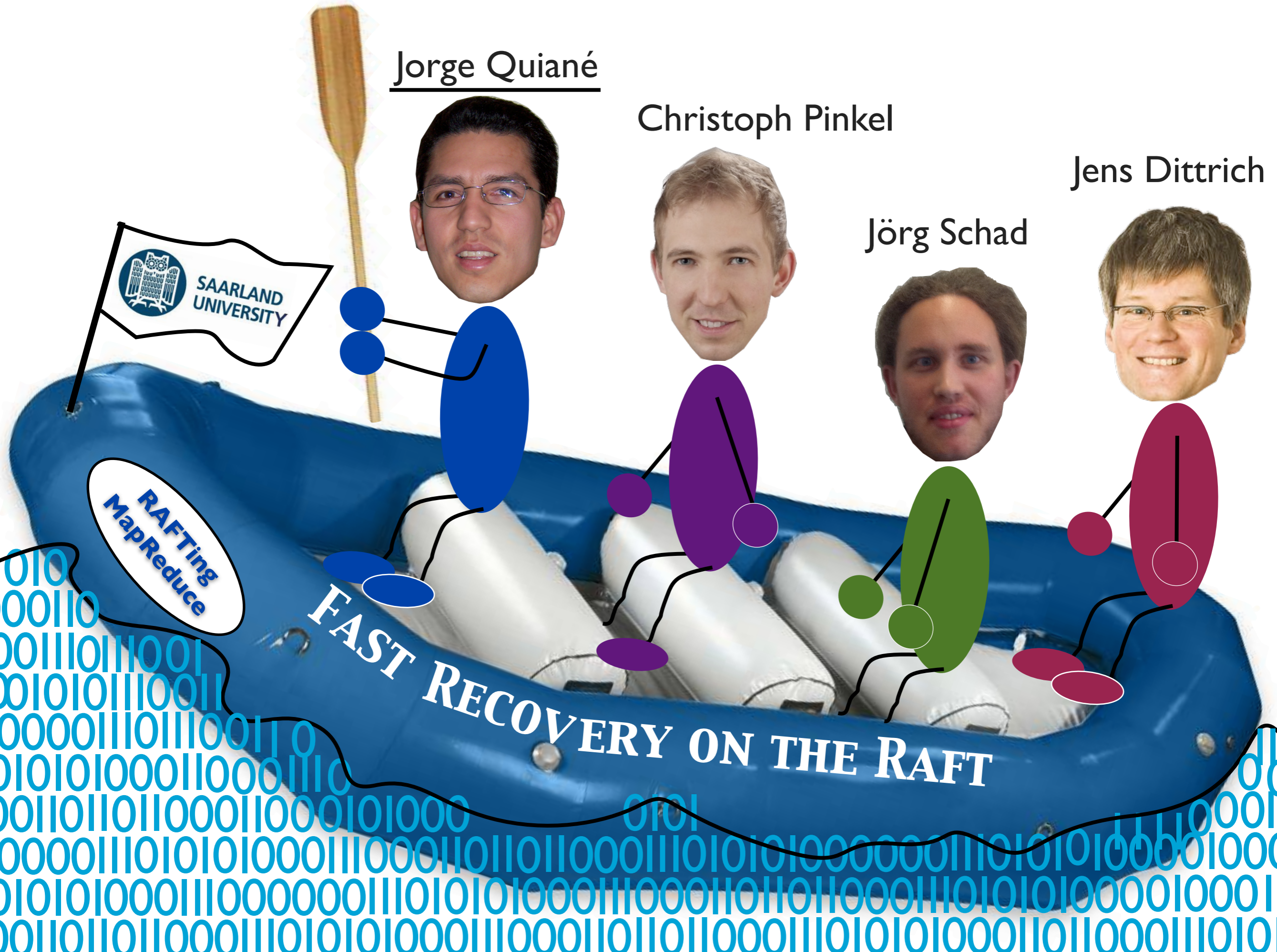


Jorge Quiané

Christoph Pinkel

Jens Dittrich

Jörg Schad

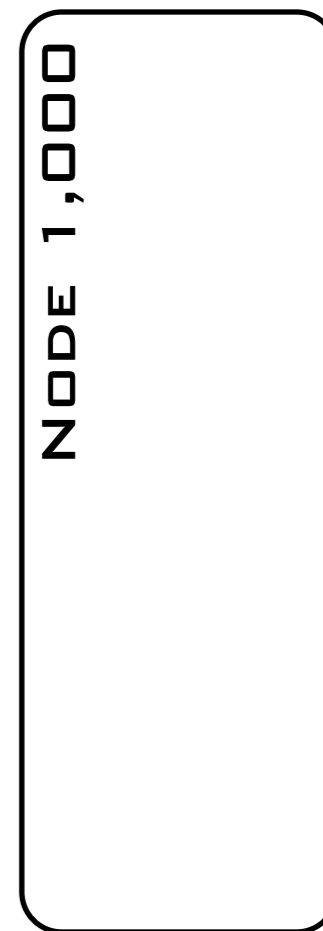
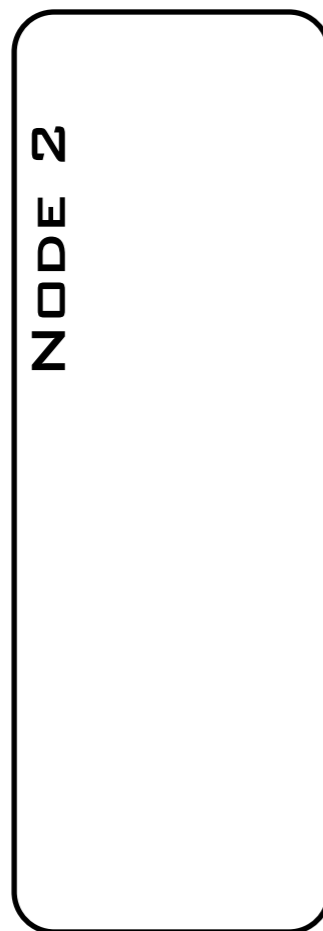
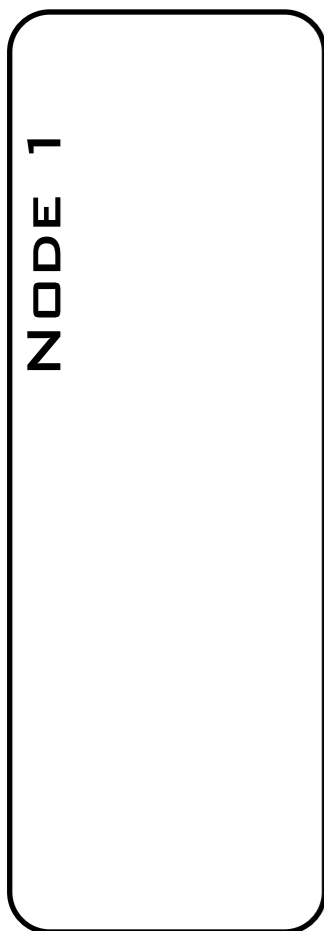
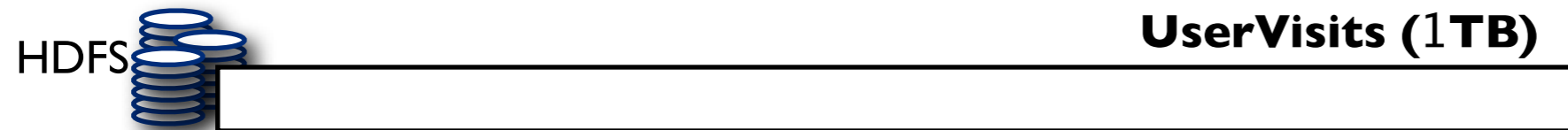


RAFTING  
MapReduce

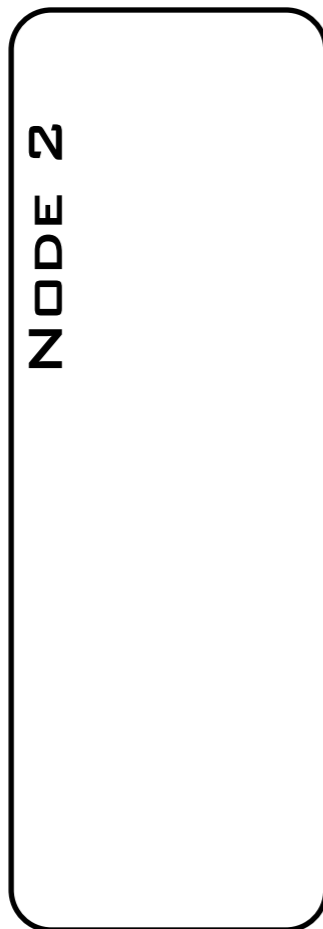
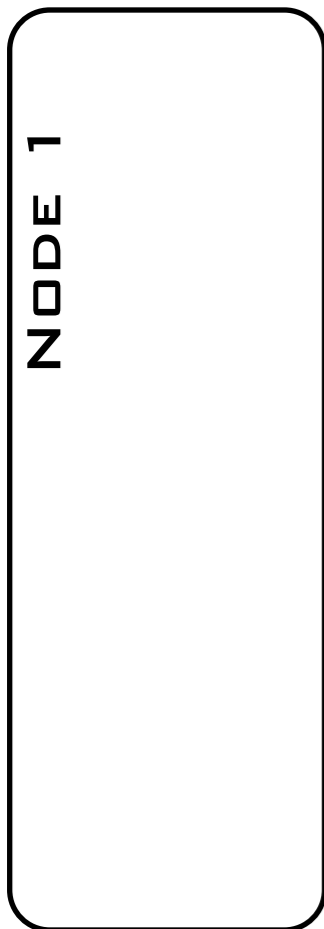
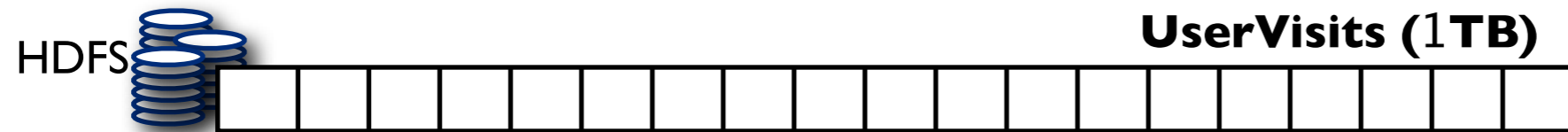
FAST RECOVERY ON THE RAFT

