

## The Streaming API: Motivation

---

- **Many organizations have developers skilled in languages other than Java, such as**
  - Ruby
  - Python
  - Perl
- **The Streaming API allows developers to use any language they wish to write Mappers and Reducers**
  - As long as the language can read from standard input and write to standard output

# The Streaming API: Advantages and Disadvantages

---

- **Advantages of the Streaming API:**

- No need for non-Java coders to learn Java
- Fast development time
- Ability to use existing code libraries

- **Disadvantages of the Streaming API:**

- Performance
- Primarily suited for handling data that can be represented as text
- Streaming jobs can use excessive amounts of RAM or fork excessive numbers of processes
- Although Mappers and Reducers can be written using the Streaming API, Partitioners, InputFormats etc. must still be written in Java

## How Streaming Works

---

- **To implement streaming, write separate Mapper and Reducer programs in the language(s) of your choice**
  - They will receive input via `stdin`
  - They should write their output to `stdout`
- **If `TextInputFormat` (the default) is used, the streaming Mapper just receives each line from the file on `stdin`**
  - No key is passed
- **Mapper and Reducer output should be sent to `stdout` as**
  - `key [tab] value [newline]`
- **Separators other than tab can be specified**



## Streaming: Example Mapper

---

- Example streaming wordcount Mapper:

```
#!/usr/bin/env perl
while (<>) {
    chomp;
    (@words) = split /\W+/;
    foreach $w (@words) {
        print "$w\t1\n";
    }
}
```

# Read lines from stdin  
# Get rid of the trailing newline  
# Create an array of words  
# Loop through the array  
# Print out the key and value

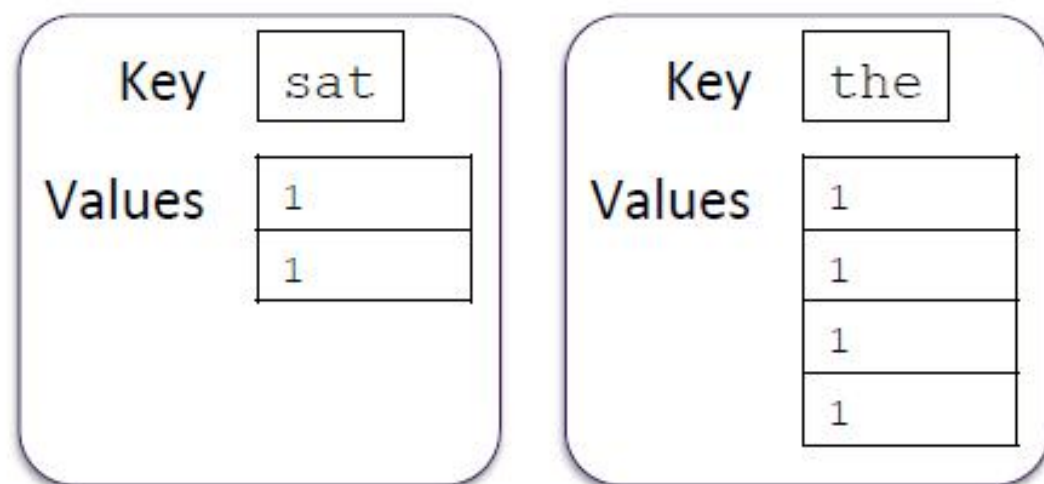
## Python MapReduce Code:

```
1 mapper.py
2 #!/usr/bin/python
3 import sys
4 #Word Count Example
5 # input comes from standard input STDIN
6 for line in sys.stdin:
7     line = line.strip() #remove leading and trailing whitespaces
8     words = line.split() #split the line into words and returns as a list
9     for word in words:
10         #write the results to standard output STDOUT
11         print '%s\t%s' % (word,1) #Emit the word
```

