

Parallel Programming Concepts

MapReduce

Frank Feinbube

Source:

MapReduce: Simplified Data Processing on Large Clusters; Dean et. Al.

Examples for Parallel Programming Support

2

	Task-Parallel Programming Model	Data-Parallel Programming Model	Actor Programming Model	Functional Programming Model	PGAS / DSM Programming Model
Shared Memory System	OpenMP, Threading Libs, Linda, Ada, Cilk	OpenMP, PLINQ, HPF	Scala, Erlang	Lisp, Clojure, Haskell, Scala, Erlang	-
Distributed Memory System	Socket communication, MPI, PVM, JXTA, MapReduce, CSP channels				-
Hybrid System	-	OpenCL	-	-	Unified Parallel C, Titanium, Fortress, X10, Chapel

MapReduce

3

- Programming model + associated implementation
- Processing and generating large data sets
- **Map:**
 - key/value pair → intermediate key/value pairs
- **Reduce:**
 - merge all intermediate values associated with the same intermediate key
- Origin: Lisp

Run-time system

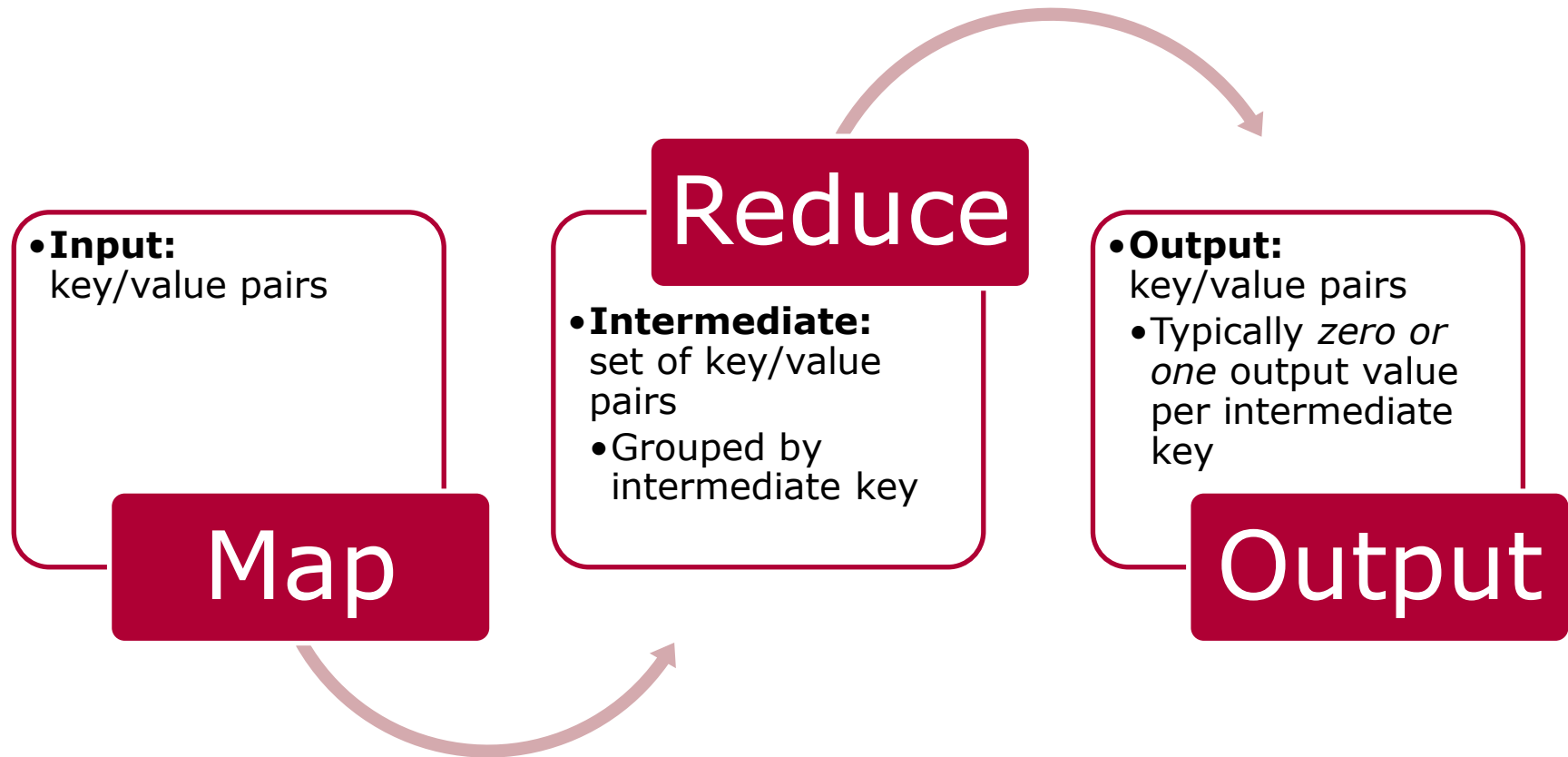
4

Automated parallelization and distribution

- Partitioning the input data
 - Scheduling the program's execution across a set of machines
 - Handling machine failures
 - Managing the required inter-machine communication
- Programmers do not have to think about parallel and distributed system specifics

Programming Model

5



Example: Wordcount

6

Counting the number of occurrences of each word in a large collection of documents:

