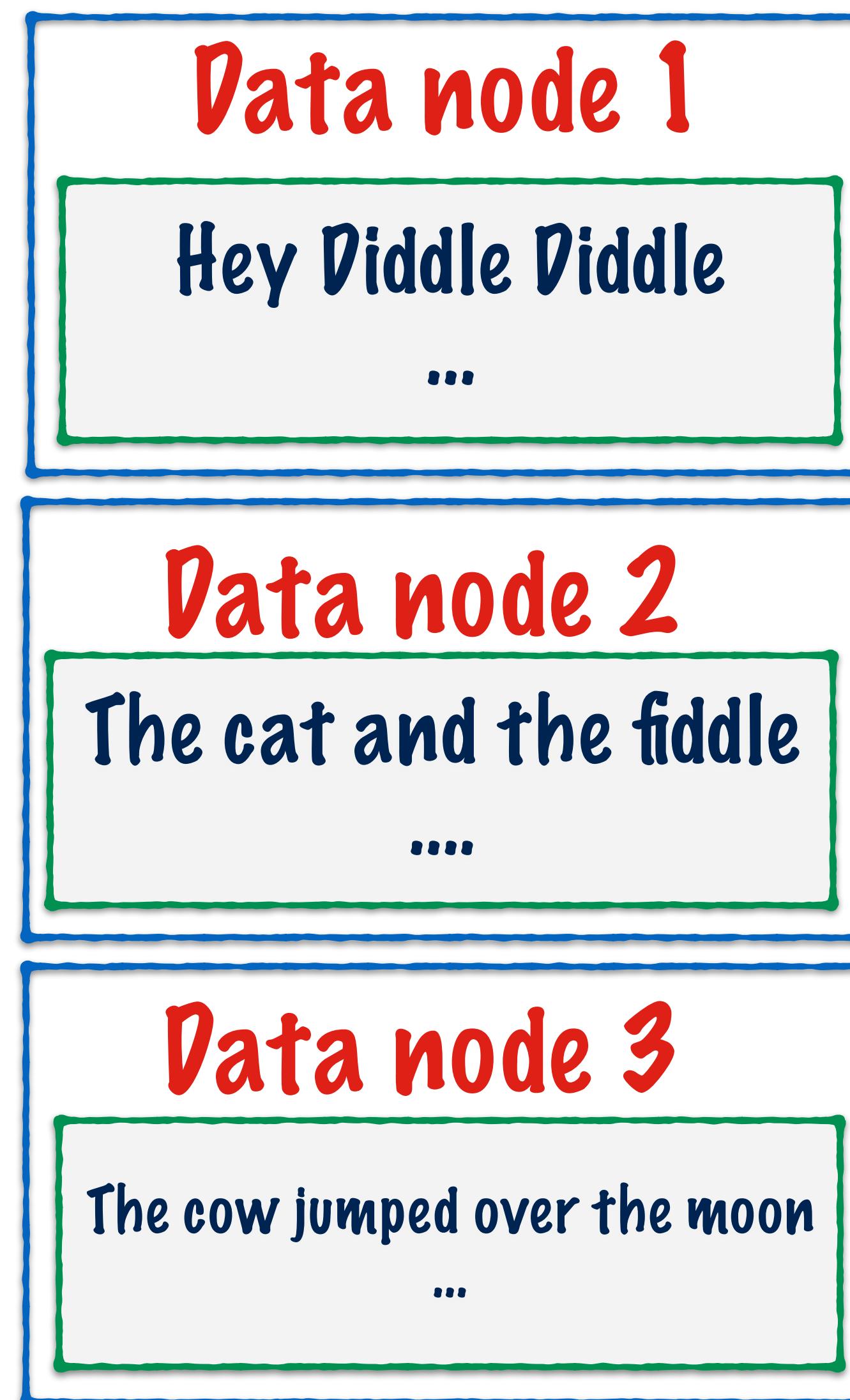
The map step will generate a list of key-value pairs on each node

From now on all the inputs and outputs are formatted as <key,value> pairs



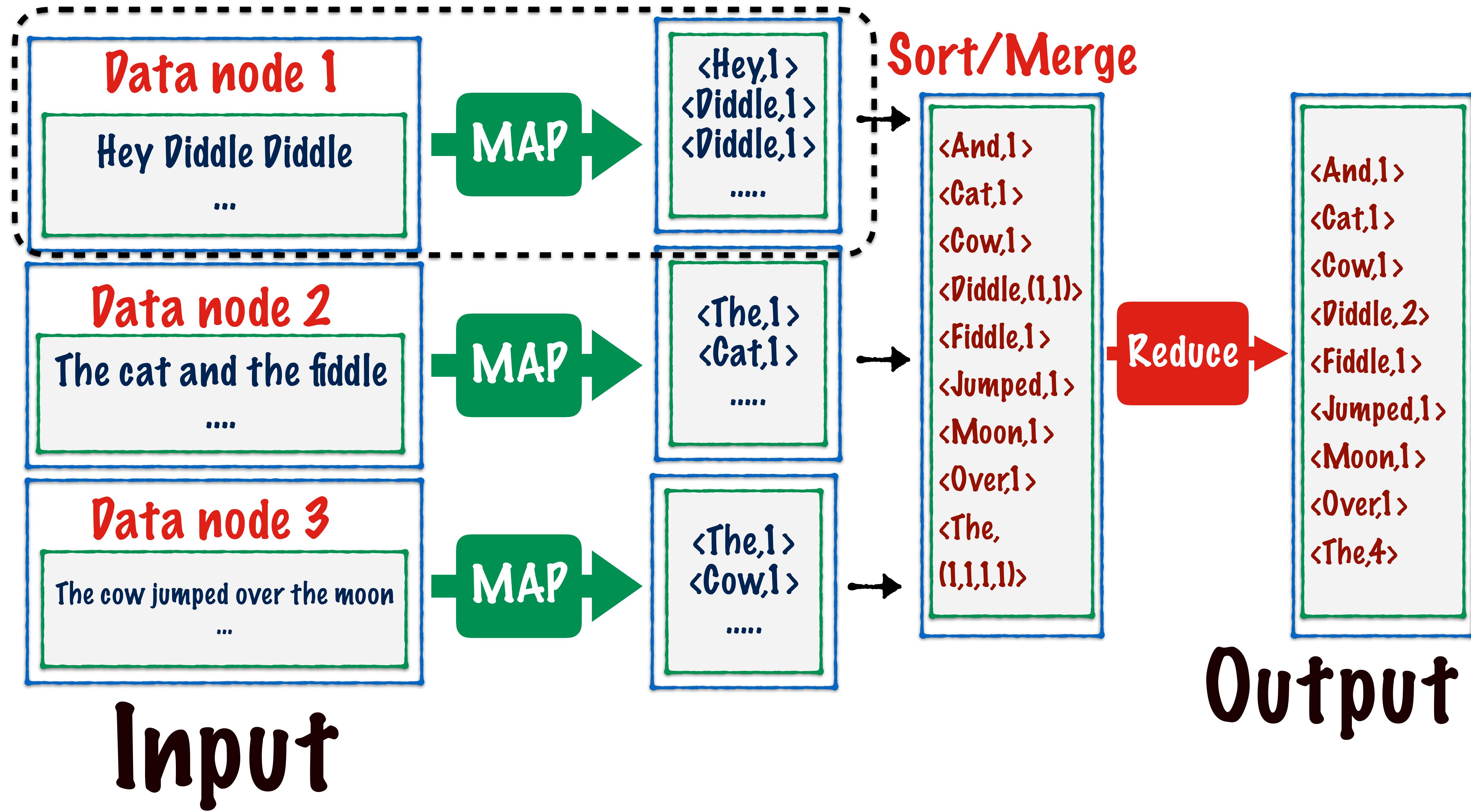
These are all copied over to one single node

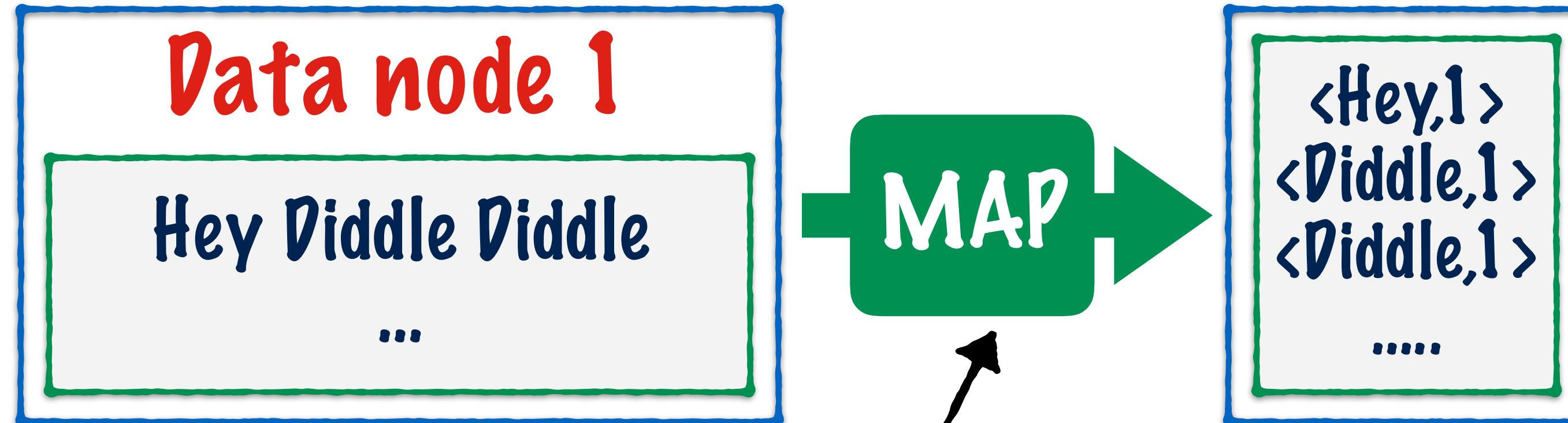
On that node an operation called **Sort/Merge** occurs



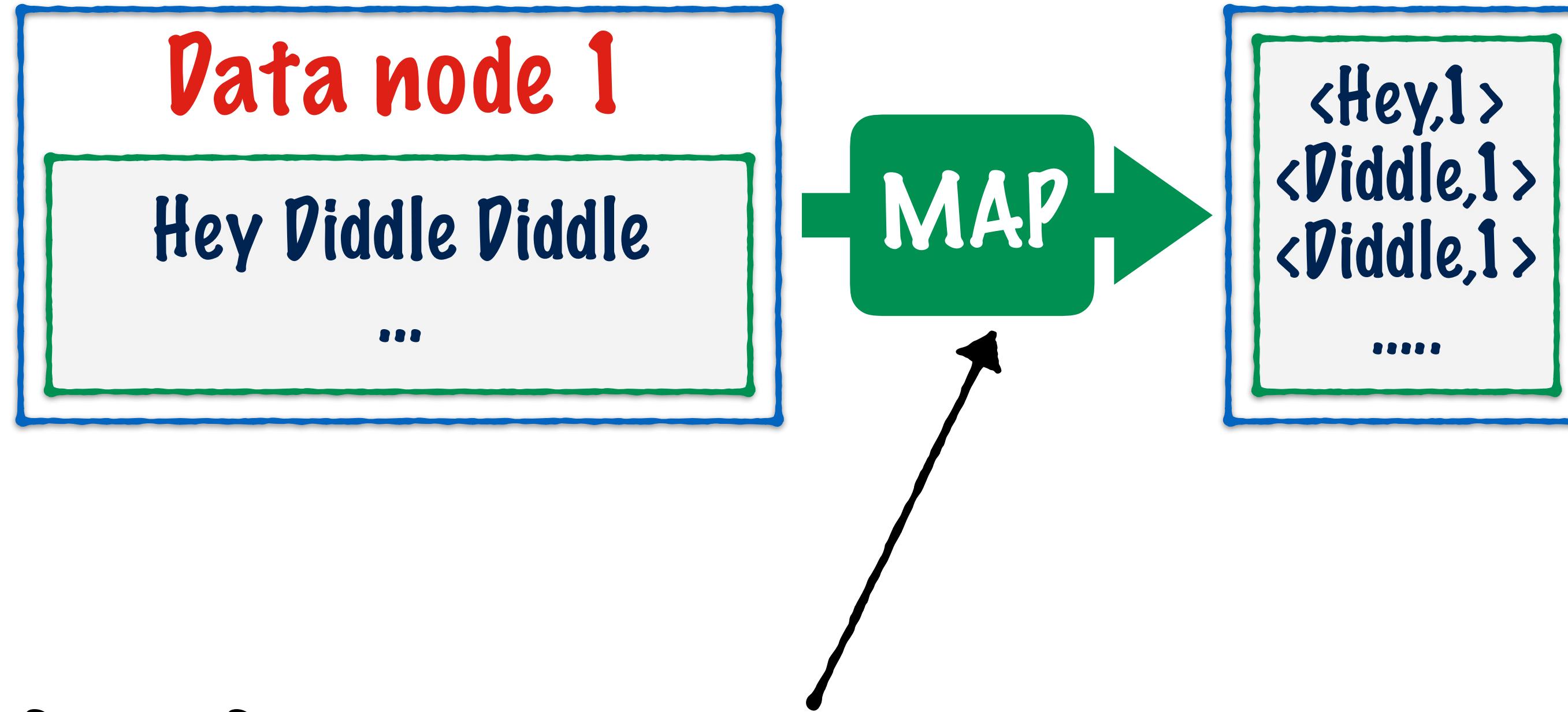
Input

Output

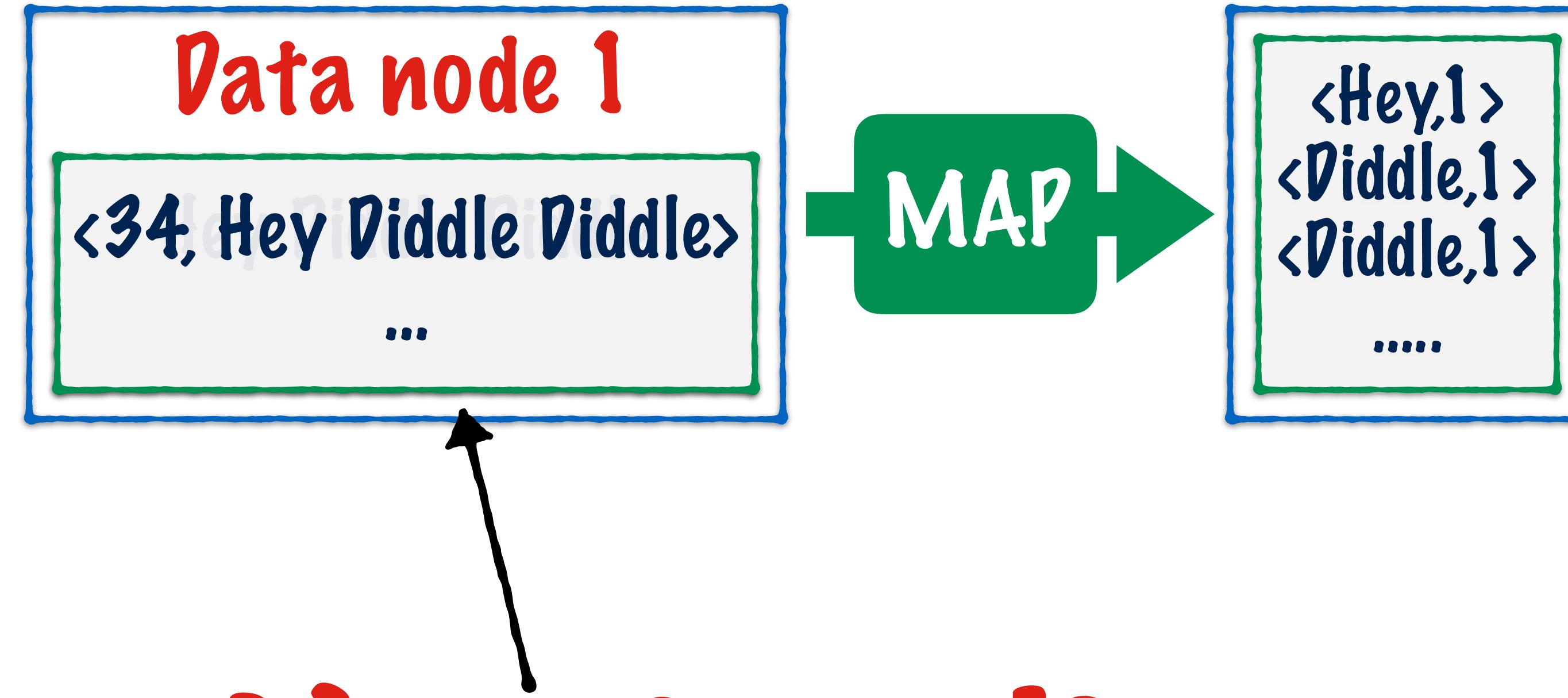




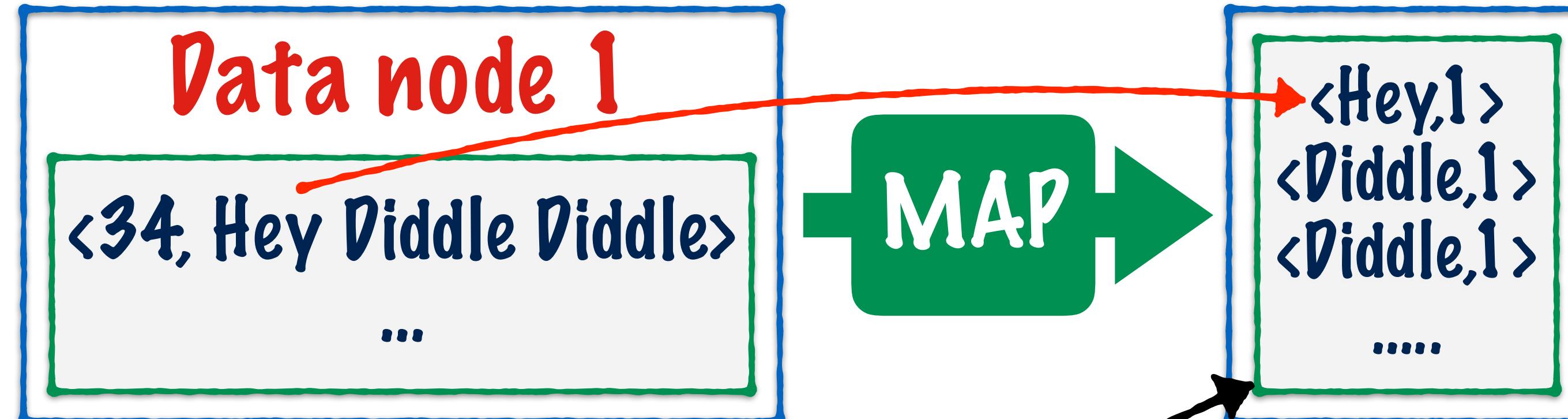
This function will run once
for each line of the text file



Both the input and output need
to be formatted as **<key, value>**
pairs

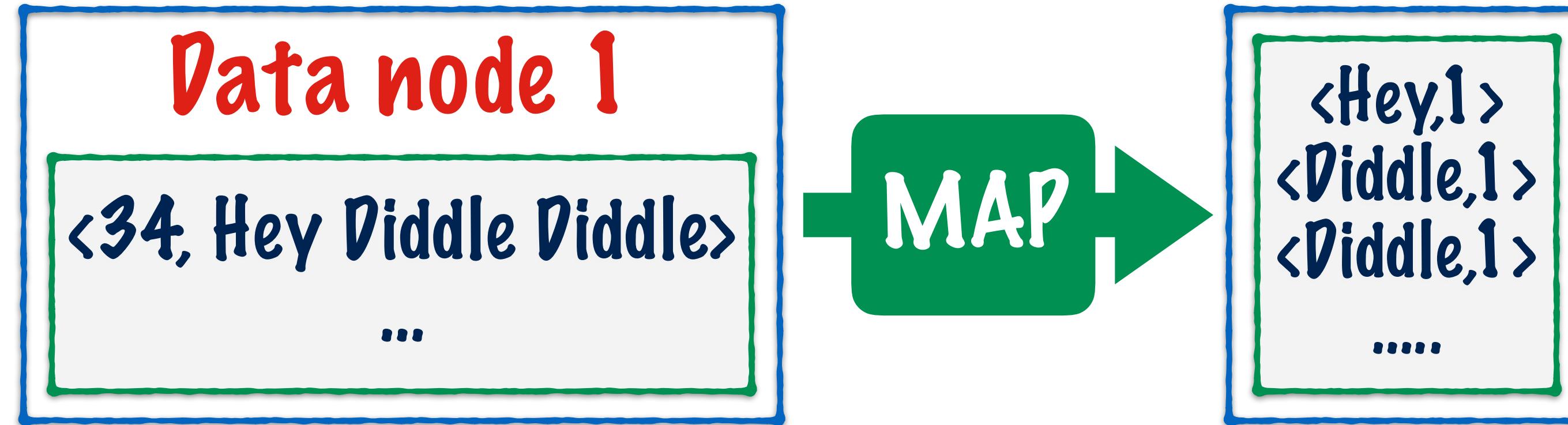


<lineNumber, line>

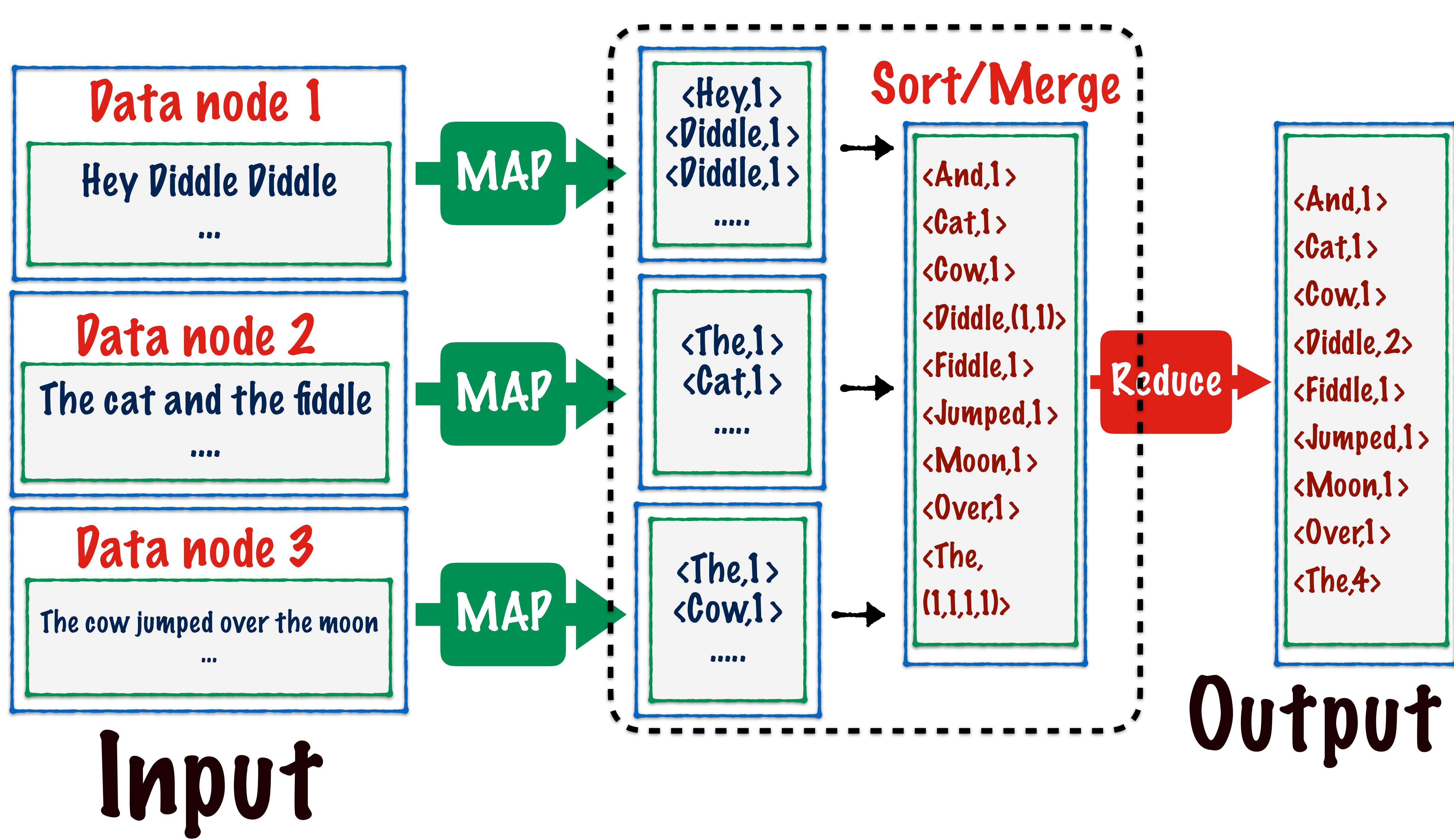


<word, 1>

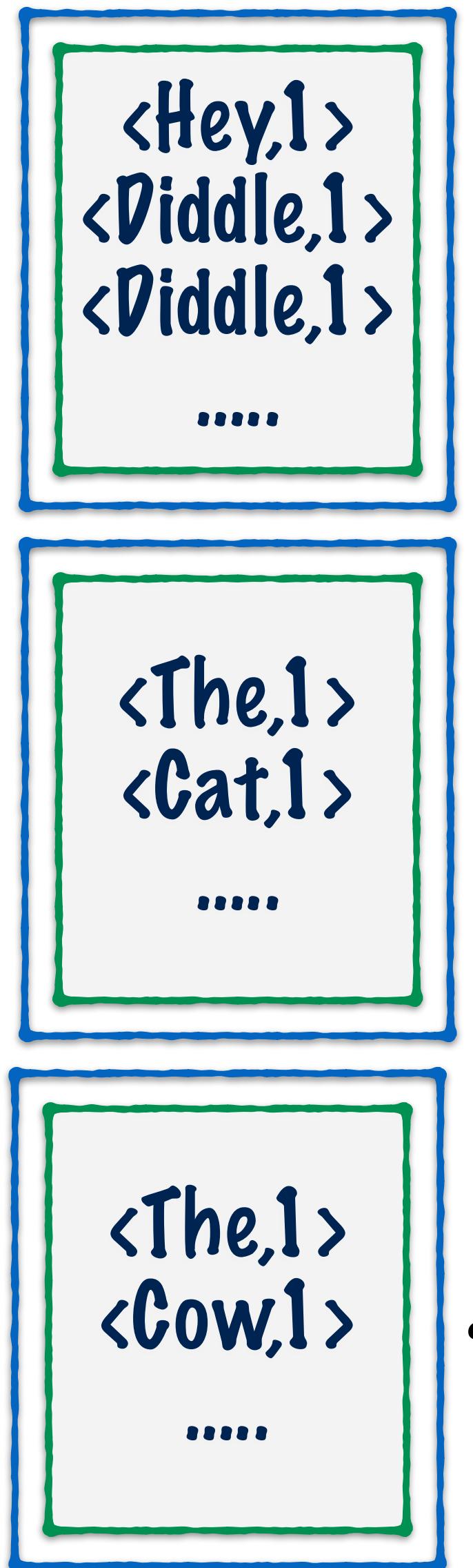
There will be one such pair for
each word in the line