1010  
000110  
001110111001  
001010111001  
00000111011100110  
010101000110001110  
001101101100011000101000  
00000111010101000111000110110110001110101010000011101010100001000  
0101010001100000011101010100011000110110110001110101010000100011  
00110110110001110101010001100011011011000111010101000110110001110101

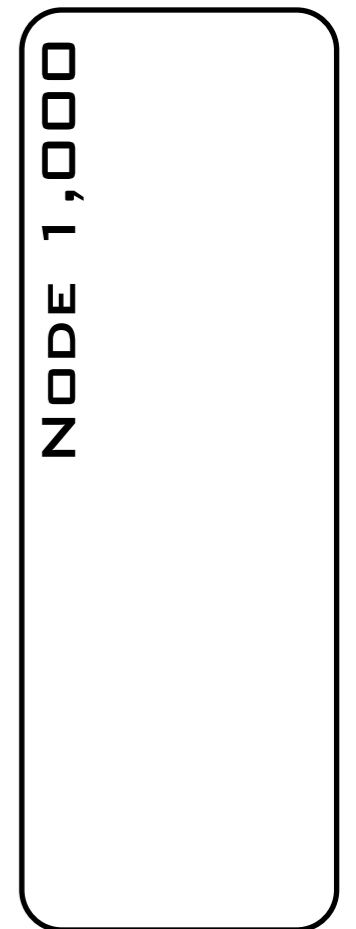
# MapReduce (Recall)



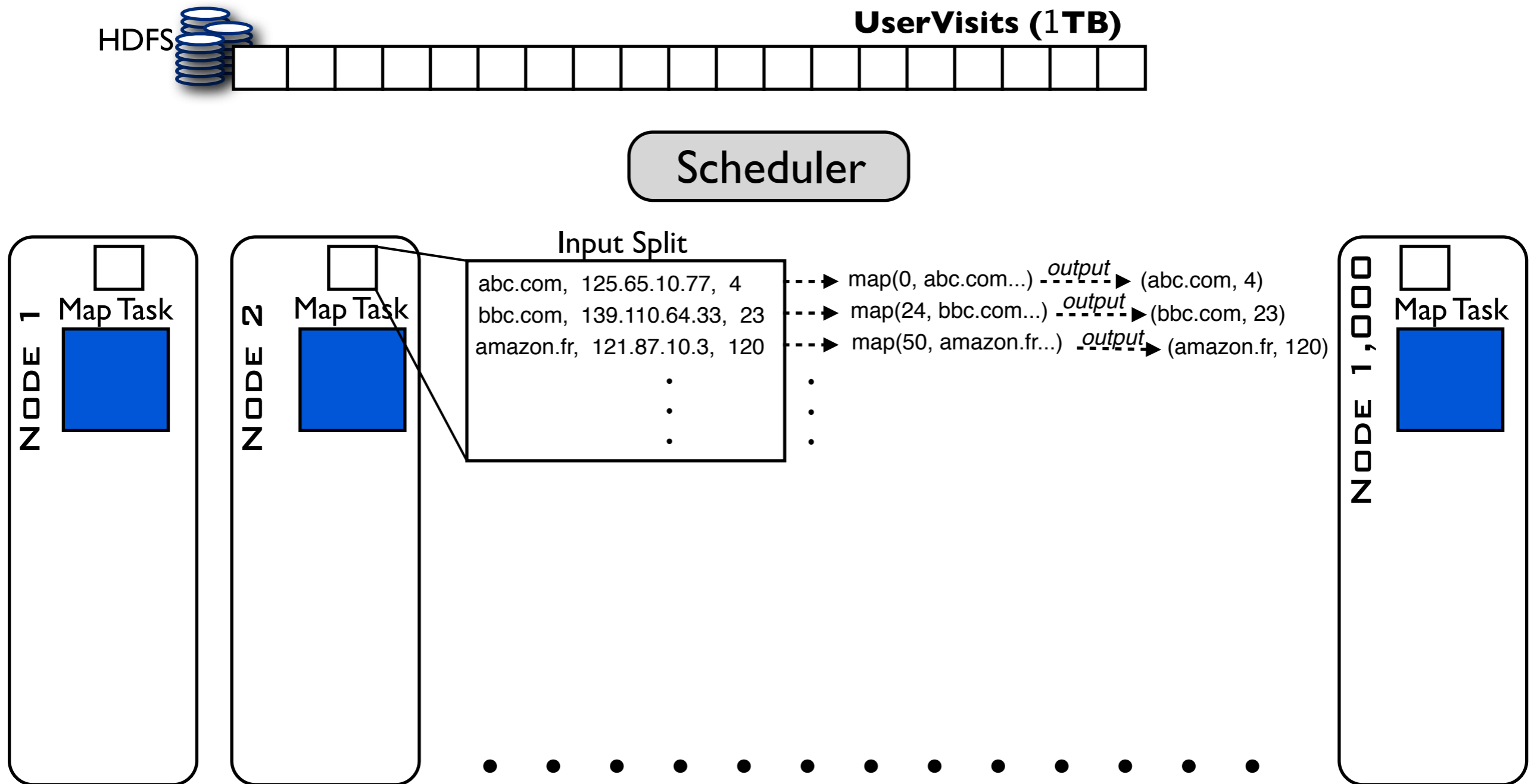
# MapReduce (Recall)



