

The background image shows a wide-angle view of a university campus. On the left, there is a long, modern building with a white facade and vertical slats. In the center, a wide, paved walkway leads through a green lawn with several trees. On the right, a tall, dark glass building with a clock tower is visible. The sky is blue with scattered white clouds.

TI2736-B

Big Data Processing

Claudia Hauff
ti2736b-ewi@tudelft.nl

Learning objectives

- **Exploit** Hadoop's Counters and setup/cleanup efficiently
- **Explain** how Hadoop addresses the problem of job scheduling
- **Explain** Hadoop's shuffle & sort phase and **use** that knowledge to improve your Hadoop code
- **Implement** strategies for efficient data input

Hadoop Programming

Revisited: setup and cleanup

Setup & cleanup

Programmer “hints”
the number of mappers
to use

- One **MAPPER object** for each map task
 - Associated with a sequence of key/value pairs (the “input split”)
 - `map()` is called **for each key/value pair** by the execution framework
- One **REDUCER object** for each reduce task
 - `reduce()` is called once per intermediate key
- MAPPER/REDUCER are **Java objects** -> **allows side effects**
 - Preserving state across multiple inputs
 - Initialise additional resources
 - Emit (intermediate) key/value pairs in one go

Programmer can
set the number of
reducers

Setup

Setup useful for one-off operations:

- opening an SQL connection
- loading a dictionary
- etc.

WordCount - count only valid dictionary terms*

```
1 public class MyMapper extends
2     Mapper<Text, IntWritable, Text, IntWritable> {
3
4     private Set<String> dictionary;//all valid words
5
6     public void setup(Context context) throws IOException {
7         dictionary = Sets.newHashSet();
8         loadDictionary();//defined elsewhere, reads file from HDFS
9     }
10
11    public void map(Text key, IntWritable val, Context context)
12        throws IOException, InterruptedException {
13        if(!dictionary.contains(key.toString()))
14            return;
15        context.write(key, new IntWritable(1));
16    }
17 }
```

Setup

Setup useful for one-off operations:

- opening an SQL connection
- loading a dictionary
- etc.

WordCount - count only valid dictionary terms*

```
1 public class MyMapper extends
2     Mapper<Text, IntWritable, Text, IntWritable> {
3
4     private Set<String> dictionary; //all valid words
5
6     public void setup(Context context) throws
7         dictionary = Sets.newHashSet();
8         loadDictionary(); //defined elsewhere
9     }
10
11     public void map(Text key, IntWritable val
12         throws IOException, I
13         if(!dictionary.contains(key.toString)
14             return;
15         context.write(key, new IntWritable(1));
16     }
17 }
```

Called once in the life cycle of a Mapper object: before any calls to map()

Called once for each key/value pair that appears in the input split

Cleanup

WordCount** - how many words start with the same letter?

```
1 public class MyReducer extends
2   Reducer<PairOfIntString, FloatWritable, NullWritable, Text> {
3   private Map<Character, Integer> cache;
4
5   public void setup(Context context) throws IOException {
6     cache = Maps.newHashMap();
7   }
8   public void reduce(PairOfIntString key, Iterable<IntWritable>
9     values, Context context) throws
10    IOException, InterruptedException {
11     char c = key.toString().charAt(0);
12     for(IntWritable iw : values){
13       //add iw to the current value of key c in cache
14     }
15   }
17   public void cleanup(Context context) throws IOException,
18     InterruptedException {
19     for (Character c : cache.keySet()) {
20       context.write(new Text(c), new IntWritable(cache.get(c)));
21     }
22   }
23 }
```


Cleanup

WordCount** - how many words start with the same letter?

```

1 public class MyReducer extends
2   Reducer<PairOfIntString, FloatWritable, NullWritable, Text> {
3   private Map<Character, Integer> cache;
4
5   public void setup(Context context) throws IOException, InterruptedException {
6       cache = Maps.newHashMap();
7   }
8   public void reduce(PairOfIntString key, Iterable<IntWritable>
9       values, Context context) throws
10      IOException, InterruptedException {
11       char c = key.toString().charAt(0);
12       for(IntWritable iw : values){
13           //add iw to the current value
14       }
15   }
16
17   public void cleanup(Context context) throws IOException,
18       InterruptedException {
19       for (Character c : cache.keySet())
20           context.write(new Text(c), new IntWritable(cache.get(c)));
21   }
22 }
23 }

```

Called once in the life cycle of a Reducer object: before any calls to `reduce()`

Called once for each key that was assigned to the reducer

Called once in the life cycle of a Reducer object: after all calls to `reduce()`

Hadoop Programming Revisited: Counters

Counter basics

- **Gathering data about the data** we are analysing, e.g.
 - Number of key/value pairs processed in map
 - Number of empty lines/invalid lines
- Wanted:
 - **Easy** to collect
 - **Estimates are viewable during job execution** (e.g. to stop a Hadoop job early at too many invalid key/value pairs)
- Why not use log messages instead?
 - Write to the error log when an invalid line occurs
 - Hadoop's logs are huge, you need to know where to look
 - Aggregating stats from the logs requires another pass over it

Counter basics

- **Gathering data about the data** we are analysing, e.g.
 - Number of key/value pairs processed in map
 - Number of empty lines/invalid lines
- Wanted:
 - **Easy** to collect
 - **Viewable during job execution** (stop Hadoop job early at too many invalid key/value pairs)
- What about log messages?
 - Write to the error log when an invalid line occurs
 - Hadoop's logs are huge, you need to know where to look
 - Aggregating stats from the logs requires another pass over it