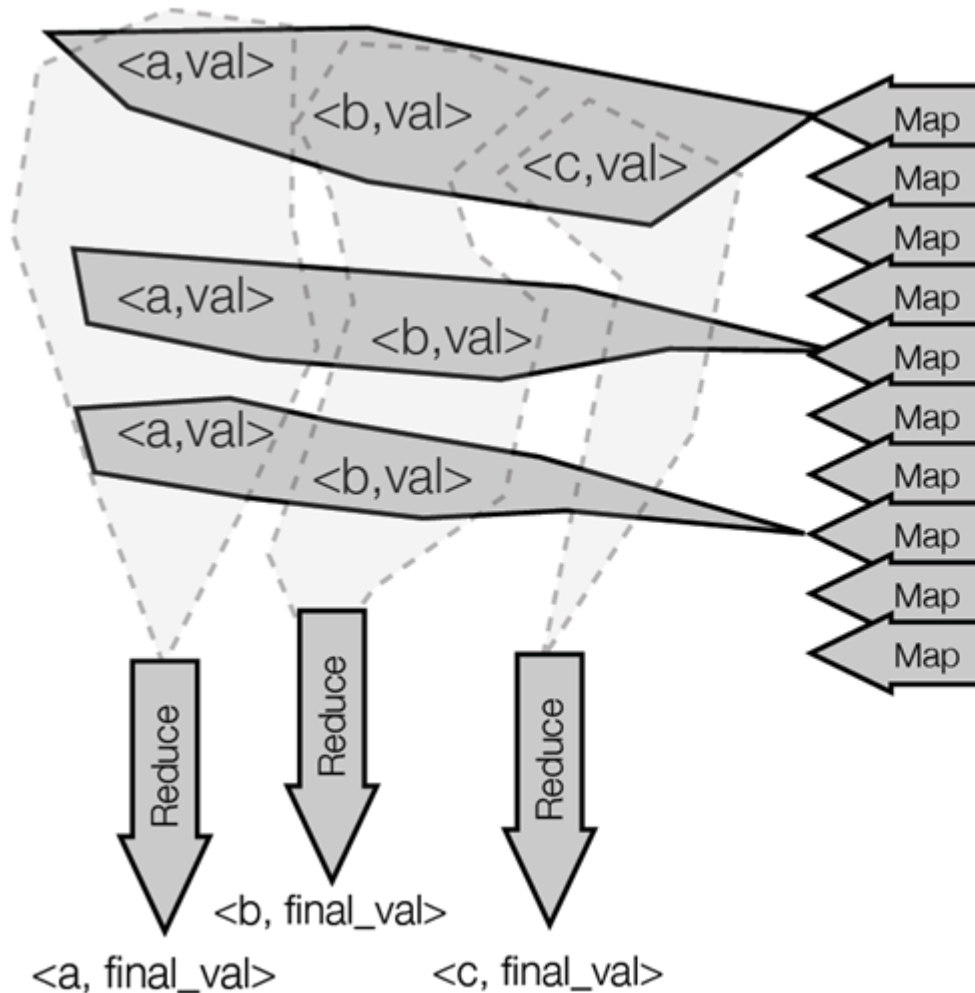
```
map(String key, String value):
  for each word w in value:
    EmitIntermediate(w, "1");
```

```
reduce(String key, Iterator values):
  int result = 0;
  for each v in values:
    result += ParseInt(v);
  Emit(AsString(result));
```

+ mapreduce specification object

Example: Wordcount

7



NAME	DATA
doc1	abcbcab
doc2	acbcabac
doc3	aaa
doc4	bcabab
doc5	aababaa
doc6	cabcbca
doc7	cbabcabc
doc8	aabbabab
doc9	bcabcbac
doc10	bcbac

Types

8

Type Specification

```
map      (k1, v1)          → list(k2, v2)
reduce  (k2, list(v2)) → list(v2)
```

Example

Work items: ID, binary content → character, number of occurrences

```
map      (long, byte[])      → list(char, int)
reduce  (char, list(int)) → list(int)
```


