



# Why Spark Is the Next Top (Compute) Model

Philly ETE 2014  
April 22-23, 2014  
@deanwampler  
[polyglotprogramming.com/talks](http://polyglotprogramming.com/talks)

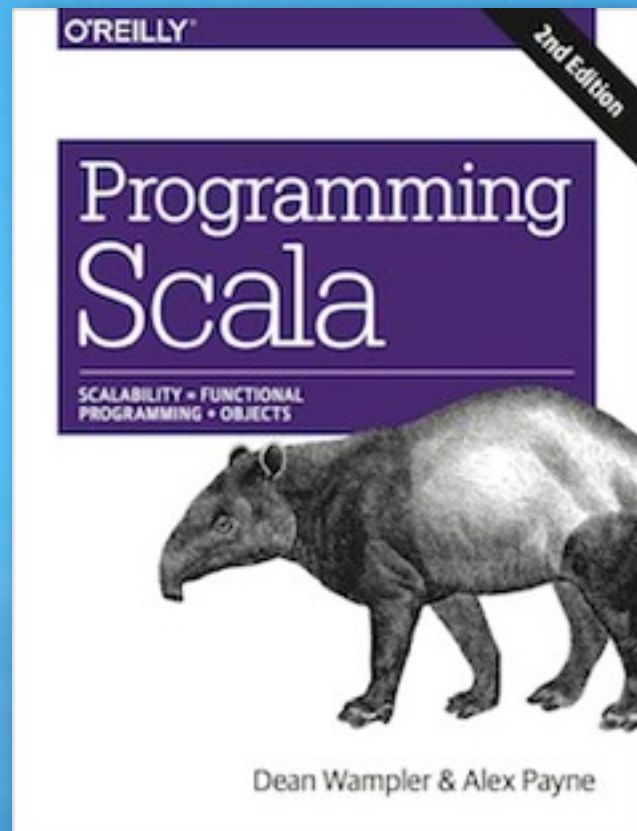
Saturday, April 19, 14

Copyright (c) 2014, Dean Wampler, All Rights Reserved, unless otherwise noted.

Image: Detail of the London Eye



# Dean Wampler



dean.wampler@typesafe.com  
polyglotprogramming.com/talks  
@deanwampler

Saturday, April 19, 14

About me. You can find this presentation and others on Big Data and Scala at [polyglotprogramming.com](http://polyglotprogramming.com).

Programming Scala, 2nd Edition is forthcoming.

photo: Dusk at 30,000 ft above the Central Plains of the U.S. on a Winter's Day.



Or  
this?

THE  
Compleat Troller,  
OR,  
THE ART  
OF  
TROLLING.  
WITH  
A Description of all the Utensils,  
Instruments, Tackling, and Mate-  
rials requisite thereto : With Rules  
and Directions how to use them





