

MR - A MapReduce Programming Language

W4115 Programming Language & Translator
Fall 2011

Siyang Dai
Zhi Zhang
Zeyang Yu
Jinxiong Tan
Shuai Yuan

Motivation

```
class HelloWorld {  
    public static void main ( String argv[] ) {  
        System.out.println("Hello World!");  
    }  
}
```

All we need actually is one line:

```
System.out.println("Hello World!");
```

Motivation

Word Count

Hello world for Mapreduce

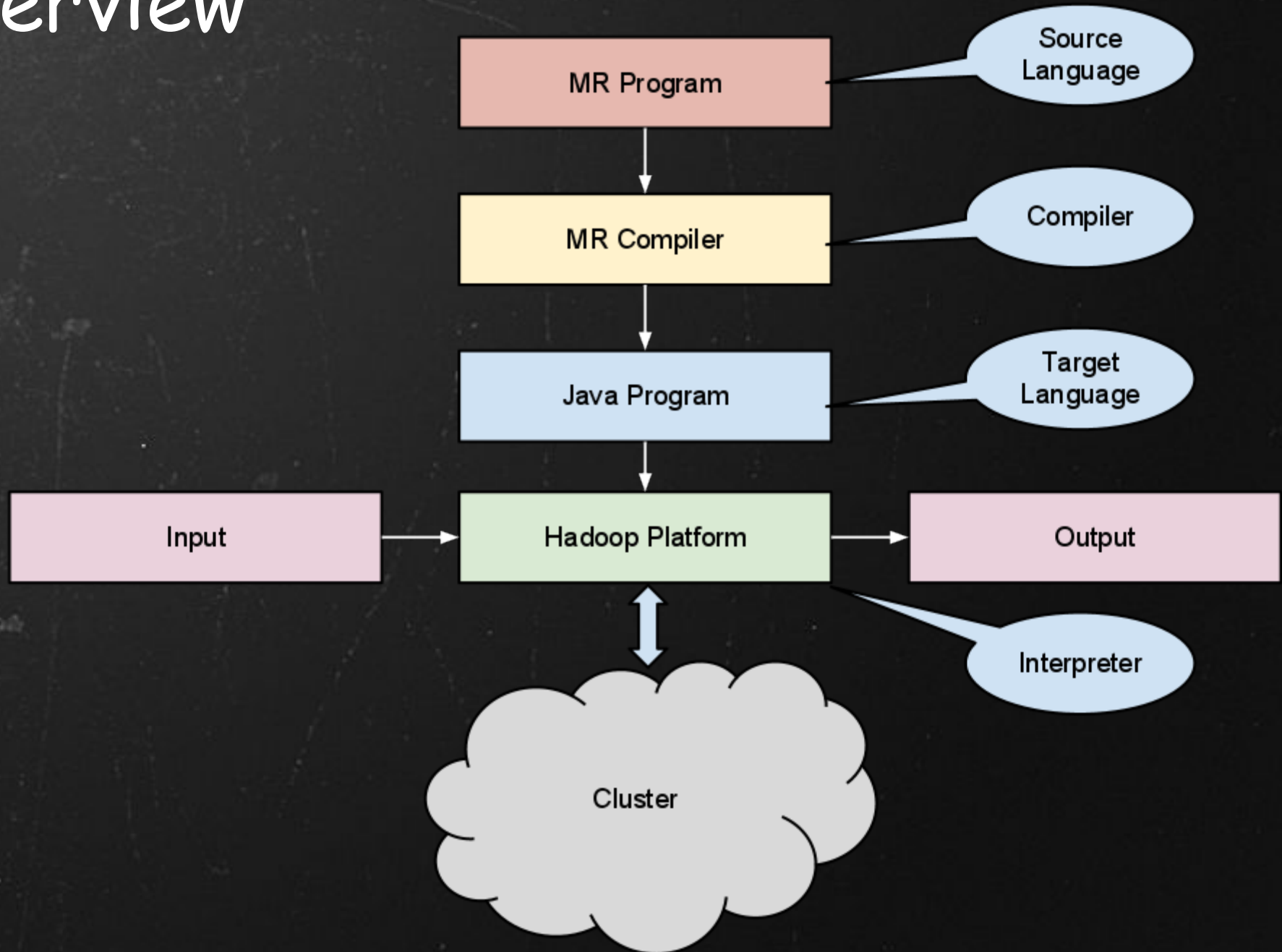
```
1 import java.io.IOException;
2 import java.lang.InterruptedException;
3 import java.util.StringTokenizer;
4
5 import org.apache.hadoop.io.*;
6 import org.apache.hadoop.conf.Configuration;
7 import org.apache.hadoop.mapreduce.*;
8 import org.apache.hadoop.fs.Path;
9 import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
10 import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
11 import org.apache.hadoop.util.GenericOptionsParser;
12
13 public class WordCount {
14     /**
15      * The map class of WordCount.
16      */
17     public static class TokenCounterMapper extends
18         Mapper<Object, Text, Text, IntWritable> {
19
20         private final static IntWritable one = new IntWritable(1);
21         private Text word = new Text();
22
23         public void map(Object key, Text value, Context context)
24             throws IOException, InterruptedException {
25             StringTokenizer itr = new StringTokenizer(value.toString());
26             while (itr.hasMoreTokens()) {
27                 word.set(itr.nextToken());
28                 context.write(word, one);
29             }
30         }
31     }
32
33     /**
34      * The reducer class of WordCount
35      */
36     public static class TokenCounterReducer extends
37         Reducer<Text, IntWritable, Text, IntWritable> {
38         public void reduce(Text key, Iterable<IntWritable> values,
39             Context context) throws IOException, InterruptedException {
40             int sum = 0;
41             for (IntWritable value : values) {
42                 sum += value.get();
43             }
44             context.write(key, new IntWritable(sum));
45         }
46     }
47
48     /**
49      * The main entry point.
50      */
51     public static void main(String[] args) throws Exception {
52         Configuration conf = new Configuration();
53         String[] otherArgs = new GenericOptionsParser(conf, args)
54             .getRemainingArgs();
55         Job job = new Job(conf, "Example Hadoop 0.20.1 WordCount");
56         job.setJarByClass(WordCount.class);
57         job.setMapperClass(TokenCounterMapper.class);
58         job.setReducerClass(TokenCounterReducer.class);
59         job.setOutputKeyClass(Text.class);
60         job.setOutputValueClass(IntWritable.class);
61         FileInputFormat.addInputPath(job, new Path(otherArgs[0]));
62         FileOutputFormat.setOutputPath(job, new Path(otherArgs[1]));
63         System.exit(job.waitForCompletion(true) ? 0 : 1);
64     }
65 }
```