More Examples

9

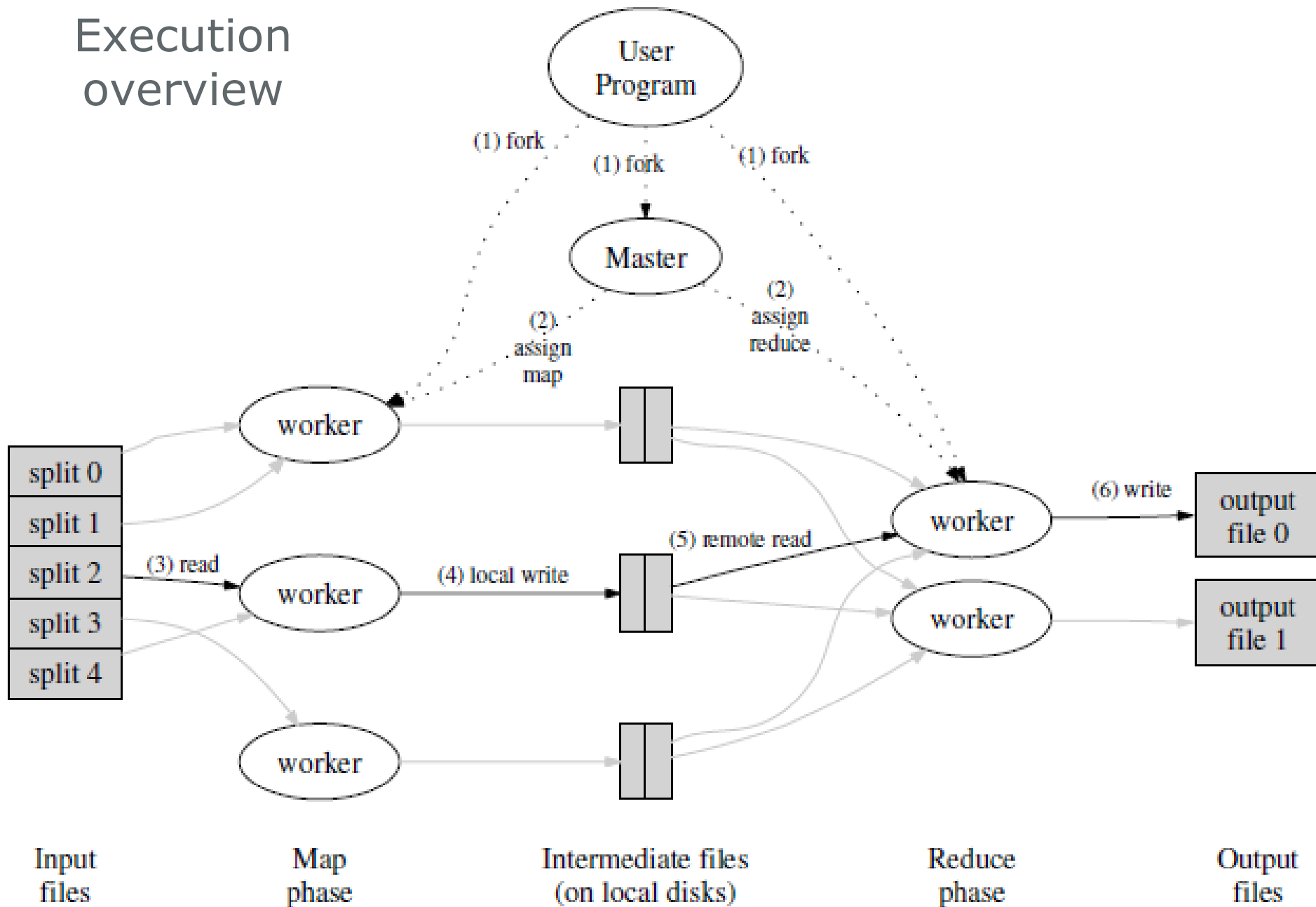
Example	Map	Reduce
Distributed Grep	Emits a line, if it matches the pattern	Emit unchanged
Count of URL access frequency	Processes logs of requests: <code><URL, 1></code>	Add values per URL: <code><URL, total count></code>
Reverse web-link graph	<code><target, source></code> , if link is found in source	<code><target, list(source)></code>
Term-vector per host (list of most important words)	<code><hostname, term vector></code> for each input document	Add all term vectors together: <code><hostname, term vector></code>