This operation can run in parallel on each data node, there is no interdependency in the inputs and outputs

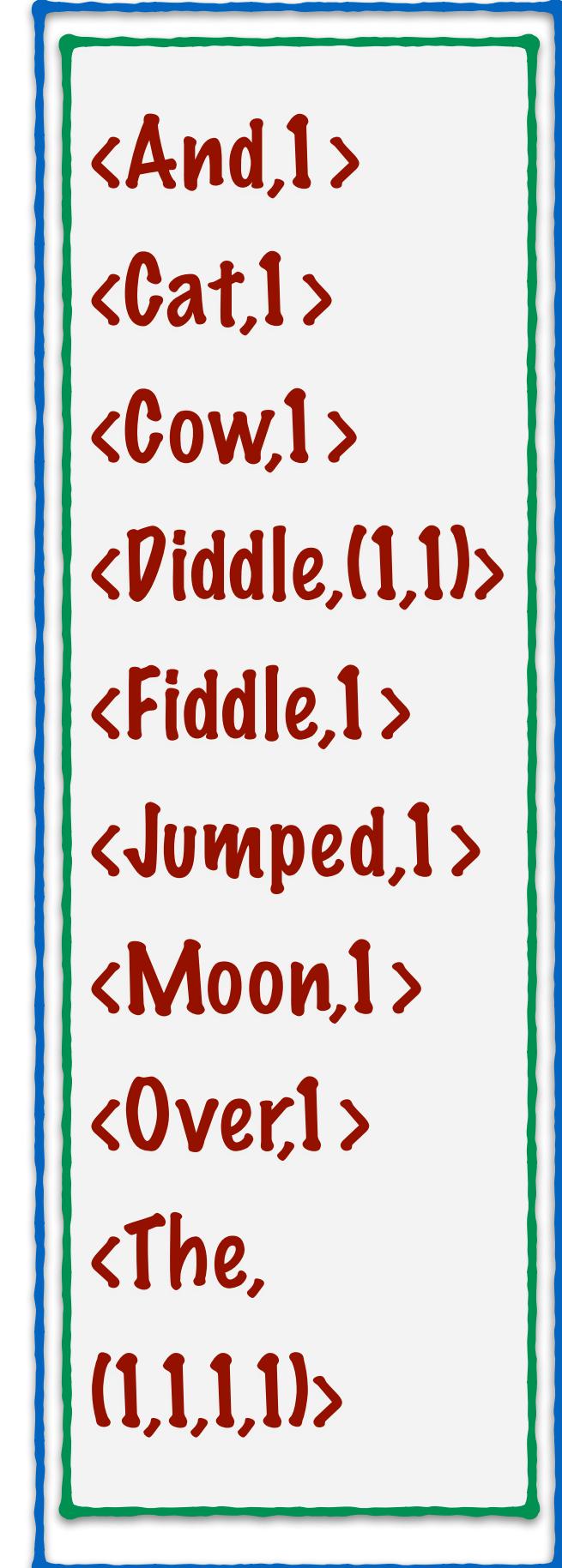


Sort/Merge

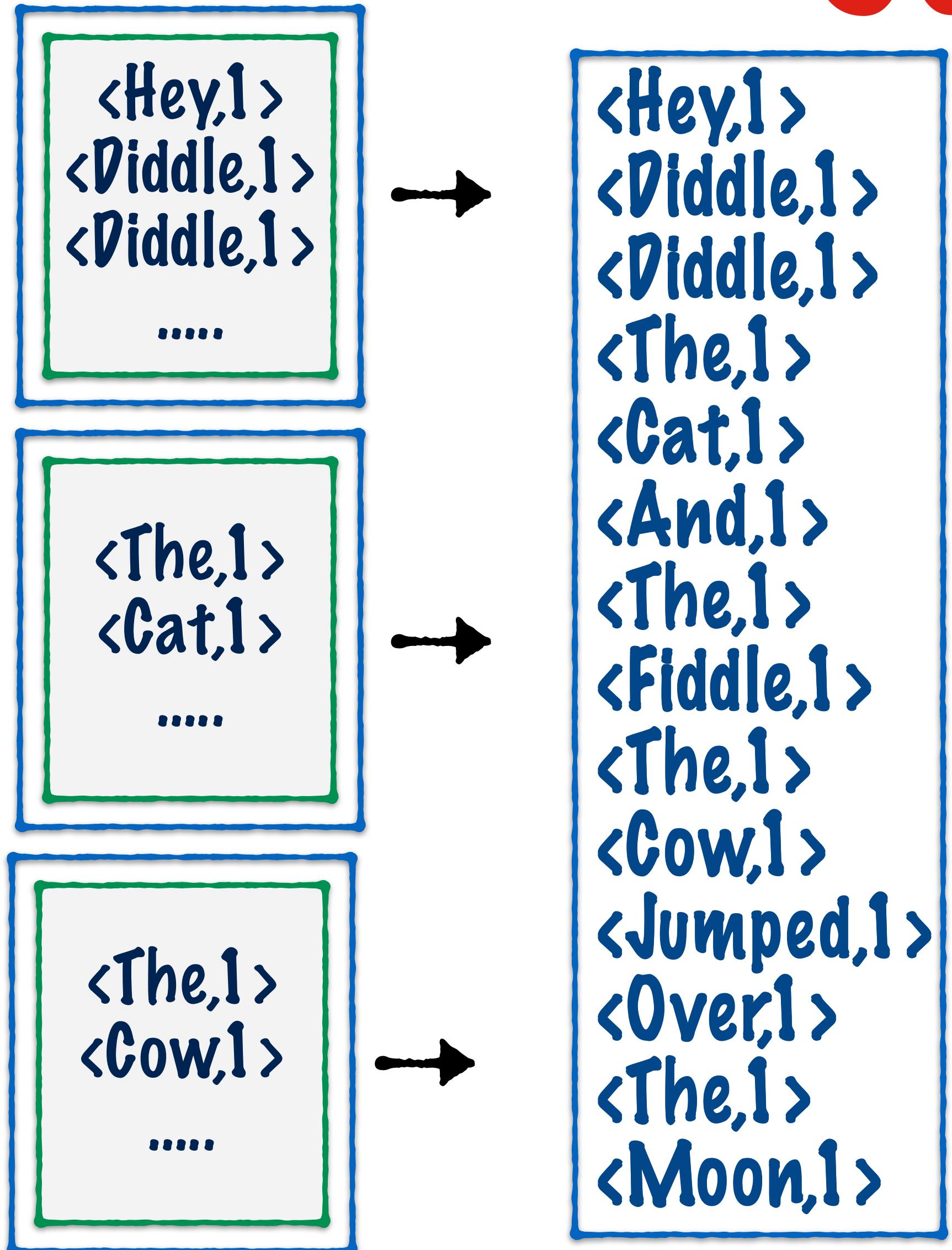


At the end of the map phase, we have **a set of key-value pairs** from each data node

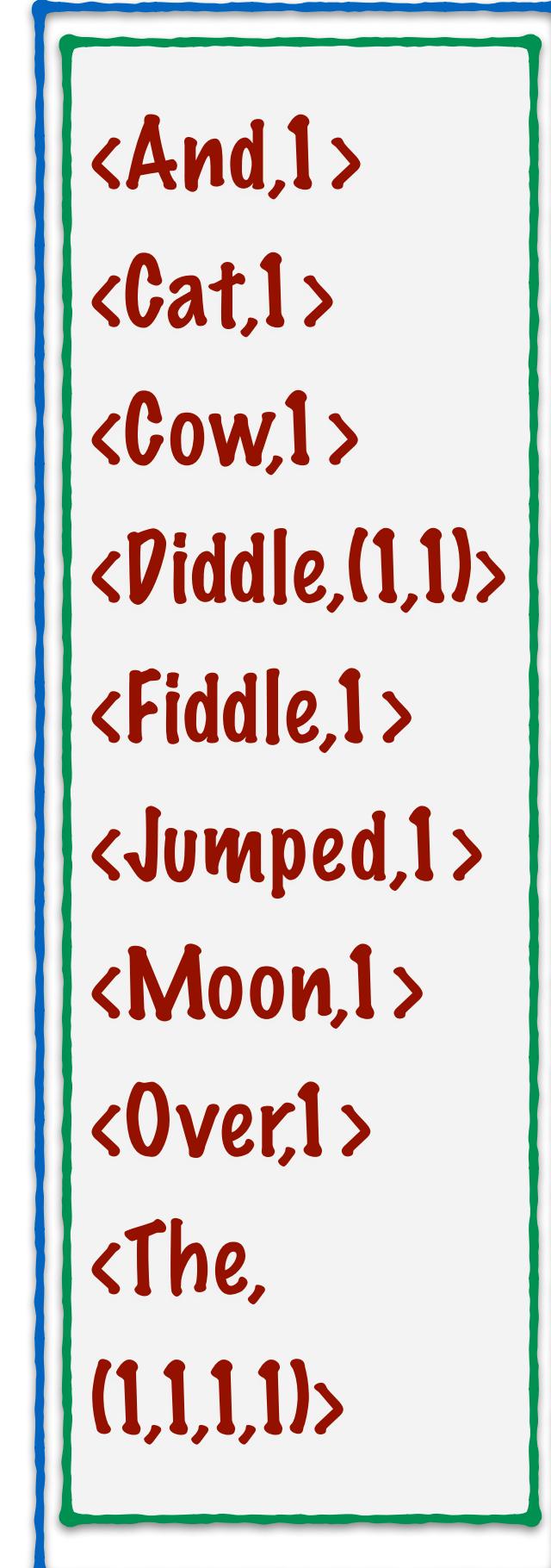
All of these results are **first copied over to a single node**



