



An Introduction

Kai Voigt, Cloudera Inc
Berlin Buzzwords, June 6th 2011

cloudera

Big Data

- Capture
- Storage
- Search
- Analytics



hadoop.apache.org



Google Filesystem
(GFS)

Hadoop Distributed
Filesystem (HDFS)

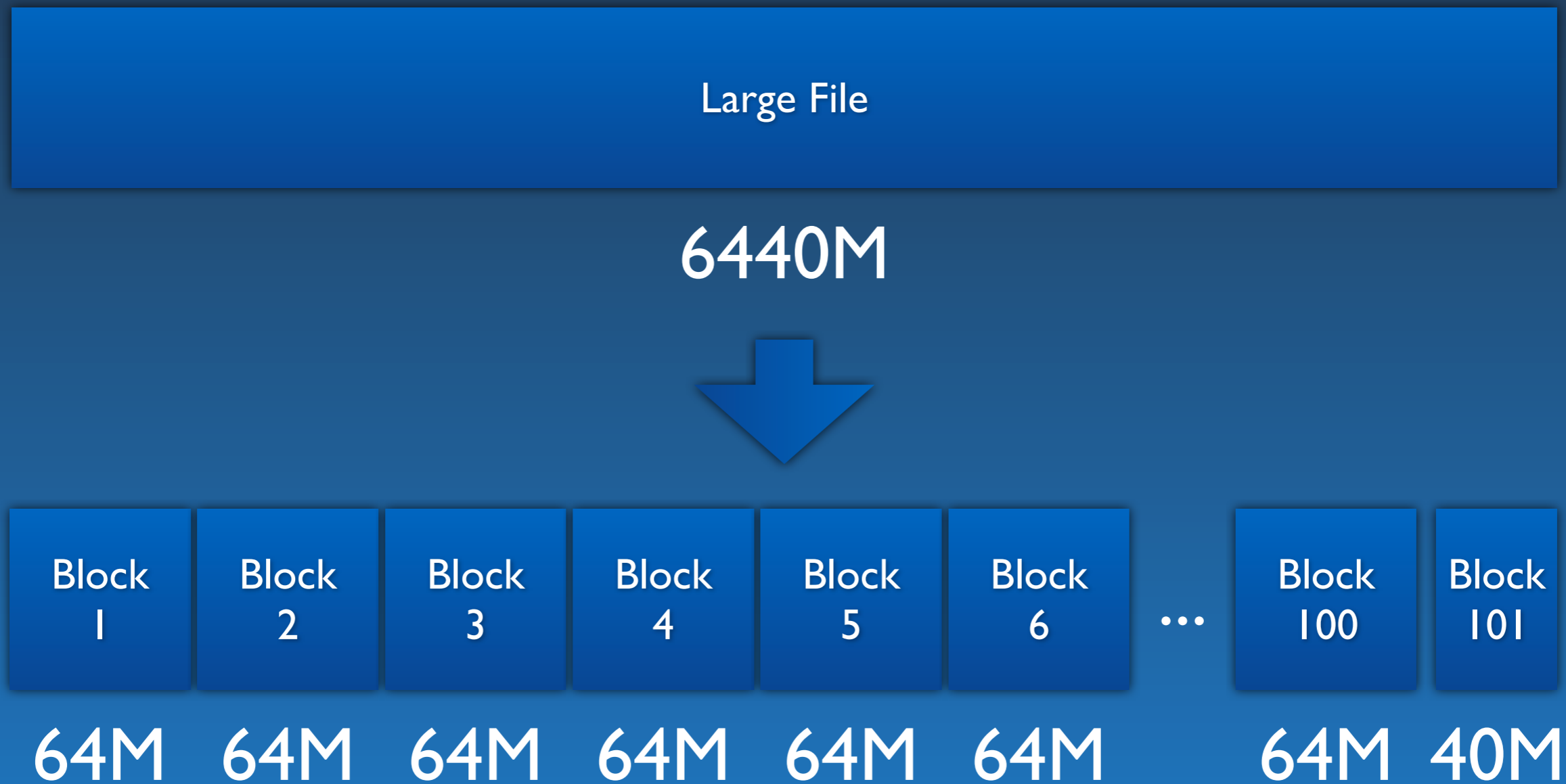
MapReduce

MapReduce

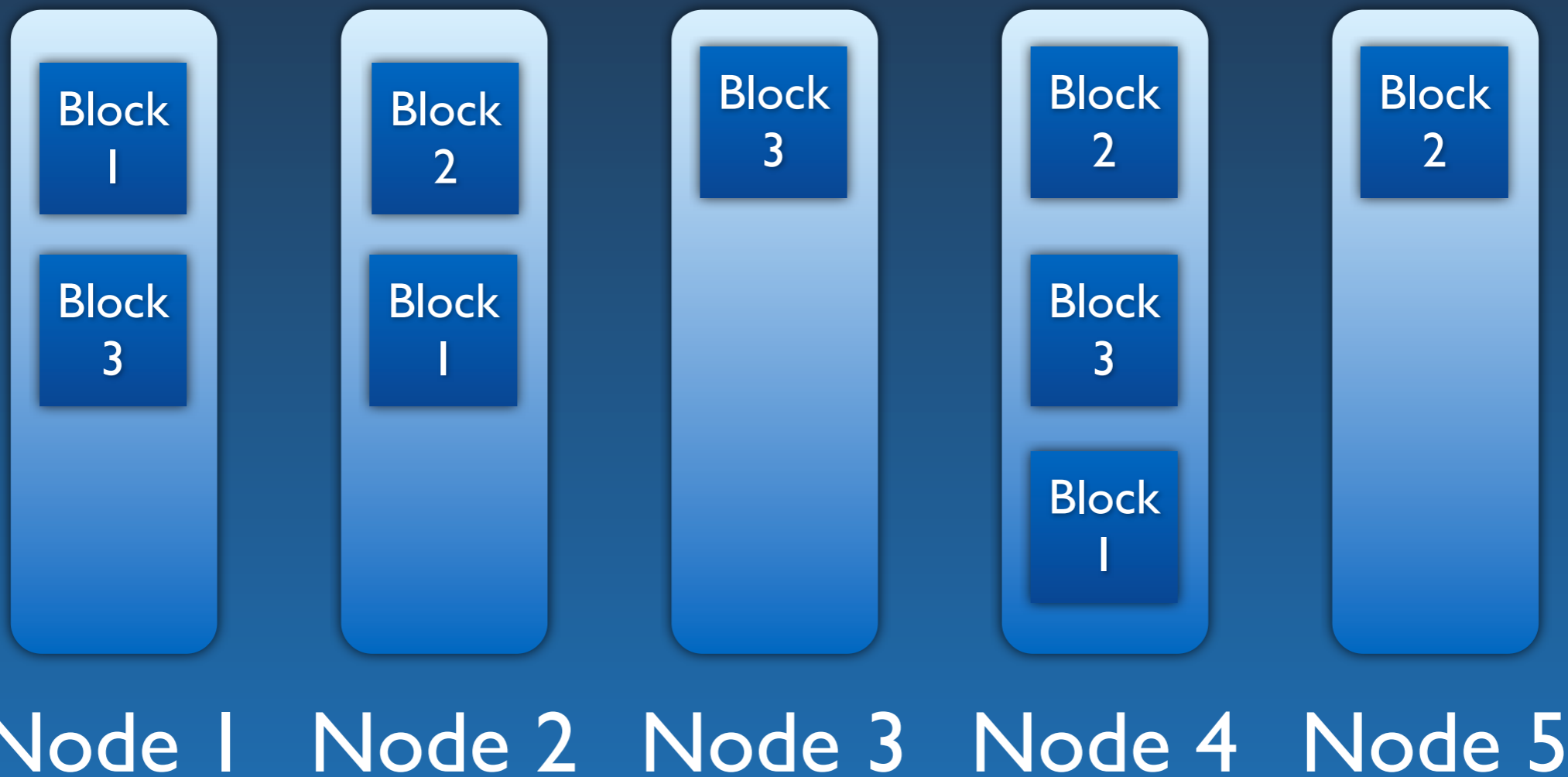