Inverted index	Parse document, emit <code><word, document ID></code>	Sort and emit <code><word, list(document ID)></code>
Distributed sort	Extract keys from records: <code><key, record></code>	Emit unchanged (done by ordering properties)

Google Implementation of MapReduce

10

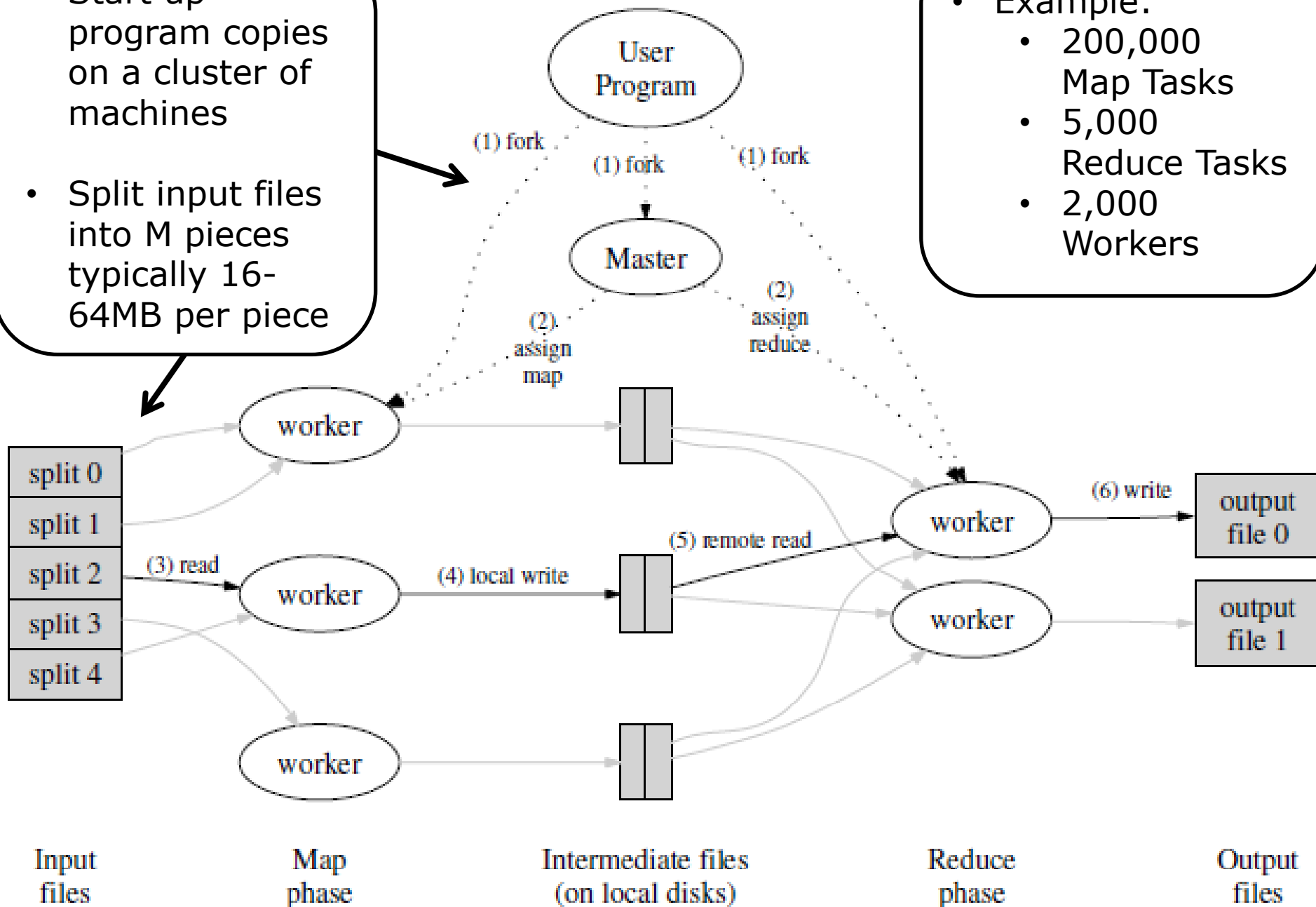
- Large cluster of standard PCs with local disks
 - x86, Ethernet: 100 Mbit/s to 1 Gbit/s, 2-4GB RAM, IDE
 - Custom global file system
 - ◇ Replication for availability and reliability
 - Job scheduling system
 - ◇ Set of tasks to set of machines
 - Machine failures are common (large number of machines)