Counter basics

- Counters: Hadoop's way of **aggregating** statistics
- Counters **count** (increment)
- **Built-in counters** maintain **metrics** of the job
 - MapReduce counters (e.g. #skipped records by all maps)
 - File system counters (e.g. #bytes read from HDFS)
 - Job counters (e.g. #launched map tasks)
- You have already seen them

Counter basics

- Counters: Hadoop's way of **aggregating** statistics
- Counters **count** (increment)
- **Built-in counters** maintain **metrics** of the job
 - MapReduce counters (e.g. #skipped records by all maps)
 - File system counters (e.g. #bytes read from HDFS)
 - Job counters (e.g. #launched map tasks)
- You have already seen them

Counter basics

- Counters: Hadoop's way of **aggregating** statistics

- Map-Reduce Framework

Map input records=5903

Map output records=47102

Combine input records=47102

Combine output records=8380

Reduce output records=5934

File System Counters

FILE: Number of bytes read=118124

FILE: Number of bytes written=1075029

HDFS: Number of bytes read=996209

HDFS: Number of bytes written=59194

Built-in vs. user-defined

- **Built-in counters**: exist for each Hadoop job
- **User-defined Counters** are maintained by the application they are associated with
 - Periodically sent to the Tasktracker and then the Jobtracker for global aggregation (pre-YARN setup)
 - Aggregated per job by the ResourceManager (YARN)

Counter values are only definite once the job has completed!
Counters may go down if a task fails!

Code example

WordCount* - count words and chars

```
1 enum Records {
2     WORDS, CHARS;
3 };
4 public class WordCount {
5     public static class MyMapper extends
6         Mapper<LongWritable, Text, Text, IntWritable> {
7
8     public void map(LongWritable key, Text value,
9         Context context) throws IOException {
10        String[] tokens = value.toString().split(" ");
11
12        for (String s : tokens) {
13            context.write(new Text(s), new IntWritable(1));
14            context.getCounter(Records.WORDS).increment(1);
15            context.getCounter(Records.CHARS).increment(s.length());
16        }
17    }
18 }
```

several enum's possible:
used to group counters

user-defined counters appear
automatically in the final status output

Code example

WordCount* - count words and chars

```
1 enum Records {
2     WORDS, CHARS;
3 };
4 public class WordC
5     public static cl
6         Mapper
7
8     public void ma
9
10        String[] t
11
12        for (Strin
13            context.
14            context.
15            context.
16        }
17    }
18 }
```

Map-Reduce Framework

Map input records=5903

Map output records=47102

Combine input records=47102

Combine output records=8380

Reduce output records=5934

...

Records

CHARS=220986

WORDS=47102

user-defined counters appear automatically in the final status output

Code example II

```
1 enum Records { MAP_WORDS, REDUCE_WORDS; };
2
3 public class WordCount {
4     --> MAPPER
5     public void map(LongWritable key, Text value,
6                     Context context)
7                     throws IOException {
8
9         String[] tokens = value.toString().split(" ");
10        for (String s : tokens) {
11            context.write(new Text(s), new IntWritable(1));
12            context.getCounter(Records.MAP_WORDS).increment(1);
13        }
14    }
15    --> REDUCER (Combiner is a copy of the Reducer)
16    public void reduce(Text key, Iterator<IntWritable> values,
17                      Context context) throws IOException {
18        int sum = 0;
19        while (values.hasNext())
20            sum += values.next().get();
21        context.getCounter(Records.REDUCE_WORDS).increment(sum);
22    }
23 }
```

...