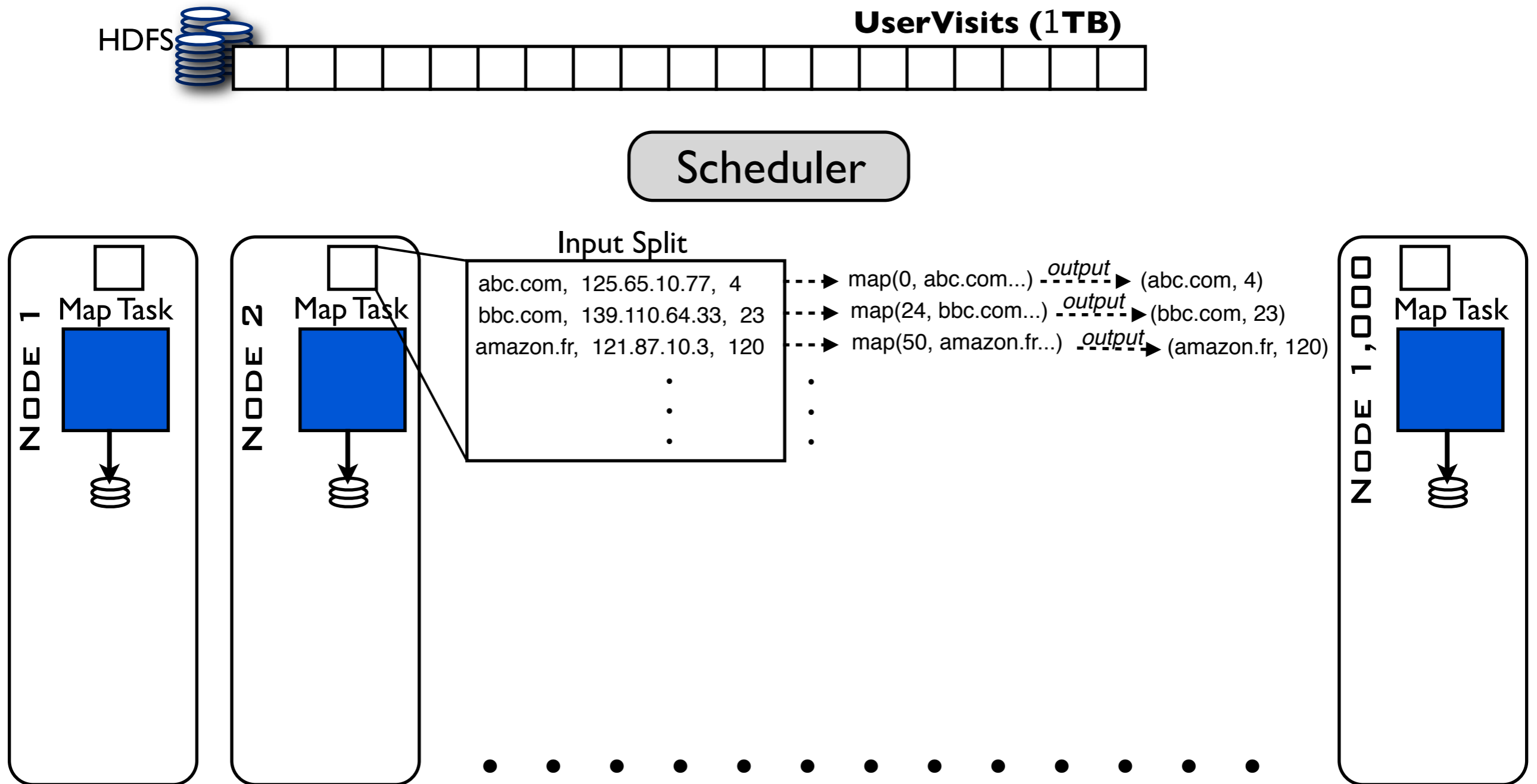
• • • • • • • • • • • • • •



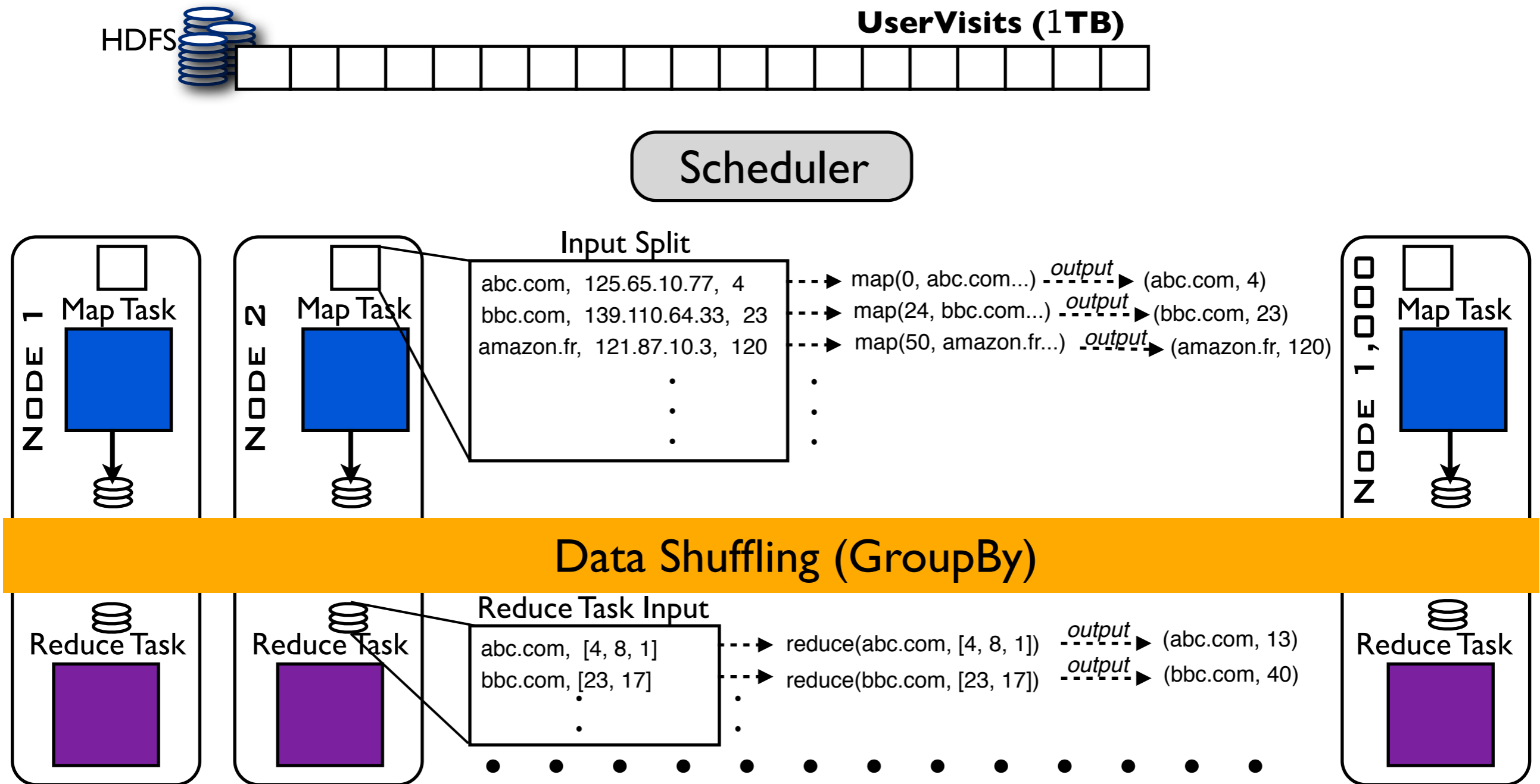
# MapReduce (Recall)



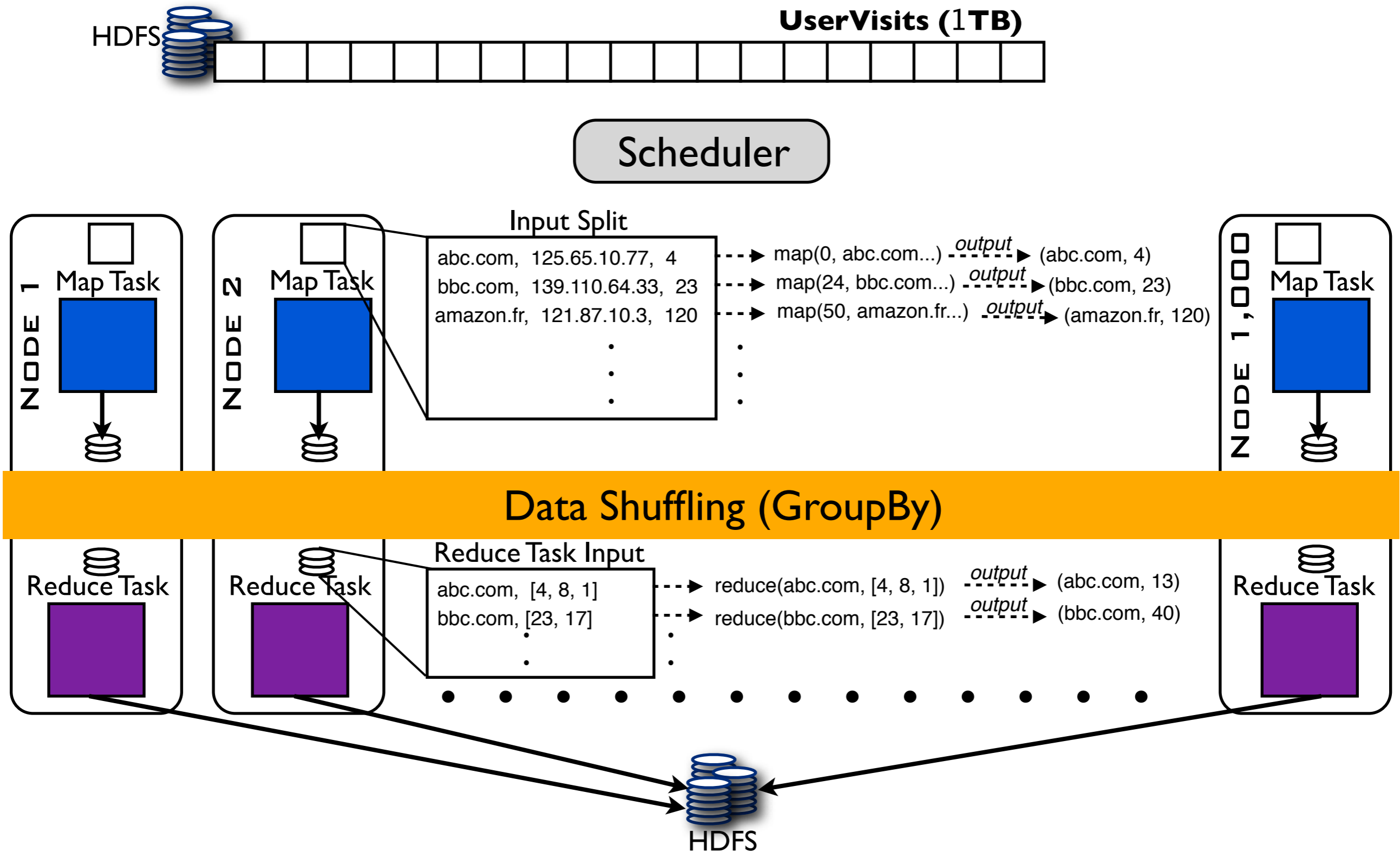
# MapReduce (Recall)