# Hadoop circa 2013

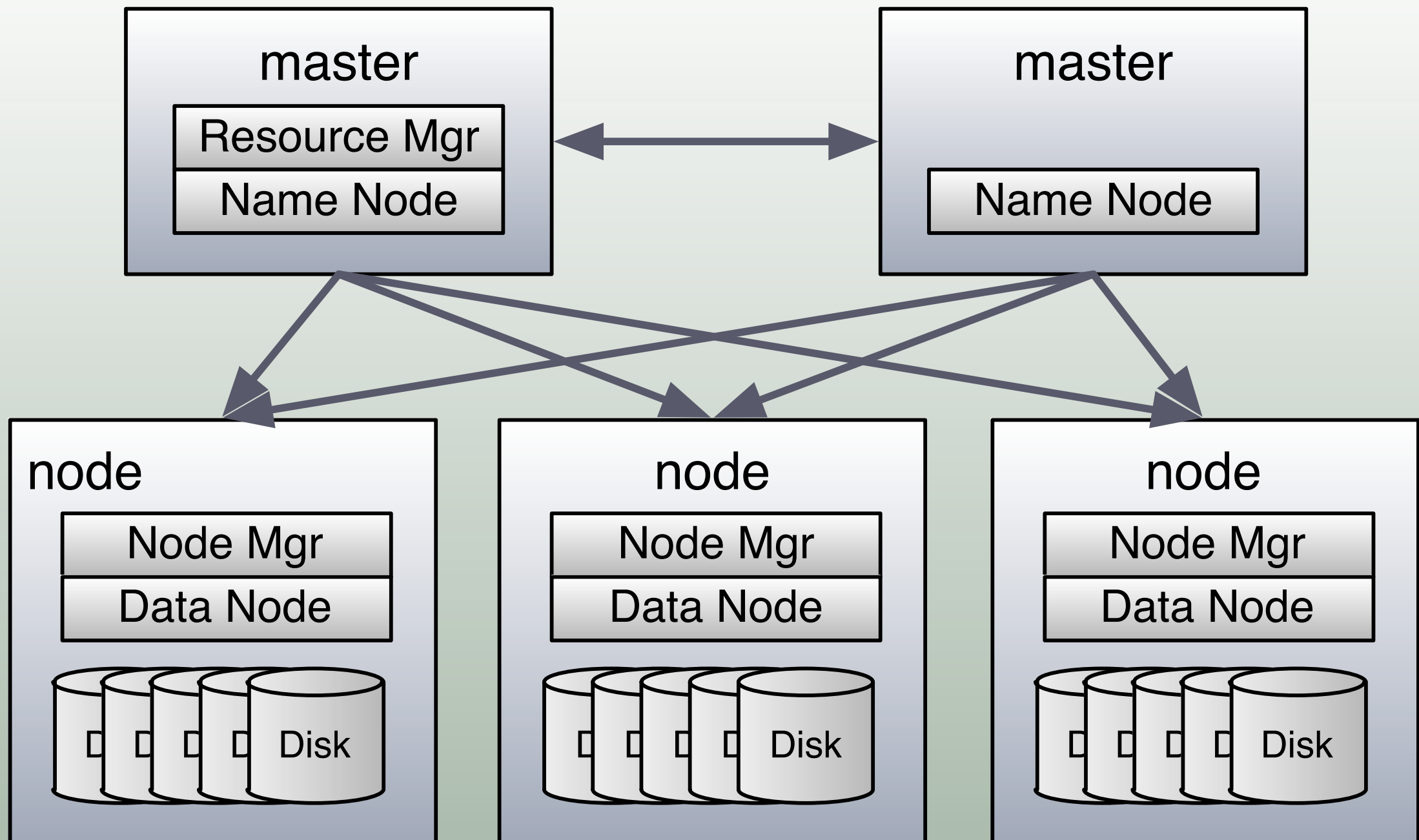
Saturday, April 19, 14

The state of Hadoop as of last year.