## Streaming Reducers: Caution

- Recall that in Java, all the values associated with a key are passed to the Reducer as an `Iterable`
- Using Hadoop Streaming, the Reducer receives its input as one key/value pair per line
- Your code will have to keep track of the key so that it can detect when values from a new key start

### Java: Iterable Objects



### Streaming: stdin

```
sat 1
sat 1
the 1
the 1
the 1
the 1
```

## Streaming: Example Reducer

---

- **Example streaming wordcount Reducer:**

```
#!/usr/bin/env perl
$sum = 0;
$last = "";
while (<>) {
    ($key, $value) = split /\t/;
    $last = $key if $last eq "";
    if ($last ne $key) {
        print "$last\t$sum\n";
        $last = $key;
        $sum = 0;
    }
    $sum += $value;
}
print "$key\t$sum\n";
```

# read lines from stdin  
# obtain the key and value  
# first time through  
# has key has changed?  
# if so output last key/value  
# start with the new key  
# reset sum for the new key  
  
# add value to tally sum for key  
  
# print the final pair

## reducer.py

```
1  #!/usr/bin/python
2  import sys
3  from operator import itemgetter
4  # using a dictionary to map words to their counts
5  current_word = None
6  current_count = 0
7  word = None
8  # input comes from STDIN
9  for line in sys.stdin:
10     line = line.strip()
11     word,count = line.split(' ',1)
12     try:
13         count = int(count)
14     except ValueError:
15         continue
16     if current_word == word:
17         current_count += count
18     else:
19         if current_word:
20             print '%s\t%s' % (current_word, current_count)
21         current_count = count
22         current_word = word
23     if current_word == word:
24         print '%s\t%s' % (current_word,current_count)
```

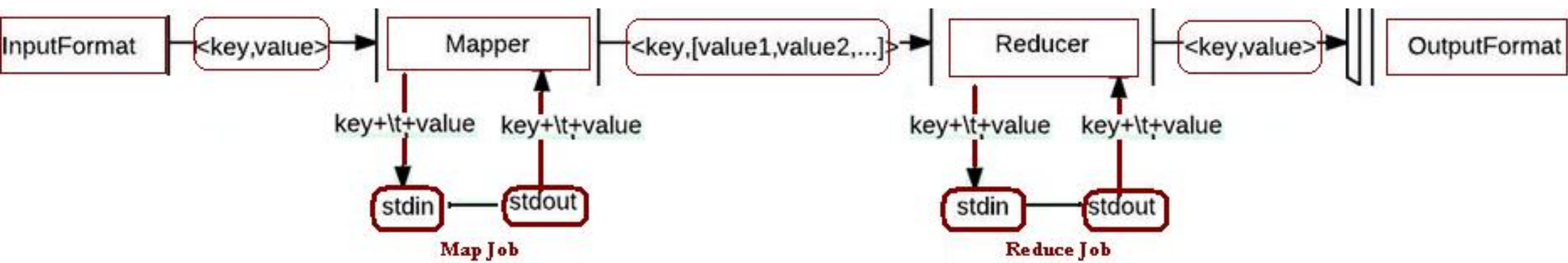


## Launching a Streaming Job

- **To launch a Streaming job, use e.g.:**

```
$ hadoop jar /usr/lib/hadoop-0.20-mapreduce/contrib/\
streaming/hadoop-streaming*.jar \  
-input myInputDirs \  
-output myOutputDir \  
-mapper myMapScript.pl \  
-reducer myReduceScript.pl \  
-file mycode/myMapScript.pl \  
-file mycode/myReduceScript.pl
```

- **Many other command-line options are available**
  - See the documentation for full details
- **Note that system commands can be used as a Streaming Mapper or Reducer**
  - For example: `awk`, `grep`, `sed`, or `wc`



```
$ bin/hadoop jar contrib/streaming-hadoop-0.18.0-streaming.jar \  
-mapper      streamMapProgram \  
-reducer     streamReduceProgram \  
-input       dfs/path \  
-output      dfs/path|
```