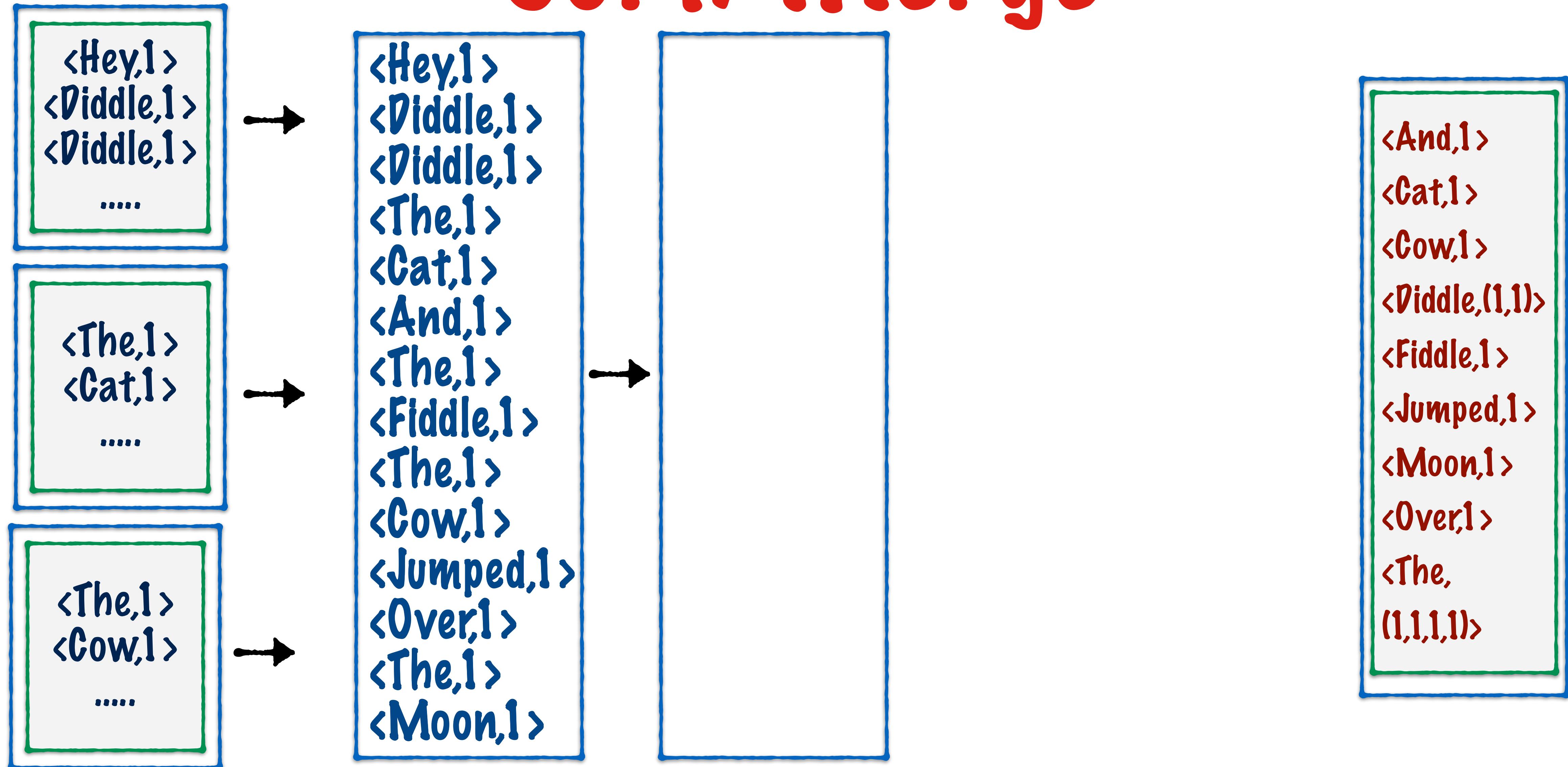
Sort/Merge



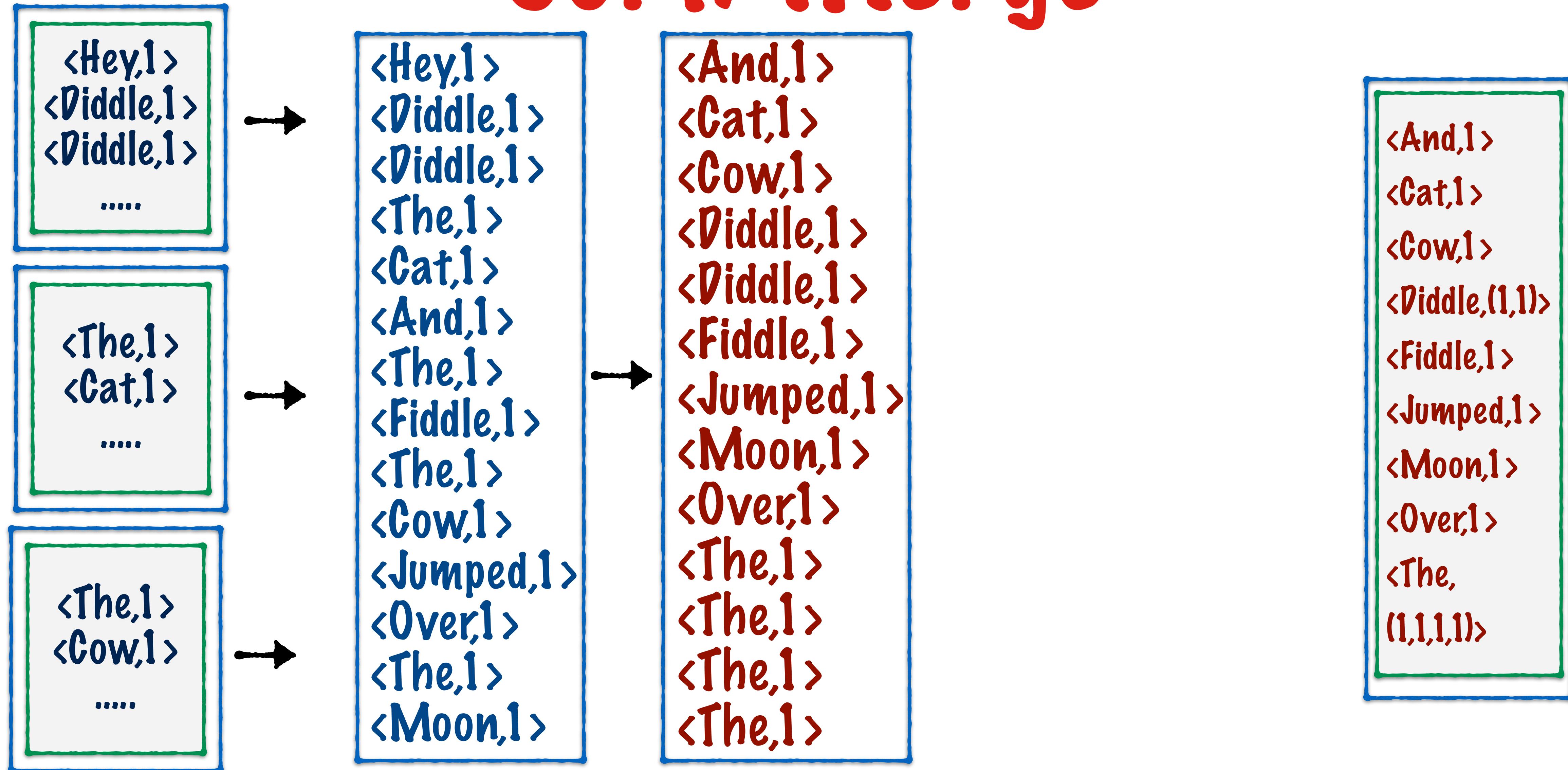
The data is
sorted by **key**



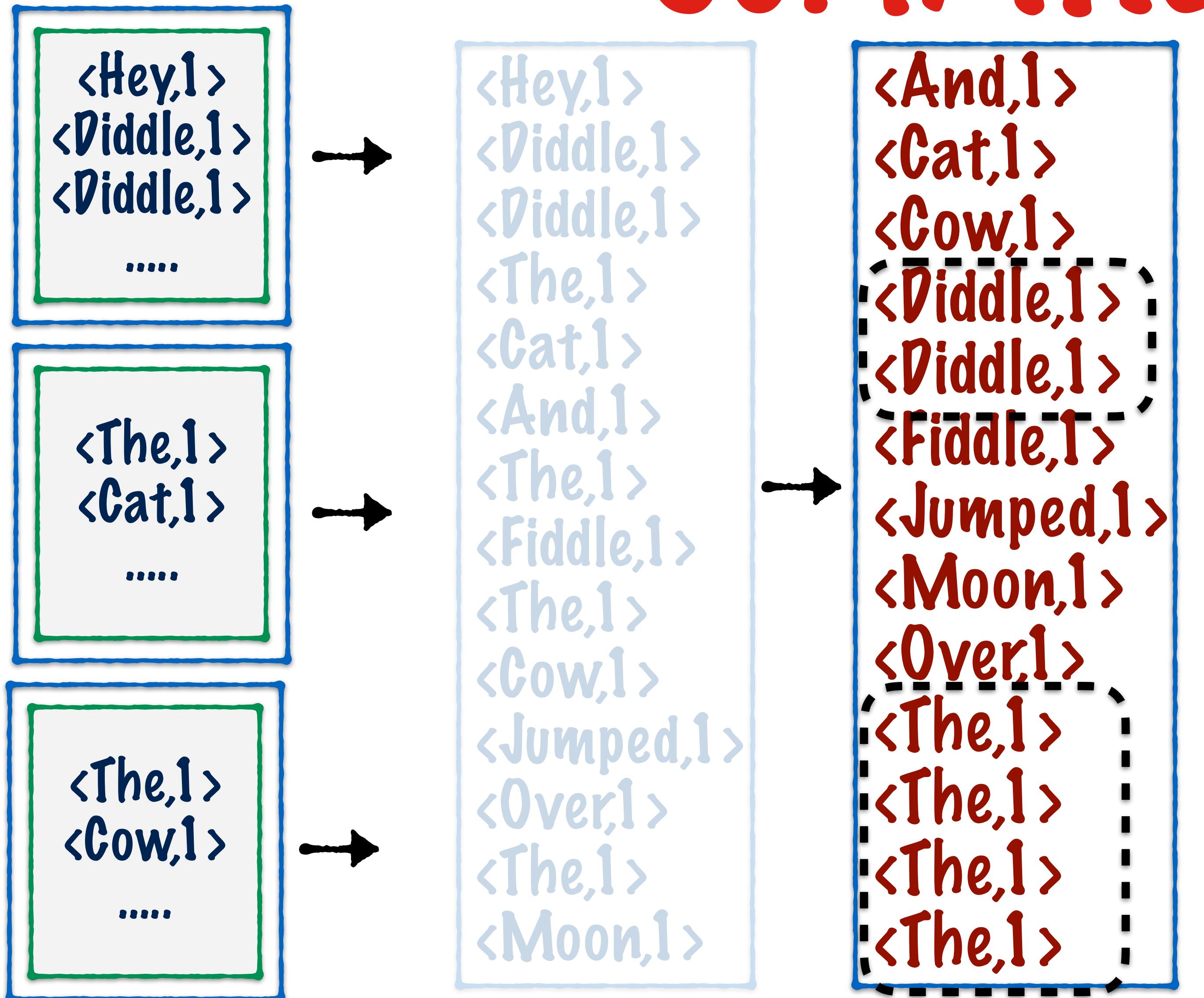
Sort/Merge



Sort/Merge

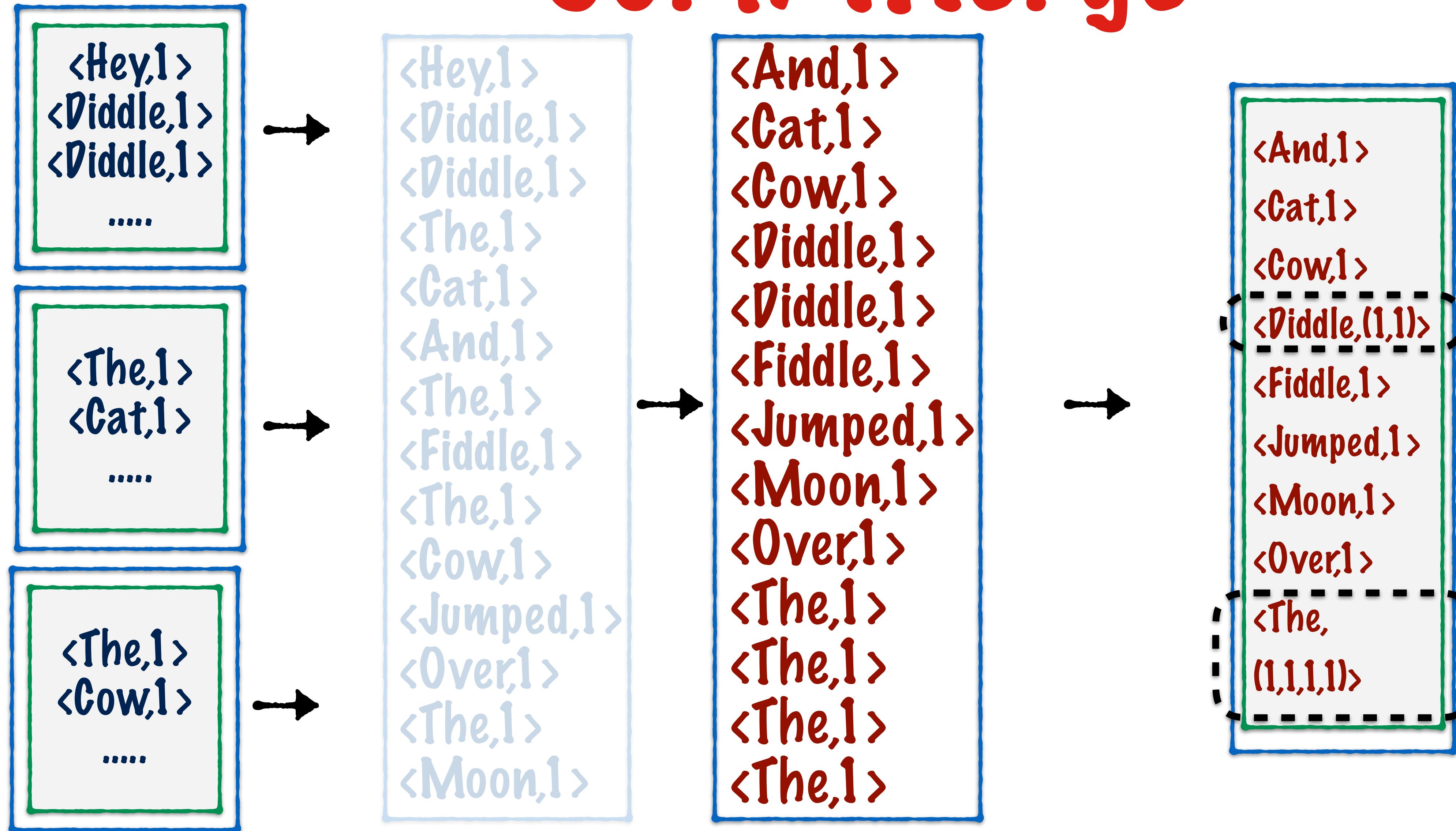


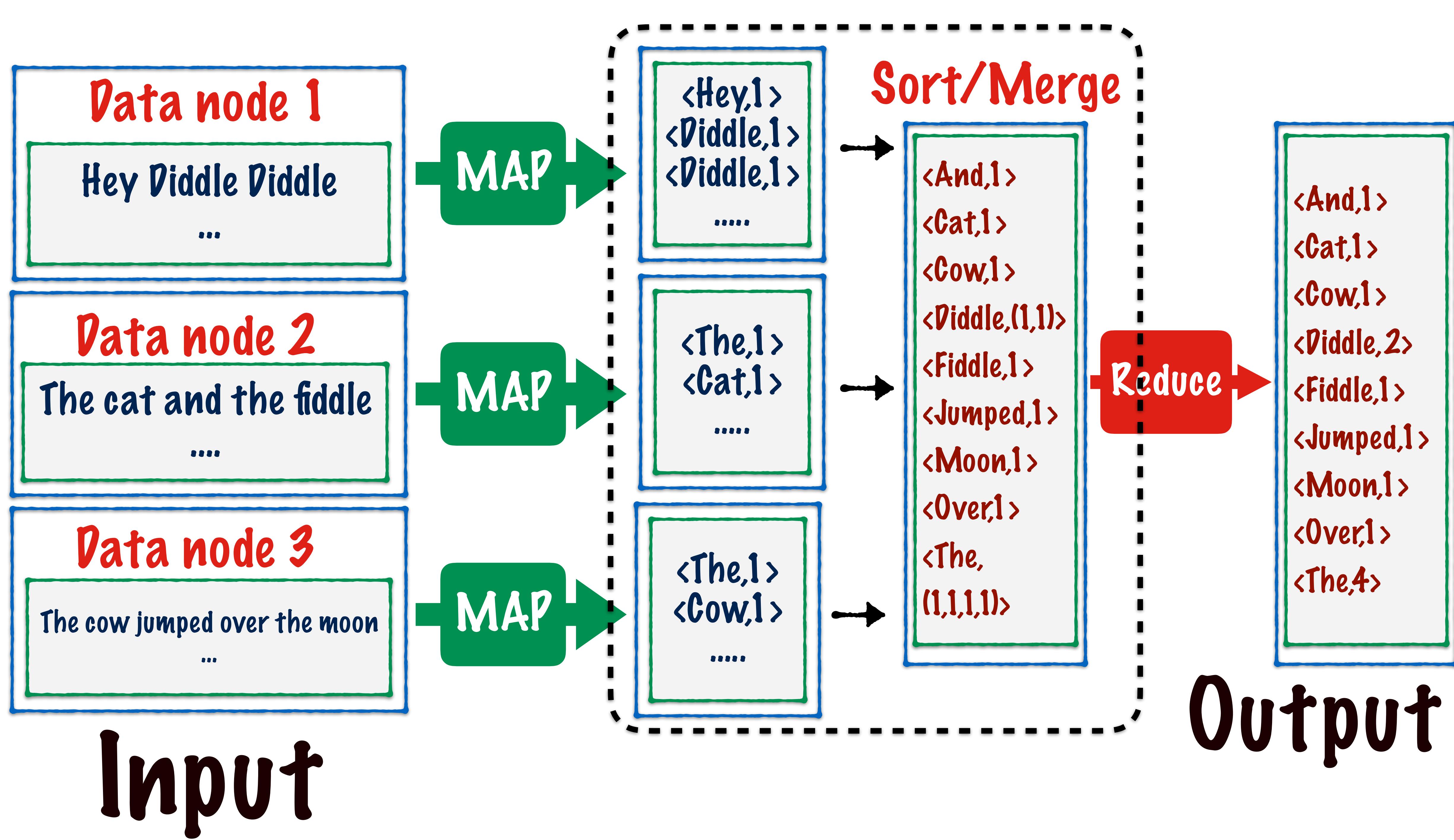
Sort/Merge

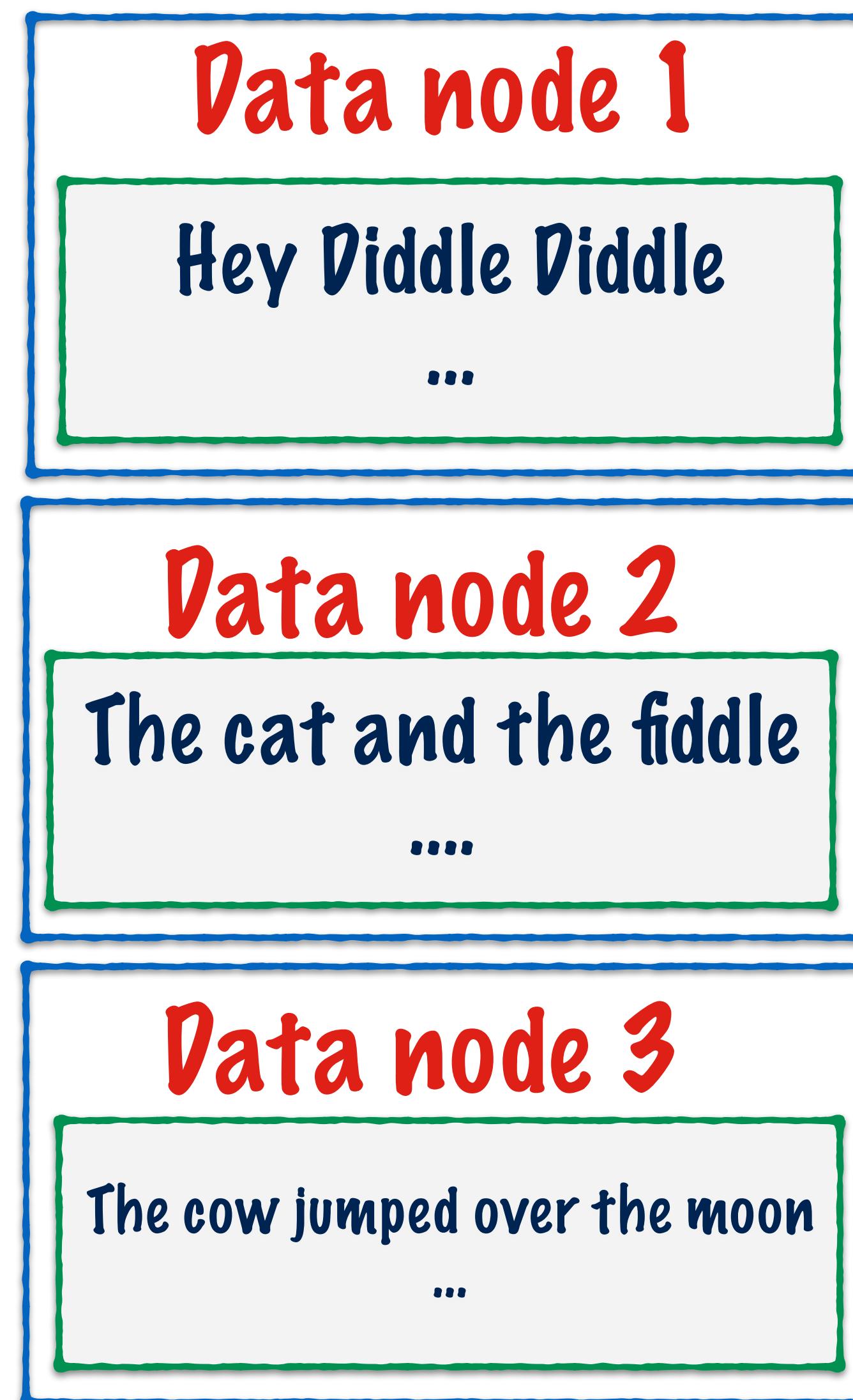


Key-Value
pairs with the
same key are
merged

Sort/Merge





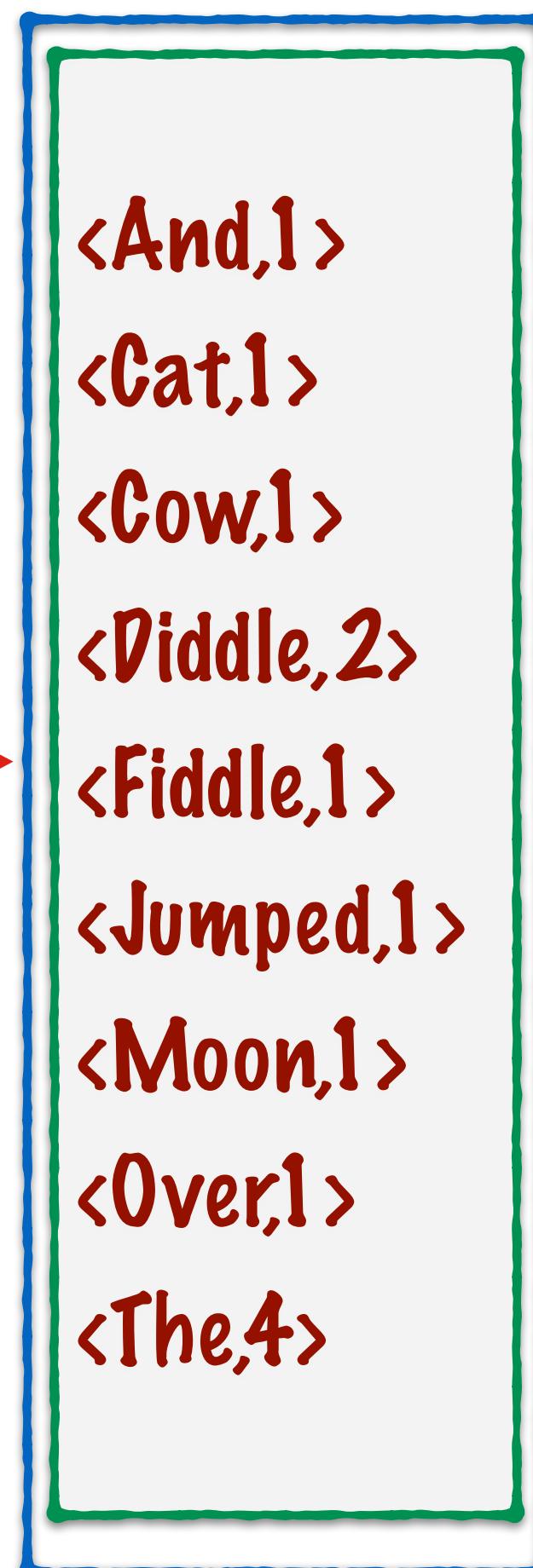
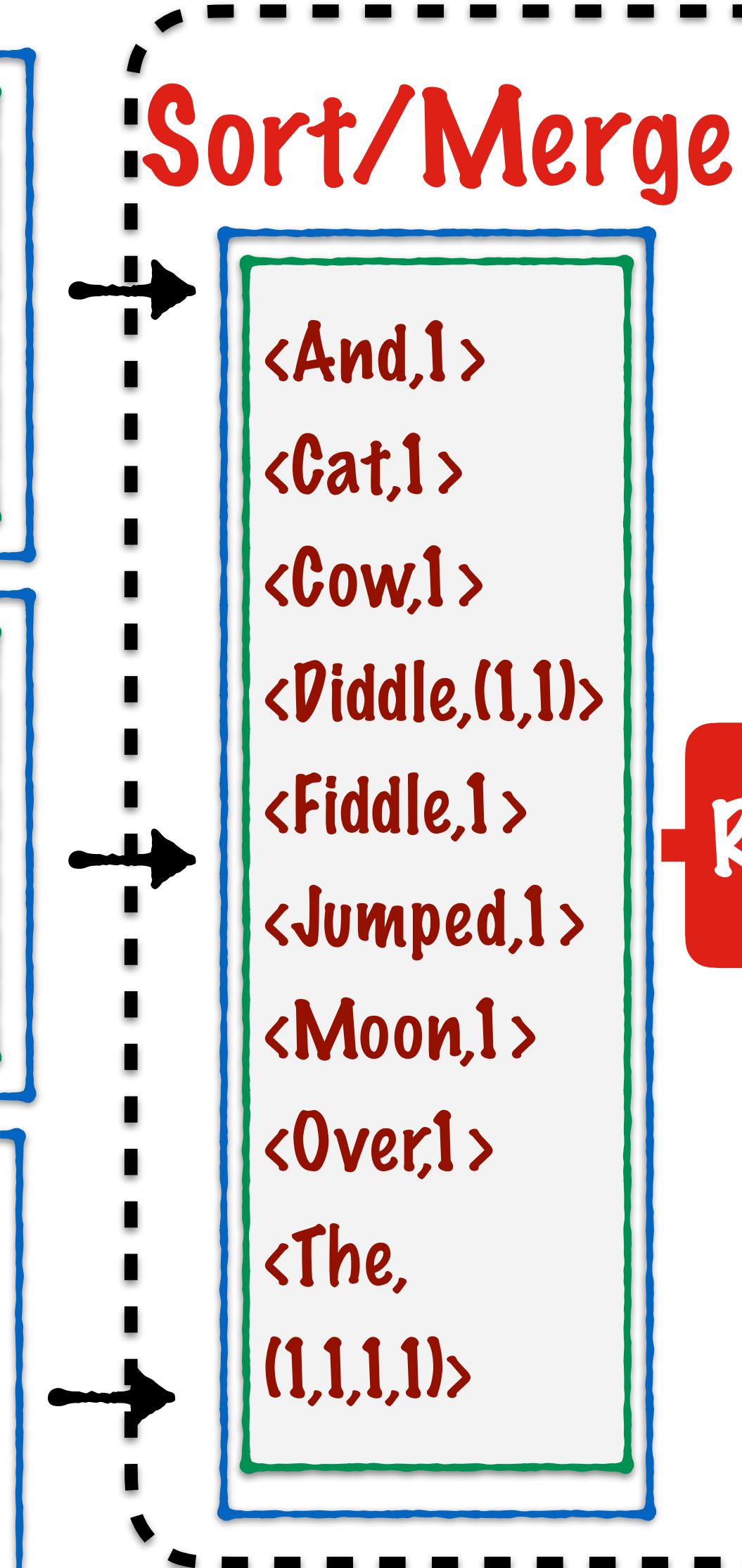


MAP

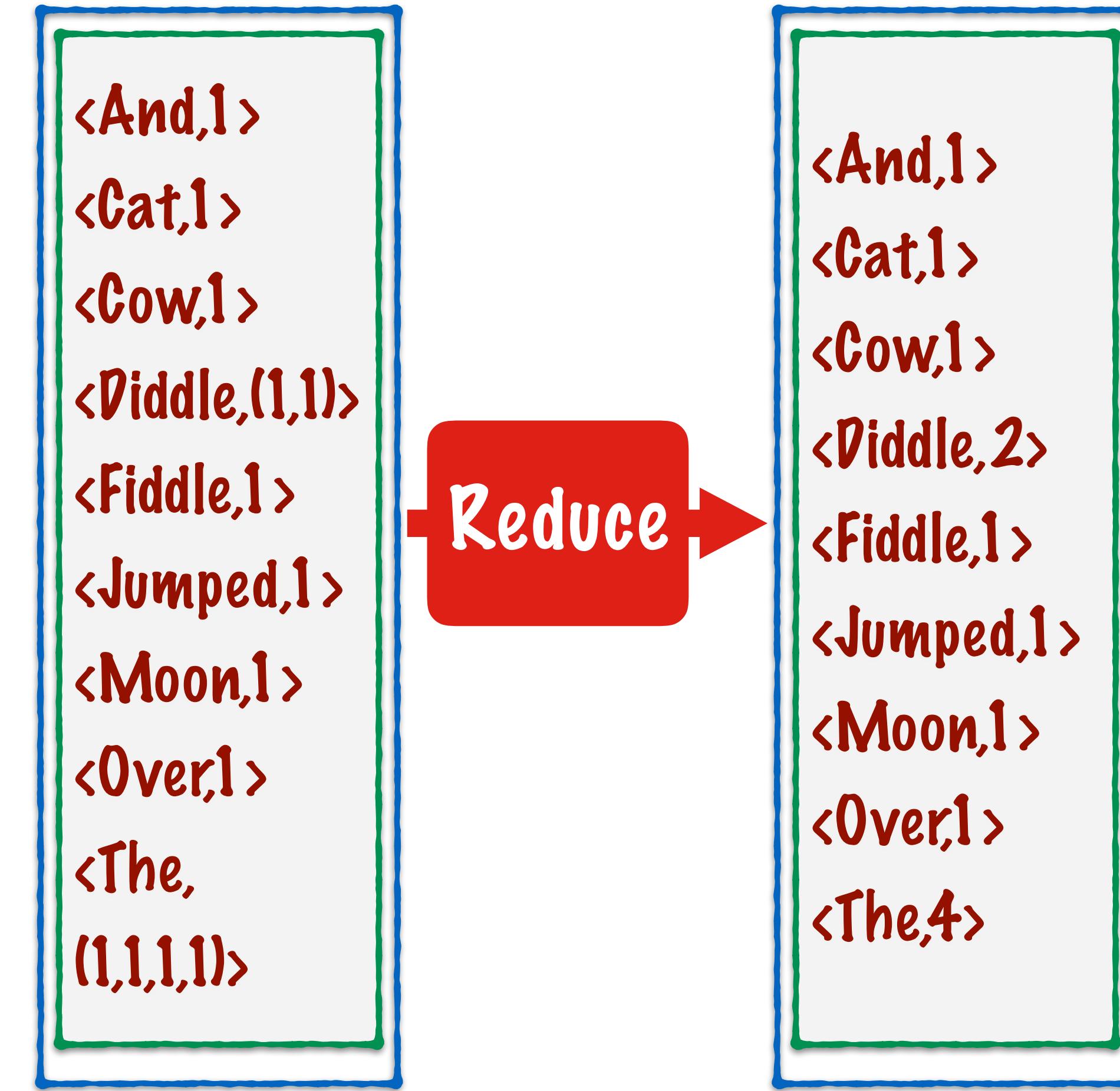
MAP

MAP

Input

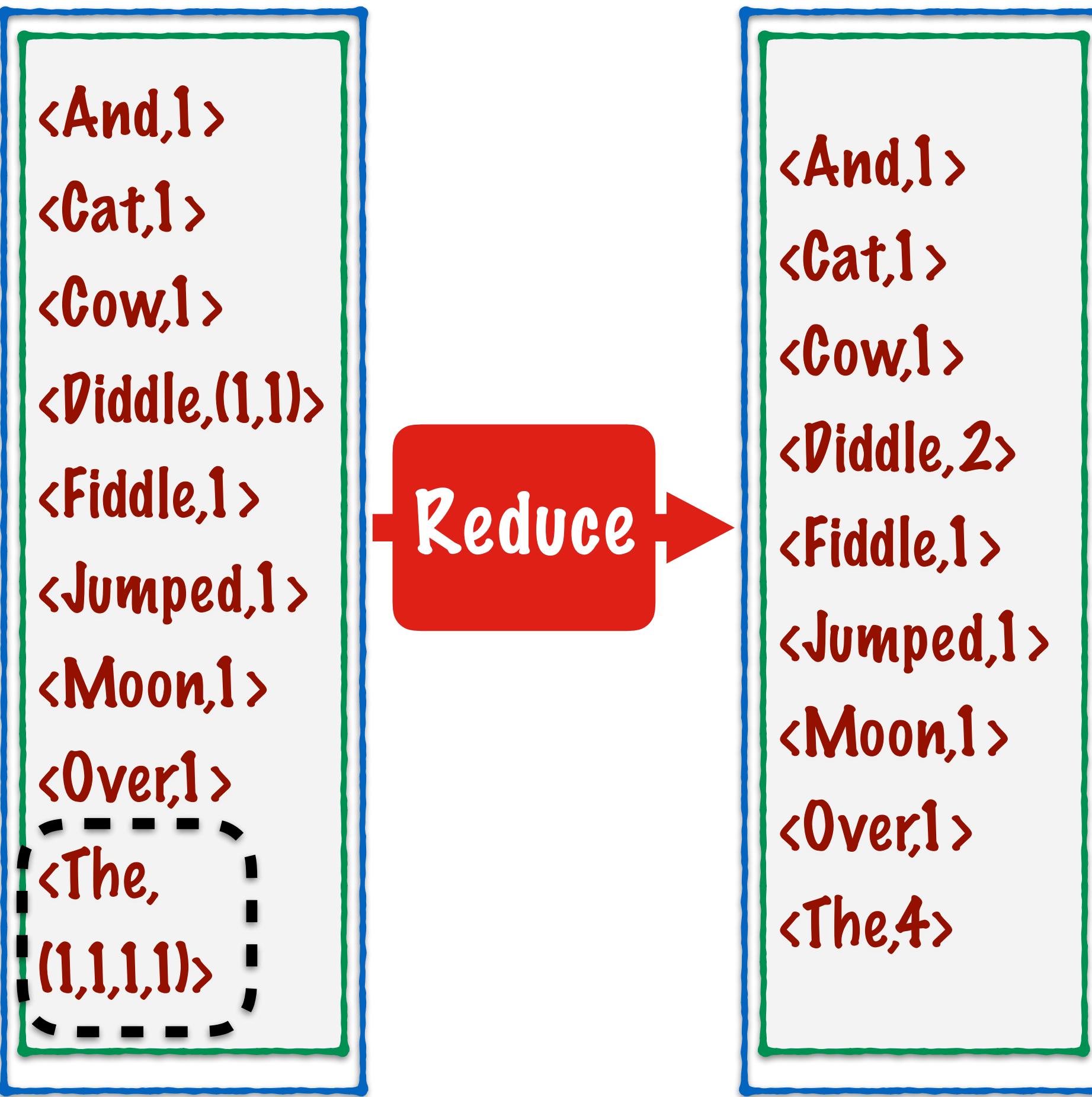


Output



reduce() will run on each pair
generated by the Sort/Merge step

<word, list>

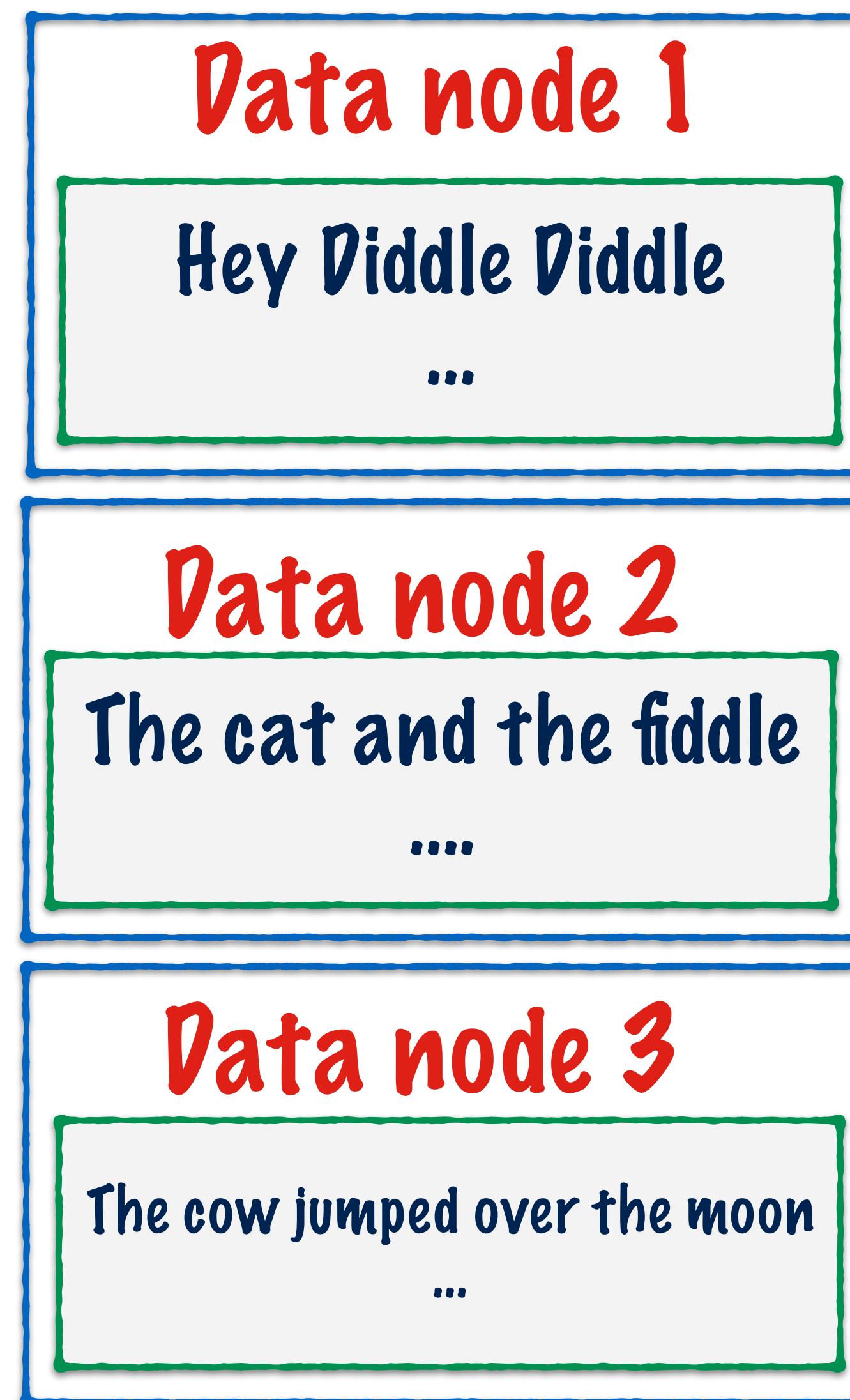


<word, list> **<word, count>**

<And,1>
<Cat,1>
<Cow,1>
<Diddle,(1,1)>
<Fiddle,1>
<Jumped,1>
<Moon,1>
<Over,1>
<The,
(1,1,1,1)>

Reduce →

<And,1>
<Cat,1>
<Cow,1>
<Diddle,2>
<Fiddle,1>
<Jumped,1>
<Moon,1>
~~<Over,1>~~
~~<The,4>~~
.....