Image: Detail of the London Eye



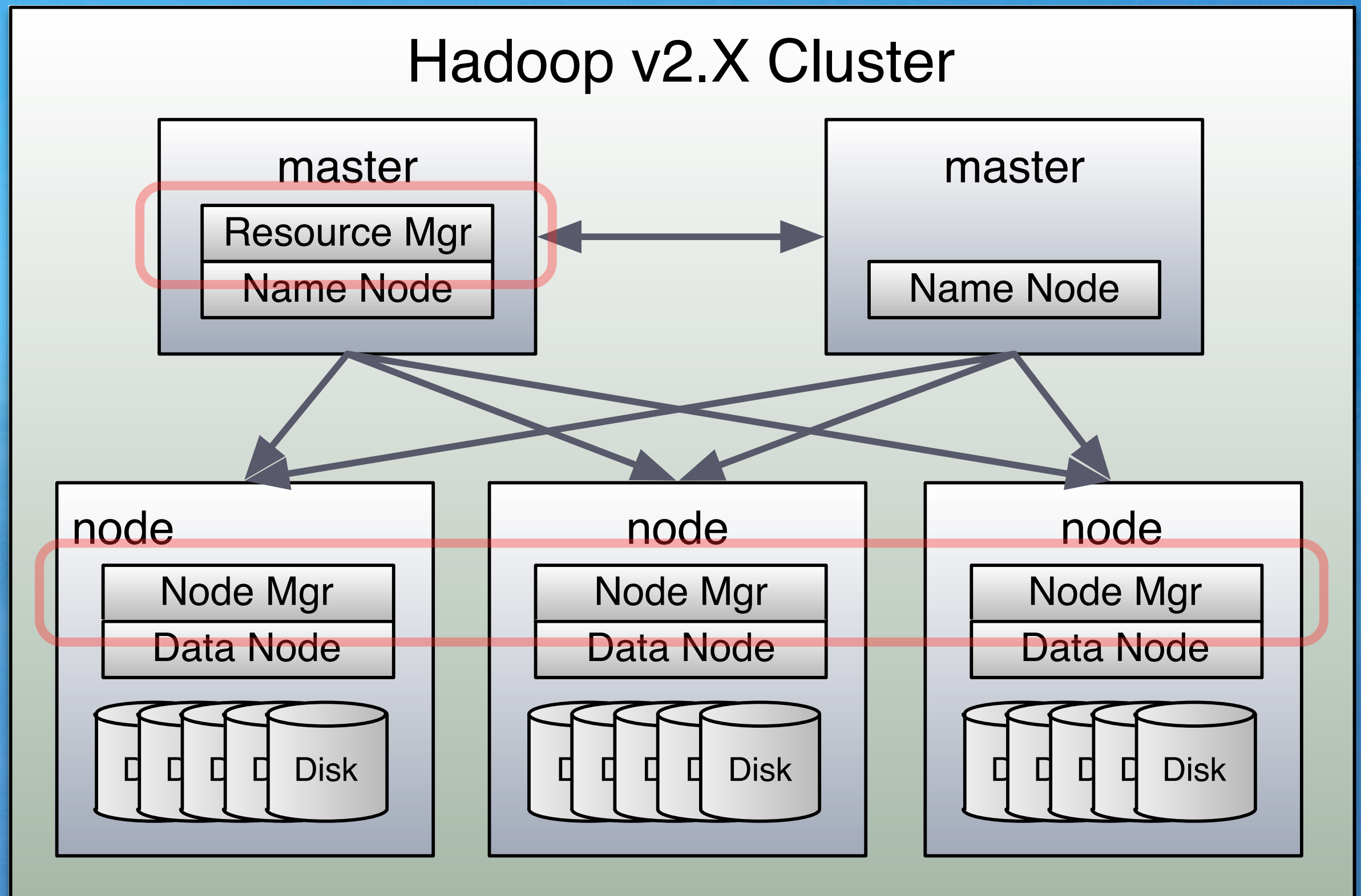
# Hadoop v2.X Cluster



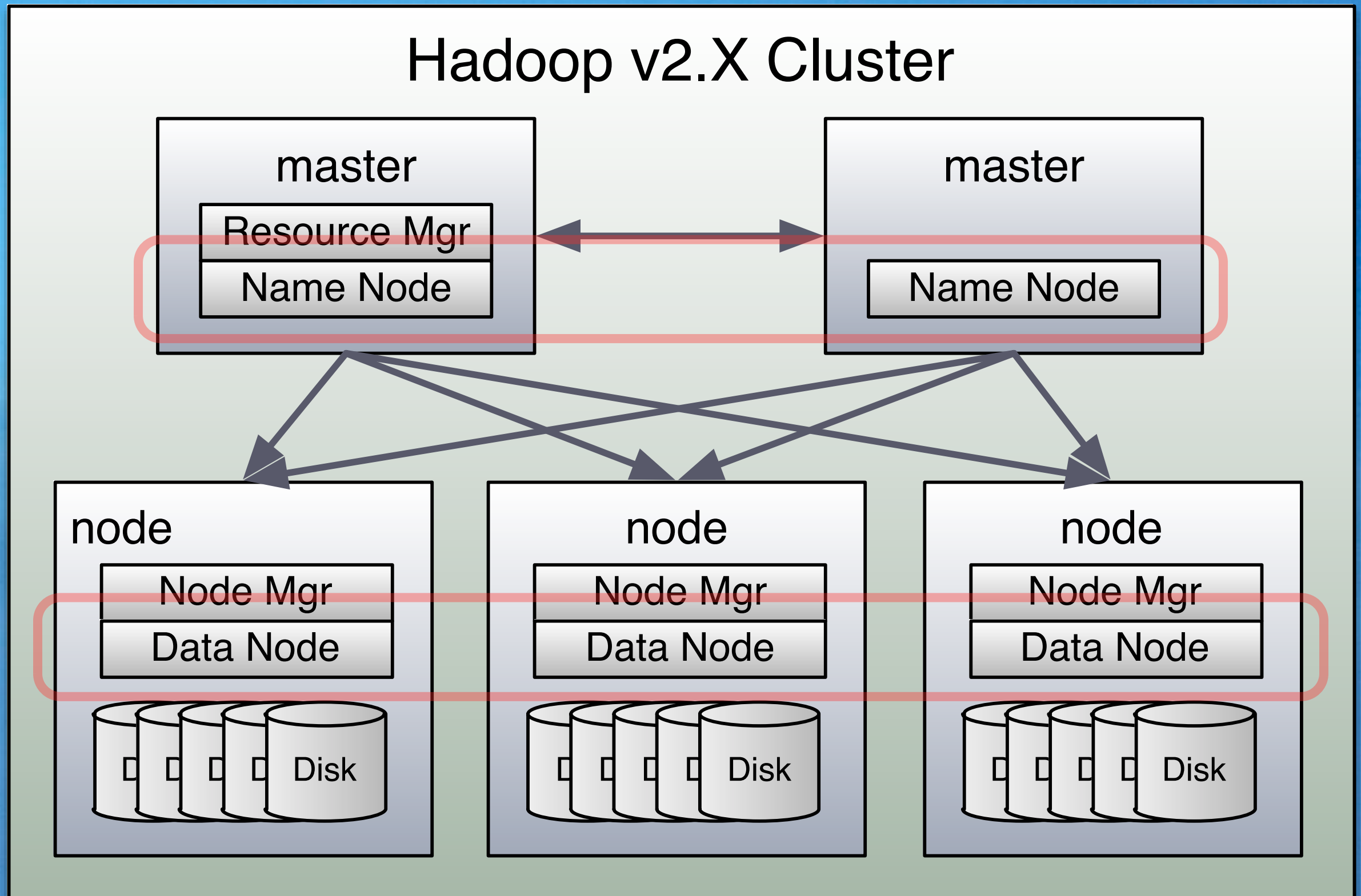
Saturday, April 19, 14

Schematic view of a Hadoop 2 cluster. For a more precise definition of the services and what they do, see e.g., <http://hadoop.apache.org/docs/r2.3.0/hadoop-yarn/hadoop-yarn-site/YARN.html> We