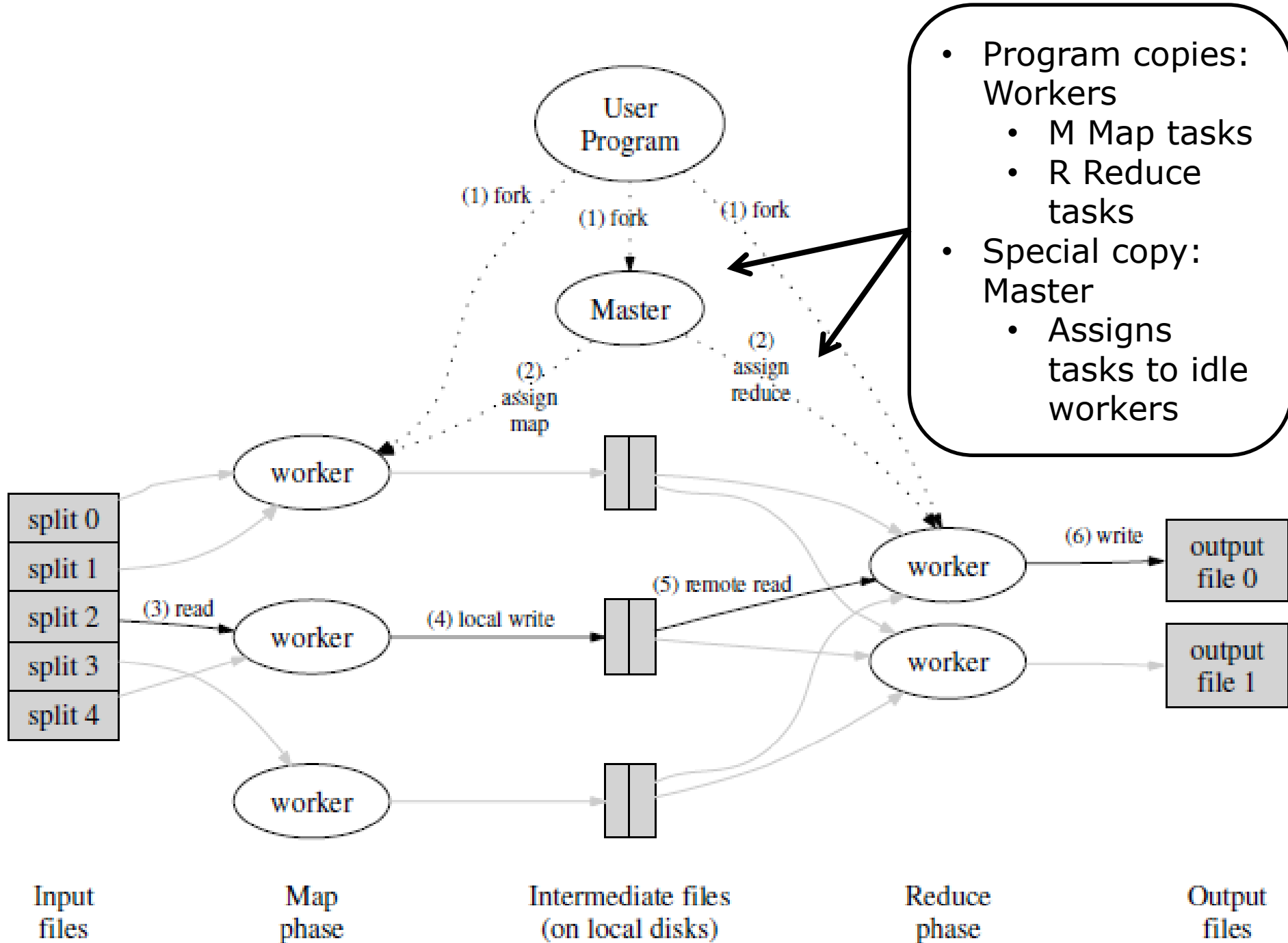
Execution overview



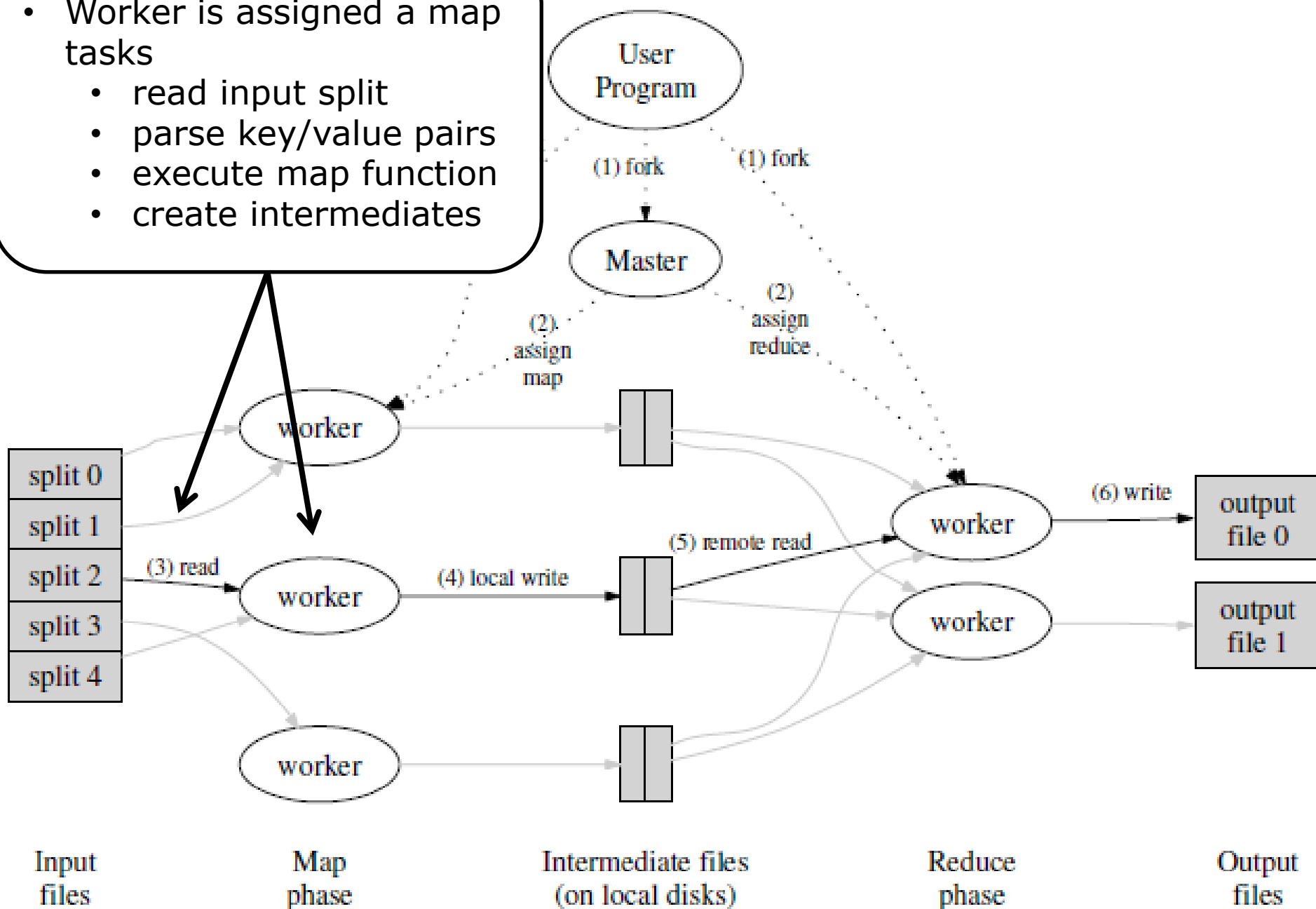
- Start up program copies on a cluster of machines
- Split input files into M pieces typically 16-64MB per piece