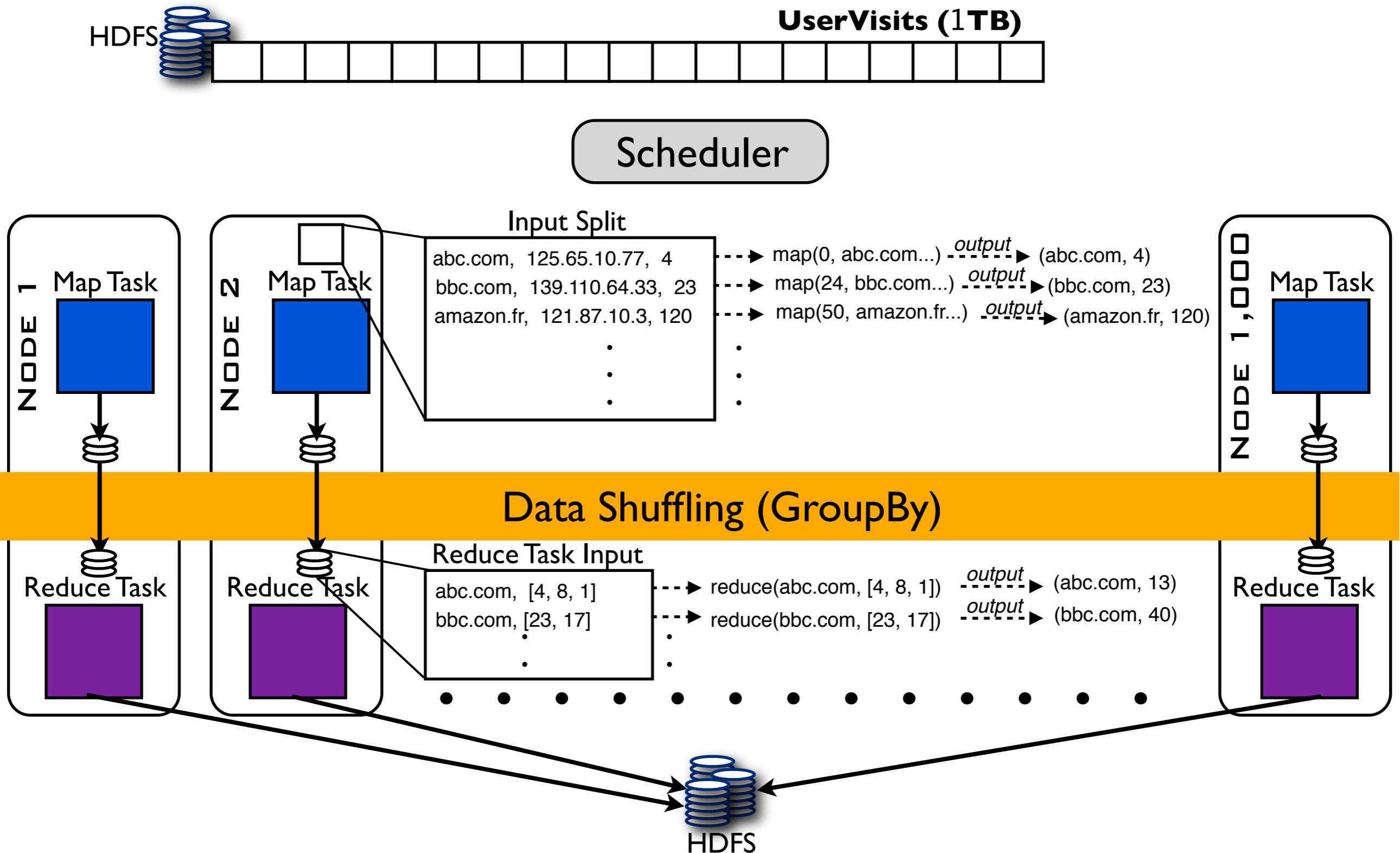
# MapReduce (Recall)