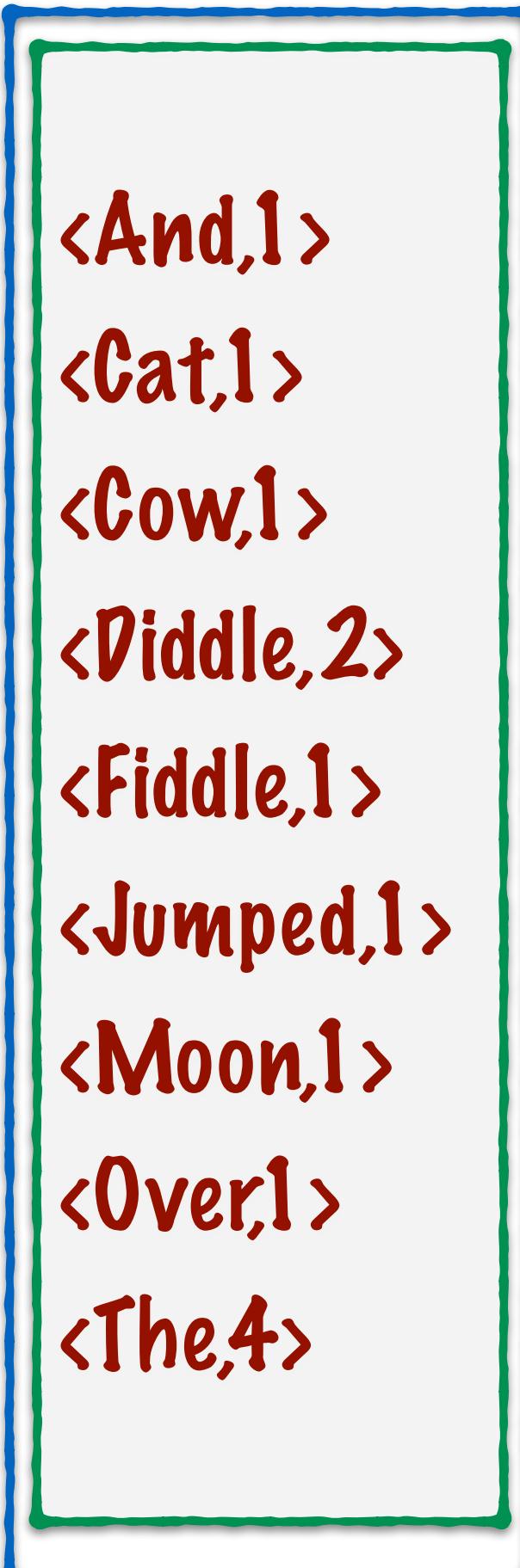
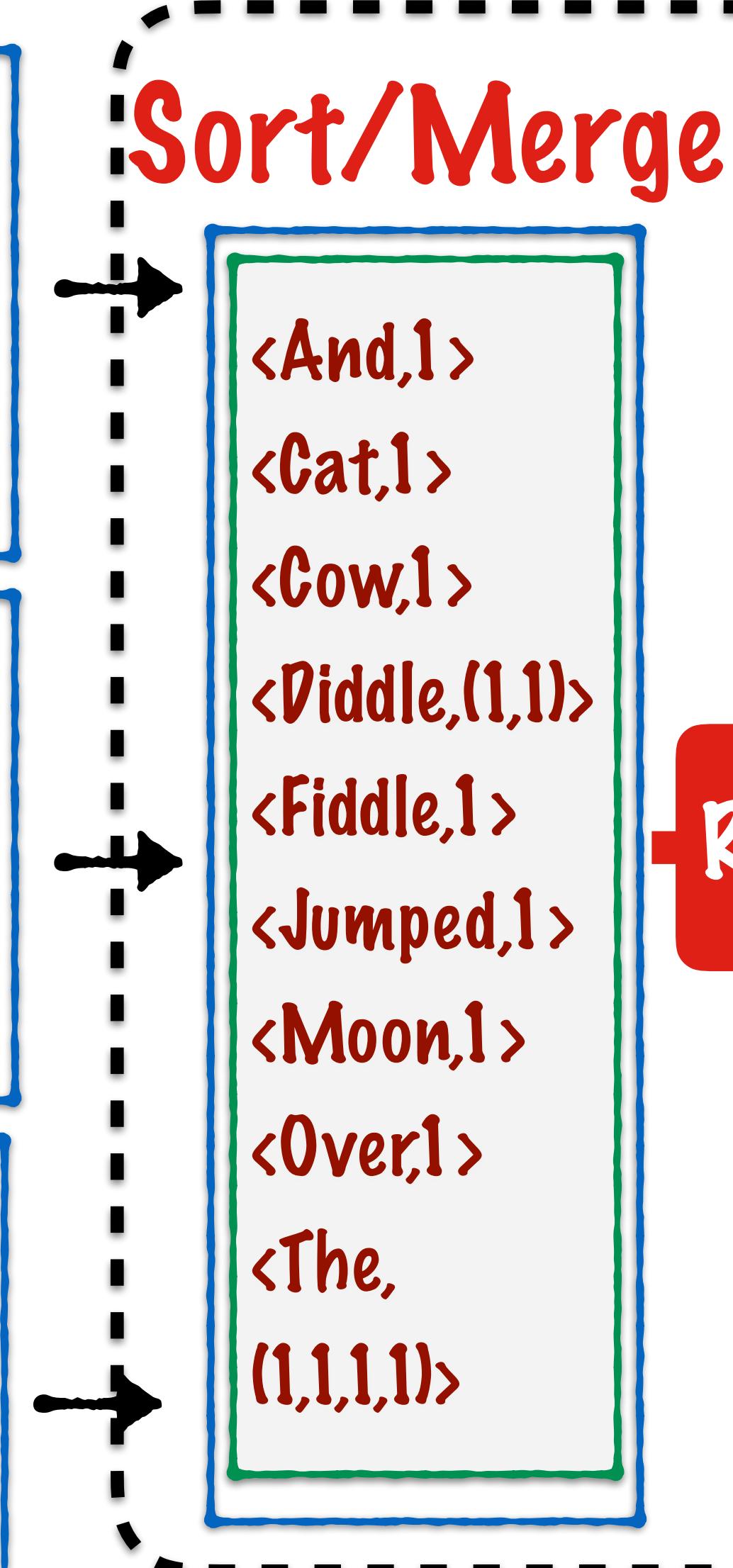
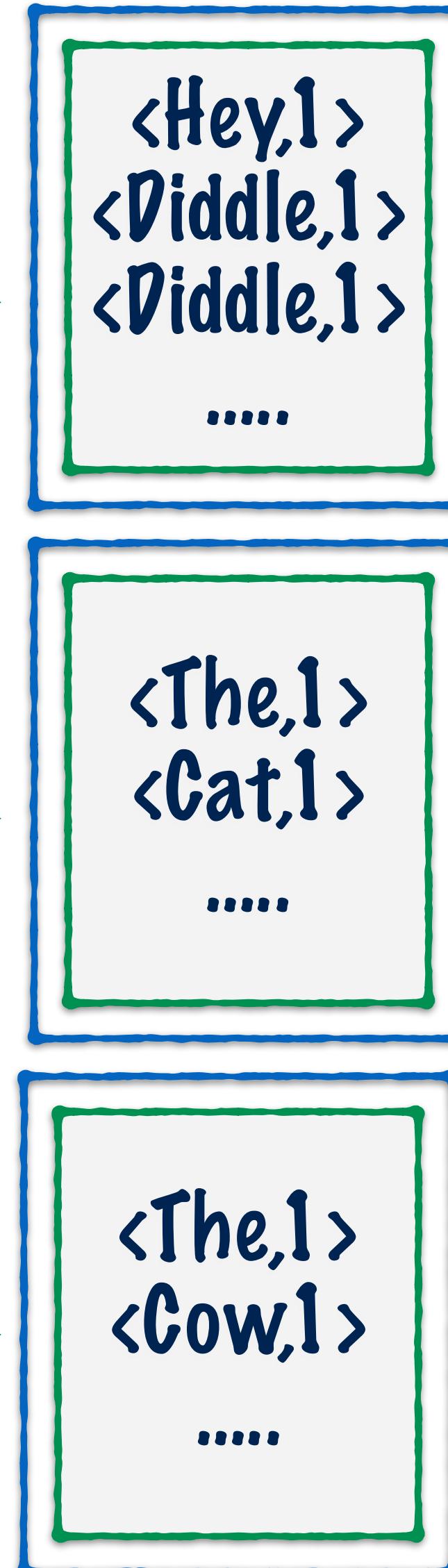


MAP

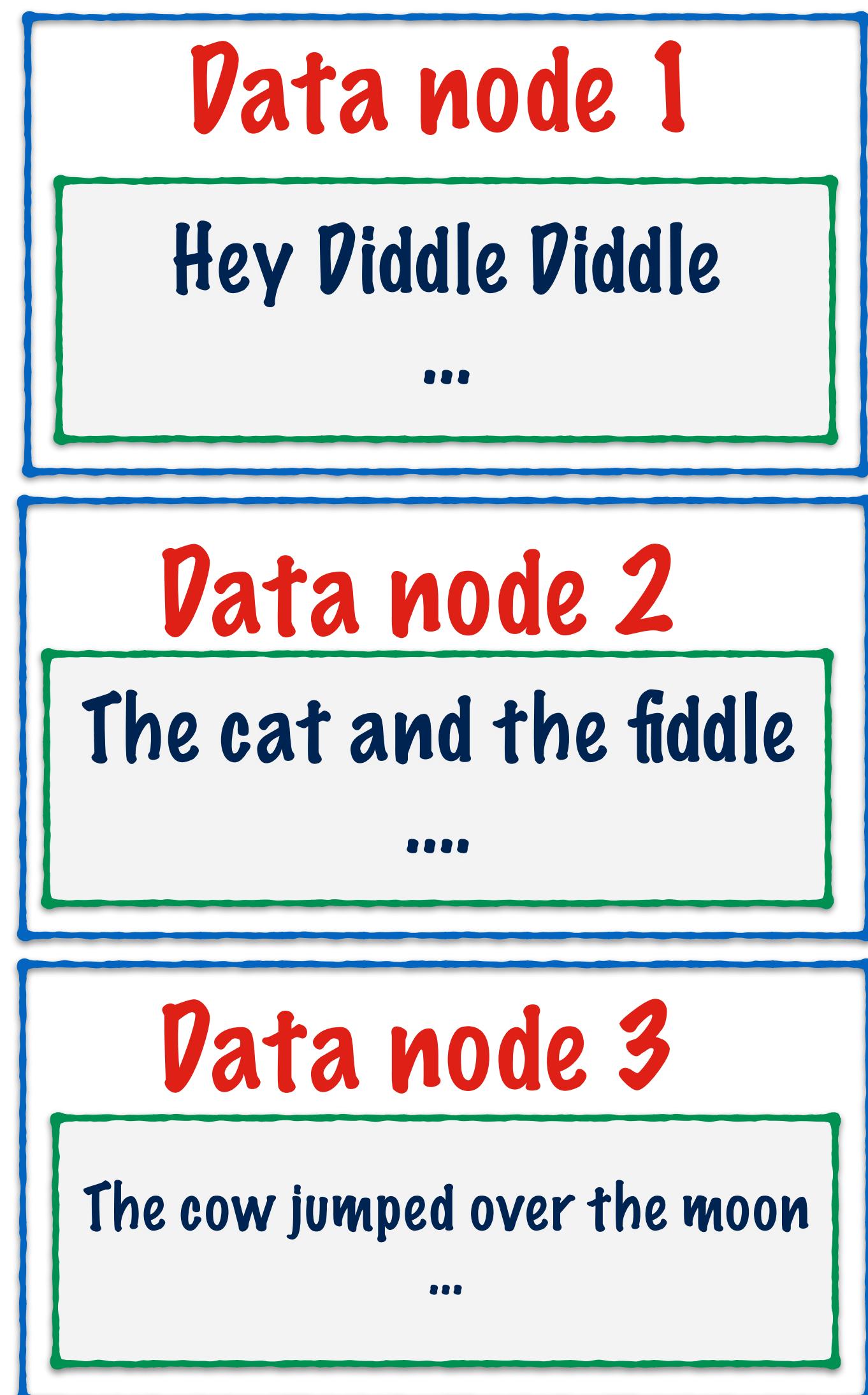
MAP

MAP

Input



Output



MAP →

MAP →

MAP →

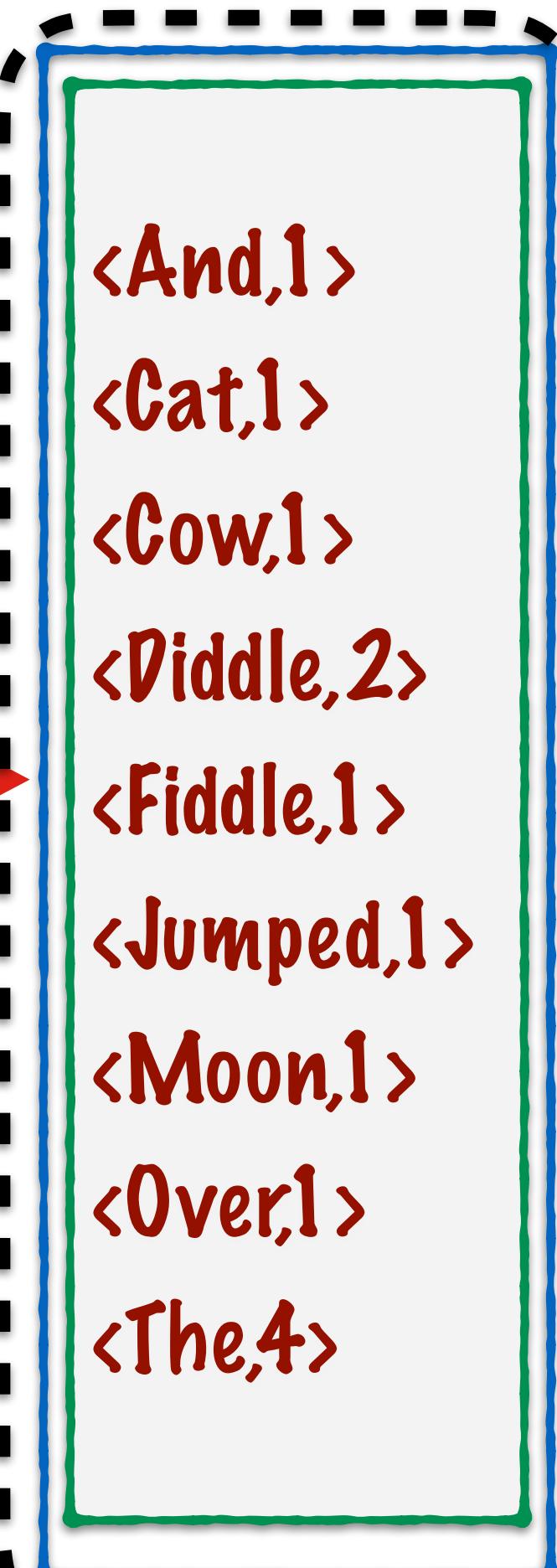
Input



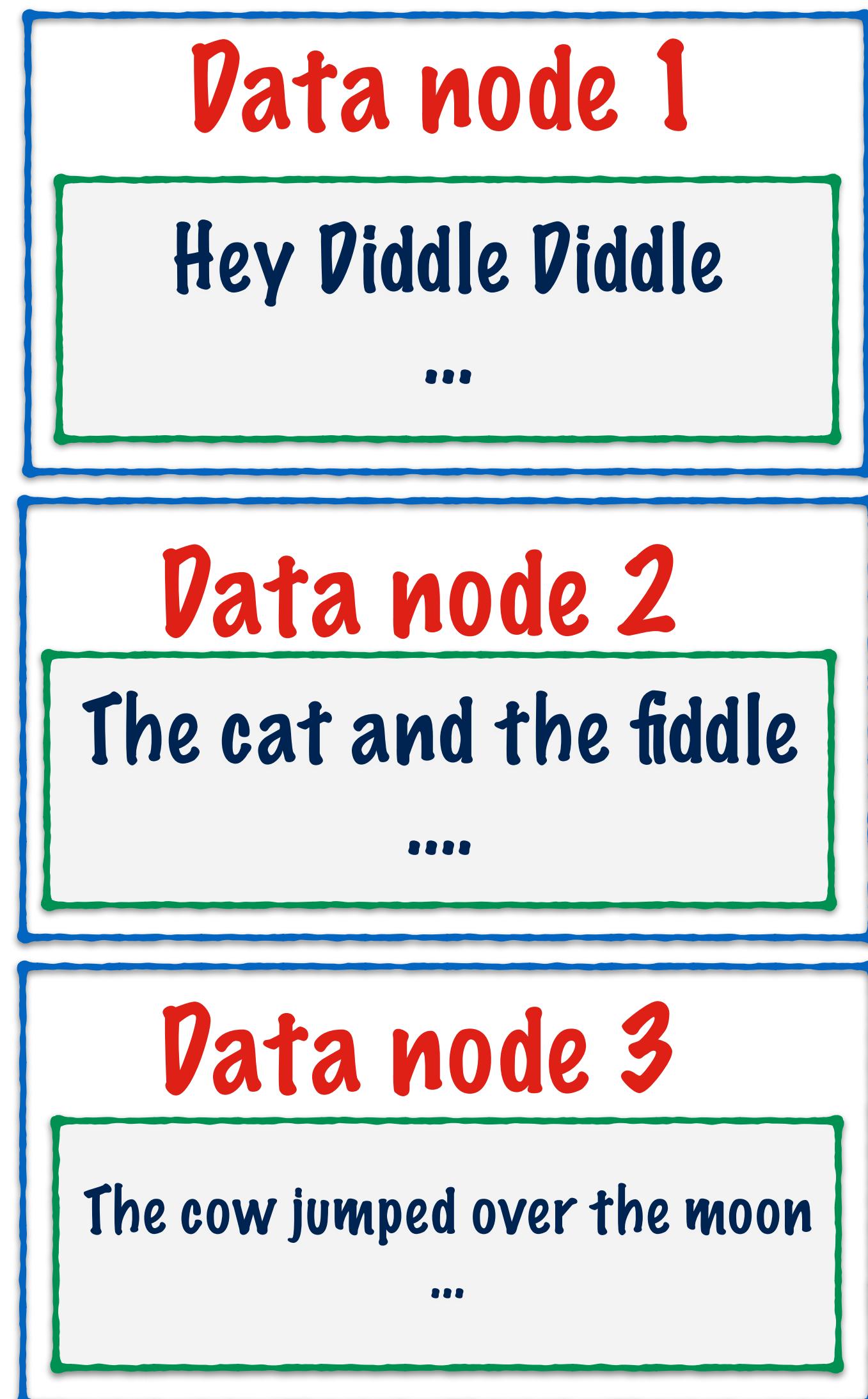
Sort/Merge



The output is sorted!!



Output



MAP →

<Hey,1>
<Diddle,1>
<Diddle,1>

MAP →

<And,1>
<Cat,1>

MAP →

<Cow,1>
....

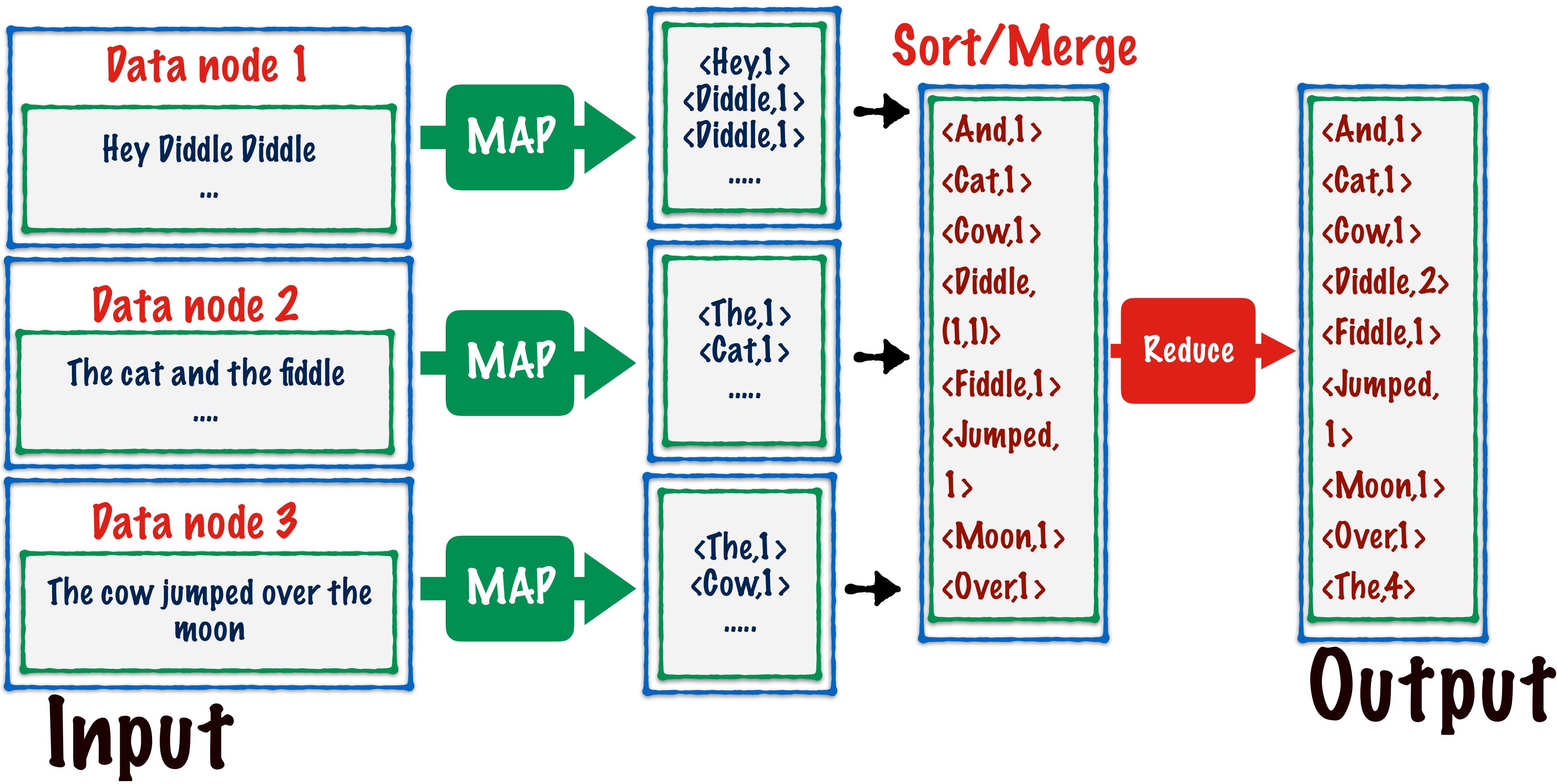
Input

Sort/Merge

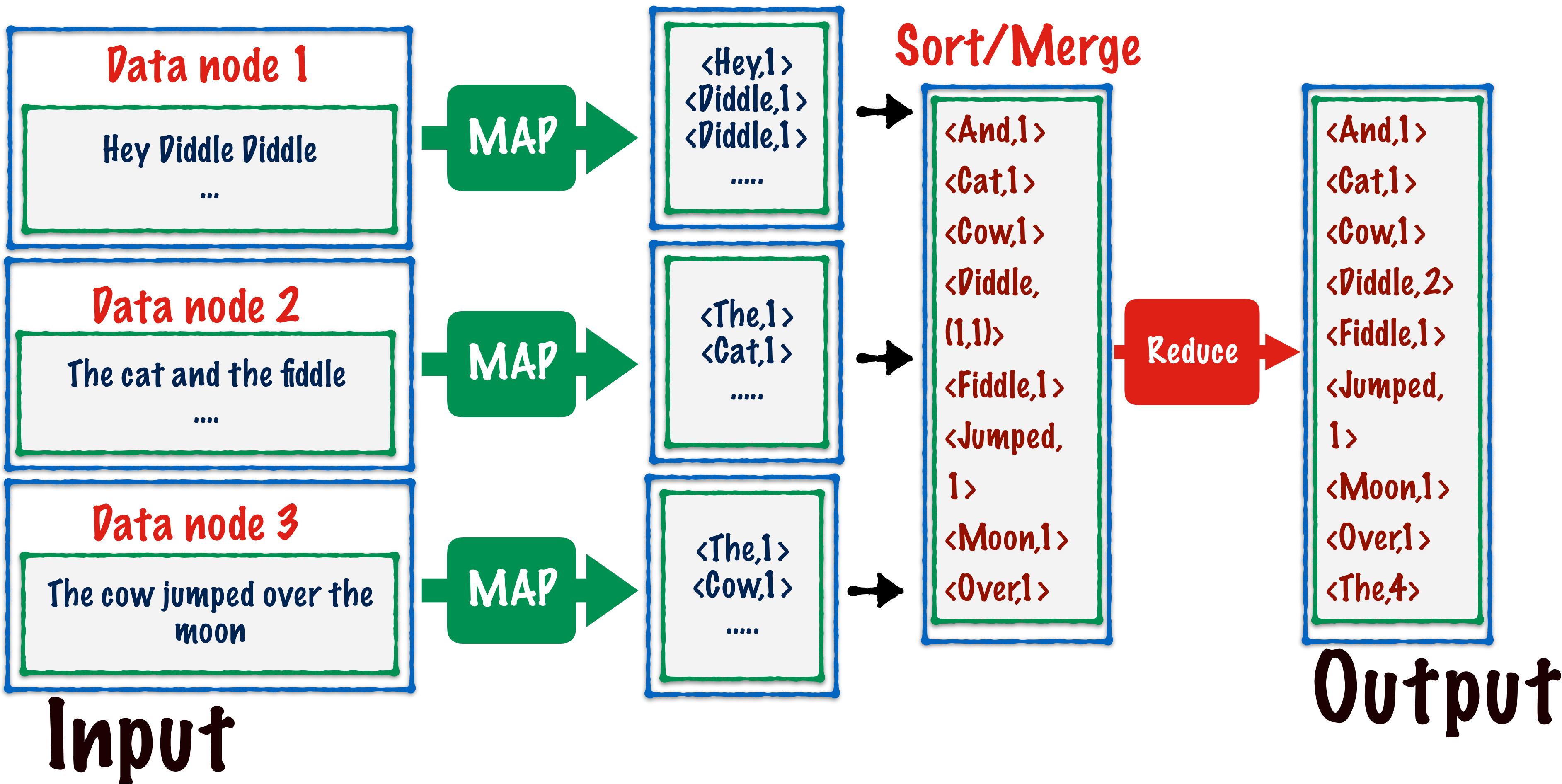
This is a happy
byproduct of
the Sort/Merge
operation