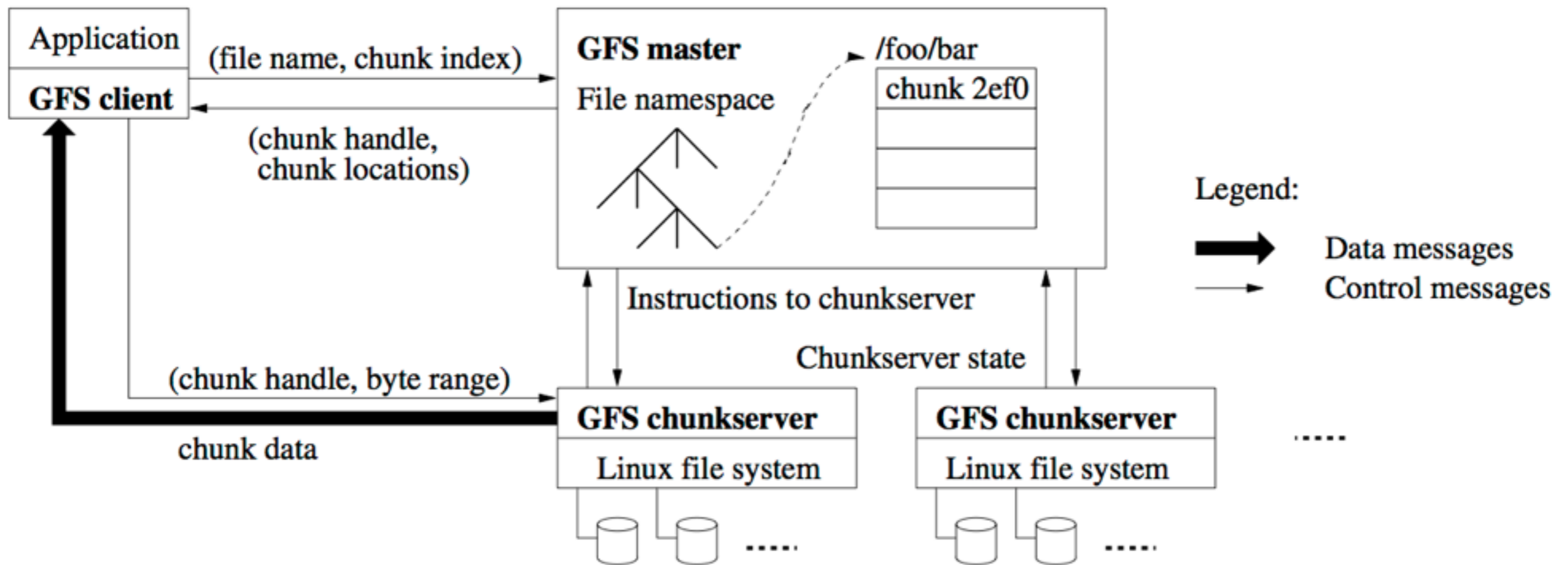
Records

MAP_WORDS=47102

REDUCE_WORDS=47102

Job Scheduling

Last time ... GFS/HDFS

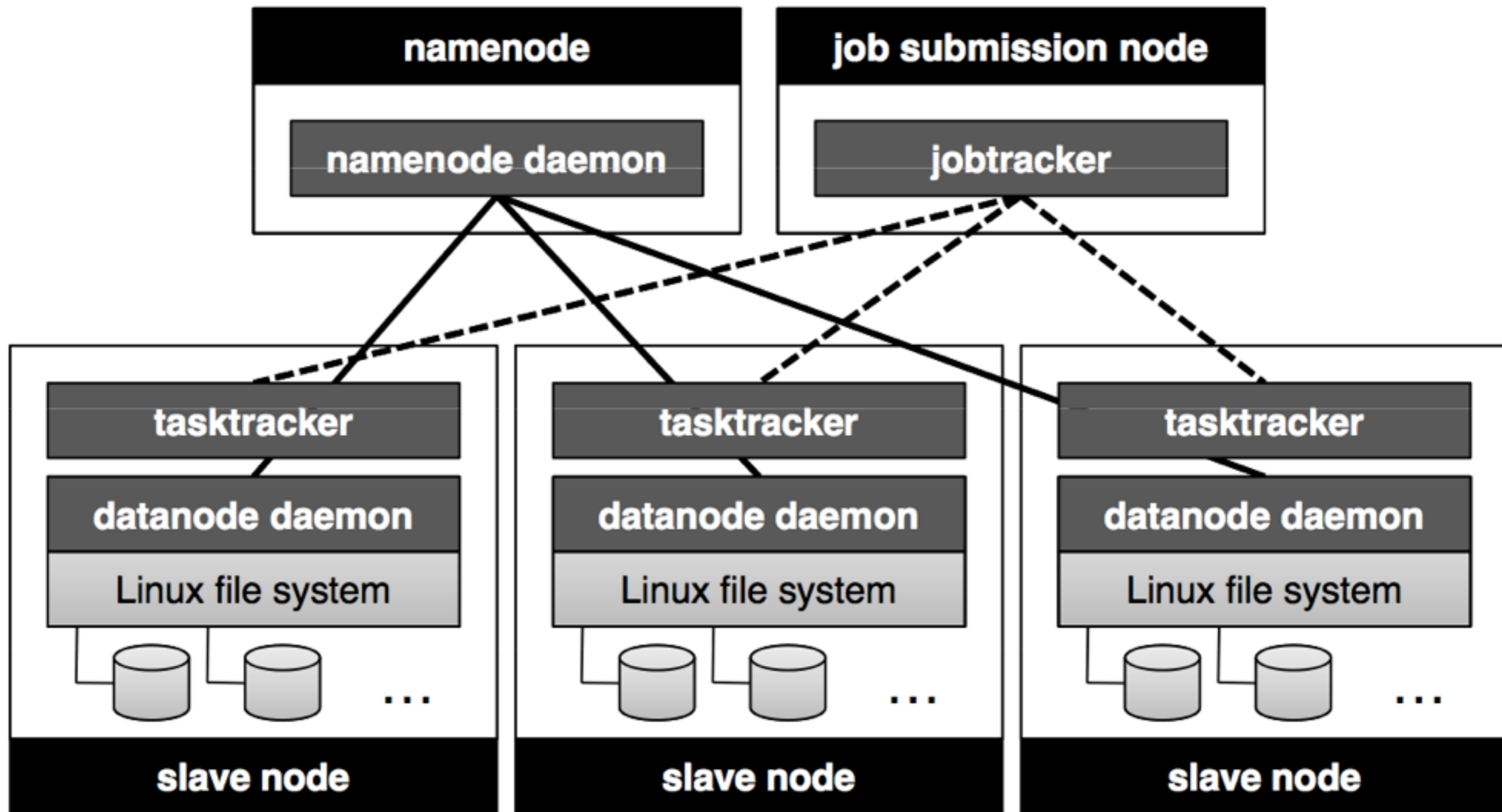


distributed file system: file systems that manage the storage across a network of machines.

What about the jobs?