Hadoop Distributed Filesystem

- *Easy Access*
- Distributed
- Redundant
- Scalable

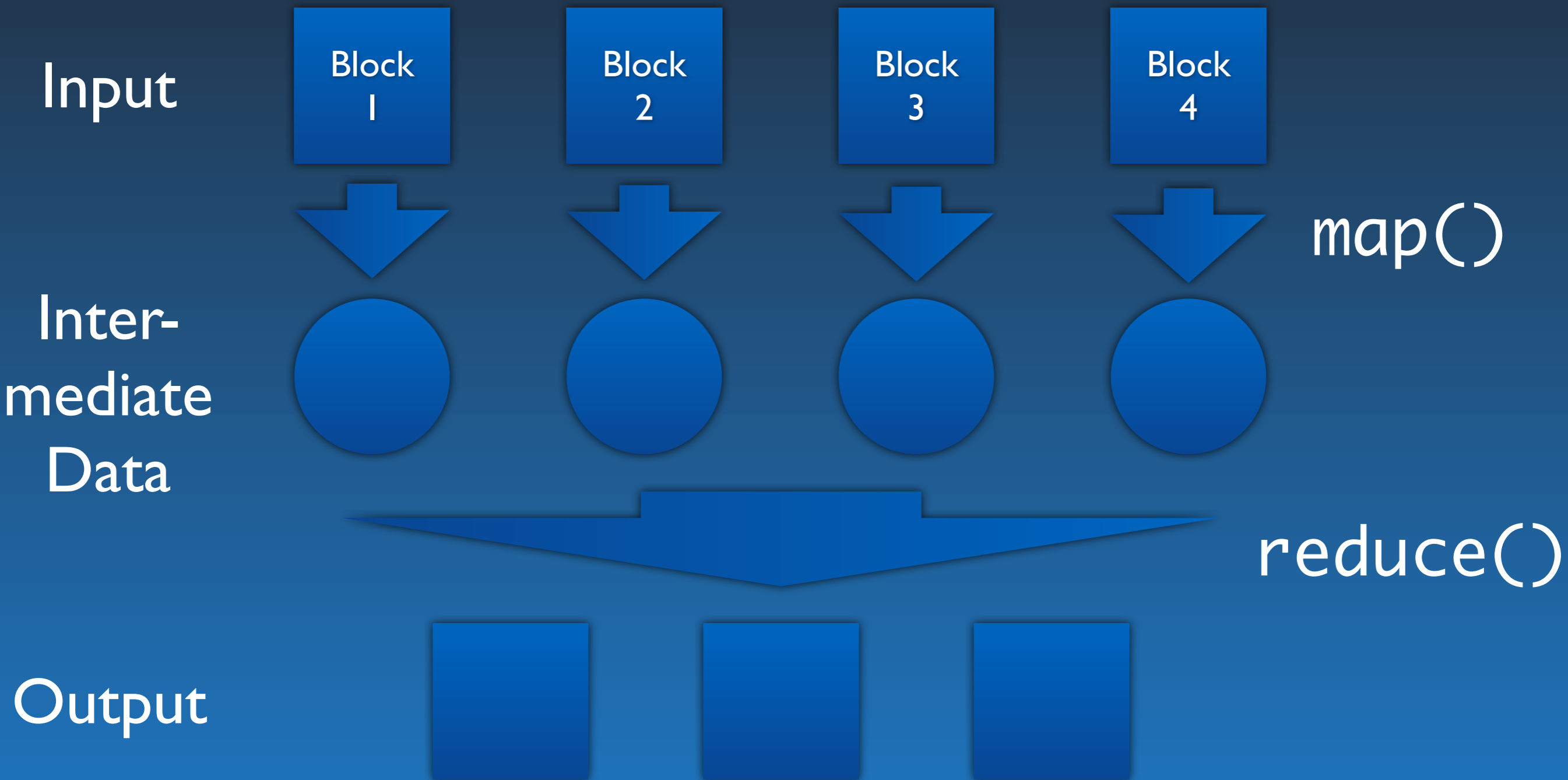
File Splits



Block Placement



MapReduce



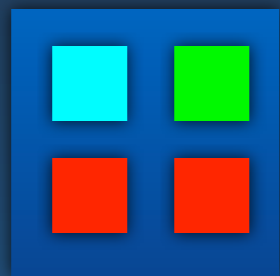
WordCount Example





```
map (offset, line) {  
  foreach word in line {  
    emit (word, 1)  
  }  
}
```


WordCount Example





```
reduce (word, count[]) {  
    total = 0;  
    foreach number in count[] {  
        total += number  
    }  
    emit (word, total)  
}
```

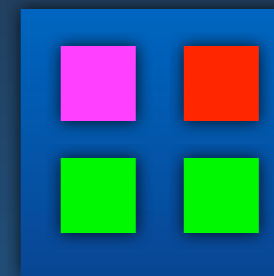
map()

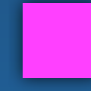
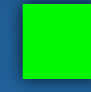




(, 1)
(, 1)
(, 1)
(, 1)