aren't interested in great details at this point, but we'll call out a few useful things to know.

# Resource and Node Managers



# Name Node and Data Nodes





# MapReduce

*The classic compute model  
for Hadoop*



# MapReduce

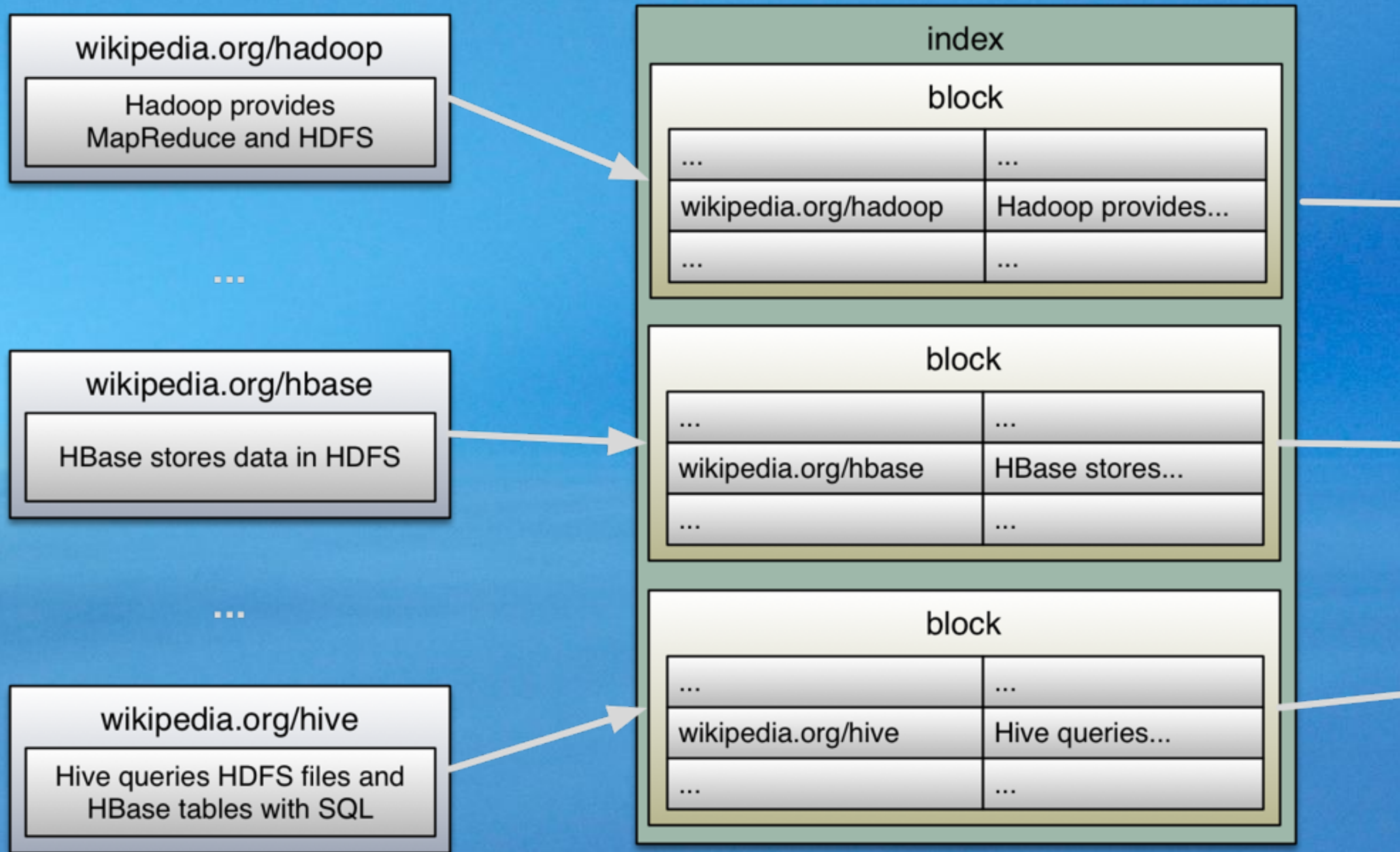
*1 map step + 1 reduce step*  
*(wash, rinse, repeat)*