- **Hadoop job**: unit of work to be performed
 - Input data
 - MapReduce program
 - Configuration information
- Hadoop divides input data into **fixed size input splits**
 - **One map task per split**
 - One map function call for each **record** in the split
 - Splits are processed in parallel (if enough DataNodes exist)

JobTracker and TaskTracker



Hadoop in practice: Yahoo! (2010)

- **40 nodes/rack** sharing one IP switch
- **16GB RAM** per cluster node, 1-gigabit Ethernet
- 70% of disk space allocated to HDFS
 - Remainder: operating system, data emitted by Mappers (not in HDFS)
- **NameNode**: up to **64GB RAM**
- **Total storage**: 9.8PB -> **3.3PB** net storage (replication: 3)
- **60 million files**, 63 million blocks
- 54,000 blocks hosted per DataNode
- **1-2 nodes lost per day**
- **Time for cluster to re-replicate lost blocks: 2 minutes**

HDFS cluster with 3,500 nodes

YARN (MapReduce 2)

- JobTracker/TaskTrackers setup becomes a **bottleneck** in clusters with thousands of nodes
- As answer YARN has been developed (**Yet Another Resource Negotiator**)
- YARN splits the JobTracker's tasks (job scheduling and task progress monitoring) into two daemons:
 - **Resource manager** (RM)
 - **Application master** (negotiates with RM for cluster resources; each Hadoop job has a dedicated master)

Job scheduling

- **Thousands of tasks** may make up **one job**
- Number of tasks can exceed number of tasks that can run concurrently
 - **Scheduler** maintains **task queue** and tracks progress of running tasks
 - Waiting tasks are assigned nodes as they become available
- “**Move code to data**”
 - Scheduler starts tasks on node that holds a particular block of data needed by the task if possible

Job scheduling

FIFO scheduler

Priority scheduler

Fair scheduler

Capacity scheduler

Basic schedulers

- Early on: **FIFO scheduler**
 - Job occupies the whole cluster while the rest waits
 - **Not feasible in larger clusters**
- Improvement: different **job priorities** VERY_HIGH, HIGH, NORMAL, LOW, or VERY_LOW
 - Next job is the one with the highest priority
 - **No pre-emption**: if a low priority job is occupying the cluster, the high priority job still has to wait

Fair Scheduler I

- Goal: every user receives a **fair share** of the cluster capacity over time
- If a **single job** runs, it uses the entire cluster
 - As **more jobs** are submitted, free task slots are given away such that each user receives a “fair share”
 - **Short jobs** complete in reasonable time, long jobs keep progressing
- A user who submits more jobs than a second user will not get more cluster resources on average

Fair Scheduler II

- Jobs are placed in **pools**, default: one pool per user
- **Pre-emption**: if a pool has not received its fair share for a period of time, the scheduler will **kill tasks** in pools running over capacity to give more slots to the pool running under capacity
 - **Task kill != Job kill**
 - Scheduler needs to keep track of all users, resources used

Capacity Scheduler

- Cluster is made up of a **number of queues** (similar to the Fair Scheduler pools)
- Each queue has an allocated **capacity**
- Within each queue, jobs are scheduled using FIFO with priorities
- Idea: users (defined using queues) **simulate** a **separate MapReduce cluster** with FIFO scheduling for each user

Speculative execution

- **Map phase is only as fast as slowest MAPPER**
- **Reduce phase is only as fast as slowest REDUCER**
- Hadoop job is sensitive to **stragglers** (tasks that take unusually long to complete)
- Idea: **identical copy** of task executed on a second node; the output of whichever node finishes first is used (improvements up to 40%)
- Can be done for both MAPPER/REDUCER
- Strategy does not help if straggler due to **skewed data distribution**

Shuffle & Sort

Shuffle & sort phase

- **Hadoop guarantee**: the input to every reducer is sorted by key
- **Shuffle**: sorting of **intermediate key/value pairs** and transferring them to the reducers (as input)
- “**Shuffle is the heart of MapReduce**”
- Understanding shuffle & sort is vital to recognise job bottlenecks
- Disclaimer: constantly evolving (*again*)

A high-level view

MAP TASK

input split

map()

in-memory buffer

partition, sort, and spill to disk

merge (disk)

reduce tasks

input split

reduce()