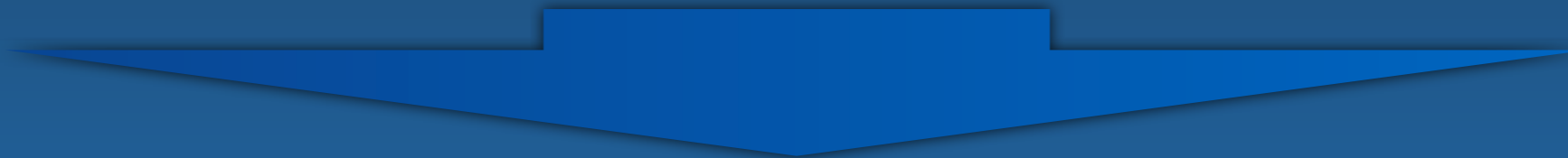


(, 1)
(, 1)
(, 1)
(, 1)


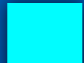





(, 1)
(, 1)
(, 1)
(, 1)

Sort & Shuffle






reduce()

(, (1,1))
(, (1,1))
(, (1))

(, (1,1,1,1))
(, (1,1))



(, 2)
(, 2)
(, 1)

(, 4)
(, 2)

Use Case: Recommendations

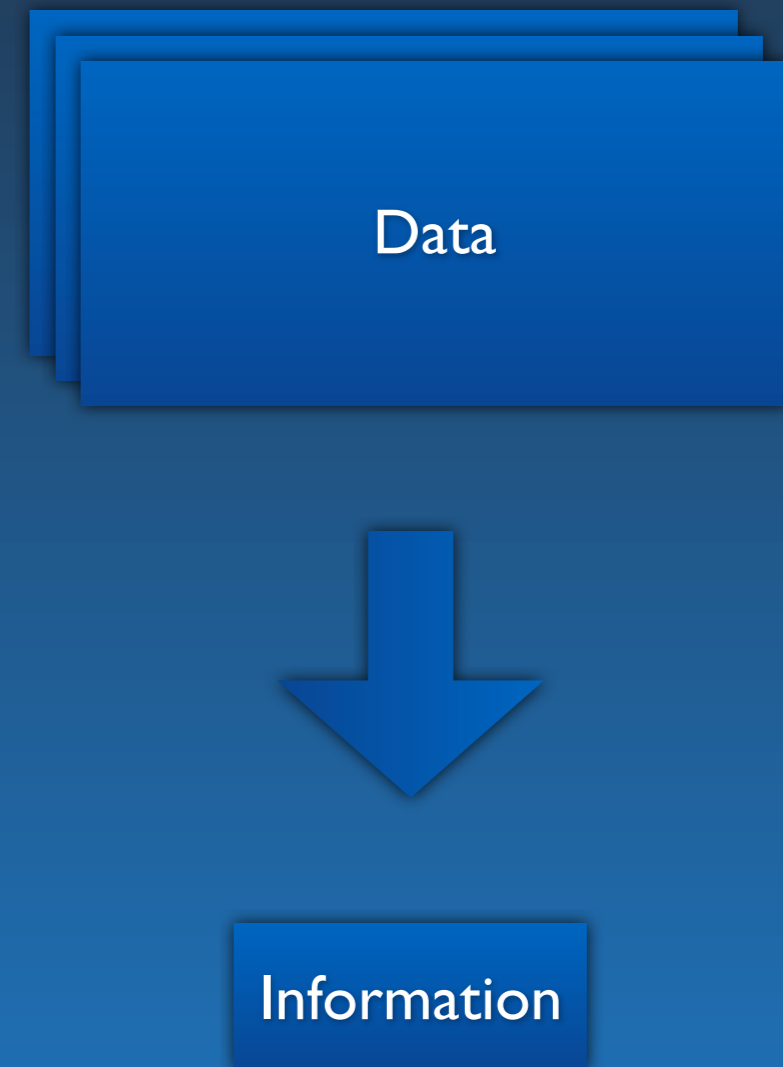
- "People looking at this article also looked at these articles"
- "You might also know these people"
- iTunes Genius Playlist
- Banner Placement