# MapReduce

## ***Example:*** *Inverted Index*

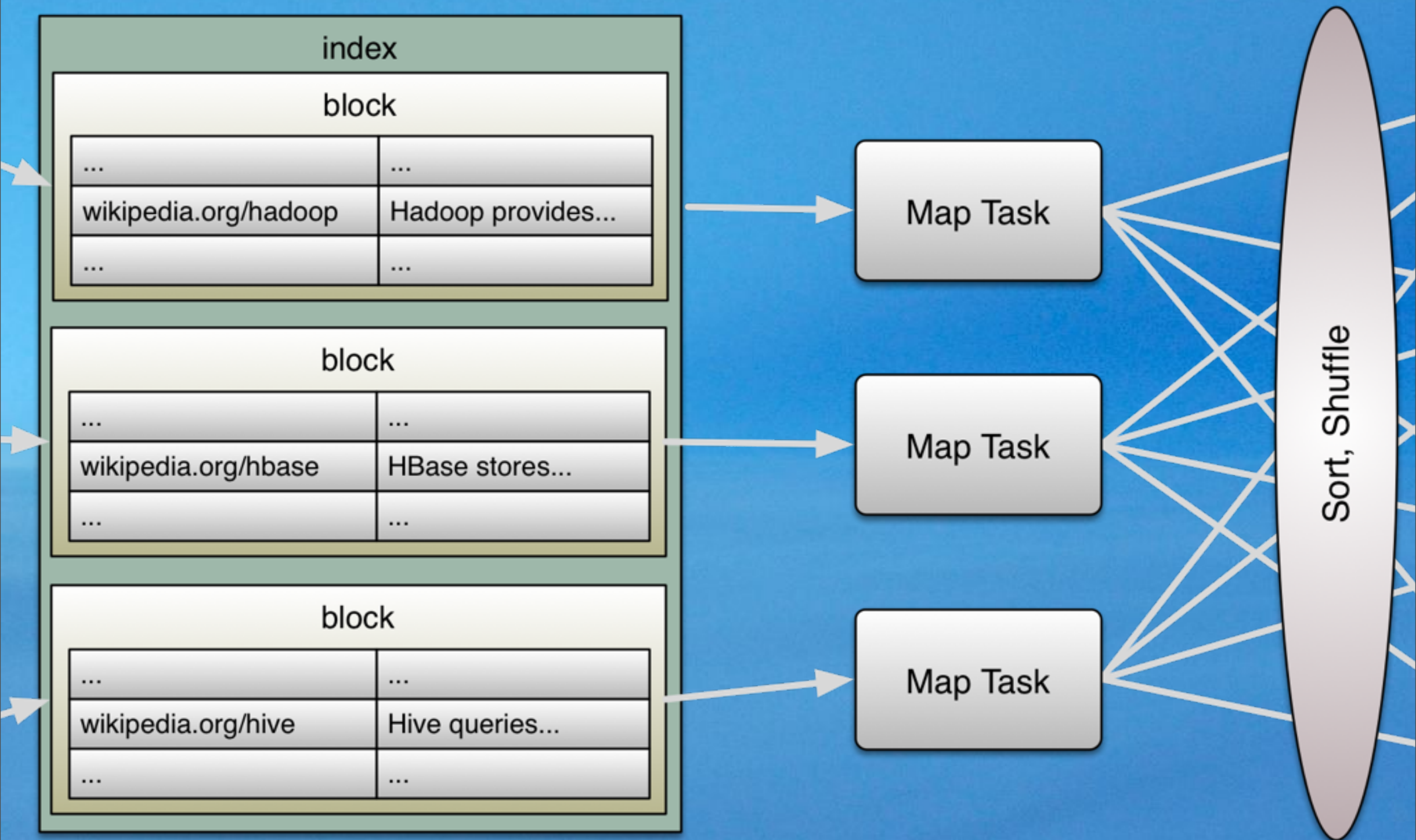


# Web Crawl



Saturday, April 19, 14

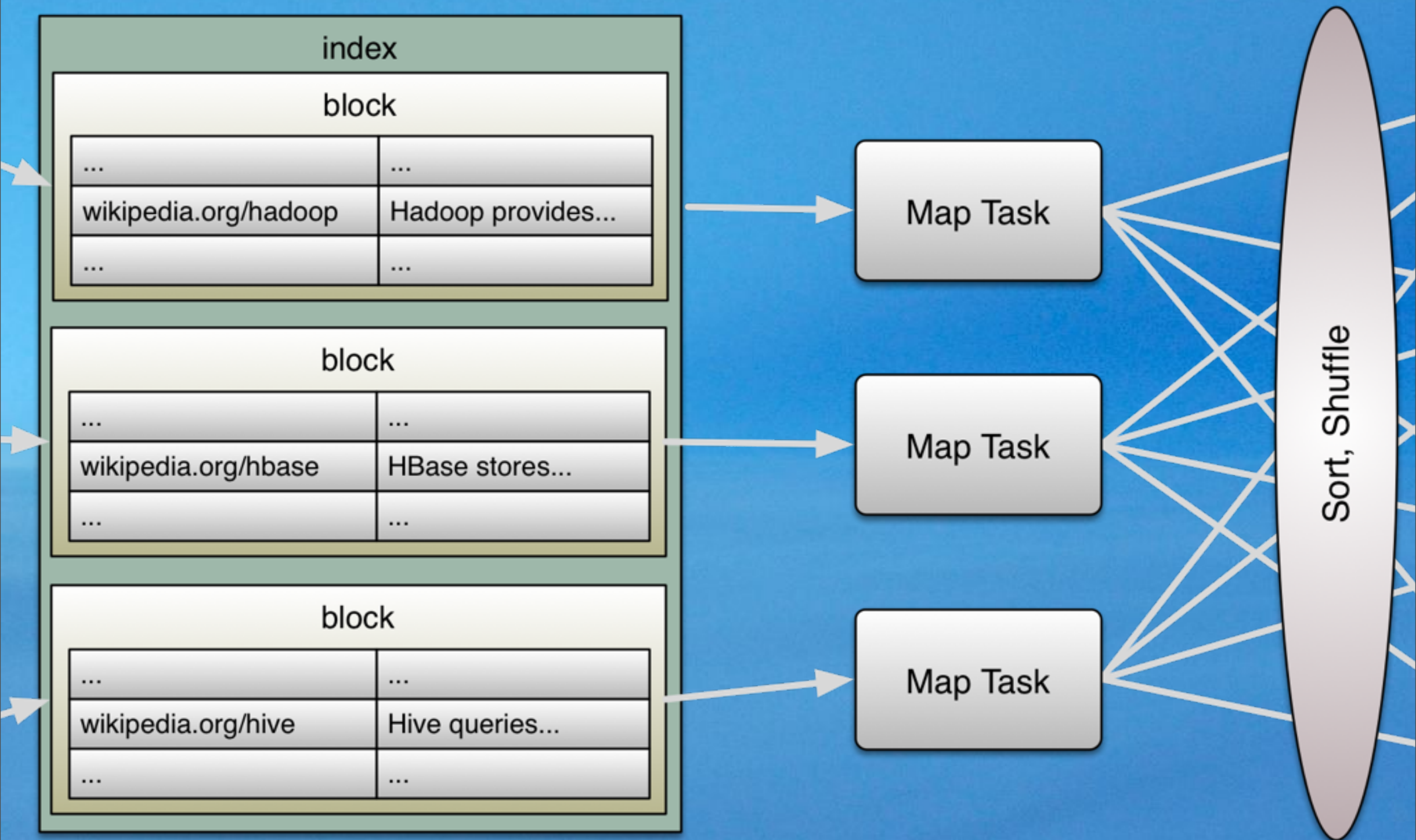
Before running MapReduce, crawl teh interwebs, find all the pages, and build a data set of URLs -> doc contents, written to flat files in HDFS or one of the more “sophisticated” formats.