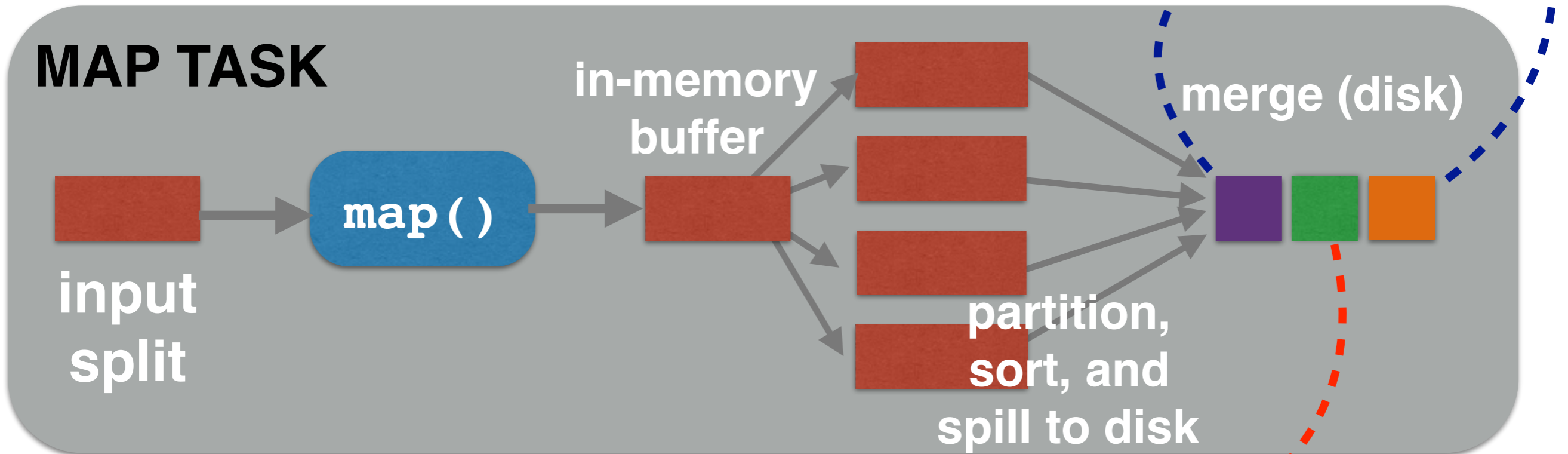
sort phase

copy across the network

map tasks

REDUCE TASK

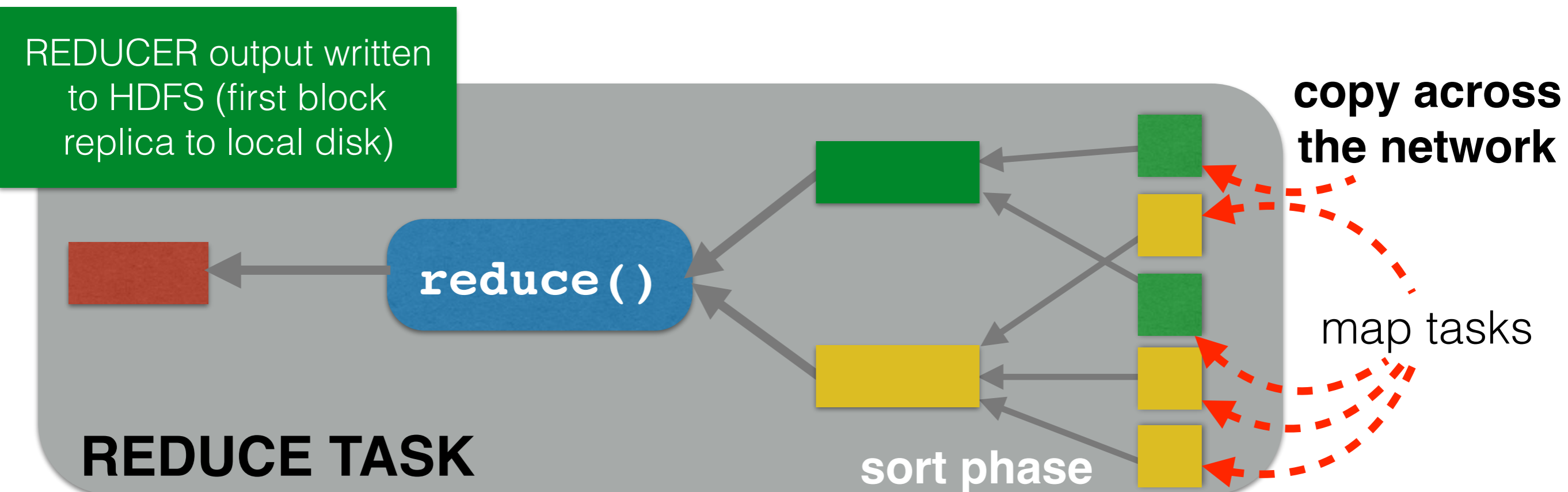
Map side



- Map task writes output to **memory buffer**
- Once the buffer is full, a background thread **spills** the content to **disk** (spill file)
 - Data is partitioned corresponding to reducers they will be send to
 - Within partition, in-memory sort by key [combiner runs on the output of sort]
- After last `map()` call, the spill files are merged [combiner may run again]

Reduce side

- Reducer requires the map output for its partition from **all map tasks of the cluster**
- Reducer **starts copying data** as soon as a map task completes ("copy phase")
- Direct copy to reducer's memory if the output is small, otherwise copy to disk
- In-memory buffer is merged and spilled to disk once it grows too large
- **Combiner may run again**
- Once **all** intermediate keys are copied the "sort phase" begins: merge of map outputs, maintaining their sort ordering



A few more details

MAP TASK



What happens to the data written to local disk by the Mapper?

Deleted after successful completion of the job.

General rule for memory usage: map/reduce/shuffle

Shuffle should get as much memory as possible; write map/reduce with low memory usage (single spill would be best)

REDUCE TASK



How does the Reducer know where to get the data from?

- Successful map task informs task tracker which informs the job tracker (via heartbeat)
- Reducer periodically queries the job tracker for map output hosts until it has retrieved all of data

reduce tasks

copy across the network

Sort phase recap

- Involves **all nodes** that **executed map tasks** and **will execute reduce tasks**
 - Job with m mappers and r reducers involves up to $m*r$ distinct copy operations
- Reducers can only start calling `reduce ()` after all mappers are finished
 - **Key/value guarantee**: one key has all values “attached”
- Copying can start earlier for intermediate keys

Data input

Input splits and logical bounds

- One **MAPPER** object for each map task
 - Associated with a sequence of key/value pairs (the “input split”)
 - **map ()** is called for each key/value pair by the execution framework

Input split	record
(part of) a text file	line of text
(range of) database table rows	a single row
(part of) an XML file	XML element
(part of) a video stream	keyframe

Input split

```
1 public abstract class InputSplit {  
2  
3     public abstract long getLength() throws  
4         IOException, InterruptedException;  
5  
6     public abstract String[] getLocations() throws  
7         IOException, InterruptedException;  
8 }
```

Scenario: There are less free map slots than input splits.

Questions: Given the input splits and their sizes, what are possible strategies of how to pick the next input split to process by a map task?

Given 3 free map slots and 7 input splits of sizes **{10, 20, 30, 100, 200, 300, 400}**, which strategy works best?