Use Case: Text Processing

- Document Indexing
- Semantic Analytics

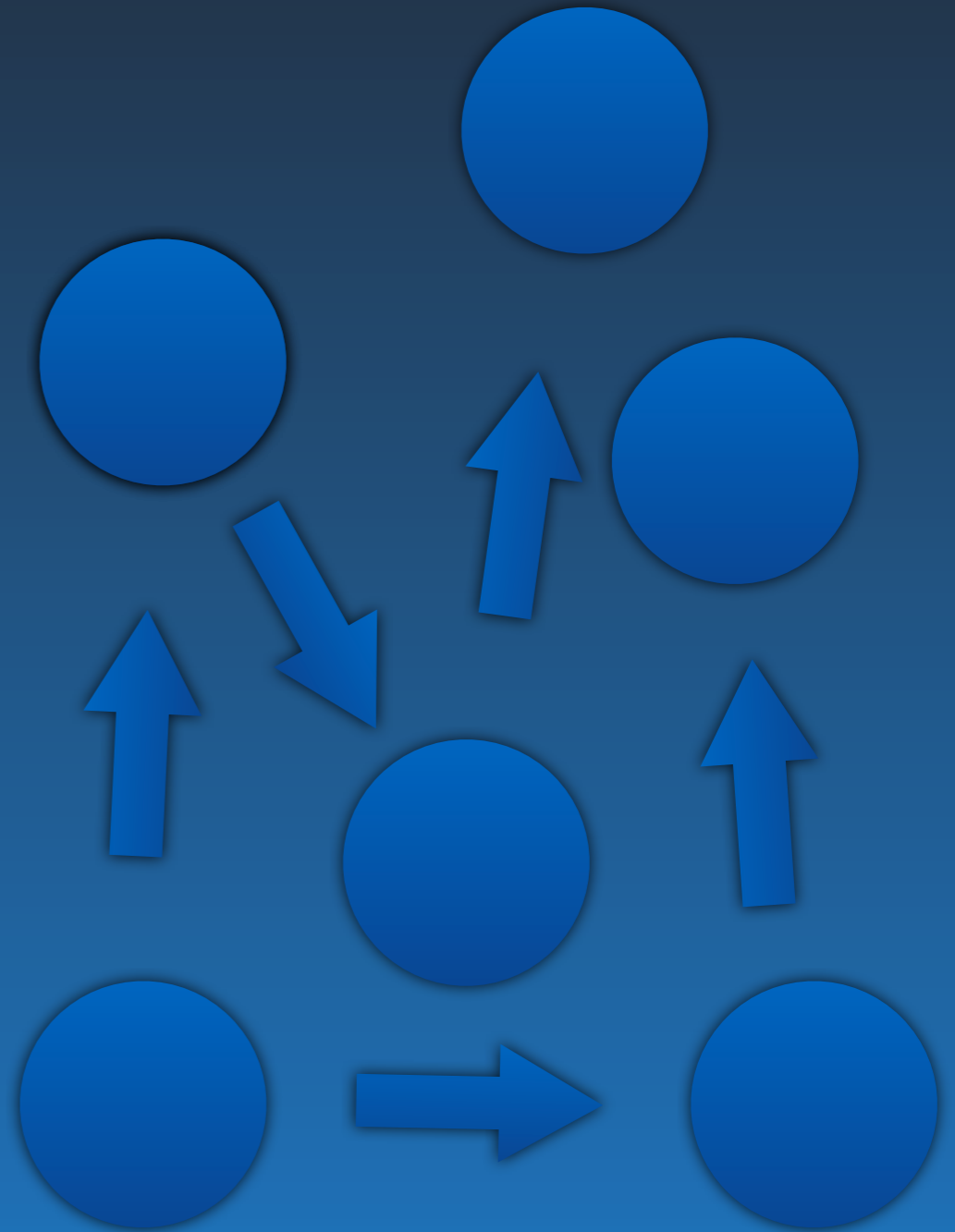
Use Case: Machine Learning

- Spam vs No Spam
- Credit Card Fraud
- "People of Interest"



Use Case: Graphs

- Shortest Paths
- Bottleneck Nodes
- Flow Optimization
- Spanning Trees
- Spanning Routes



Hadoop Ecosystem

Hive & Pig	High Level Language Access
HBase	Real Time Access
Sqoop	SQL to/from Hadoop
Flume	Distributed Data Collection
Oozie	Job Workflow
Mahout	Machine Learning Library

many more

Homework

- Cloudera's Distribution including Hadoop (CDH3)
- Online Tutorials
- WordCount Example
- Conference Demo Cluster

Quick Demo!

Thank you!

- Kai Voigt
- kai@cloudera.com
- <http://www.cloudera.com/>
- <http://apache.hadoop.org/>