<And,1>
<Cat,1>
<Cow,1>
<Diddle,2>
<Fiddle,1>
<Jumped,1>
<Moon,1>
<Over,1>
<The,4>

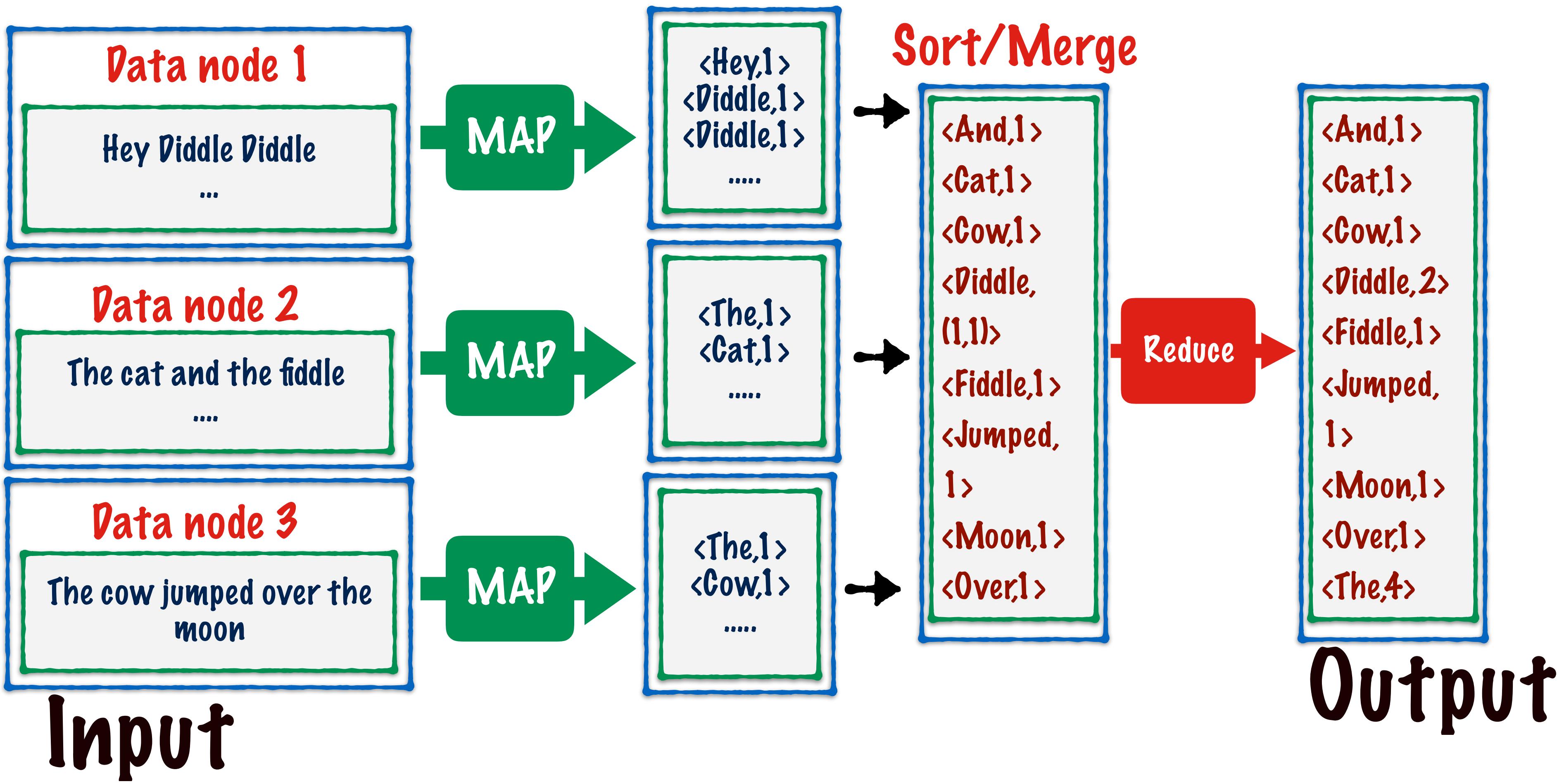
Output



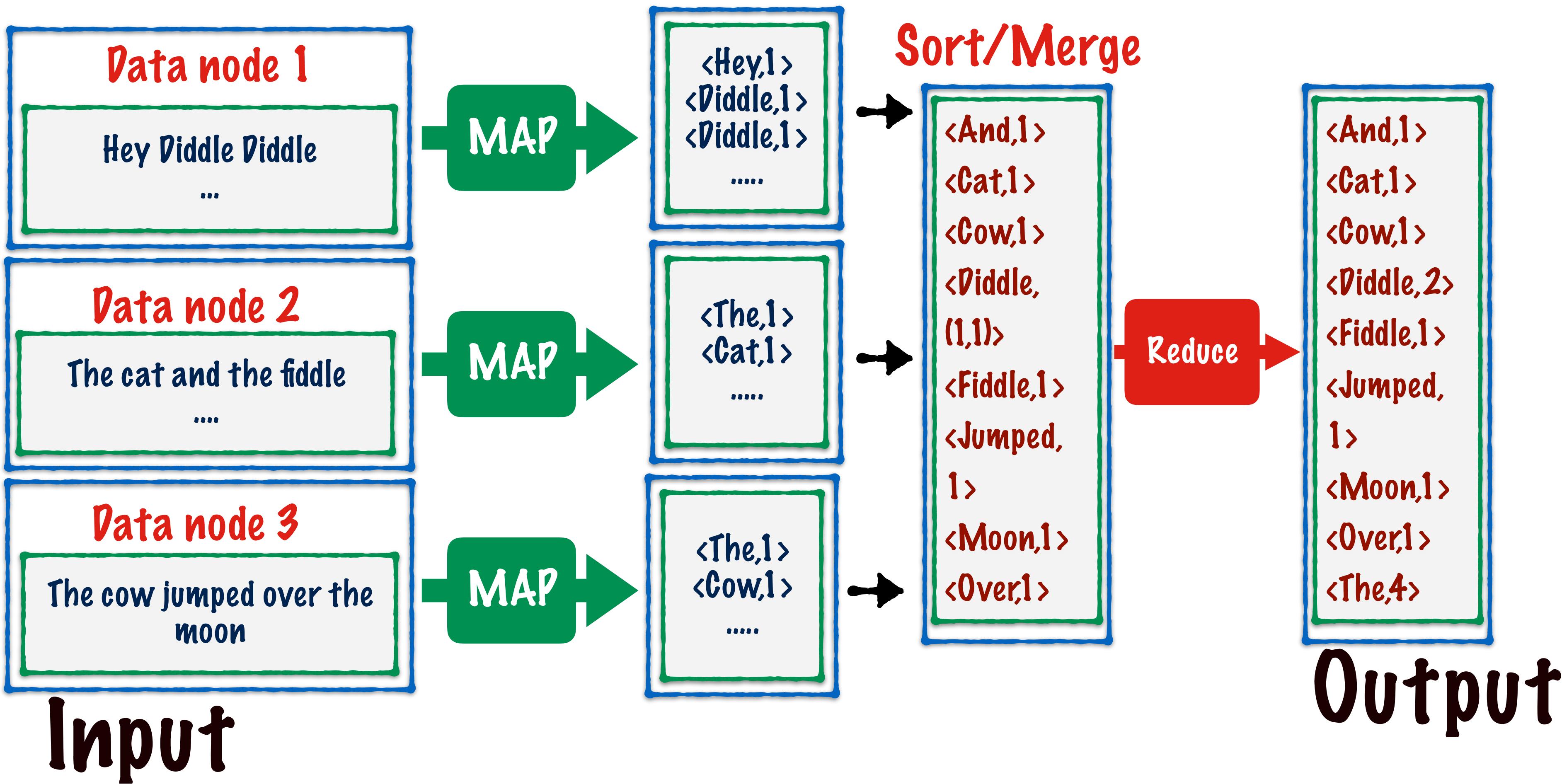
This is the general idea behind every
MapReduce operation



The map function is chosen such that it can run in parallel on all the nodes



The output of map should be key, value pairs



The `reduce()` function will combine all the values for the same key in some way

MapReduce Hello World

Let's write our first
MapReduce program

MapReduce Programs are
usually written in Java

Hadoop has an API called Streaming
that can be used for Ruby, Python etc
We'll get to that a little later

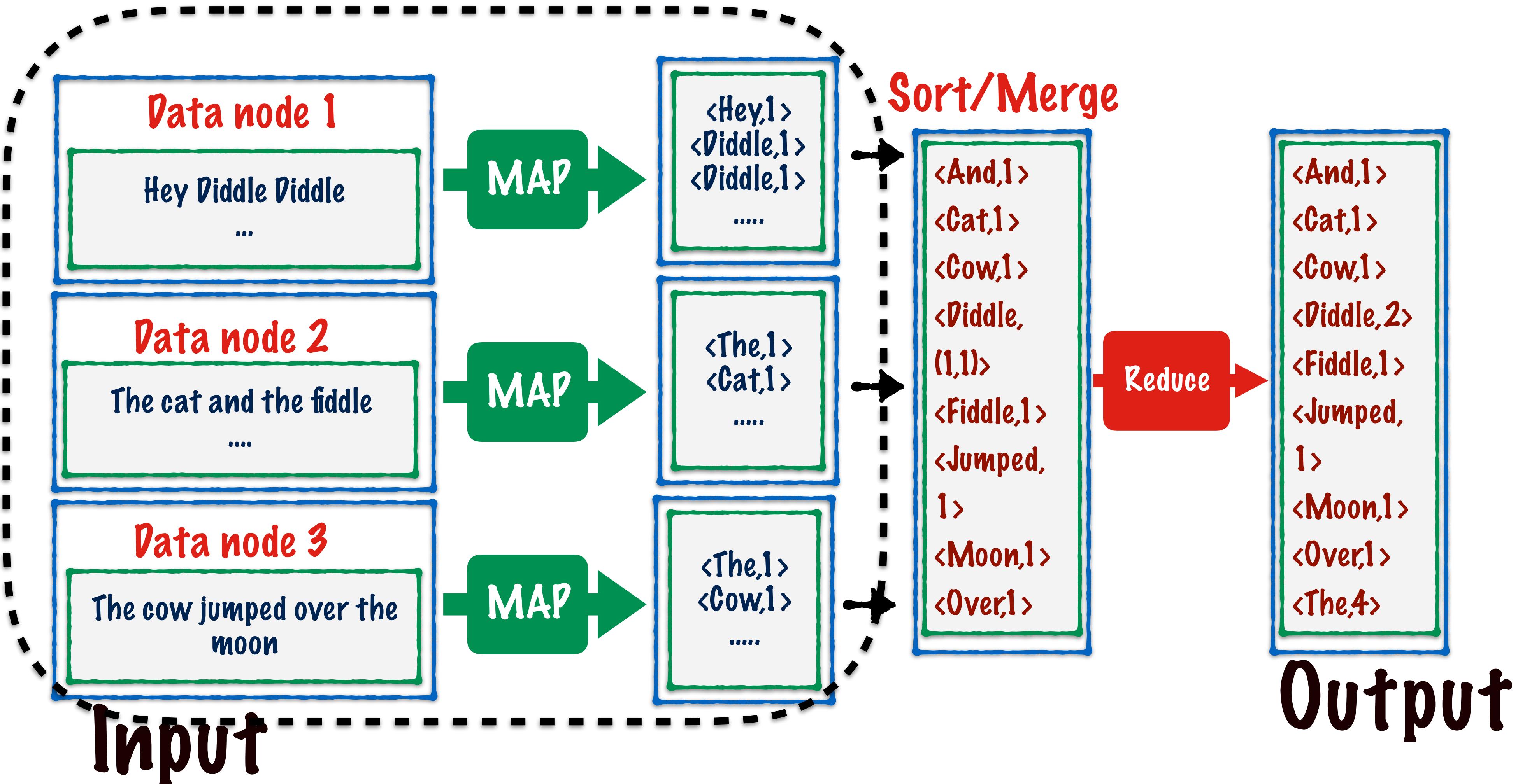
Objective: Create a **Frequency Distribution** of words in a text file

Step 1: Write a map() function

Step 2: Write a reduce() function

Step 3: Setup a driver that points to our map and reduce implementations

Objective: Create a Frequency Distribution of words in a text file



WordMapper Class

Mapper Class

Sort/Merge

<And,1>
<Cat,1>
<Cow,1>
<Diddle,
(1,1)>
<Fiddle,1>
<Jumped,
1>
<Moon,1>
<Over,1>
<Wh,1>

<And,1>
<Cat,1>
<Cow,1>
<Diddle,2>
<Fiddle,1>
<Jumped,
1>
<Moon,1>
<Over,1>
<Wh,4>

The map() function
is implemented in a
class that extends
the Mapper Class

WordMapper Class

`<input key type,
input value type,
output key type,
output value type>`

Mapper Class

Sort/Merge

| |
|-------------------------------|
| <code><And,1></code> |
| <code><Cat,1></code> |
| <code><Cow,1></code> |
| <code><Diddle,</code> |
| <code>(1,1)></code> |
| <code><Fiddle,1></code> |
| <code><Jumped,</code> |
| <code>1></code> |
| <code><Moon,1></code> |
| <code><Over,1></code> |
| <code><Over,1></code> |

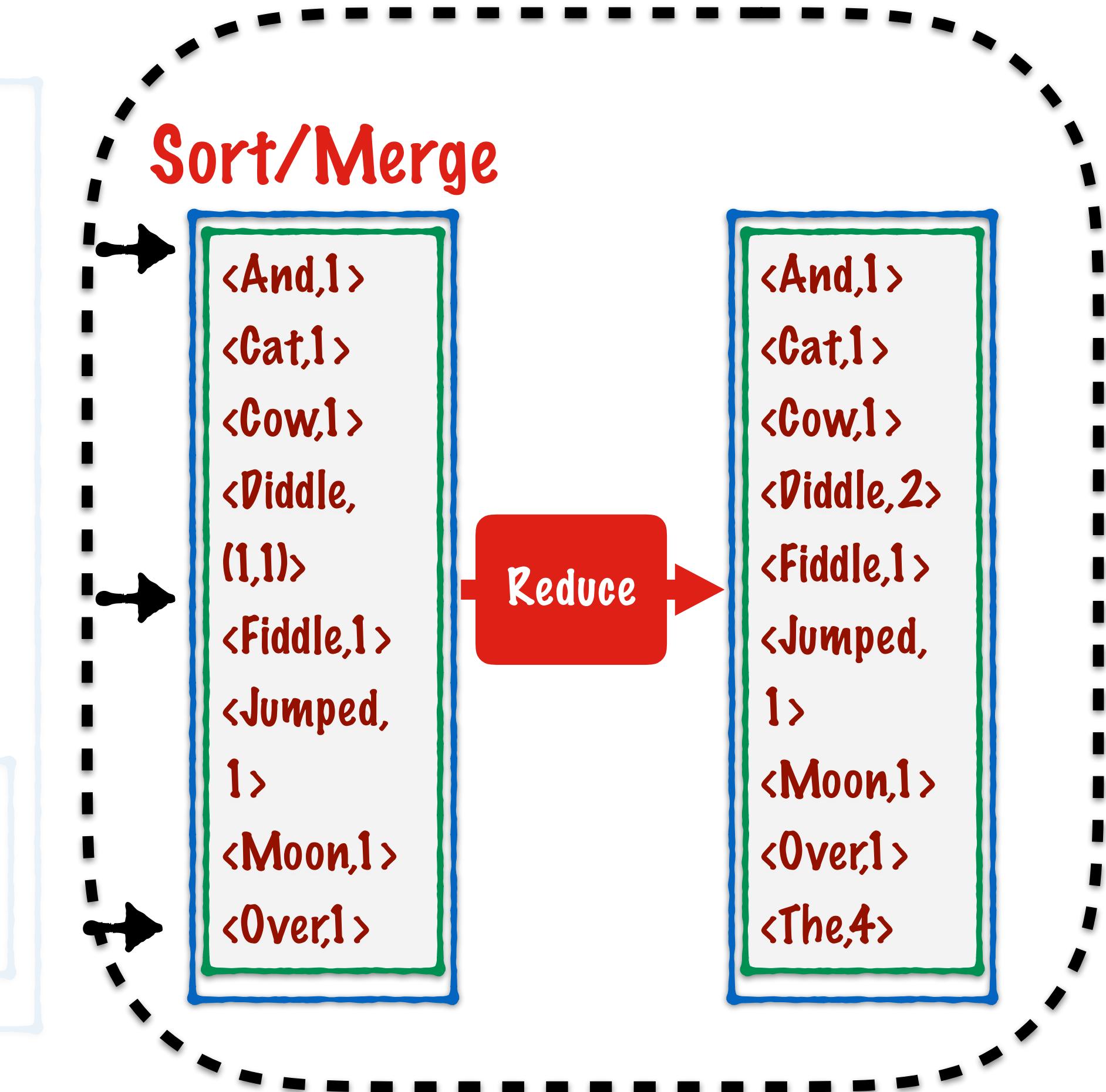


| |
|-------------------------------|
| <code><And,1></code> |
| <code><Cat,1></code> |
| <code><Cow,1></code> |
| <code><Diddle,2></code> |
| <code><Fiddle,1></code> |
| <code><Jumped,</code> |
| <code>1></code> |
| <code><Moon,1></code> |
| <code><Over,1></code> |
| <code><The,4></code> |

This is a generic
class, with 4
type parameters

WordMapper Class
`<input key type,
input value type,
output key type,
output value type>`

Mapper Class



WordMapper Class

```
<input key type,  
input value type,  
output key type,  
output value type>
```

Mapper Class

The reduce() function is implemented in a class that extends the Reducer Class

WordReducer Class

Reducer Class

WordMapper Class

`<input key type,
input value type,
output key type,
output value type>`

Mapper Class

WordReducer Class

`<input key type,
input value type,
output key type,
output value type>`

Reducer Class

This is a generic class too,
with 4 type parameters

WordMapper Class

`<input key type,
input value type,
output key type,
output value type>`

Mapper Class

WordReducer Class

`<input key type,
input value type,
output key type,
output value type>`

Reducer Class

The output types of the
Mapper should match the
input types of the Reducer

WordMapper Class

`<input key type,
input value type,
output key type,
output value type>`

Mapper Class

WordReducer Class

`<input key type,
input value type,
output key type,
output value type>`

Reducer Class

The output types of the
Mapper should match the
input types of the Reducer

WordMapper Class

WordReducer Class

These 2 classes are used
by a Job that is configured
in the Main Class

WordCount Class

Job Object

WordMapper Class

The Job has a bunch of properties that need to be configured

WordReducer Class

WordCount Class

Job Object

Input filepath

Output filepath

Mapper class

Reducer class

Output data types

WordMapper Class

The Mapper and Reducer will point to the classes with our implementation

WordReducer Class

WordCount Class

Job Object

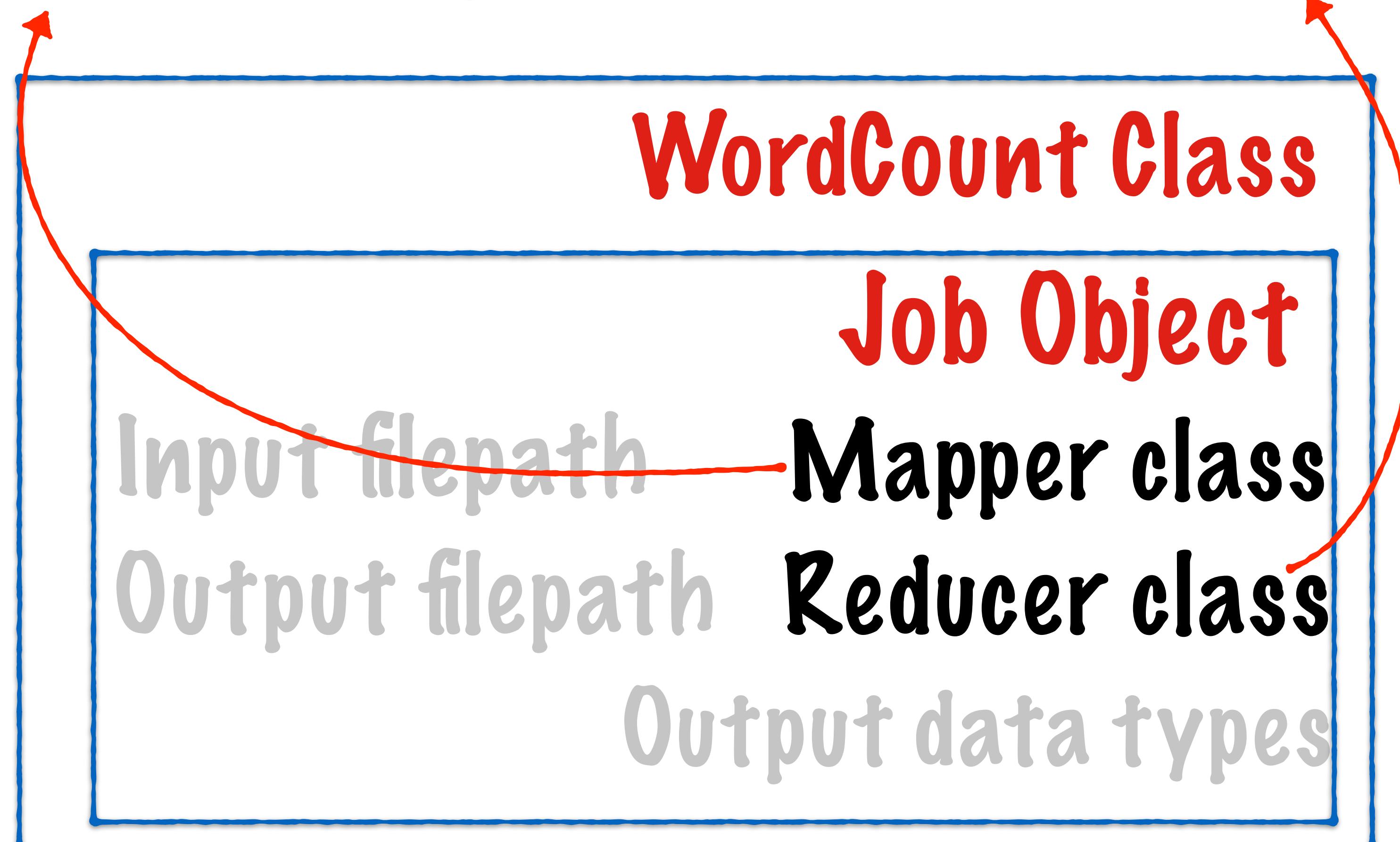
Input filepath

Output filepath

Output data types

Mapper class

Reducer class



WordMapper Class

The output data types need to match the type parameters of the Mapper and Reducer Classes

WordReducer Class

WordCount Class

Job Object

Input filepath

Output filepath

Mapper class

Reducer class

Output data types

All these classes
are put into a JAR
file which contains
all the Hadoop
JARS as well

WordMapper Class

WordReducer Class

WordCount Class

Job Object

When you **run** a job,
this JAR file **gets**
distributed to all the
nodes where the
computation is run

WordCount.JAR
WordMapper Class
WordReducer Class
WordCount Class

WordMapper Class

In Java the **map()** function is represented by the **Mapper Class**

```
package com.company;

import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;

import java.io.IOException;

public class WordCountMapper extends Mapper<LongWritable, Text, Text, IntWritable> {
    @Override
    public void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException {
        String line = value.toString();
        for (String word : line.split(" ")){
            if(word.length()>0){
                context.write(new Text(word), new IntWritable(1));
            }
        }
    }
}
```