Input splits and logical bounds

3 free map slots and 7 input splits of sizes **{10, 20, 30, 100, 200, 300, 400}**

Random selection: worst case scenario

=530

100
10
20
400

Shortest input split first

=510

10
100
400

20
300

30
300

Longest input split first

=400

400

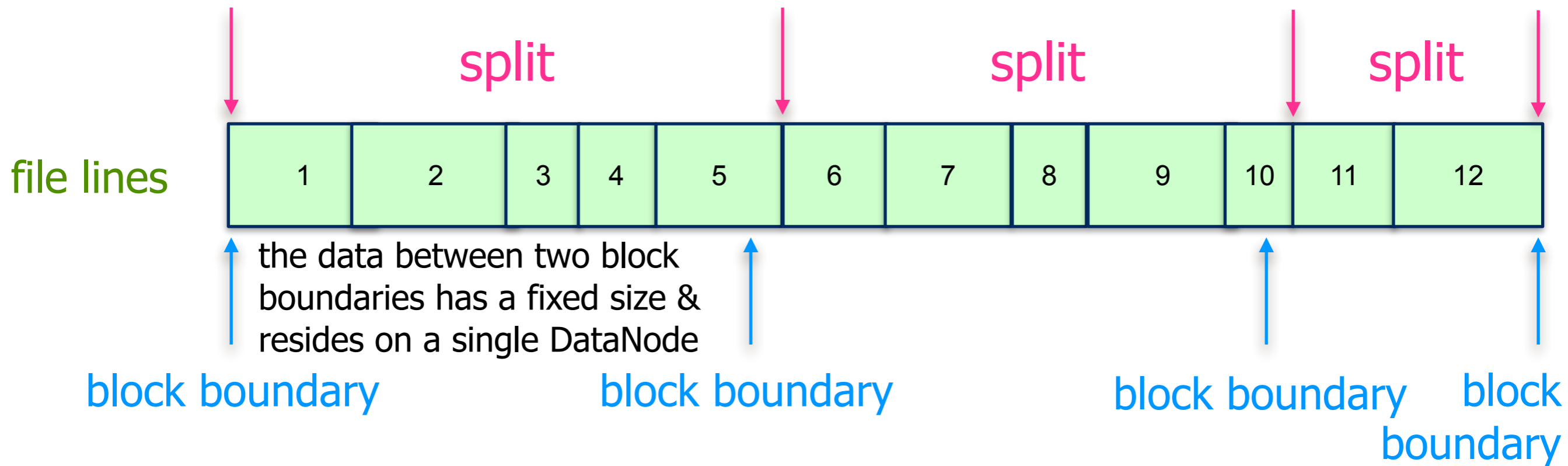
300
20
10

200
100
30

greedy approximation of optimal approach; strategy used by Hadoop

Input splits vs. HDFS blocks

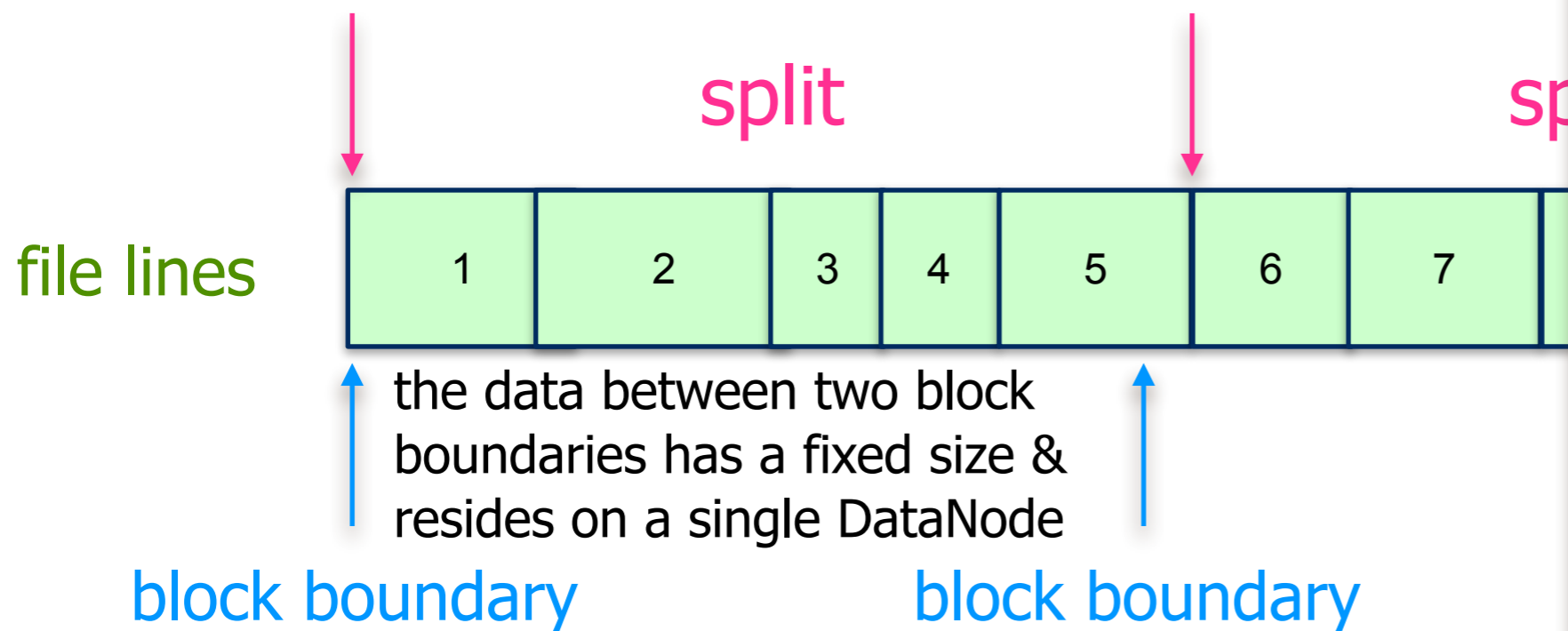
an input split ends at a logical record; it can cross the block boundary; requires additional remote read of the missing data