- Example:
 - 200,000 Map Tasks
 - 5,000 Reduce Tasks
 - 2,000 Workers

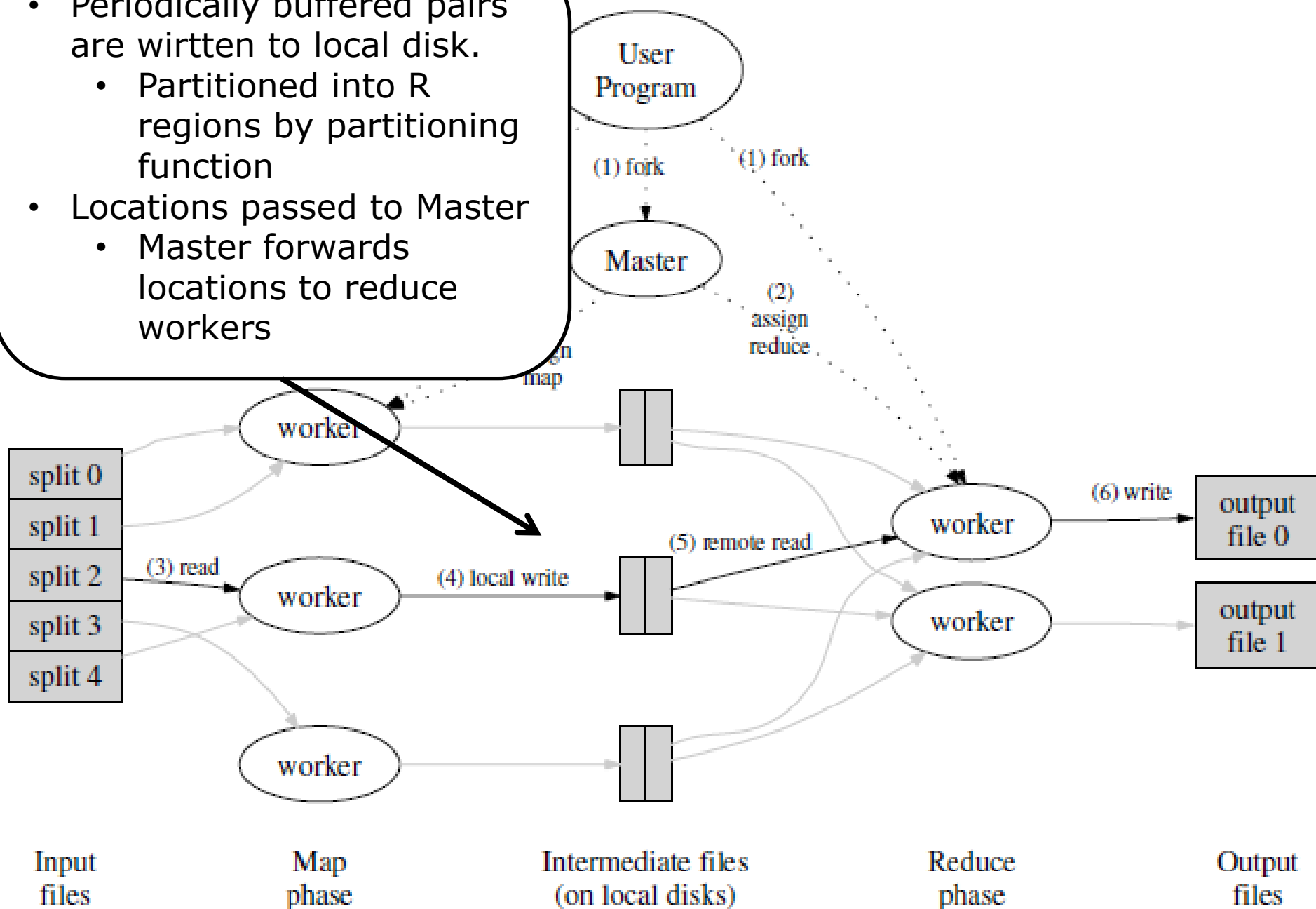


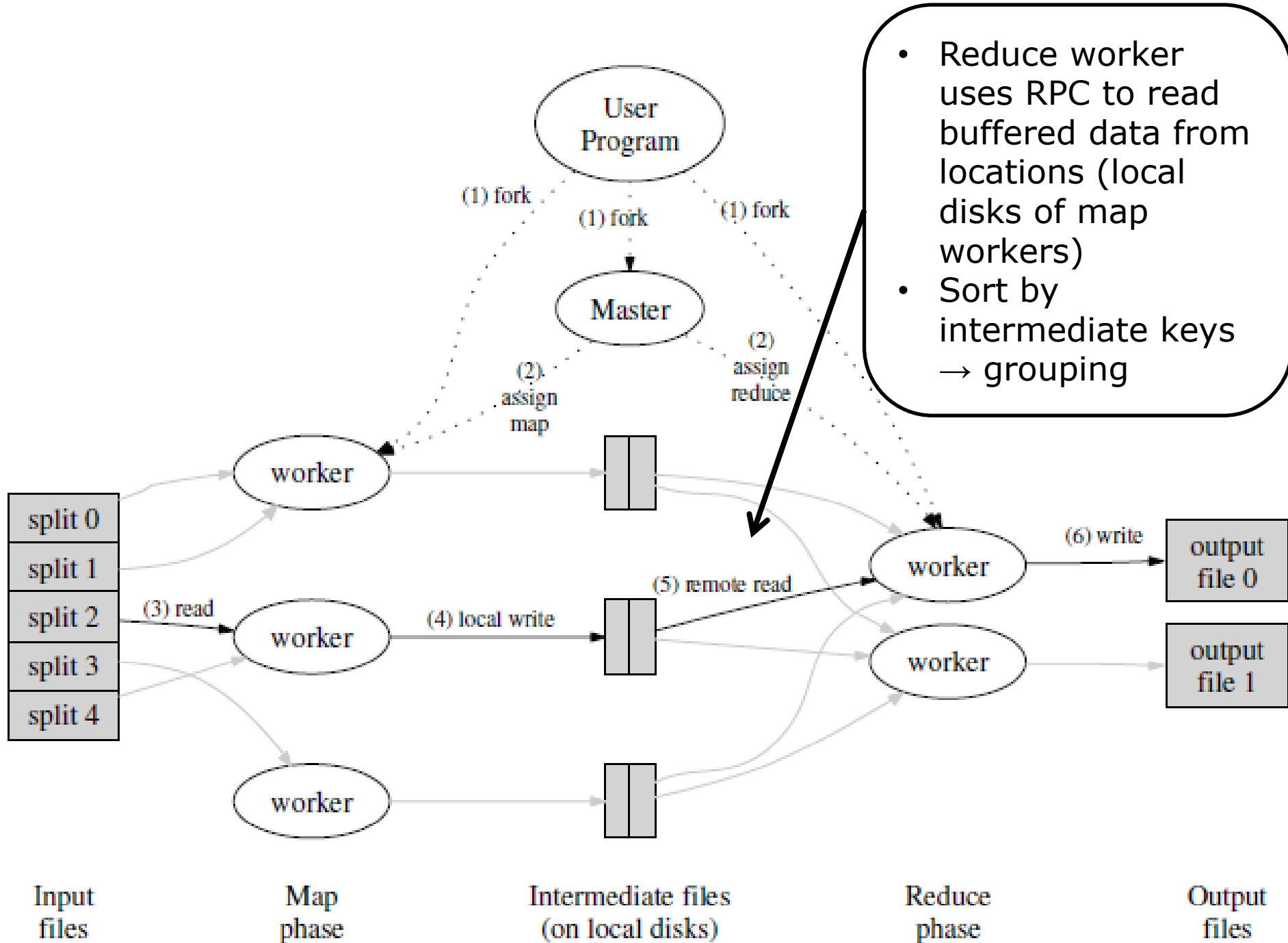