Motivation

Why not like this:

```
1 #JobName = "WordCount"
2
3 def wordcount_map(offset, line) <(Int, Text) -> (Text , Int)> : Mapper
4 {
5     List<Text> words = split line by " ";
6     foreach word in words
7         emit(word, 1);
8 }
9
10 def wordcount_reduce (word, counts) <(Text , Int) -> (Text, Int)> : Reducer
11 {
12     Int total = 0;
13     foreach count in counts
14         total = total + count;
15
16     emit(word, total);
17 }
```


Overview



Tutorial: wordcount.mr

```
#JobName = "WordCount"
//map function definition
def wordcount_map <(Int, Text) -> (Text , Int)> (offset, line): Mapper
{
  List<Text> words;
    words = split line by " ";
  foreach Text word in words {
    emit(word, 1);
  }
}
//reduce function definition
def wordcount_reduce <(Text , Int) -> (Text, Int)> (word, counts): Reducer
{
  Int total;
  total = 0;
  foreach Int count in counts {
    total = total + count;
  }
  emit(word, total);
}
```

Map function

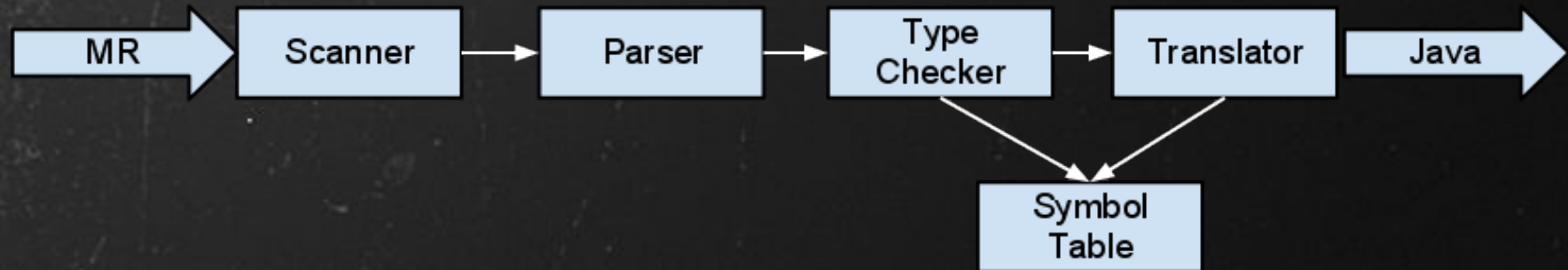
```
//map function definition
def wordcount_map <(Int, Text) -> (Text , Int)> (offset, line): Mapper
{
  List<Text> words;
  words = split line by " ";
  foreach Text word in words {
    emit(word, 1);
  }
}
```

Reduce function

//reduce function definition

```
def wordcount_reduce <(Text , Int) -> (Text, Int)> (word, counts): Reducer
{
  Int total;
  total = 0;
  foreach Int count in counts {
    total = total + count;
  }
  emit(word, total);
}
```