TextInputFormat

Input splits vs. HDFS blocks

an input split ends at a logical record; it can cross the block boundary; requires additional remote read of the missing data



What about:

```
<wikipedia>
<page></page>
. . .
10GB later
. . .
<page></page>
</wikipedia>
```

TextInputFormat

HDFS: Compression and Small Files

Splittable compression

Compression format	Tool	Algorithm	Filename extension	Splittable?
DEFLATE ^[a]	N/A	DEFLATE	<i>.deflate</i>	No
gzip	<i>gzip</i>	DEFLATE	<i>.gz</i>	No
bzip2	<i>bzip2</i>	bzip2	<i>.bz2</i>	Yes
LZO	<i>lzop</i>	LZO	<i>.lzo</i>	No ^[b]
LZ4	N/A	LZ4	<i>.lz4</i>	No
Snappy	N/A	Snappy	<i>.snappy</i>	No

space/time tradeoff: faster (de)compression means less space savings

Splittable compression

Compression format	Tool	Algorithm	Filename extension	Splittable?
DEFLATE ^[a]	N/A	DEFLATE	<i>.deflate</i>	No
gzip	<i>gzip</i>	middle ground		No
bzip2	<i>bzip2</i>	better compression		Yes
LZO	<i>lzop</i>	optimized for speed, less effective compression		No ^[b]
LZ4	N/A		No	
Snappy	N/A		No	

space/time tradeoff: faster (de)compression means less space savings

Splittable is an important attribute

- 1GB uncompressed file
 - Stored within 16 blocks on HDFS (block size 64MB)
 - Hadoop job creates 16 input splits, **each processed by one map task**
- 1GB gzip-compressed file
 - Stored within 16 blocks on HDFS
 - Hadoop job cannot create 16 input splits (reading at an arbitrary point does not work)
 - A **single map task** will process the 16 HDFS blocks

Hadoop Archives

- Storing a large number of small files is inefficient
 - **But**: not all files can be **easily** converted to blocks (e.g. millions of images)
- Files and blocks **occupy namespace** which is limited by the physical memory in the NameNode
 - Small files take up large portion of namespace but not the disk space
 - Rule of thumb: 150 bytes per file/directory/block (1 million files of one block each: 300MB of memory)
- Hadoop Archive (***.har**) is a solution

small=substantially less
than the block size
(64MB/128MB)

A Web special: WARC

- **Web ARCHive format**: aggregates digital resources in an archive and keeps track of related information
 - Per resource: text header and arbitrary data
- Extension of the **Internet Archive's ARC format**
- Commonly used to store **Web crawls**

A Web special: WARC

WARC/0.17

WARC-Type: response

WARC-Target-URI: http://www.archive.org/robots.txt

WARC-Date: 2008-04-30T20:48:25Z

WARC-Payload-Digest: sha1:SUCGMUVXDKVB5CS2NL4R4JABNX7K466U

WARC-IP-Address: 207.241.229.39

WARC-Record-ID: <urn:uuid:e7c9eff8-f5bc-4aeb-b3d2-9d3df99afb30>

Content-Type: application/http; msgtype=response

Content-Length: 782

HTTP/1.1 200 OK

Date: Wed, 30 Apr 2008 20:48:24 GMT

Server: Apache/2.0.54 (Ubuntu) PHP/5.0.5-2ubuntu1.4 mod_ssl/2.0.54 OpenSSL/
0.9.7g

Last-Modified: Sat, 02 Feb 2008 19:40:44 GMT

ETag: "47c3-1d3-11134700"

Accept-Ranges: bytes

Content-Length: 467

Connection: close

Content-Type: text/plain; charset=UTF-8

#####

Welcome to the Archive!

Hadoop's SequenceFile format

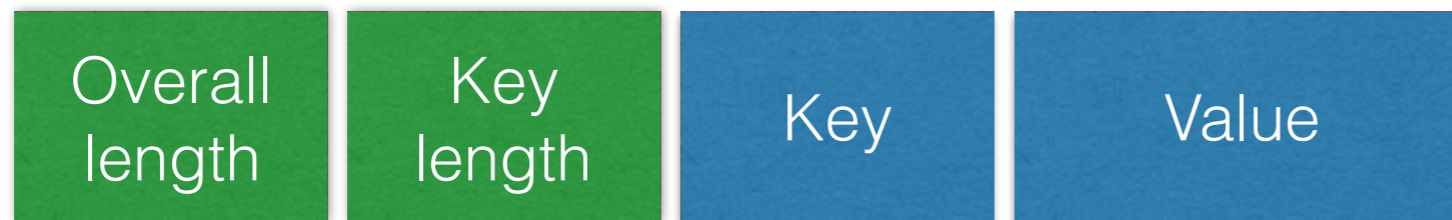
- Main usage: **intermediate** output of Mappers written in this format
- Flat file, consisting of **binary key-value pairs**
- Defines a `Reader`, `Writer` and `Sorter`
- **Three types:**
 - Uncompressed key-value pairs
 - Compressed values (“record compressed”)
 - Keys and values compressed (“block compressed”)
- From small files to SequenceFile:
(`some_key, file_content`)