# MapReduce (Recall)

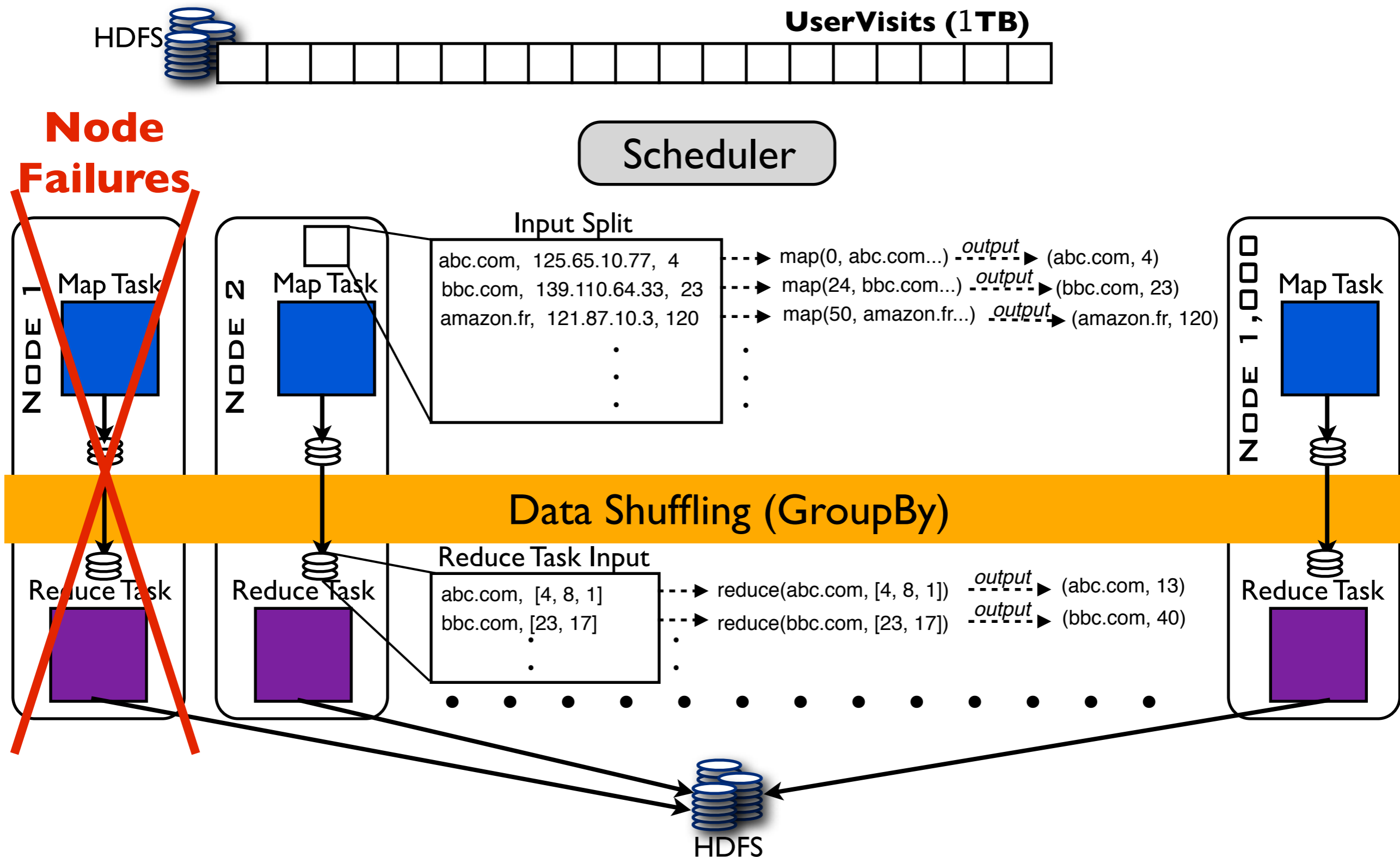


# Failures are the Rule!

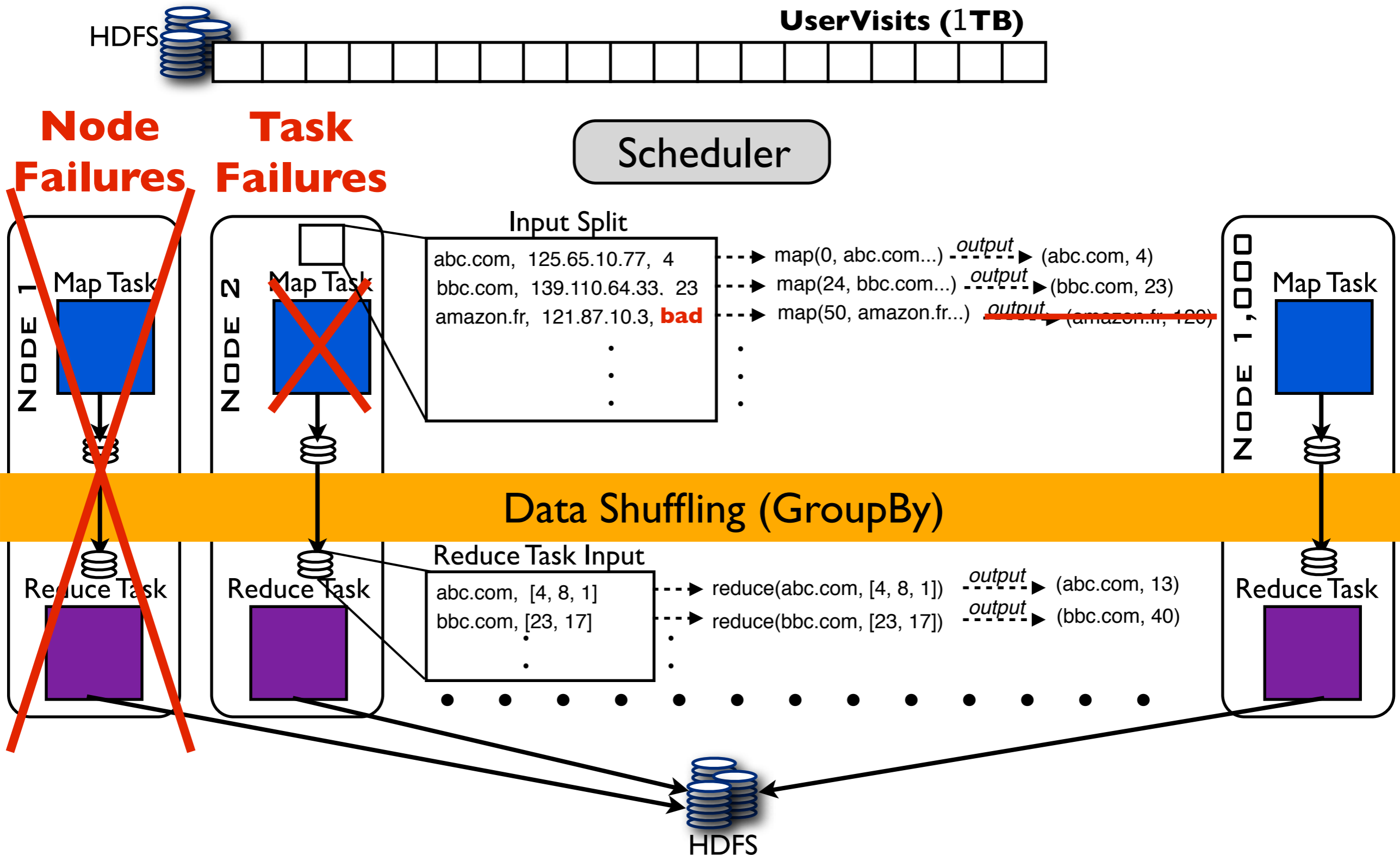




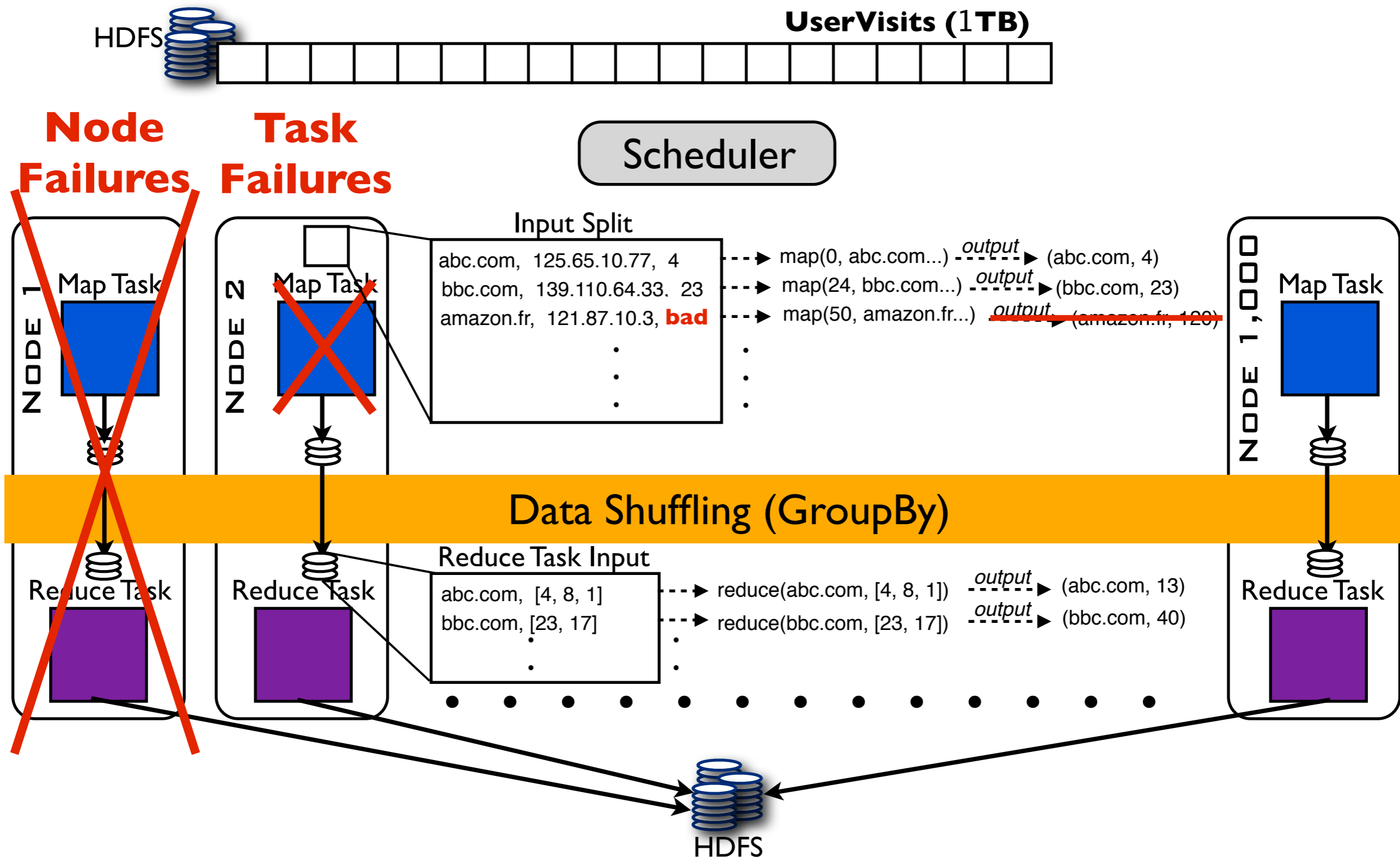
# Failures are the Rule!



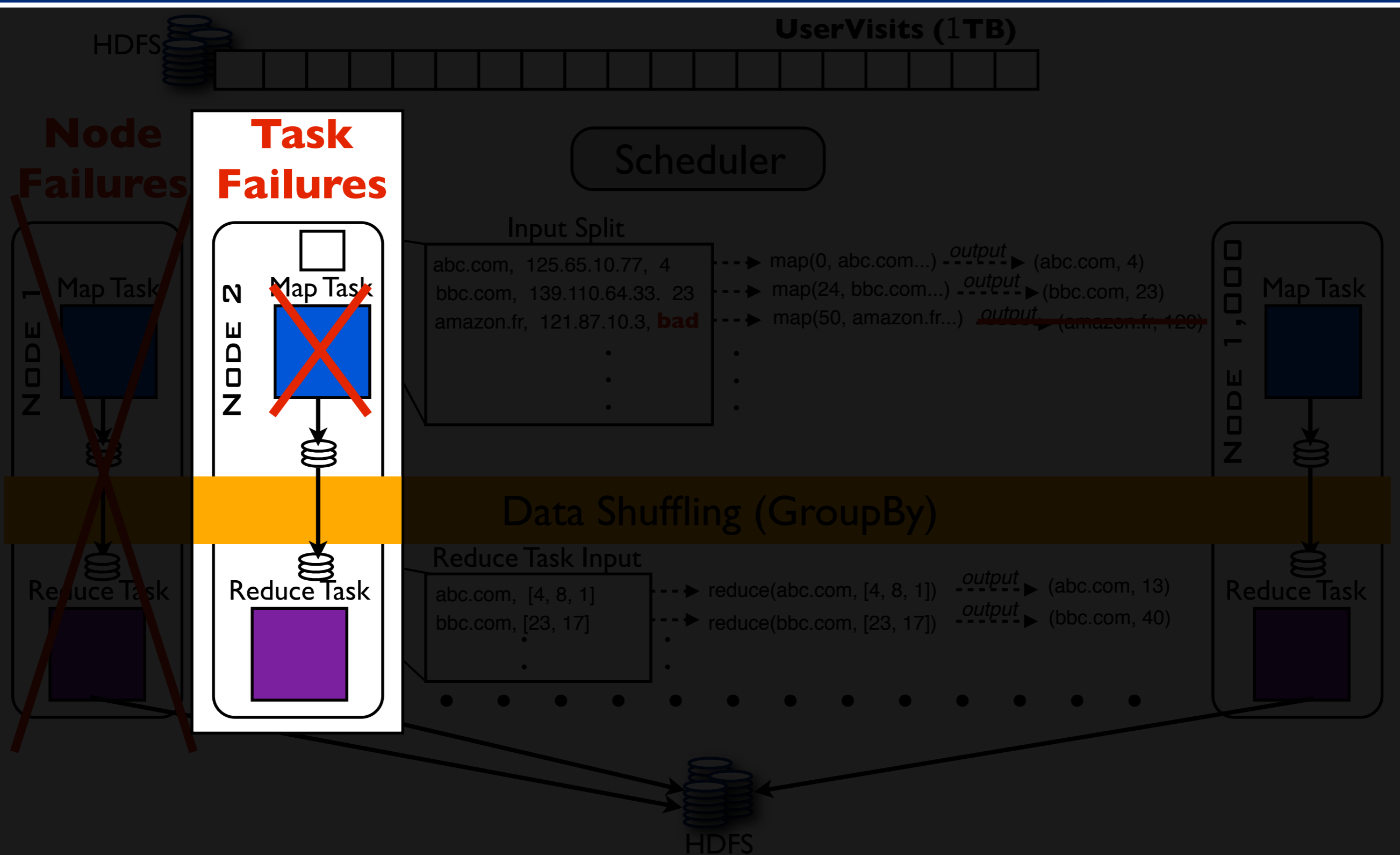
# Failures are the Rule!