WordMapper Class

package

```
import  
import  
import  
import  
  
import  
  
public class  
@Override
```

```
{
```

```
    String line = value.toString();
```

```
<1 ,Hey diddle diddle> → map() → <Hey, 1>  
context.write(
```

```
<Linenum, Line of text>
```

```
}
```

In Java the **map()** function is represented by the **Mapper Class**

This class represents the code for the step

```
<Hey, 1>  
<Diddle, 1>  
<Diddle, 1>  
<Word, 1>
```

WordMapper Class

```
package com.company;

import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;

import java.io.IOException;

public class WordCountMapper extends Mapper<LongWritable, Text, Text, IntWritable>
{
    @Override
    public void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException
    {
        String line = value.toString();
        for (String word : line.split(" ")){
            if(word.length()>0){
                context.write(new Text(word),new IntWritable(1));
            }
        }
    }
}
```

Mapper class is a generic It has 4 type parameters

Input
Key Type

Output
key type

Input
Value Type

Output
Value type

WordMapper Class

```
package com.company;

import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;

import java.io.IOException;

public class WordCountMapper extends Mapper<LongWritable, Text, Text, IntWritable>
{
    @Override
    public void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException{

        String line = value.toString();
        for (String word : line.split(" ")){
            if(word.length()>0){
                context.write(new Text(word),new IntWritable(1));
            }
        }
    }
}
```

Hadoop has it's own set of basic types optimized for network serialization

WordMapper Class

```
package com.company;

import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;

import java.io.IOException;

public class WordCountMapper extends Mapper<LongWritable, Text, Text, IntWritable> {
    @Override
    public void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException{
        String line = value.toString();
        for (String word : line.split(" ")){
            if(word.length()>0){
                context.write(new Text(word),new IntWritable(1));
            }
        }
    }
}
```

These are **wrappers** around
Java primitive types

WordMapper Class

```
package com.company;

import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;

import java.io.IOException;

public class WordCountMapper extends Mapper<LongWritable, Text, Text, IntWritable> {
    @Override
    public void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException{
        String line = value.toString();
        for (String word : line.split(" ")){
            if(word.length()>0){
                context.write(new Text(word),new IntWritable(1));
            }
        }
    }
}
```

These types are implemented using the **Writable Interface**

WordMapper Class

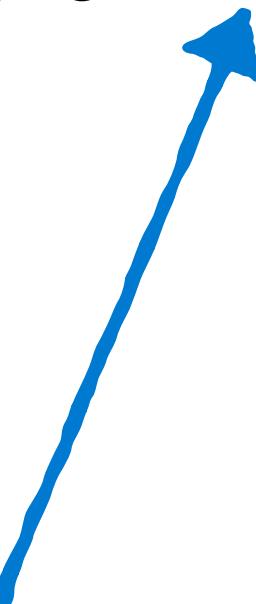
```
package com.company;

import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;

import java.io.IOException;

public class WordCountMapper extends Mapper<LongWritable, Text, Text, IntWritable> {
    @Override
    public void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException{
        String line = value.toString();
        for (String word : line.split(" ")){
            if(word.length() > 0){
                context.write(new Text(word), new IntWritable(1));
            }
        }
    }
}
```

Corresponding Hadoop type
for Java Long



WordMapper Class

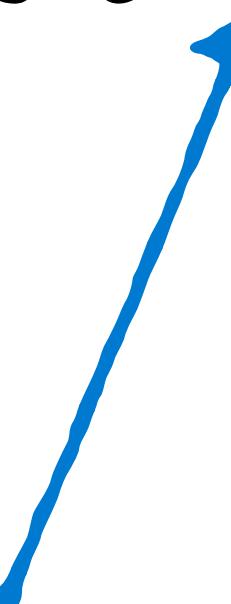
```
package com.company;

import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;

import java.io.IOException;

public class WordCountMapper extends Mapper<LongWritable, Text, Text, IntWritable> {
    @Override
    public void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException{
        String line = value.toString();
        for (String word : line.split(" ")){
            if(word.length()>0){
                context.write(new Text(word),new IntWritable(1));
            }
        }
    }
}
```

Corresponding Hadoop type
for Java String



Corresponding Hadoop type for Java Integer

```
package com.company;

import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;

import java.io.IOException;

public class WordCountMapper extends Mapper<LongWritable, Text, Text, IntWritable> {
    @Override
    public void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException{
        String line = value.toString();
        for (String word : line.split(" ")){
            if(word.length()>0){
                context.write(new Text(word), new IntWritable(1));
            }
        }
    }
}
```

WordMapper Class

```
package com.company;

import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;

import java.io.IOException;

public class WordCountMapper extends Mapper<LongWritable, Text, Text, IntWritable>
{
    @Override
    public void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException{

        String line = value.toString();
        for (String word : line.split(" ")){
            if(word.length()>0){
                context.write(new Text(word),new IntWritable(1));
            }
        }
    }
}
```

In our example, the input to map has the form
<lineNum , line of Text>



WordMapper Class

```
package com.company;

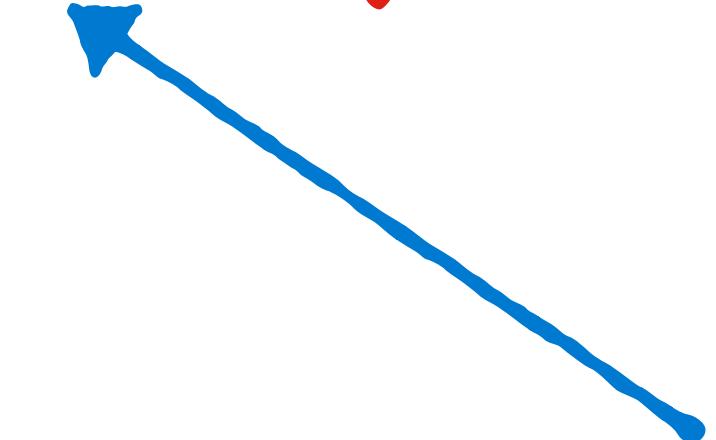
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;

import java.io.IOException;

public class WordCountMapper extends Mapper<LongWritable, Text, Text, IntWritable> {
    @Override
    public void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException{
        String line = value.toString();
        for (String word : line.split(" ")){
            if(word.length()>0){
                context.write(new Text(word),new IntWritable(1));
            }
        }
    }
}
```

In our example, the output of map has the form

<Word,1>



WordMapper Class

```
package com.company;

import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;

import java.io.IOException;

public class WordCountMapper extends Mapper<LongWritable, Text, Text, IntWritable> {
    @Override
    public void map(LongWritable key, Text value, Context context) throws
IOException, InterruptedException {

        String line = value.toString();
        for (String word : line.split(" ")){
            if(word.length()>0){
                context.write(new Text(word),new IntWritable(1));
            }
        }
    }
}
```

The Mapper class has an abstract method called **map()**



WordMapper Class

```
package com.company;

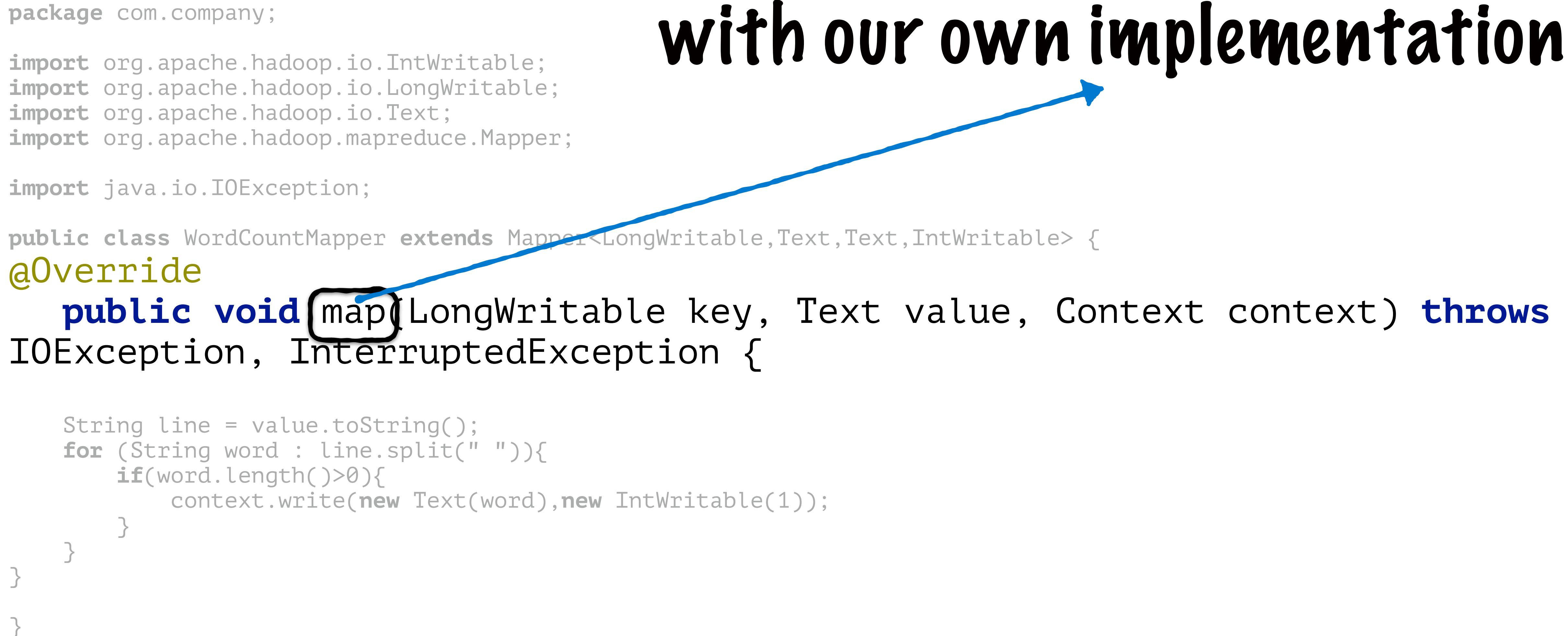
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;

import java.io.IOException;

public class WordCountMapper extends Mapper<LongWritable, Text, Text, IntWritable> {
    @Override
    public void map(LongWritable key, Text value, Context context) throws
    IOException, InterruptedException {

        String line = value.toString();
        for (String word : line.split(" ")){
            if(word.length()>0){
                context.write(new Text(word),new IntWritable(1));
            }
        }
    }
}
```

We'll override this method with our own implementation



WordMapper Class

```
package com.company;

import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;

import java.io.IOException;

public class WordCountMapper extends Mapper<LongWritable, Text, Text, IntWritable> {
    @Override
    public void map(LongWritable key, Text value, Context context) throws
IOException, InterruptedException {

        String line = value.toString();
        for (String word : line.split(" ")){
            if(word.length()>0){
                context.write(new Text(word),new IntWritable(1));
            }
        }
    }
}
```

The `map()` method
takes a Key and a Value

WordMapper Class

```
package com.company;

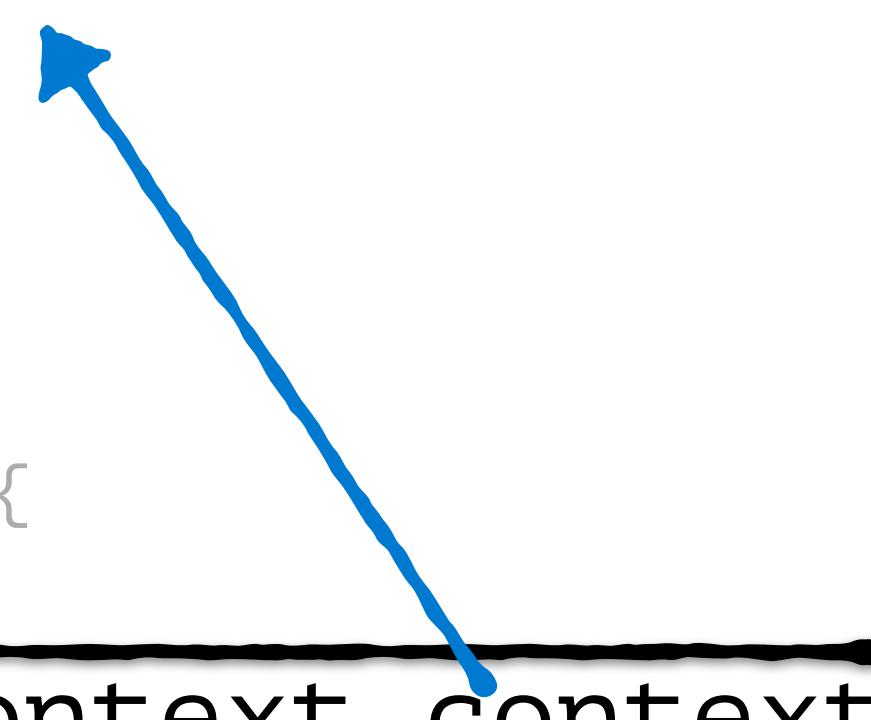
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;

import java.io.IOException;

public class WordCountMapper extends Mapper<LongWritable, Text, Text, IntWritable> {
    @Override
    public void map(LongWritable key, Text value, Context context) throws
IOException, InterruptedException {

        String line = value.toString();
        for (String word : line.split(" ")){
            if(word.length()>0){
                context.write(new Text(word),new IntWritable(1));
            }
        }
    }
}
```

It processes the input and writes the output to a Context object



WordMapper Class

```
package com.company;

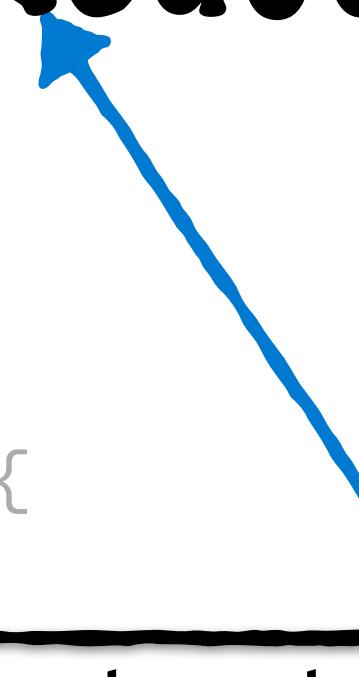
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;

import java.io.IOException;

public class WordCountMapper extends Mapper<LongWritable, Text, Text, IntWritable> {
    @Override
    public void map(LongWritable key, Text value, Context context) throws
IOException, InterruptedException {

        String line = value.toString();
        for (String word : line.split(" ")){
            if(word.length()>0){
                context.write(new Text(word),new IntWritable(1));
            }
        }
    }
}
```

The **Context object** stores the output and is accessed by the rest of the MapReduce system



throws

WordMapper Class

```
package com.company;

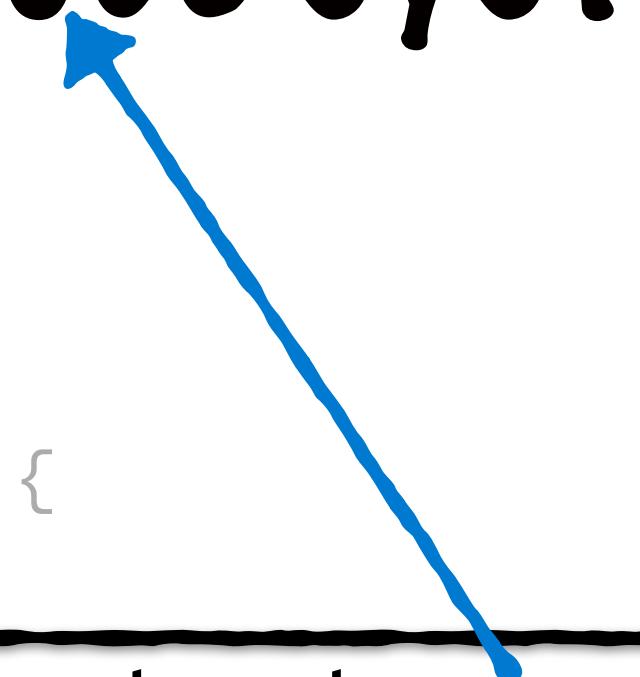
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;

import java.io.IOException;

public class WordCountMapper extends Mapper<LongWritable, Text, Text, IntWritable> {
    @Override
    public void map(LongWritable key, Text value, Context context) throws
IOException, InterruptedException {

        String line = value.toString();
        for (String word : line.split(" ")){
            if(word.length()>0){
                context.write(new Text(word),new IntWritable(1));
            }
        }
    }
}
```

It's how the Mapper / Reducer
communicate with the rest of the
MapReduce system



throws

WordMapper Class

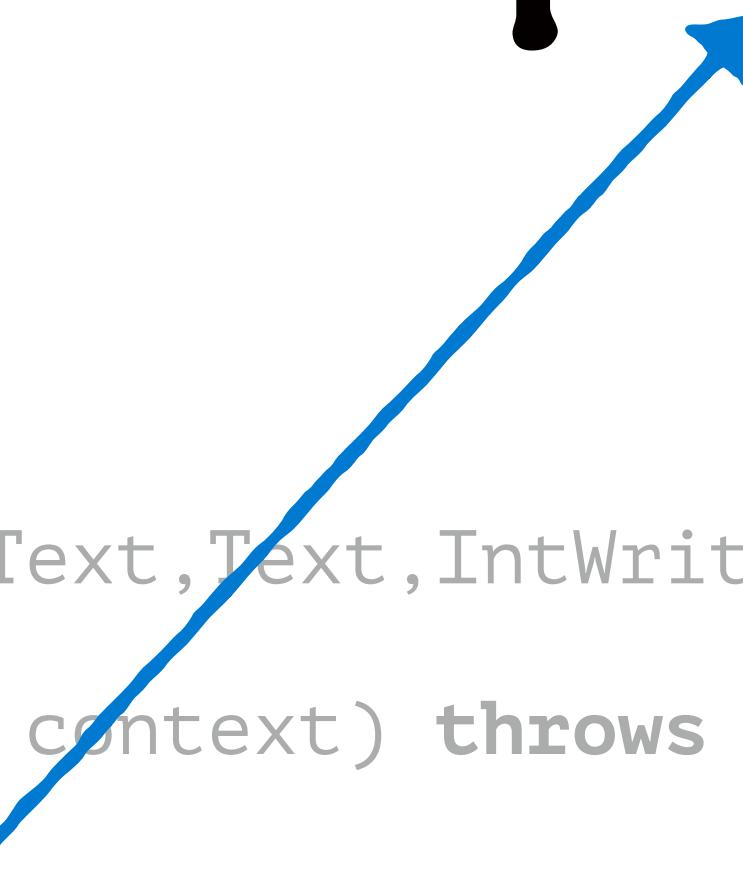
```
package com.company;

import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;

import java.io.IOException;

public class WordCountMapper extends Mapper<LongWritable, Text, Text, IntWritable> {
    @Override
    public void map(LongWritable key, Text value, Context context) throws IOException,
    InterruptedException {
        String line = value.toString();
        for (String word : line.split(" ")){
            if(word.length()>0){
                context.write(new Text(word), new IntWritable(1));
            }
        }
    }
}
```

Here we've implemented
the `map()` function



WordMapper Class

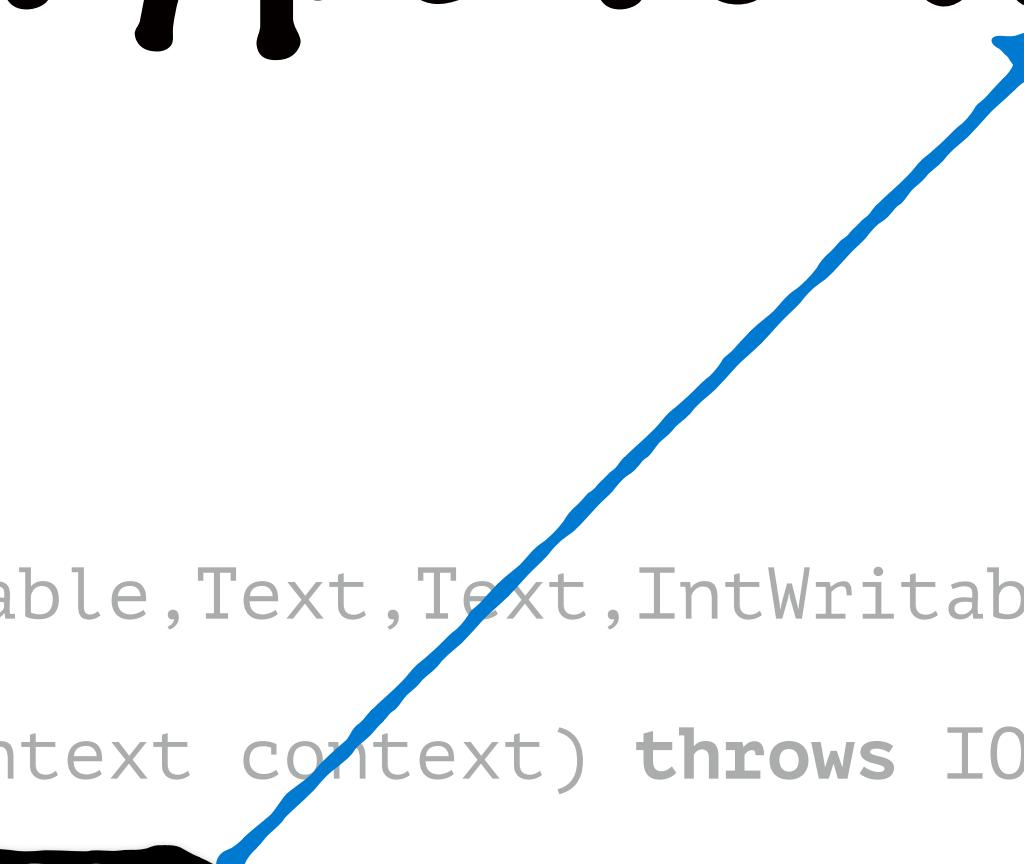
Convert the Hadoop Text type to Java String

```
package com.company;

import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;

import java.io.IOException;

public class WordCountMapper extends Mapper<LongWritable, Text, Text, IntWritable> {
    @Override
    public void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException {
        String line = value.toString();
        for (String word : line.split(" ")){
            if(word.length()>0){
                context.write(new Text(word), new IntWritable(1));
            }
        }
    }
}
```



WordMapper Class

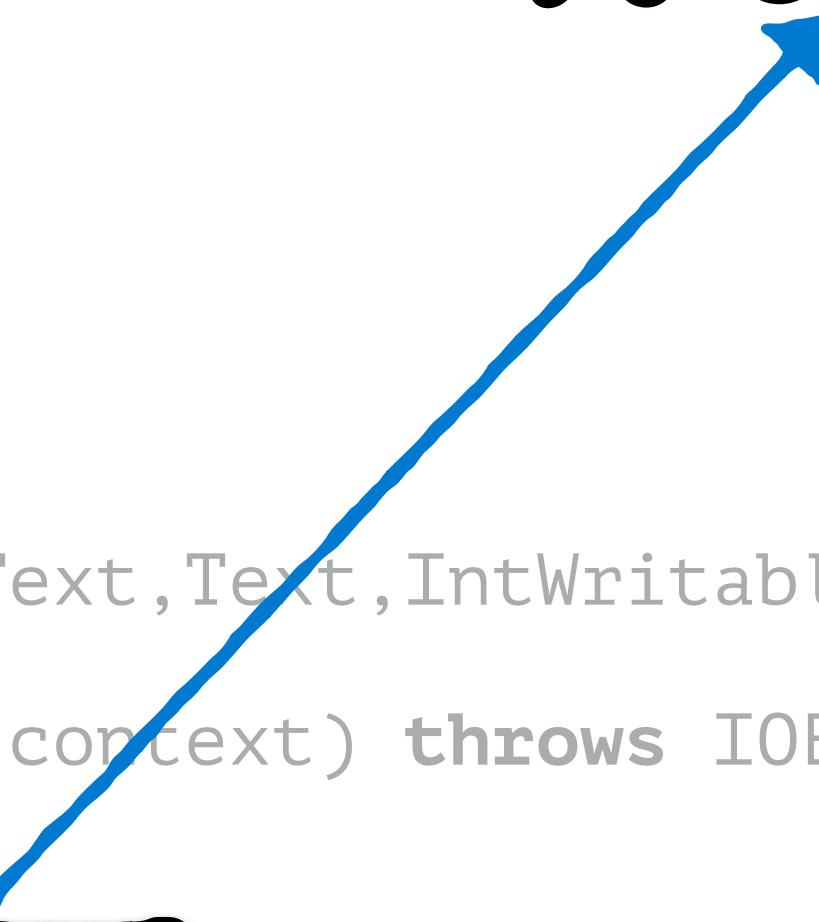
```
package com.company;

import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;

import java.io.IOException;

public class WordCountMapper extends Mapper<LongWritable, Text, Text, IntWritable> {
    @Override
    public void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException {
        String line = value.toString();
        for (String word : line.split(" ")){
            if(word.length() > 0){
                context.write(new Text(word), new IntWritable(1));
            }
        }
    }
}
```