Hadoop's SequenceFile format

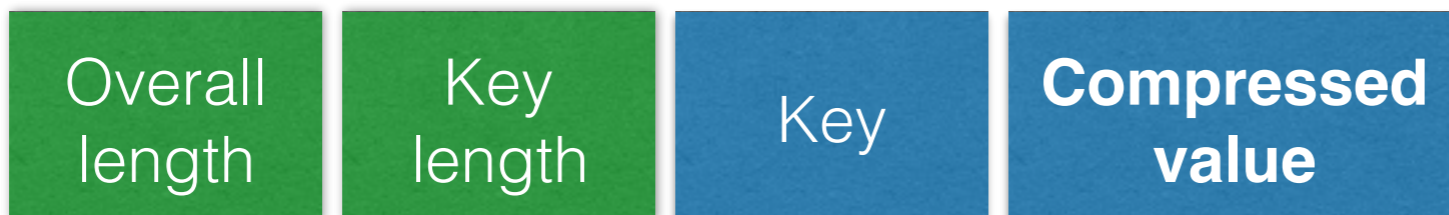
compression details, file meta-data, etc.



used to synchronise to a record boundary

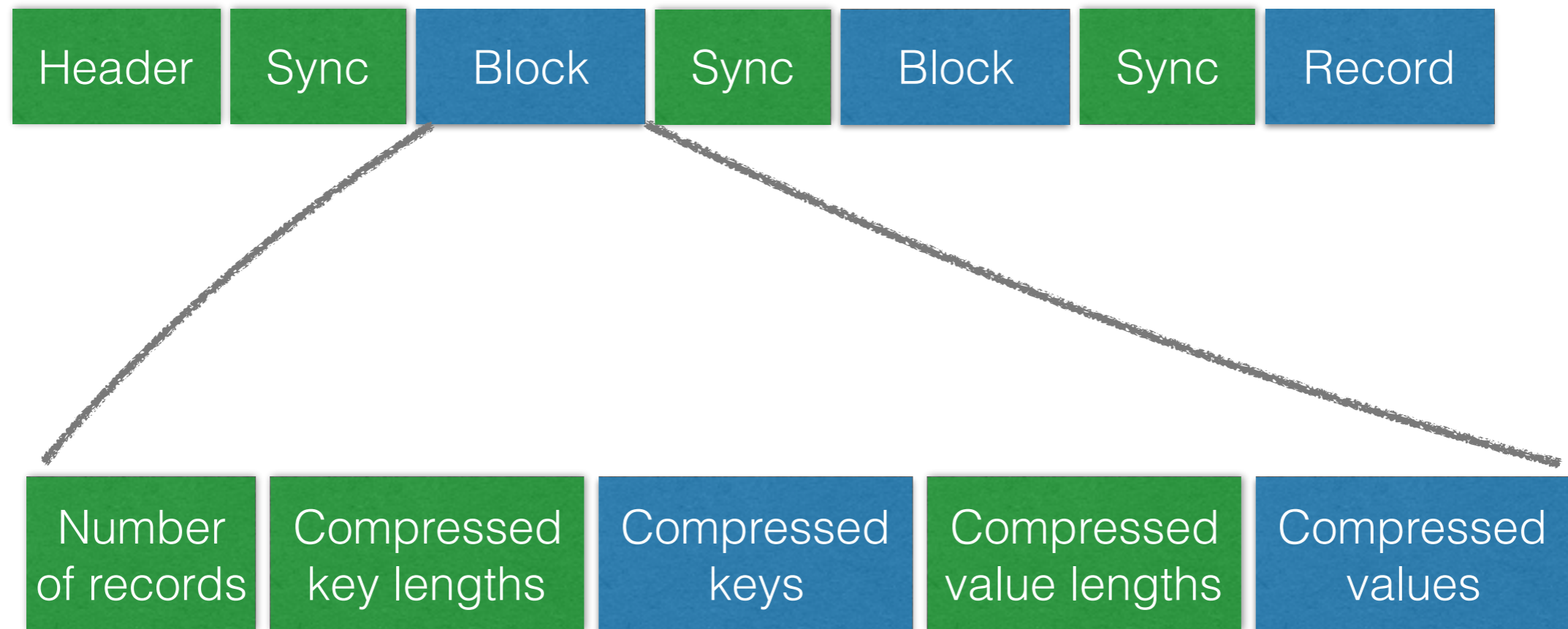


No compression



Record compression

Hadoop's SequenceFile format



Block compression (allows most compression)

HDFS: the rest

HDFS is just one possible implementation

Filesystem	URI scheme	Java implementation (all under <code>org.apache.hadoop</code>)	Description
Local	<code>file</code>	<code>fs.LocalFileSystem</code>	A filesystem for a locally connected disk with client-side checksums. Use <code>RawLocalFileSystem</code> for a local filesystem without checksums. See LocalFileSystem .
HDFS	<code>hdfs</code>	<code>fs.DistributedFileSystem</code>	Hadoop's distributed filesystem. HDFS is designed to work efficiently in conjunction with MapReduce.
FTP	<code>ftp</code>	<code>fs.ftp.FTPFileSystem</code>	A filesystem backed by an FTP server.
S3 (native)	<code>s3n</code>	<code>fs.s3native.NativeS3FileSystem</code>	A filesystem backed by Amazon S3. See http://wiki.apache.org/hadoop/AmazonS3 .
S3 (block-based)	<code>s3</code>	<code>fs.s3.S3FileSystem</code>	A filesystem backed by Amazon S3, which stores files in blocks (much like HDFS) to overcome S3's 5 GB file size limit.
HAR	<code>har</code>	<code>fs.HarFileSystem</code>	A filesystem layered on another filesystem for archiving files. Hadoop Archives are typically used for archiving files in HDFS to reduce the namenode's memory usage. See Hadoop Archives .

we can also use the local filesystem

Summary

- Hadoop Counters, setup/cleanup
- Job scheduling
- Shuffle & sort
- Data input

THE END