## Map Task

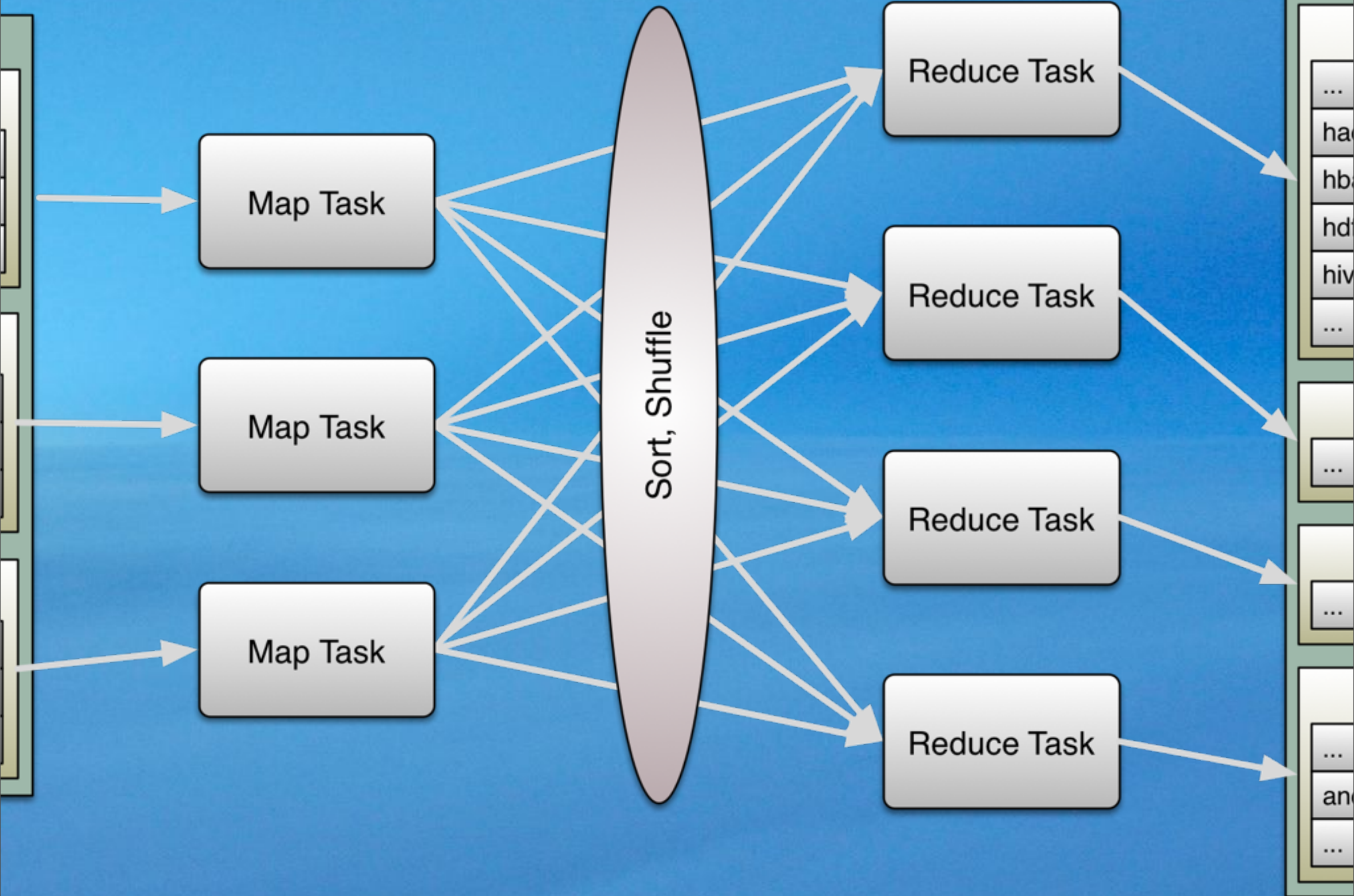
(hadoop,(wikipedia.org/hadoop,1))  
(provides,(wikipedia.org/hadoop,1))  
(mapreduce,(wikipedia.org/hadoop, 1))  
(and,(wikipedia.org/hadoop,1))  
(hdfs,(wikipedia.org/hadoop, 1))