Split the line into words



WordMapper Class

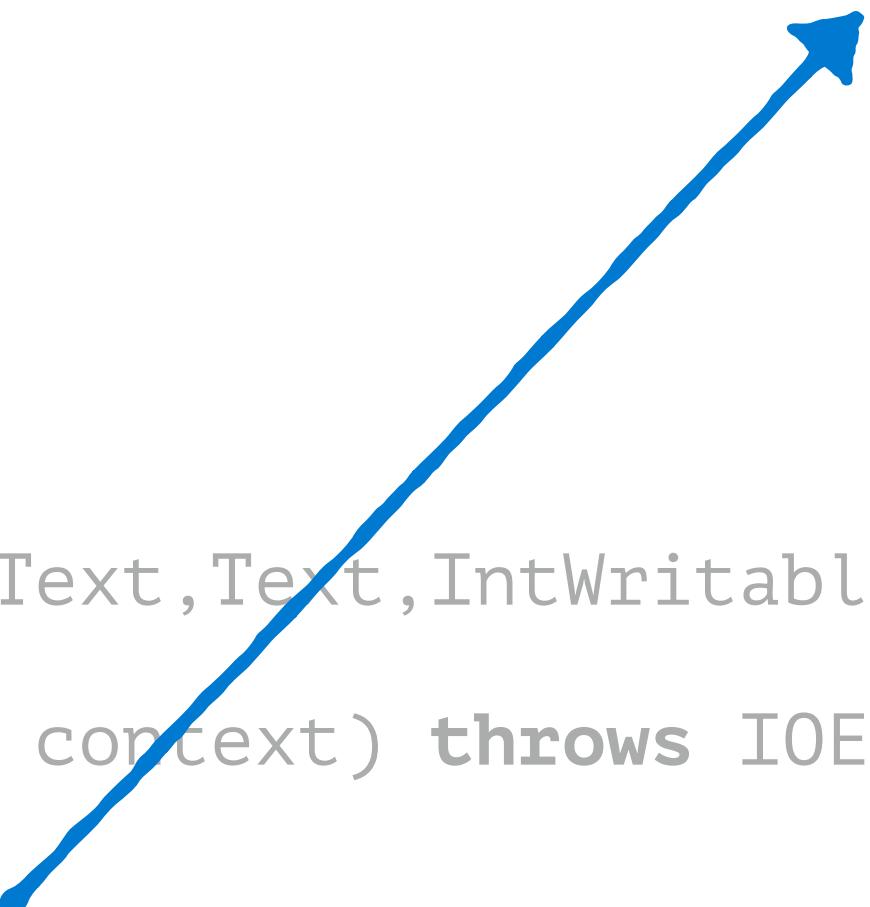
```
package com.company;

import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;

import java.io.IOException;

public class WordCountMapper extends Mapper<LongWritable, Text, Text, IntWritable> {
    @Override
    public void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException {
        String line = value.toString();
        for (String word : line.split(" ")){
            if(word.length() > 0){
                context.write(new Text(word), new IntWritable(1));
            }
        }
    }
}
```

For each of the words, we'll write out a key-value pair



WordMapper Class

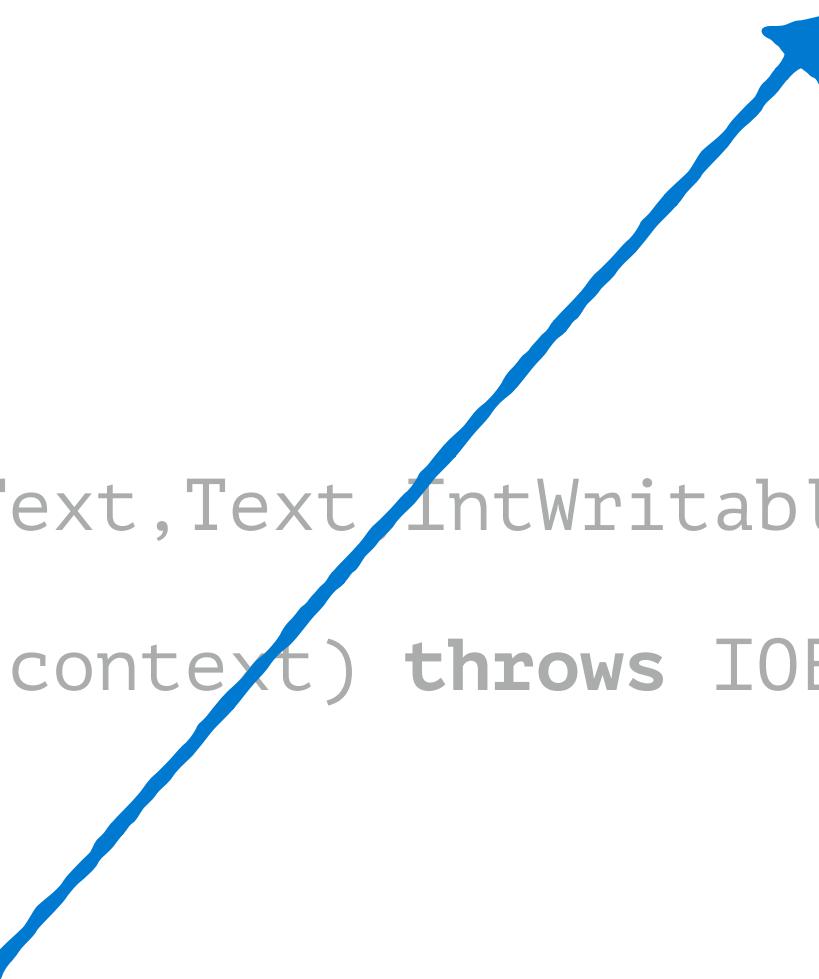
```
package com.company;

import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;

import java.io.IOException;

public class WordCountMapper extends Mapper<LongWritable, Text, Text, IntWritable> {
    @Override
    public void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException {
        String line = value.toString();
        for (String word : line.split(" ")){
            if(word.length() > 0){
                context.write(new Text(word), new IntWritable(1));
            }
        }
    }
}
```

For each of the words, we'll write out a key-value pair



WordMapper Class

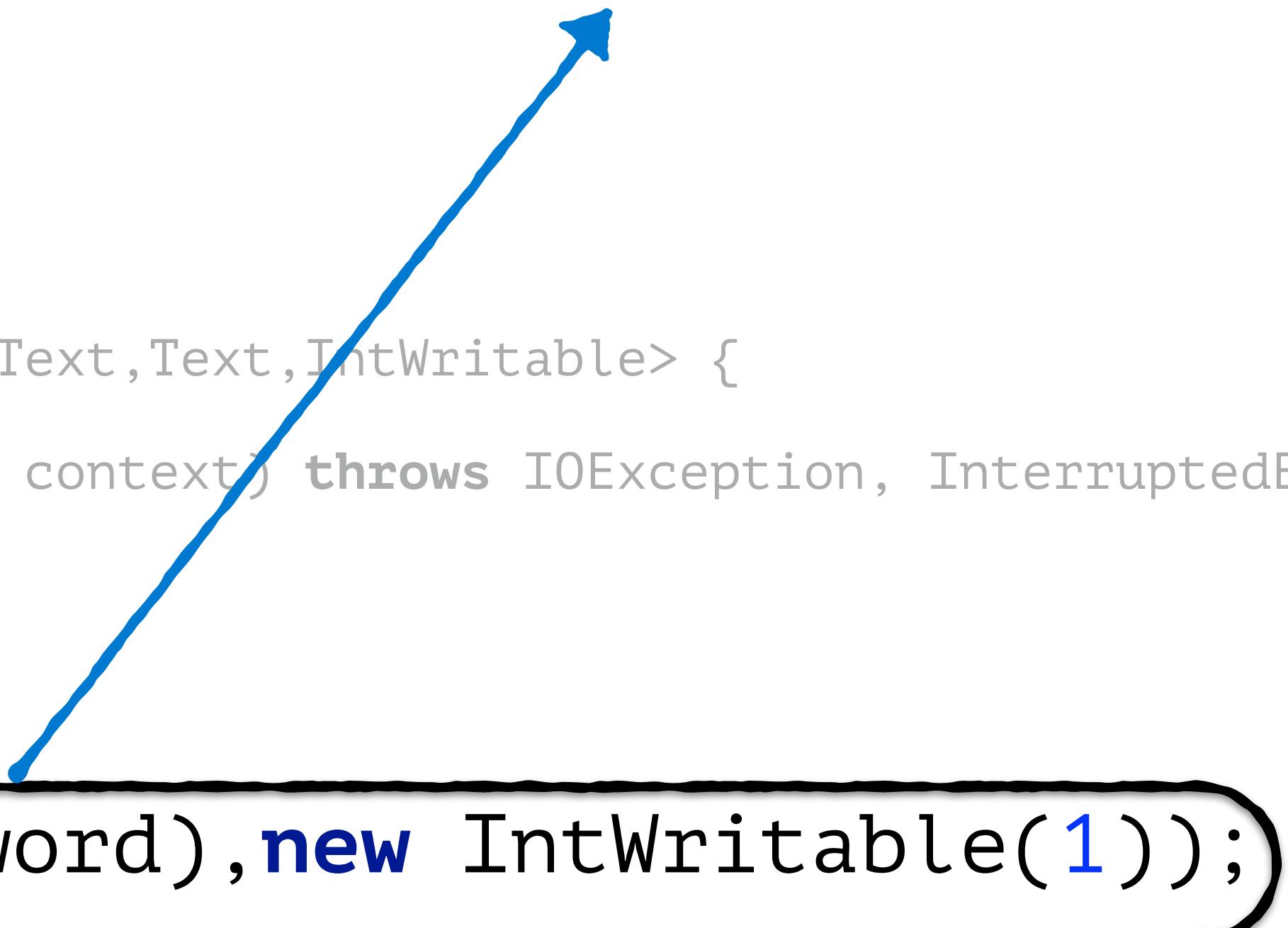
```
package com.company;

import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;

import java.io.IOException;

public class WordCountMapper extends Mapper<LongWritable, Text, Text, IntWritable> {
    @Override
    public void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException {
        String line = value.toString();
        for (String word : line.split(" ")){
            if(word.length()>0){
                context.write(new Text(word), new IntWritable(1));
            }
        }
    }
}
```

The pair **(word, 1)** is written to the context object



When you **run** a job,
this JAR file **gets**
distributed to all the
nodes where the
computation is run

WordCount.JAR
✓ WordMapper Class
WordReducer Class
WordCount Class

WordReducer Class

In Java the `reduce()` function
is represented by the
Reducer Class

```
package com.company;

import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Reducer;

import java.io.IOException;

public class WordCountReducer extends Reducer<Text, IntWritable, Text, IntWritable> {
    @Override
    public void reduce(Text key, Iterable<IntWritable> values, Context context) throws IOException,
    InterruptedException{
        int count = 0;
        for (IntWritable value:values)
        {
            count += value.get();
        }

        context.write(key, new IntWritable(count));
    }
}
```

WordReducer Class

package

```
import  
import  
import  
import
```

```
public class
```

```
InterruptedException{
```

```
<Diddle, (1, 1) > → reduce() → <Diddle, 2>  
<word, list> → <word, count>
```

This is created by taking the map output
and merging all the pairs with same key

In Java the **reduce()** function is
represented by the
Reducer Class

This class represents the
code for the step

WordReducer Class

```
package com.company;

import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Reducer;

import java.io.IOException;

public class WordCountReducer extends Reducer<Text, IntWritable, Text, IntWritable> {
    @Override
    public void reduce(Text key, Iterable<IntWritable> values, Context context) throws IOException,
    InterruptedException{
        int count = 0;
        for (IntWritable value:values)
        {
            count += value.get();
        }

        context.write(key, new IntWritable(count));
    }
}
```

Reducer class is a generic It has 4 type parameters

**Input
Key Type**

Type of each
element in
the list of
Input Values

**Output
key type**

Output
Value type

WordReducer Class

```
package com.company;

import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Reducer;

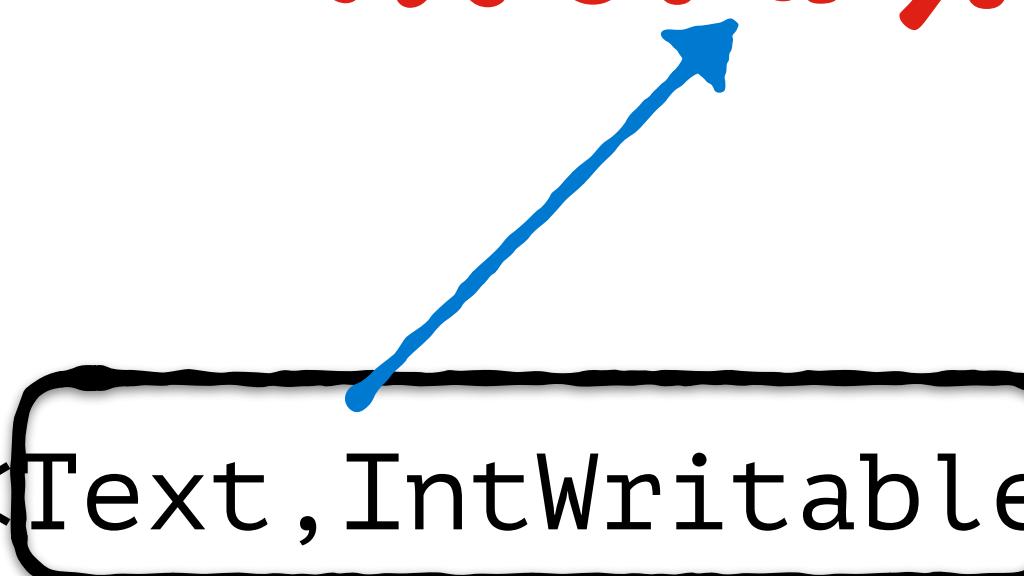
import java.io.IOException;

public class WordCountReducer extends Reducer<Text, IntWritable, Text, IntWritable> {
    @Override
    public void reduce(Text key, Iterable<IntWritable> values, Context context) throws IOException,
InterruptedException{
        int count = 0;
        for (IntWritable value:values)
        {
            count += value.get();
        }

        context.write(key, new IntWritable(count));
    }
}
```

In our example, the input to map has the form

<word ,list of ls>



WordReducer Class

```
package com.company;

import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Reducer;

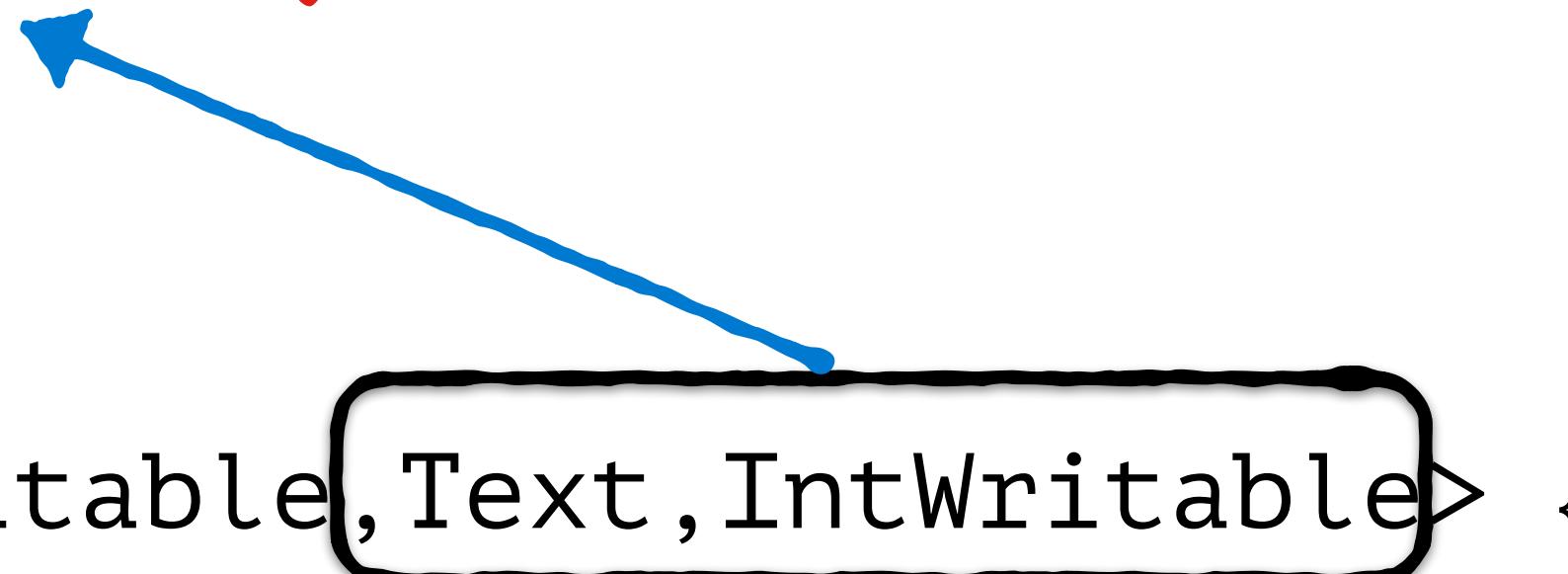
import java.io.IOException;

public class WordCountReducer extends Reducer<Text, IntWritable, Text, IntWritable> {
    @Override
    public void reduce(Text key, Iterable<IntWritable> values, Context context) throws IOException,
    InterruptedException{
        int count = 0;
        for (IntWritable value:values)
        {
            count += value.get();
        }

        context.write(key, new IntWritable(count));
    }
}
```

In our example, the output
of reduce has the form

<word , count>



WordReducer Class

The Reducer class has an abstract method called **reduce()**

```
package com.company;

import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Reducer;

import java.io.IOException;

public class WordCountReducer extends Reducer<Text, IntWritable, Text, IntWritable> {
    @Override
    public void reduce(Text key, Iterable<IntWritable> values, Context context)
throws IOException, InterruptedException{
        int count = 0;
        for (IntWritable value:values)
        {
            count += value.get();
        }

        context.write(key, new IntWritable(count));
    }
}
```

WordReducer Class

We'll override this method
with our own implementation

```
package com.company;

import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Reducer;

import java.io.IOException;

public class WordCountReducer extends Reducer<Text, IntWritable, Text, IntWritable> {
    @Override
    public void reduce(Text key, Iterable<IntWritable> values, Context context)
throws IOException, InterruptedException{
        int count = 0;
        for (IntWritable value:values)
        {
            count += value.get();
        }

        context.write(key, new IntWritable(count));
    }
}
```

WordReducer Class

The `reduce()` method takes a Key and an Iterable of Values

```
package com.company;

import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Reducer;

import java.io.IOException;

public class WordCountReducer extends Reducer<Text, IntWritable, Text, IntWritable> {
    @Override
    public void reduce(Text key, Iterable<IntWritable> values, Context context)
        throws IOException, InterruptedException{
        int count = 0;
        for (IntWritable value:values)
        {
            count += value.get();
        }

        context.write(key, new IntWritable(count));
    }
}
```

WordReducer Class

It processes the input and writes the output to a Context object

```
package com.company;

import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Reducer;

import java.io.IOException;

public class WordCountReducer extends Reducer<Text, IntWritable, Text, IntWritable> {
    @Override
    public void reduce(Text key, Iterable<IntWritable> values, Context context)
        throws IOException, InterruptedException{
        int count = 0;
        for (IntWritable value:values)
        {
            count += value.get();
        }

        context.write(key, new IntWritable(count));
    }
}
```

WordReducer Class

```
package com.company;

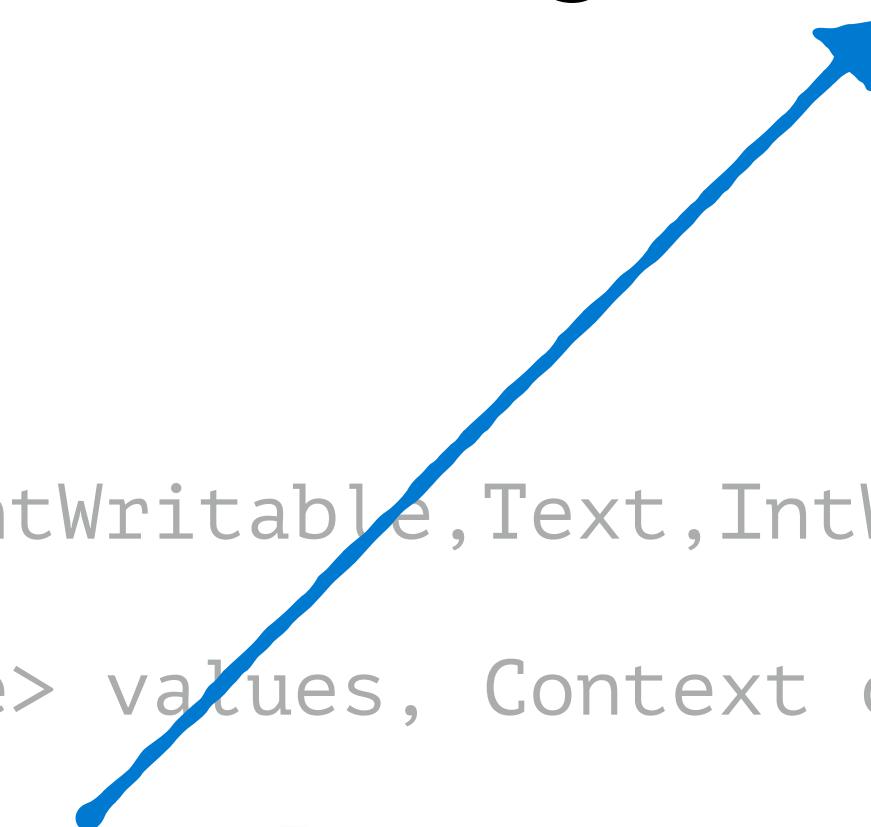
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Reducer;

import java.io.IOException;

public class WordCountReducer extends Reducer<Text, IntWritable, Text, IntWritable> {
    @Override
    public void reduce(Text key, Iterable<IntWritable> values, Context context) throws IOException,
InterruptedException{
        int count = 0;
        for (IntWritable value:values)
        {
            count += value.get();
        }

        context.write(key, new IntWritable(count));
    }
}
```

Here we've implemented
the reduce() function



WordReducer Class

```
package com.company;

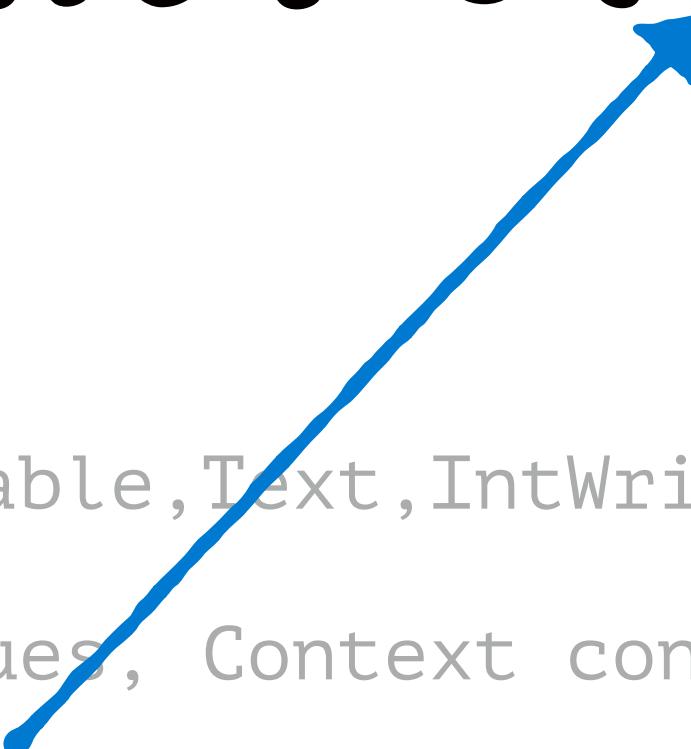
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Reducer;

import java.io.IOException;

public class WordCountReducer extends Reducer<Text, IntWritable, Text, IntWritable> {
    @Override
    public void reduce(Text key, Iterable<IntWritable> values, Context context) throws IOException,
    InterruptedException{
        int count = 0;
        for (IntWritable value:values)
        {
            count += value.get();
        }

        context.write(key, new IntWritable(count));
    }
}
```

We iterate through the
list of values



WordReducer Class

```
package com.company;

import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Reducer;

import java.io.IOException;

public class WordCountReducer extends Reducer<Text, IntWritable, Text, IntWritable> {
    @Override
    public void reduce(Text key, Iterable<IntWritable> values, Context context) throws IOException,
InterruptedException{
        int count = 0;
        for (IntWritable value:values)
        {
            count += value.get();
        }
        context.write(key, new IntWritable(count));
    }
}
```

