



1. The InputFormat generates $\langle k1, v1 \rangle$ pairs
2. Converts the $\langle k1, v1 \rangle$ pairs to lines and sends them to the stdin of the process
3. The stdout of the process is converted into $\langle k2, v2 \rangle$ pairs
4. Converts $\langle k2, (v2, v2, \dots) \rangle$ pairs to lines and send them to the stdin of the process
5. The stdout of the process is converted into $\langle k3, v3 \rangle$ pairs

Running a Hadoop Streaming Job

A Streaming job is a MapReduce job defined in the `hadoop-streaming.jar` file:

```
hadoop jar $HADOOP_HOME/lib/hadoop-streaming.jar  
  -input input_directories  
  -output output_directories  
  -mapper mapper_script  
  -reducer reducer_script
```