# Task Failures



# Task Failures



# Current Approach

## Hadoop without Failures

Scheduler

Map Task: 20 seconds

Time: 0s

NODE 1

NODE 2

## Hadoop with Task Failures

Scheduler

Time: 0s

NODE 1

NODE 2

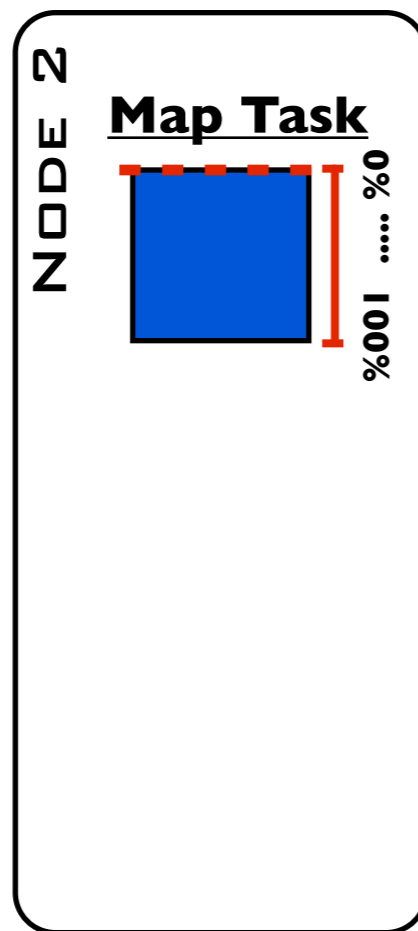
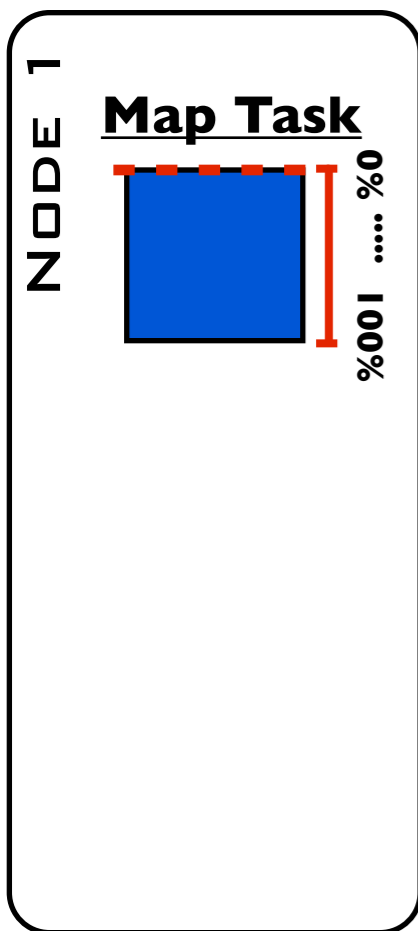
# Current Approach

## Hadoop without Failures

Map Task: 20 seconds

Scheduler

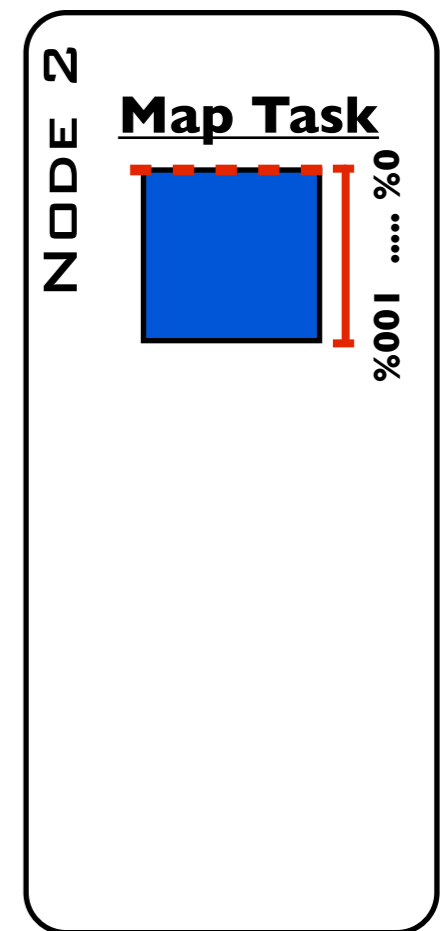
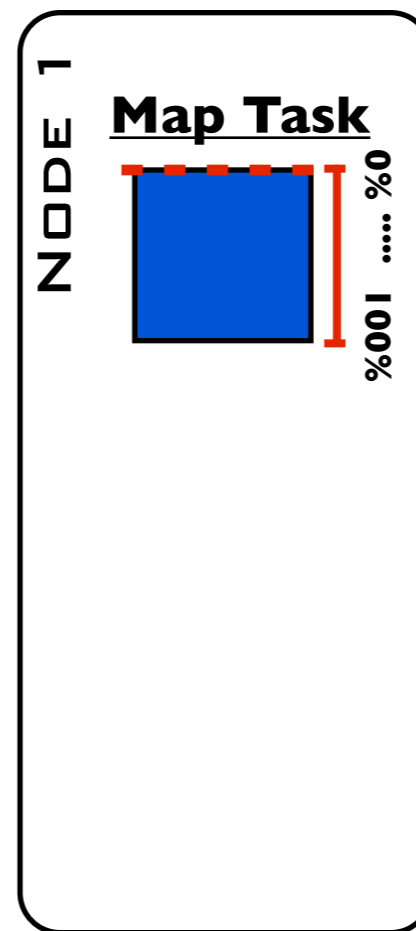
Time: 0s



## Hadoop with Task Failures

Scheduler

Time: 0s



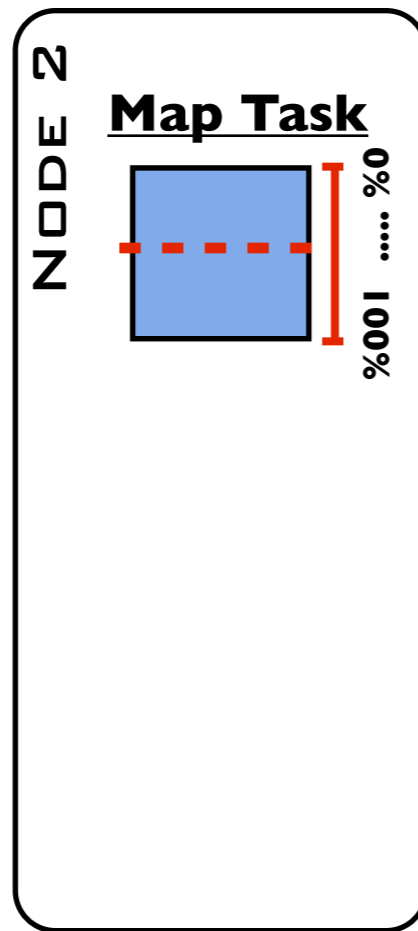
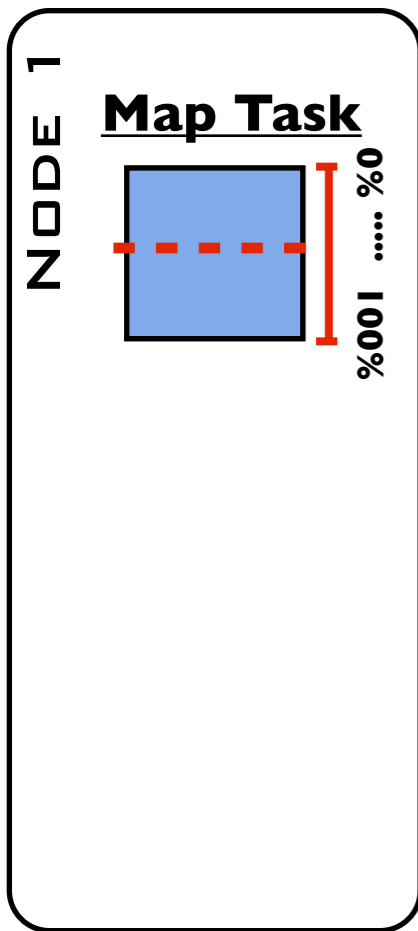
# Current Approach

