# 50002 - Software Engineering Design - Lecture 10

Oliver Killane

14/11/21

# MapReduce

A system built by Google to work on very large datasets, allowing data processing tasks to be spread over many computers in a farm.
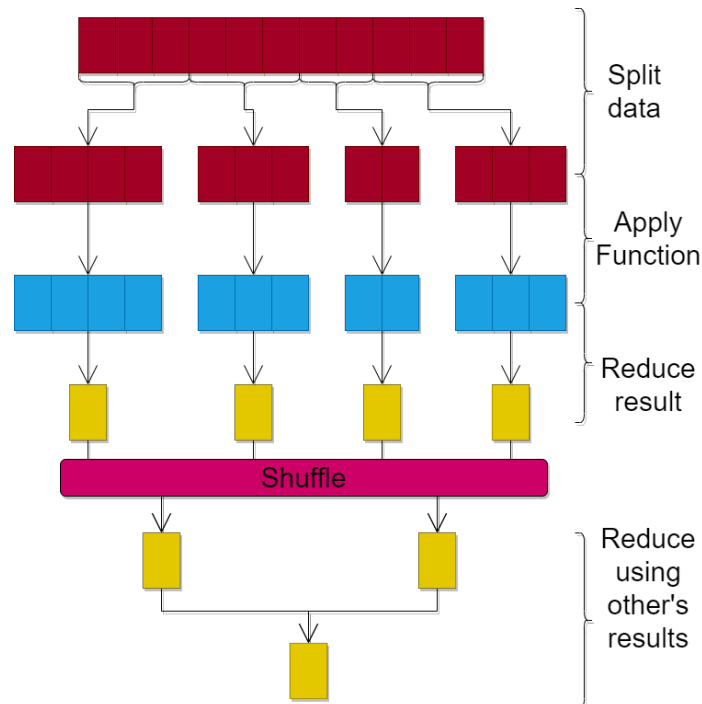
The main difficulties win distributed computing considered were:

- If many projects are changed to be distributed, then lots of code must be changed, much will be duplicated & there is potential for bugs.
- Ordinging interactions and sharing data across many networked computers is difficult.
- Probability of a hardware failure grows with the amount of hardware.

MapReduce takes tasks written by programmers, and distributises it across a number of systems, handling failures & can be easily used through an **API**. Furthermore any improvements to the system will immediately be reflected in all projects using it.

Google wanted to index the world wide web, this is a large problem that is best suited to highly distributed programming.

```
1   // Calculate the sum of the squares of numbers in a list.
2
3   var square = function (x) { return x * x }
4   var sum = function (x,y) { return x + y }
5
6   [1,2,3,4,5].map(square).reduce(sum);
```

The shuffle stage is important, to increase performance we want to transmit/copy the minimum amount of data between computers.

For **MapReduce** the signatures are different to lists & integers:

$$map : (k_1, v_1) \rightarrow list(k_2, v_2)$$

$$reduce : list(k_2, list(v_2)) \rightarrow list(k_2, v_2)$$

For example: We have $k_1 = $ line number and $v_1 = $ line of text

| Key | Value |
|---|---|
| 1 | I think that it's extraordinarily important that we in computer science keep fun in computing. |
| 2 | When it started out, it was an awful lot of fun. Of course, the paying customers got shafted every |
| 3 | now and then, and after a while we began to take their complaints seriously. We began to feel as |
| 4 | if we really were responsible for the successful, error-free perfect use of these machines. I don't |
| 5 | think we are. I think we're responsible for stretching them, setting them off in new directions, and |
| 6 | keeping fun in the house. I hope the field of computer science never loses its sense of fun |

$map(k_1, v_1) \rightarrow [(k_2, v_2)]$ where $k_2 = $ word and $v_2 = 1$.

```
1  [("I",1),("think",1),("that",1),("it's",1),("extraordinarily",1),("important",1),("
   ↪  that",1),("we",1),("in",1),("computer",1),("science",1),("keep",1),("fun",1),(
   ↪  "in",1),("computing.",1)]
```

The we reduce with + for each key:
$reduce[(k_2, v_2)] \rightarrow [(k_2, [v_2])]$:

```
1  [('I', 1),( 'think', 1),( 'that', 2),( "it's", 1),( 'extraordinarily', 1),( '
   ↪ important', 1),( 'we', 1),( 'in', 2),( 'computer', 1),( 'science', 1),( 'keep'
   ↪ , 1),( 'fun', 1),( 'computing.', 1)]
```

### Alternatives

Other projects have been developed to perform a similar role to **MapReduce** such as **Hadoop**.

## Java Example

```java
1   public class Map extends Mapper<LongWritable, Text, Text, IntWritable> {
2       private final static IntWritable one = new IntWritable(1);
3       private Text word = new Text();
4
5       public void map(LongWritable key, Text value, Context context)
6       throws IOException, InterruptedException {
7
8         String line = value.toString();
9         StringTokenizer tokenizer = new StringTokenizer(line);
10        while (tokenizer.hasMoreTokens()) {
11           word.set(tokenizer.nextToken());
12           context.write(word, one);
13        }
14      }
15  }
16
17
18  public class Reduce extends Reducer<Text, IntWritable, Text, IntWritable> {
19
20      public void reduce(Text key, Iterable<IntWritable> values, Context context)
21                                       throws IOException, InterruptedException {
22
23          int sum = 0;
24          for (IntWritable val : values) {
25              sum += val.get();
26          }
27          context.write(key, new IntWritable(sum));
28      }
29  }
30
31
32  public class WordCount {
33      public static void main(String[] args) throws Exception {
34          Configuration conf = new Configuration();
35
36          Job job = new Job(conf, "wordcount");
37
38          job.setOutputKeyClass(Text.class);
39          job.setOutputValueClass(IntWritable.class);
40
41          job.setMapperClass(Map.class);
42          job.setReducerClass(Reduce.class);
43
44          job.setInputFormatClass(TextInputFormat.class);
45          job.setOutputFormatClass(TextOutputFormat.class);
46
```

```
47          FileInputFormat.addInputPath(job, new Path("some path");
48          FileOutputFormat.setOutputPath(job, new Path("some path");
49
50          job.waitForCompletion(true);
51        }
52    }
```