And compute the sum
of the values

<Diddle , (1 , 1) >

Sum the values in this list

WordReducer Class

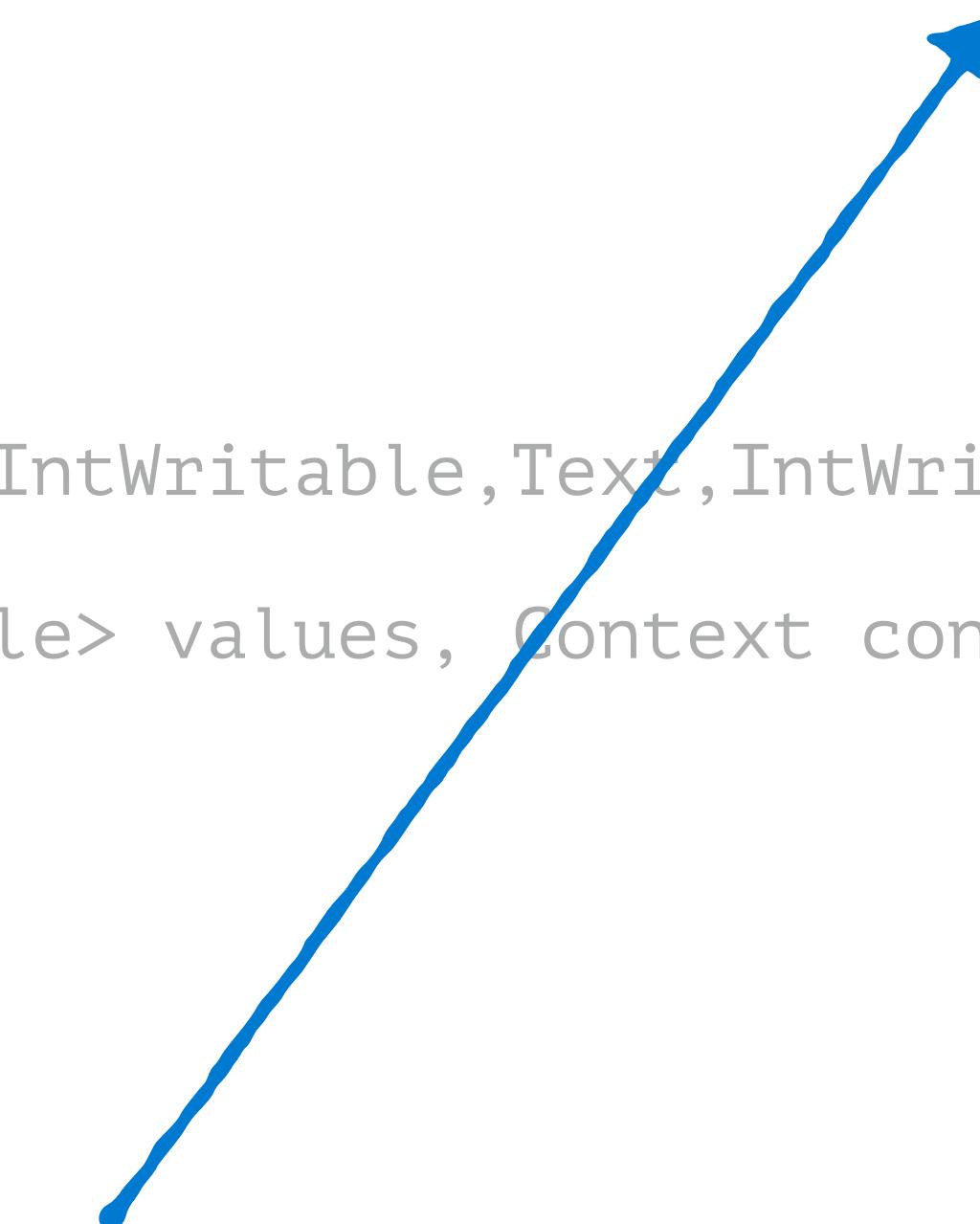
```
package com.company;

import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Reducer;

import java.io.IOException;

public class WordCountReducer extends Reducer<Text, IntWritable, Text, IntWritable> {
    @Override
    public void reduce(Text key, Iterable<IntWritable> values, Context context) throws IOException,
InterruptedException{
        int count = 0;
        for (IntWritable value:values)
        {
            count += value.get();
        }
        context.write(key, new IntWritable(count));
    }
}
```

In the end, we write out a pair
(word, count) to the Context object



When you **run** a job,
this JAR file **gets**
distributed to all the
nodes where the
computation is run

WordCount.JAR

✓ WordMapper Class

✓ WordReducer Class

WordCount Class

WordCount Class

```
package com.company;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;

public class WordCount {

    public static void main(String[] args) throws Exception {
        if(args.length !=2){
            System.err.println("Invalid Command");
            System.err.println("Usage: WordCount <input path> <output path>");
            System.exit(0);
        }
        Configuration conf = new Configuration();
        Job job = new Job(conf, "wordcount");
        job.setJarByClass(WordCount.class);

        FileInputFormat.addInputPath(job, new Path(args[0]));
        FileOutputFormat.setOutputPath(job, new Path(args[1]));
        job.setMapperClass(WordCountMapper.class);
        job.setReducerClass(WordCountReducer.class);

        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);

        System.exit(job.waitForCompletion(true)?0:1);
    }
}
```

We'll write a driver class
whose **main()** method will point
to our Mapper and Reducer

WordCount Class

```
package com.company;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
```

```
public class WordCount {
```

**public static void main(String[] args) throws
Exception {**

```
    if(args.length !=2){
        System.err.println("Invalid Command");
        System.err.println("Usage: WordCount <input path> <output path>");
        System.exit(0);
    }
    Configuration conf = new Configuration();
    Job job = new Job(conf, "wordcount");
    job.setJarByClass(WordCount.class);

    FileInputFormat.addInputPath(job, new Path(args[0]));
    FileOutputFormat.setOutputPath(job, new Path(args[1]));
    job.setMapperClass(WordCountMapper.class);
    job.setReducerClass(WordCountReducer.class);
    job.setOutputKeyClass(Text.class);
    job.setOutputValueClass(IntWritable.class);
    System.exit(job.waitForCompletion(true)?0:1);
}
```

The main() method will
take 2 String inputs

WordCount Class

```
package com.company;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
```

```
public class WordCount {
```

public static void main(String[] args) throws Exception {

```
    if(args.length !=2){
        System.err.println("Invalid Command");
        System.err.println("Usage: WordCount <input path> <output path>");
        System.exit(0);
    }
    Configuration conf = new Configuration();
    Job job = new Job(conf, "wordcount");
    job.setJarByClass(WordCount.class);

    FileInputFormat.addInputPath(job, new Path(args[0]));
    FileOutputFormat.setOutputPath(job, new Path(args[1]));
    job.setMapperClass(WordCountMapper.class);
    job.setReducerClass(WordCountReducer.class);
    job.setOutputKeyClass(Text.class);
    job.setOutputValueClass(IntWritable.class);
    System.exit(job.waitForCompletion(true)?0:1);
}
```

2 String inputs

Input text file path,
Output file path

WordCount Class

```
package com.company;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;

public class WordCount {

    public static void main(String[] args) throws Exception {
        if(args.length != 2){
            System.err.println("Invalid Command");
            System.err.println("Usage: WordCount <input path> <output path>");
            System.exit(0);
        }
        Configuration conf = new Configuration();
        Job job = new Job(conf, "wordcount");
        job.setJarByClass(WordCount.class);

        FileInputFormat.addInputPath(job, new Path(args[0]));
        FileOutputFormat.setOutputPath(job, new Path(args[1]));
        job.setMapperClass(WordCountMapper.class);
        job.setReducerClass(WordCountReducer.class);
        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);
        System.exit(job.waitForCompletion(true)?0:1);
    }
}
```

We check whether the right number of inputs are given, and if not print an error

WordCount Class

```
package com.company;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
```

```
public class WordCount {

    public static void main(String[] args) throws Exception {
        if(args.length !=2){
            System.err.println("Invalid Command");
            System.err.println("Usage: WordCount <input path> <output path>");
            System.exit(0);
    }
```

A configuration object

```
Configuration conf = new Configuration();
```

```
Job job = new Job(conf, "wordcount");

job.setJarByClass(WordCount.class);

FileInputFormat.addInputPath(job, new Path(args[0]));

FileOutputFormat.setOutputPath(job, new Path(args[1]));

job.setMapperClass(WordCountMapper.class);
job.setReducerClass(WordCountReducer.class);

job.setOutputKeyClass(Text.class);
job.setOutputValueClass(IntWritable.class);

System.exit(job.waitForCompletion(true)?0:1);
}
```

WordCount Class By default, this will pick up the configuration of your Hadoop deployment

```
package com.company;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;

public class WordCount {

    public static void main(String[] args) throws Exception {
        if(args.length !=2){
            System.err.println("Invalid Command");
            System.err.println("Usage: WordCount <input path> <output path>");
            System.exit(0);
        }

        Configuration conf = new Configuration();

        Job job = new Job(conf, "wordcount");
        job.setJarByClass(WordCount.class);

        FileInputFormat.addInputPath(job, new Path(args[0]));
        FileOutputFormat.setOutputPath(job, new Path(args[1]));

        job.setMapperClass(WordCountMapper.class);
        job.setReducerClass(WordCountReducer.class);

        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);

        System.exit(job.waitForCompletion(true)?0:1);
    }
}
```

WordCount Class

```
package com.company;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;

public class WordCount {
```

```
    public static void main(String[] args) throws Exception {
        if(args.length !=2){
            System.err.println("Invalid Command");
            System.err.println("Usage: WordCount <input path> <output path>");
            System.exit(0);
        }
    }
```

```
    Configuration conf = new Configuration();
```

Job job = new Job(conf, "wordcount");

```
    job.setJarByClass(WordCount.class);

    FileInputFormat.addInputPath(job, new Path(args[0]));
    FileOutputFormat.setOutputPath(job, new Path(args[1]));

    job.setMapperClass(WordCountMapper.class);
    job.setReducerClass(WordCountReducer.class);

    job.setOutputKeyClass(Text.class);
    job.setOutputValueClass(IntWritable.class);

    System.exit(job.waitForCompletion(true)?0:1);
}
```

We instantiate a new Job
whose name is “wordcount”

WordCount Class

```
package com.company;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;

public class WordCount {
```

```
    public static void main(String[] args) throws Exception {
        if(args.length !=2){
            System.err.println("Invalid Command");
            System.err.println("Usage: WordCount <input path> <output path>");
            System.exit(0);
        }
    }
```

```
    Configuration conf = new Configuration();
```

Job job = new Job(conf, "wordcount");

```
    job.setJarByClass(WordCount.class);

    FileInputFormat.addInputPath(job, new Path(args[0]));
    FileOutputFormat.setOutputPath(job, new Path(args[1]));

    job.setMapperClass(WordCountMapper.class);
    job.setReducerClass(WordCountReducer.class);

    job.setOutputKeyClass(Text.class);
    job.setOutputValueClass(IntWritable.class);

    System.exit(job.waitForCompletion(true)?0:1);
}
```

This object will hold all the specifications including the Mapper and Reducer to be used

WordCount Class

```
package com.company;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;

public class WordCount {

    public static void main(String[] args) throws Exception {
        if(args.length !=2){
            System.err.println("Invalid Command");
            System.err.println("Usage: WordCount <input path> <output path>");
            System.exit(0);
        }
        Configuration conf = new Configuration();
        Job job = new Job(conf, "wordcount");

        job.setJobName("Word Count");
        FileInputFormat.addInputPath(job, new Path(args[0]));
        FileOutputFormat.setOutputPath(job, new Path(args[1]));

        job.setMapperClass(WordCountMapper.class);
        job.setReducerClass(WordCountReducer.class);

        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);

        System.exit(job.waitForCompletion(true)?0:1);
    }
}
```

We specify the JAR file
that contains our code

job.setJarByClass(WordCount.class);

WordCount Class

```
package com.company;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
```

public class WordCount {

```
    public static void main(String[] args) throws Exception {
        if(args.length !=2){
            System.err.println("Invalid Command");
            System.err.println("Usage: WordCount <input path> <output path>");
            System.exit(0);
        }
        Configuration conf = new Configuration();
        Job job = new Job(conf, "wordcount");
```

job.setJarByClass(WordCount.class);

```
        job.setJobName("Word Count");
        FileInputFormat.addInputPath(job, new Path(args[0]));
        FileOutputFormat.setOutputPath(job, new Path(args[1]));
        job.setMapperClass(WordCountMapper.class);
        job.setReducerClass(WordCountReducer.class);
        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);
        System.exit(job.waitForCompletion(true)?0:1);
    }
```

The JAR is specified by the
Main class of our code

WordCount Class

```
package com.company;  
  
import org.apache.hadoop.conf.Configuration;  
import org.apache.hadoop.fs.Path;  
import org.apache.hadoop.io.IntWritable;  
import org.apache.hadoop.io.Text;  
import org.apache.hadoop.mapreduce.Job;  
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;  
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
```

```
public class WordCount {
```

```
    public static void main(String[] args) throws Exception {  
        if(args.length !=2){  
            System.err.println("Invalid Command");  
            System.err.println("Usage: WordCount <input path> <output path>");  
            System.exit(0);  
        }  
        Configuration conf = new Configuration();  
        Job job = new Job(conf, "wordcount");
```

```
        job.setJarByClass(WordCount.class);
```

```
        job.setJobName("Word Count");  
        FileInputFormat.addInputPath(job, new Path(args[0]));  
        FileOutputFormat.setOutputPath(job, new Path(args[1]));  
        job.setMapperClass(WordCountMapper.class);  
        job.setReducerClass(WordCountReducer.class);  
        job.setOutputKeyClass(Text.class);  
        job.setOutputValueClass(IntWritable.class);  
        System.exit(job.waitForCompletion(true)?0:1);  
    }
```

Hadoop will look for the
JAR file that contains this
class

WordCount Class

```
package com.company;  
  
import org.apache.hadoop.conf.Configuration;  
import org.apache.hadoop.fs.Path;  
import org.apache.hadoop.io.IntWritable;  
import org.apache.hadoop.io.Text;  
import org.apache.hadoop.mapreduce.Job;  
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;  
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
```

```
public class WordCount {
```

```
    public static void main(String[] args) throws Exception {  
        if(args.length !=2){  
            System.err.println("Invalid Command");  
            System.err.println("Usage: WordCount <input path> <output path>");  
            System.exit(0);  
        }  
        Configuration conf = new Configuration();  
        Job job = new Job(conf, "wordcount");
```

```
        job.setJarByClass(WordCount.class);
```

```
        job.setJobName("Word Count");  
        FileInputFormat.addInputPath(job, new Path(args[0]));  
        FileOutputFormat.setOutputPath(job, new Path(args[1]));  
        job.setMapperClass(WordCountMapper.class);  
        job.setReducerClass(WordCountReducer.class);  
        job.setOutputKeyClass(Text.class);  
        job.setOutputValueClass(IntWritable.class);  
        System.exit(job.waitForCompletion(true)?0:1);  
    }
```

It will then distribute this
JAR file to all the nodes
where our job will run

WordCount Class

```
package com.company;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
```

```
public class WordCount {

    public static void main(String[] args) throws Exception {
        if(args.length !=2){
            System.err.println("Invalid Command");
            System.err.println("Usage: WordCount <input path> <output path>");
            System.exit(0);
        }
        Configuration conf = new Configuration();
        Job job = new Job(conf, "wordcount");
        job.setJarByClass(WordCount.class);
```

`FileInputFormat.addInputPath(job, new Path(args[0]));`

`FileOutputFormat.setOutputPath(job, new Path(args[1]));`

```
        job.setMapperClass(WordCountMapper.class);
        job.setReducerClass(WordCountReducer.class);

        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);

        System.exit(job.waitForCompletion(true)?0:1);
    }
```

We set the **input file path** and **output file path** for the job using the arguments passed to `main()`

WordCount Class

```
package com.company;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
```

```
public class WordCount {

    public static void main(String[] args) throws Exception {
        if(args.length !=2){
            System.err.println("Invalid Command");
            System.err.println("Usage: WordCount <input path> <output path>");
            System.exit(0);
        }
        Configuration conf = new Configuration();
        Job job = new Job(conf, "wordcount");
        job.setJarByClass(WordCount.class);
```

`FileInputFormat.addInputPath(job, new Path(args[0]));`

`FileOutputFormat.setOutputPath(job, new Path(args[1]));`

```
        job.setMapperClass(WordCountMapper.class);
        job.setReducerClass(WordCountReducer.class);

        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);

        System.exit(job.waitForCompletion(true)?0:1);
    }
```

The **input path** can be a file or a directory or a file name pattern

WordCount Class

```
package com.company;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;

public class WordCount {

    public static void main(String[] args) throws Exception {
        if(args.length !=2){
            System.err.println("Invalid Command");
            System.err.println("Usage: WordCount <input path> <output path>");
            System.exit(0);
        }
        Configuration conf = new Configuration();
        Job job = new Job(conf, "wordcount");
        job.setJarByClass(WordCount.class);

        FileInputFormat.addInputPath(job, new Path(args[0]));

        FileOutputFormat.setOutputPath(job, new Path(args[1]));

        job.setMapperClass(WordCountMapper.class);
        job.setReducerClass(WordCountReducer.class);

        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);

        System.exit(job.waitForCompletion(true)?0:1);
    }
}
```

The output path must be a directory that does not exist

FileInputFormat.addInputPath(job, new Path(args[0]));

FileOutputFormat.setOutputPath(job, new Path(args[1]));

WordCount Class

```
package com.company;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;

public class WordCount {

    public static void main(String[] args) throws Exception {
        if(args.length !=2){
            System.err.println("Invalid Command");
            System.err.println("Usage: WordCount <input path> <output path>");
            System.exit(0);
        }
        Configuration conf = new Configuration();

        Job job = new Job(conf, "wordcount");

        job.setJarByClass(WordCount.class);
        job.setJobName("Word Count");
        FileInputFormat.addInputPath(job, new Path(args[0]));
        FileOutputFormat.setOutputPath(job, new Path(args[1]));

        job.setMapperClass(WordCountMapper.class);
        job.setReducerClass(WordCountReducer.class);

        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);

        System.exit(job.waitForCompletion(true)?0:1);
    }
}
```

We set the Mapper class
and Reducer class to the
classes we've written

WordCount Class

```
package com.company;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;

public class WordCount {

    public static void main(String[] args) throws Exception {
        if(args.length !=2){
            System.err.println("Invalid Command");
            System.err.println("Usage: WordCount <input path> <output path>");
            System.exit(0);
        }
        Configuration conf = new Configuration();

        Job job = new Job(conf, "wordcount");

        job.setJarByClass(WordCount.class);
        job.setJobName("Word Count");

        FileInputFormat.addInputPath(job, new Path(args[0]));
        FileOutputFormat.setOutputPath(job, new Path(args[1]));

        job.setMapperClass(WordCountMapper.class);
        job.setReducerClass(WordCountReducer.class);

        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);

        System.exit(job.waitForCompletion(true)?0:1);
    }
}
```

Explicitly set the type of
the output keys and values
of the Reducer Class

WordCount Class

```
package com.company;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;

public class WordCount {

    public static void main(String[] args) throws Exception {
        if(args.length !=2){
            System.err.println("Invalid Command");
            System.err.println("Usage: WordCount <input path> <output path>");
            System.exit(0);
        }
        Configuration conf = new Configuration();
        Job job = new Job(conf, "wordcount");
        job.setJarByClass(WordCount.class);
        job.setJobName("Word Count");
        FileInputFormat.addInputPath(job, new Path(args[0]));
        FileOutputFormat.setOutputPath(job, new Path(args[1]));
        job.setMapperClass(WordCountMapper.class);
        job.setReducerClass(WordCountReducer.class);

        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);

        System.exit(job.waitForCompletion(true)?0:1);
    }
}
```

This has to match the types
we specified in our Reducer
Class

**public class WordCountReducer extends
Reducer<Text, IntWritable, Text, IntWritable>**

**job.setOutputKeyClass(Text.class);
job.setOutputValueClass(IntWritable.class);**

WordCount Class

```
package com.company;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;

public class WordCount {

    public static void main(String[] args) throws Exception {
        if(args.length !=2){
            System.err.println("Invalid Command");
            System.err.println("Usage: WordCount <input path> <output path>");
            System.exit(0);
        }
        Configuration conf = new Configuration();
        Job job = new Job(conf, "wordcount");
        job.setJarByClass(WordCount.class);
        job.setJobName("Word Count");
        FileInputFormat.addInputPath(job, new Path(args[0]));
        FileOutputFormat.setOutputPath(job, new Path(args[1]));
        job.setMapperClass(WordCountMapper.class);
        job.setReducerClass(WordCountReducer.class);

        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);

        System.exit(job.waitForCompletion(true)?0:1);
    }
}
```

In our case the Mapper class also has the same output type

**public class WordCountMapper extends
Mapper<LongWritable,Text,Text,IntWritable>**

**job.setOutputKeyClass(Text.class);
job.setOutputValueClass(IntWritable.class);**

WordCount Class

```
package com.company;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;

public class WordCount {

    public static void main(String[] args) throws Exception {
        if(args.length !=2){
            System.err.println("Invalid Command");
            System.err.println("Usage: WordCount <input path> <output path>");
            System.exit(0);
        }
        Configuration conf = new Configuration();
        Job job = new Job(conf, "wordcount");
        job.setJarByClass(WordCount.class);
        job.setJobName("Word Count");
        FileInputFormat.addInputPath(job, new Path(args[0]));
        FileOutputFormat.setOutputPath(job, new Path(args[1]));
        job.setMapperClass(WordCountMapper.class);
        job.setReducerClass(WordCountReducer.class);

        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);

        System.exit(job.waitForCompletion(true)?0:1);
    }
}
```

Otherwise, the Mapper output types
should be set using
setMapOutputKeyClass() and
setMapOutputValueClass()

**public class WordCountMapper extends
Mapper<LongWritable,Text,Text,IntWritable>**

**job.setOutputKeyClass(Text.class);
job.setOutputValueClass(IntWritable.class);**

WordCount Class

```
package com.company;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;

public class WordCount {

    public static void main(String[] args) throws Exception {
        if(args.length !=2){
            System.err.println("Invalid Command");
            System.err.println("Usage: WordCount <input path> <output path>");
            System.exit(0);
        }
        Configuration conf = new Configuration();

        Job job = new Job(conf, "wordcount");

        job.setJarByClass(WordCount.class);
        job.setJobName("Word Count");

        FileInputFormat.addInputPath(job, new Path(args[0]));
        FileOutputFormat.setOutputPath(job, new Path(args[1]));

        job.setMapperClass(WordCountMapper.class);
        job.setReducerClass(WordCountReducer.class);

        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);

        System.exit(job.waitForCompletion(true)?0:1);
    }
}
```

The input types are set by default for text files in the way we want, so we don't need to specify that

WordCount Class

```
package com.company;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;

public class WordCount {

    public static void main(String[] args) throws Exception {
        if(args.length !=2){
            System.err.println("Invalid Command");
            System.err.println("Usage: WordCount <input path> <output path>");
            System.exit(0);
        }
        Configuration conf = new Configuration();

        Job job = new Job(conf, "wordcount");

        job.setJarByClass(WordCount.class);
        job.setJobName("Word Count");

        FileInputFormat.addInputPath(job, new Path(args[0]));
        FileOutputFormat.setOutputPath(job, new Path(args[1]));

        job.setMapperClass(WordCountMapper.class);
        job.setReducerClass(WordCountReducer.class);

        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);

        System.exit(job.waitForCompletion(true)?0:1);
    }
}
```

waitForCompletion() method will submit the Job for execution and wait for it to complete

WordCount Class

```
package com.company;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;

public class WordCount {

    public static void main(String[] args) throws Exception {
        if(args.length !=2){
            System.err.println("Invalid Command");
            System.err.println("Usage: WordCount <input path> <output path>");
            System.exit(0);
        }
        Configuration conf = new Configuration();

        Job job = new Job(conf, "wordcount");

        job.setJarByClass(WordCount.class);

        job.setJobName("Word Count");

        FileInputFormat.addInputPath(job, new Path(args[0]));

        FileOutputFormat.setOutputPath(job, new Path(args[1]));

        job.setMapperClass(WordCountMapper.class);
        job.setReducerClass(WordCountReducer.class);

        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);

        System.exit(job.waitForCompletion(true)?0:1);
    }
}
```

The method will return a Boolean indicating whether the Job ran successfully

WordCount Class

```
package com.company;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;

public class WordCount {

    public static void main(String[] args) throws Exception {
        if(args.length !=2){
            System.err.println("Invalid Command");
            System.err.println("Usage: WordCount <input path> <output path>");
            System.exit(0);
        }
        Configuration conf = new Configuration();
        Job job = new Job(conf, "wordcount");
        job.setJarByClass(WordCount.class);
        job.setJobName("Word Count");
        FileInputFormat.addInputPath(job, new Path(args[0]));
        FileOutputFormat.setOutputPath(job, new Path(args[1]));
        job.setMapperClass(WordCountMapper.class);
        job.setReducerClass(WordCountReducer.class);
        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);

        System.exit(job.waitForCompletion(true)?0:1);
    }
}
```

We translate that
Boolean to the
System exit code
true : 0 (success)
false : 1 (failure)

WordCount Class

```
package com.company;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;

public class WordCount {

    public static void main(String[] args) throws Exception {
        if(args.length !=2){
            System.err.println("Invalid Command");
            System.err.println("Usage: WordCount <input path> <output path>");
            System.exit(0);
        }
        Configuration conf = new Configuration();
        Job job = new Job(conf, "wordcount");
        job.setJarByClass(WordCount.class);
        job.setJobName("Word Count");
        FileInputFormat.addInputPath(job, new Path(args[0]));
        FileOutputFormat.setOutputPath(job, new Path(args[1]));
        job.setMapperClass(WordCountMapper.class);
        job.setReducerClass(WordCountReducer.class);
        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);

        System.exit(job.waitForCompletion(true)?0:1);
    }
}
```

true here means that
the job progress will
be printed to screen

true

When you **run** a job,
this JAR file **gets**
distributed to all the
nodes where the
computation is run

WordCount.JAR

WordMapper Class

WordReducer Class

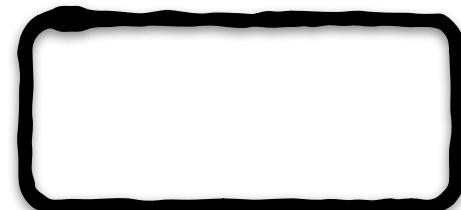
WordCount Class

Running a
MapReduce job

Step 1: Build a JAR file that contains all the code we just wrote

Step 1a: In the JAR, include all the Hadoop JARS as dependencies

Step 2: Run the job at the command line



Step 2: Run the job at the command line

```
$ hadoop jar <JARpath> <mainclass> <inputFile> <outputDir>
```

The path to the
JAR file

Step 2: Run the job at the command line

```
$ hadoop jar <JARpath> <mainclass> <inputFile> <outputDir>
```

The complete path (including package name) of the Main Class

Ex: com.company.WordCount

Step 2: Run the job at the command line

```
$ hadoop jar <JARpath> <mainclass> <inputFile> <outputDir>
```

The inputFile or a directory
containing multiple input files

Step 2: Run the job at the command line

```
$ hadoop jar <JARpath> <mainclass> <inputFile> <outputDir>
```

A directory that does not
exist