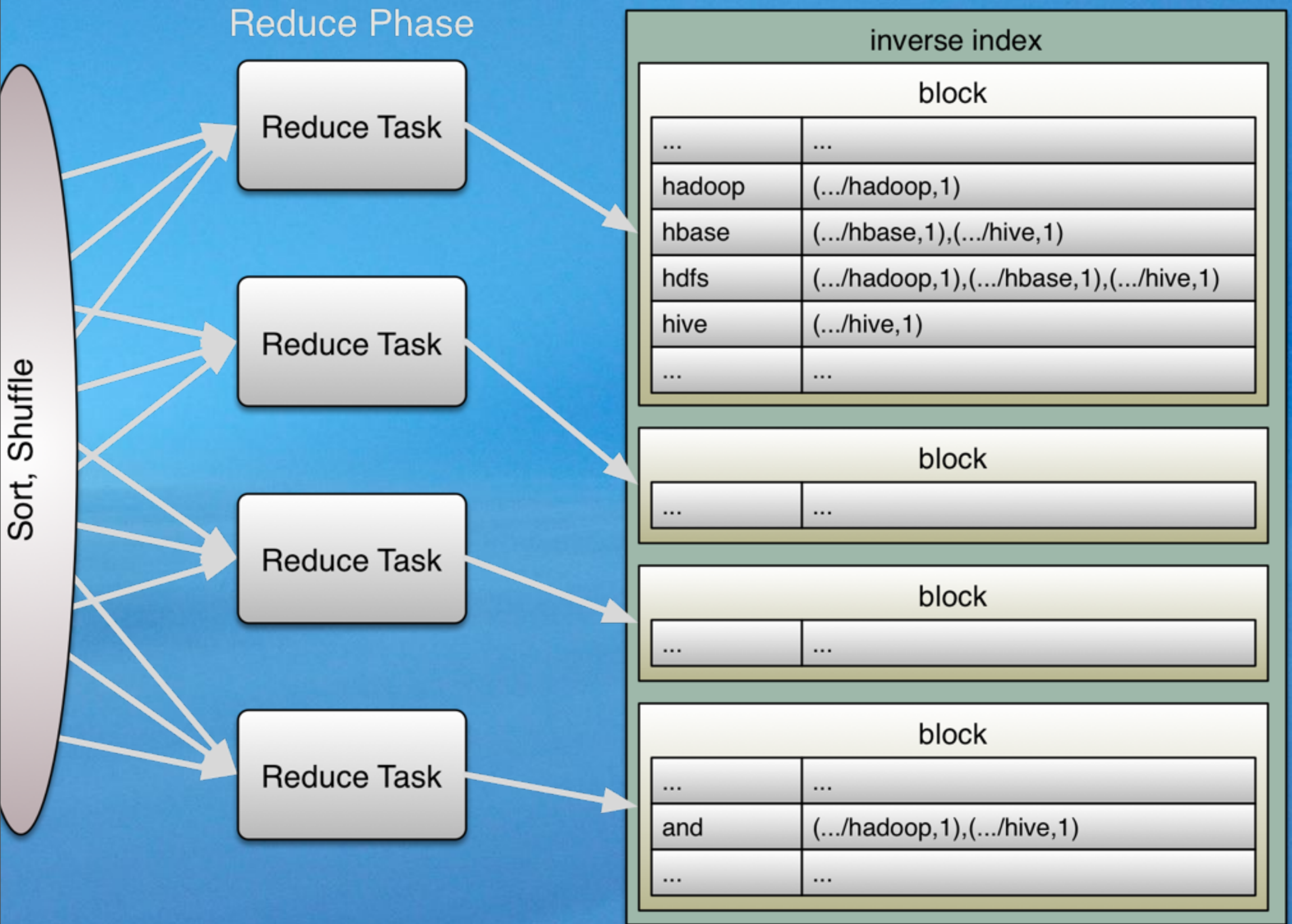
## Map Phase

## Reduce Phase



Saturday, April 19, 14

The output tuples are sorted by key locally in each map task, then “shuffled” over the cluster network to reduce tasks (each a JVM process, too), where we want all occurrences of a given key to land on the same reduce task.







# Why Spark Is the Next Top (Compute) Model

Philly ETE 2014  
April 22-23, 2014  
@deanwampler  
[polyglotprogramming.com/talks](http://polyglotprogramming.com/talks)

Saturday, April 19, 14

Copyright (c) 2014, Dean Wampler, All Rights Reserved, unless otherwise noted.

Image: Detail of the London Eye