## Hadoop without Failures

Map Task: 20 seconds

Scheduler

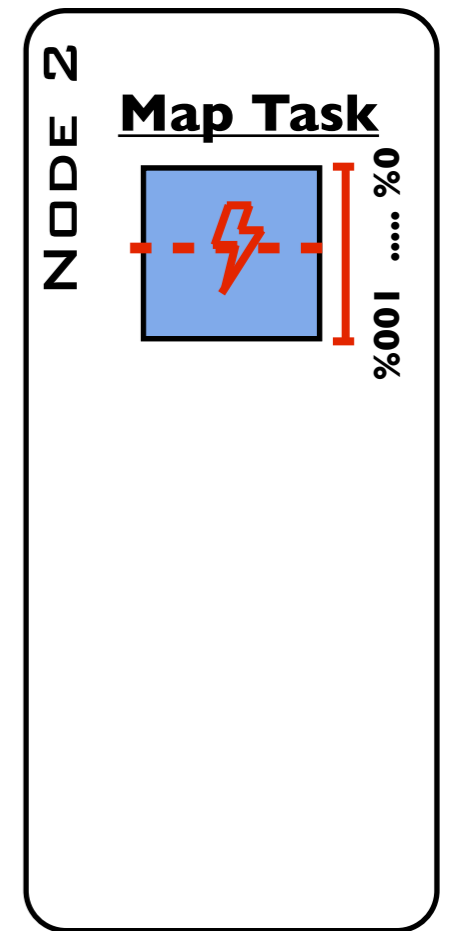
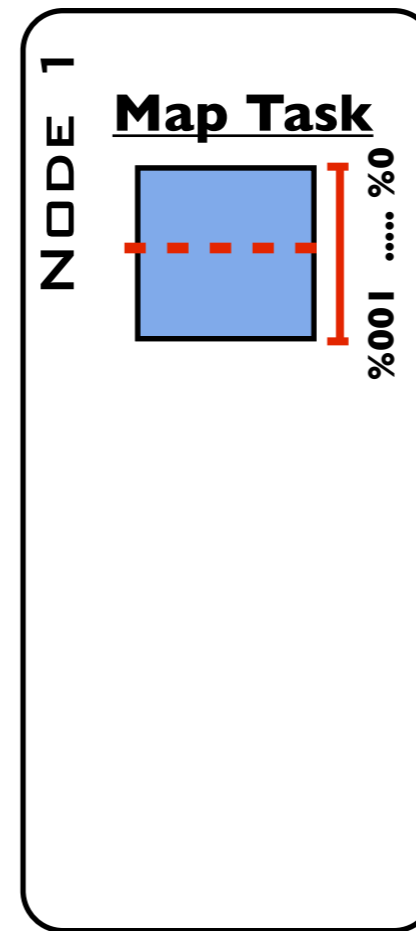
Time: 10s



## Hadoop with Task Failures

Scheduler

Time: 10s



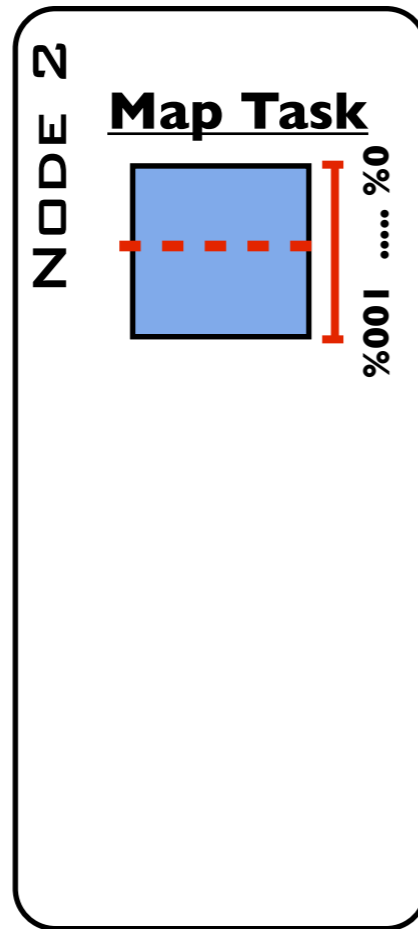
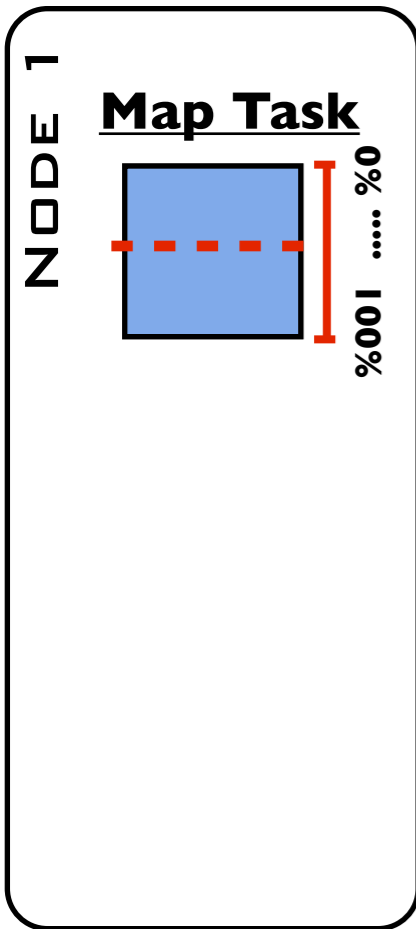
# Current Approach

## Hadoop without Failures

Map Task: 20 seconds

Scheduler

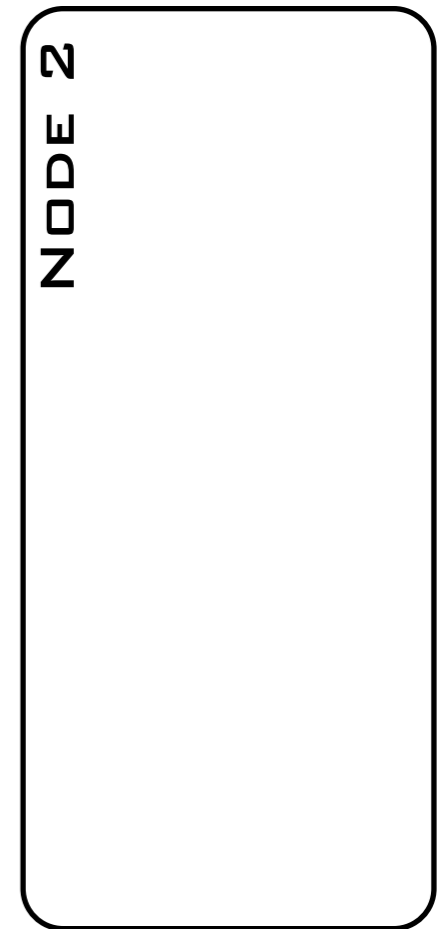
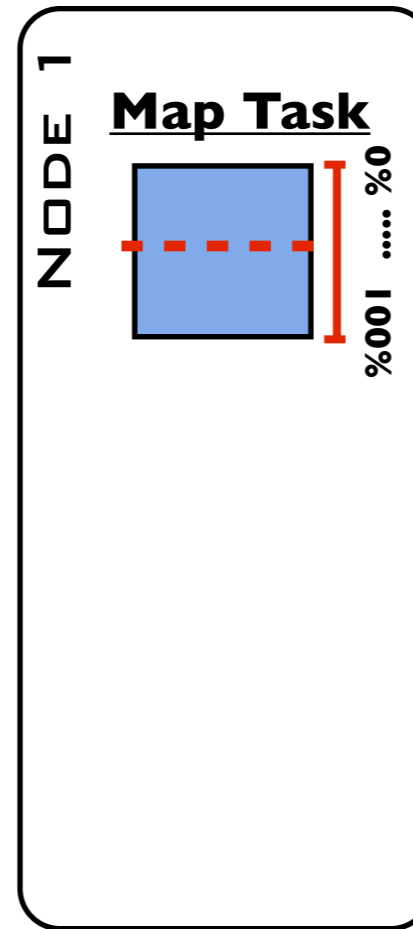
Time: 10s



## Hadoop with Task Failures

Scheduler

Time: 10s





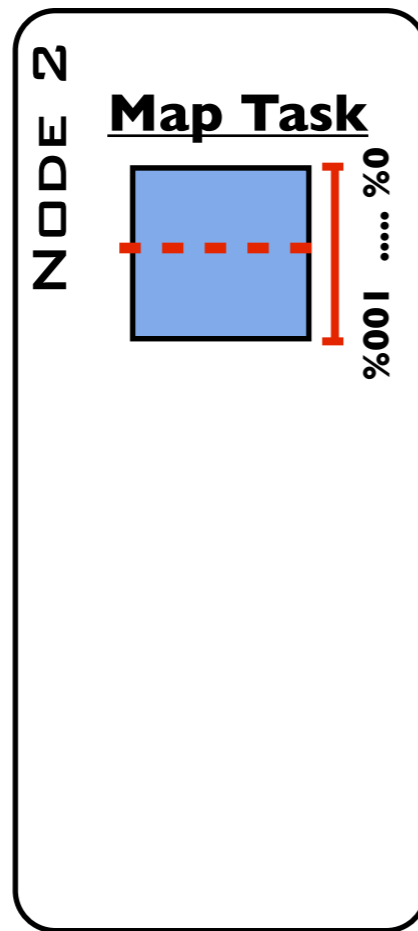
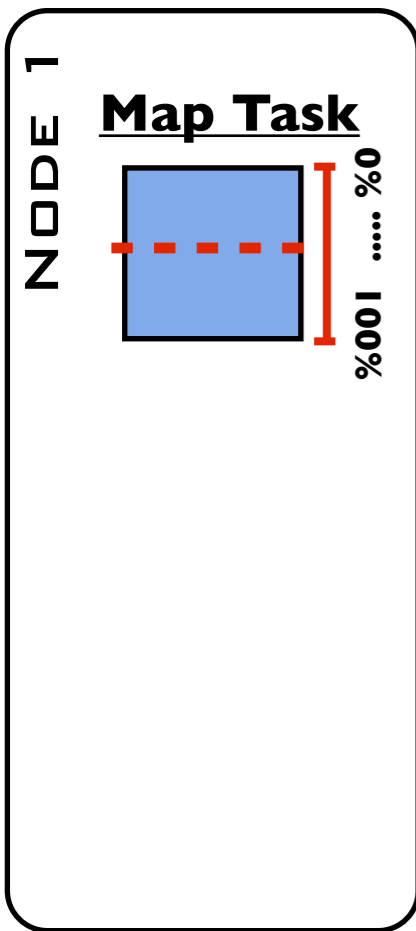
# Current Approach