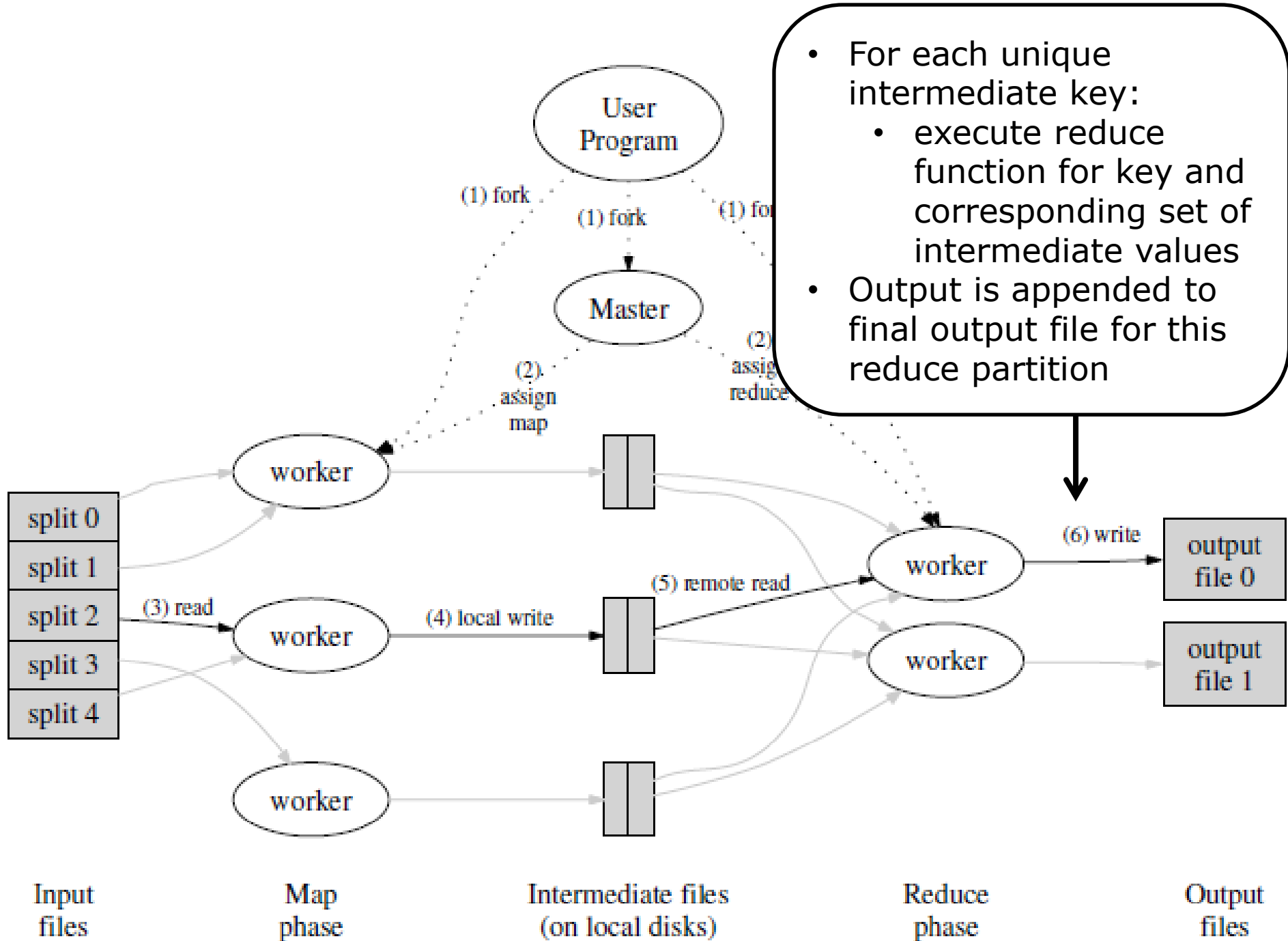
- Worker is assigned a map tasks
 - read input split
 - parse key/value pairs
 - execute map function
 - create intermediates