Compiler Structure



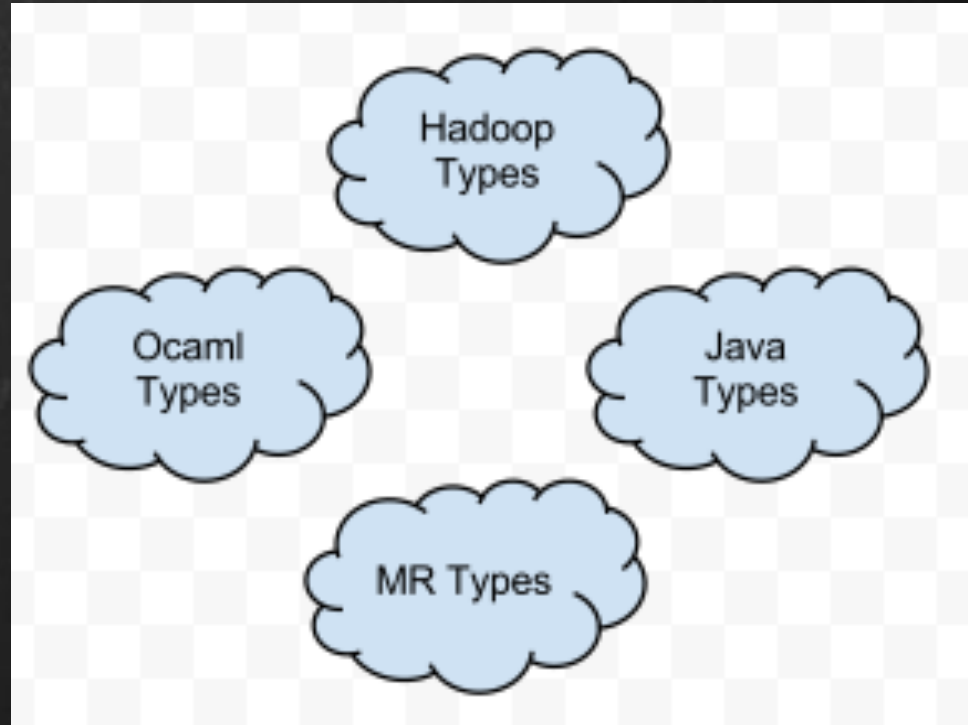
Translator

```
Mapper / Reducer (Parameters in Hadoop types) {  
  mapper / reducer logic  
  emit the values in Hadoop types  
}
```

Problem:

Hadoop types do not support basic operations.

Types Puzzle



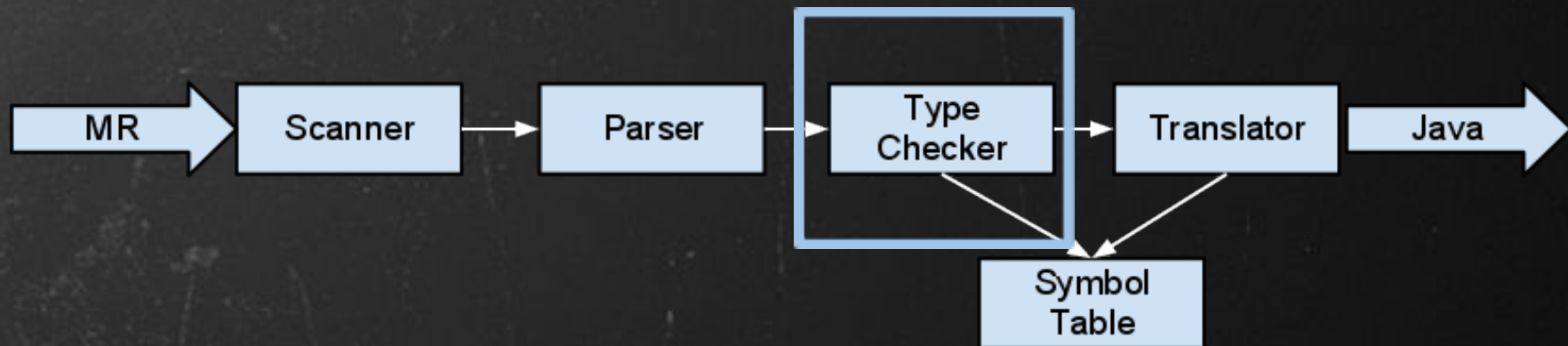
Translator

```
Mapper / Reducer (Parameters in Hadoop types) {  
    Convert from Hadoop types to Java types  
    mapper / reducer logic  
    Convert from Java types to Hadoop types  
    emit the values in Hadoop types  
}
```

Solution:

Another layer of indirection

Type Checker



Prevent Ridiculous Expression/Statement

e.g.

```
Int a = "oops"; //illegal assignment
```

```
if("oops"){ //it should be boolean type
```

```
...
```

```
}
```


Type Checker

Variable Definition

Assignment Operation

Binary Operation (Arithmetic, Relational, Logical)

Split Expression

If-Expression

Foreach Statement

Emit Statement

Lessons Learned

1. Parametric Polymorphism
2. Think Recursively
3. Universal Solution - Indirection
4. Teamwork

Demo