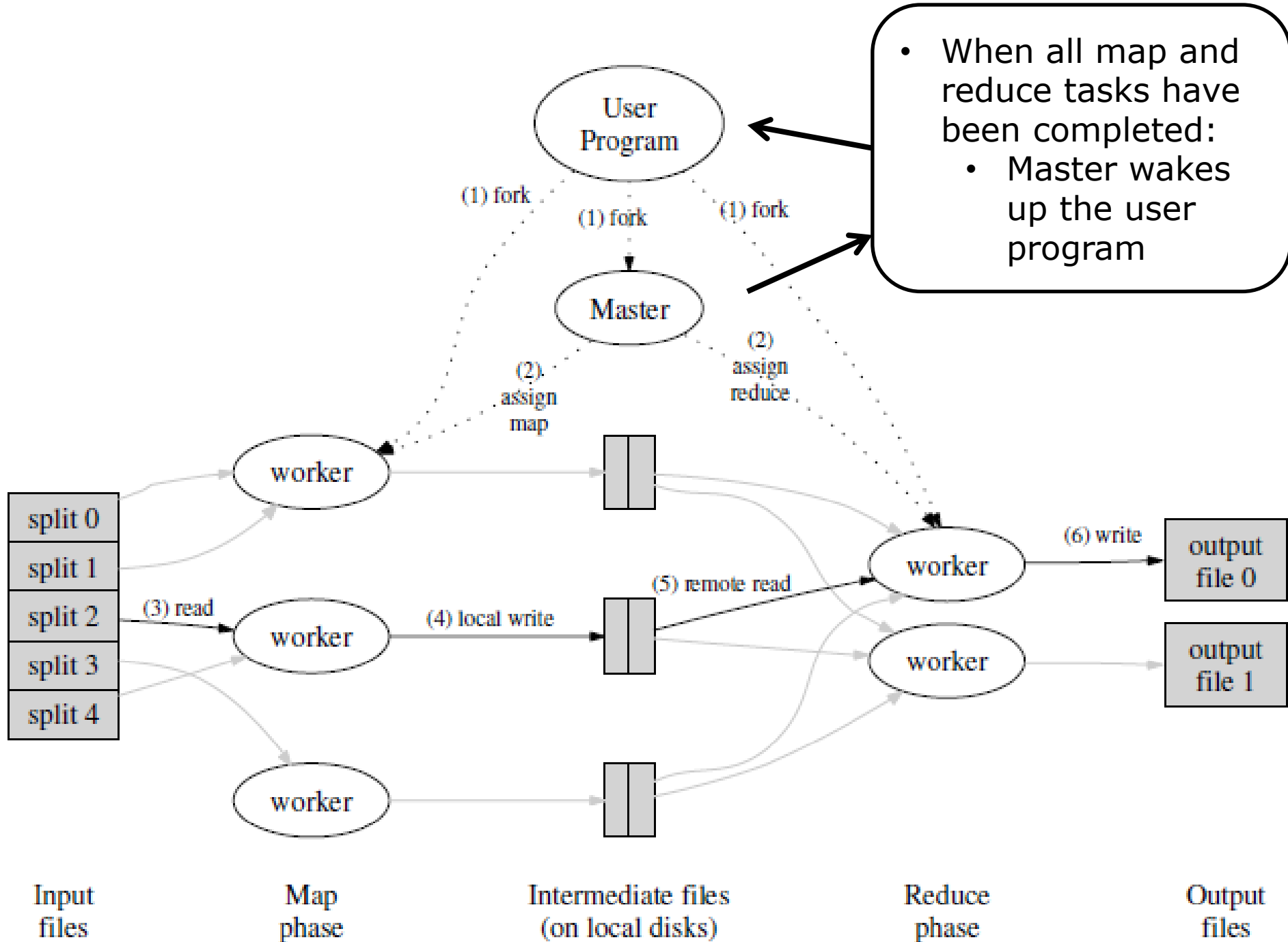


- Periodically buffered pairs are written to local disk.
 - Partitioned into R regions by partitioning function
- Locations passed to Master
 - Master forwards locations to reduce workers









Specific Google properties

19

- Network bottleneck in Google cluster
 - Master tries to use locality information about the input data, which is stored in the distributed file system
 - For large MapReduce tasks, most input data is read locally
- **Fault tolerance**
 - Periodic heartbeat between master and workers
 - For a failed worker, re-execute completed and in-progress map tasks
 - For a failed master, MapReduce is aborted → user has to re-execute
 - Span backup tasks (cloned workers, same task) when MapReduce is close to completion, to compensate faulty (delaying) workers

Refinements

20

Refinement	Description
Partitioning Function	User functions for data partitioning are possible (<code>hash(key) mod R</code> is default)
Ordering Guarantees	Intermediate key/value pairs are ordered inc.
Combiner Function	Partial merging of local data (like reduce)
Input and Output Types	Some standard formats; user can specify more
Side-Effects	Additional files have to be addressed by the user
Skipping Bad Records	Ignore records with deterministic crashes (optional)
Local Execution	Special MapReduce library for sequential execution
Status Information	Master runs an internal HTTP server for diagnosis
Counters	Count occurrences of various events; user defined