## Hadoop without Failures

Map Task: 20 seconds

Scheduler

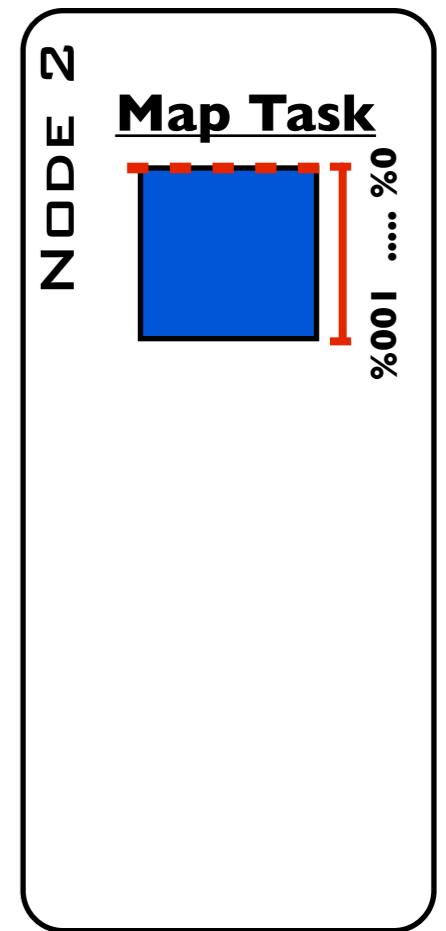
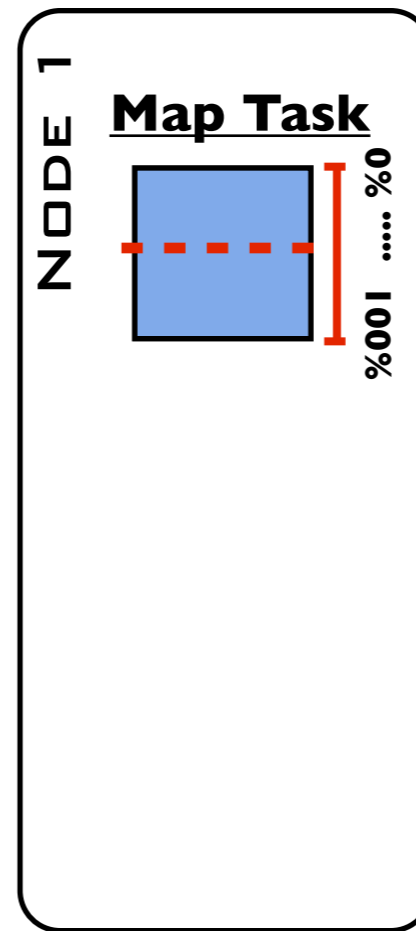
Time: 10s



## Hadoop with Task Failures

Scheduler

Time: 10s



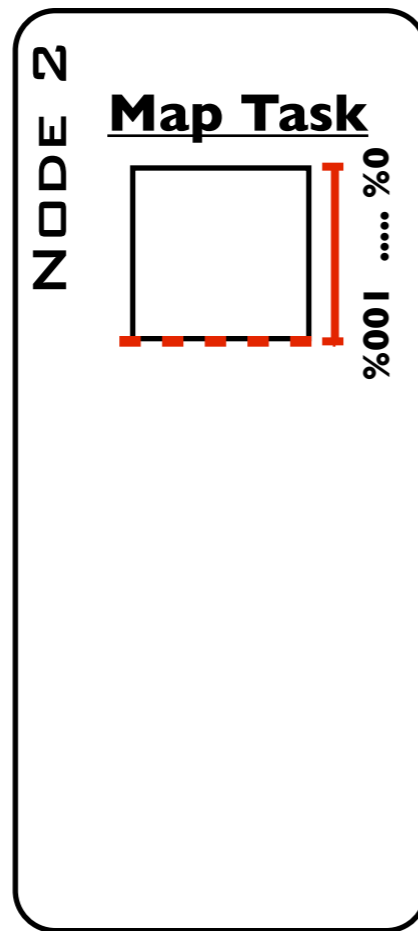
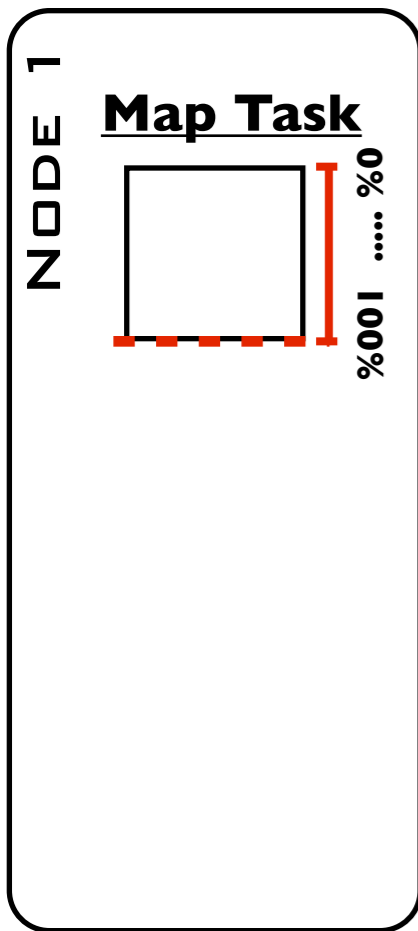
# Current Approach

## Hadoop without Failures

Map Task: 20 seconds

Scheduler

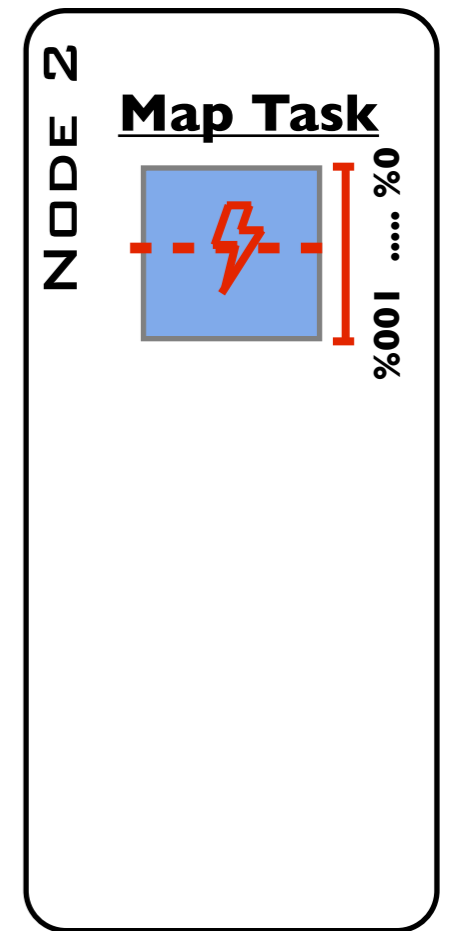
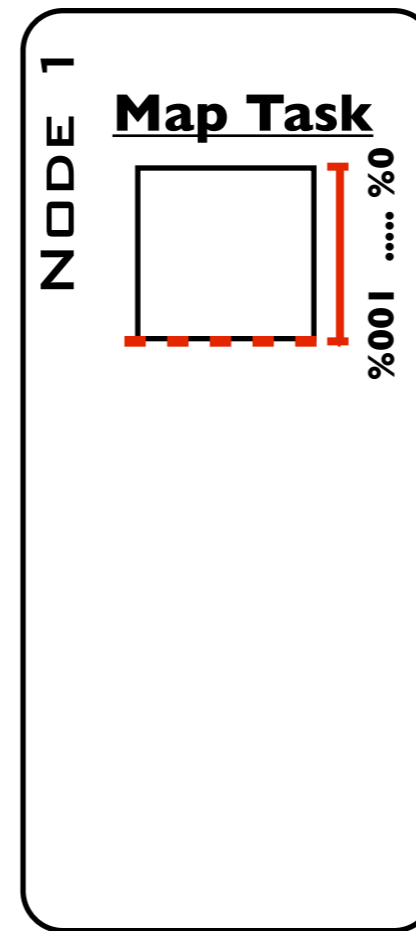
Time:20s



## Hadoop with Task Failures

Scheduler

Time:20s



# Current Approach

## Hadoop without Failures

Map Task: 20 seconds

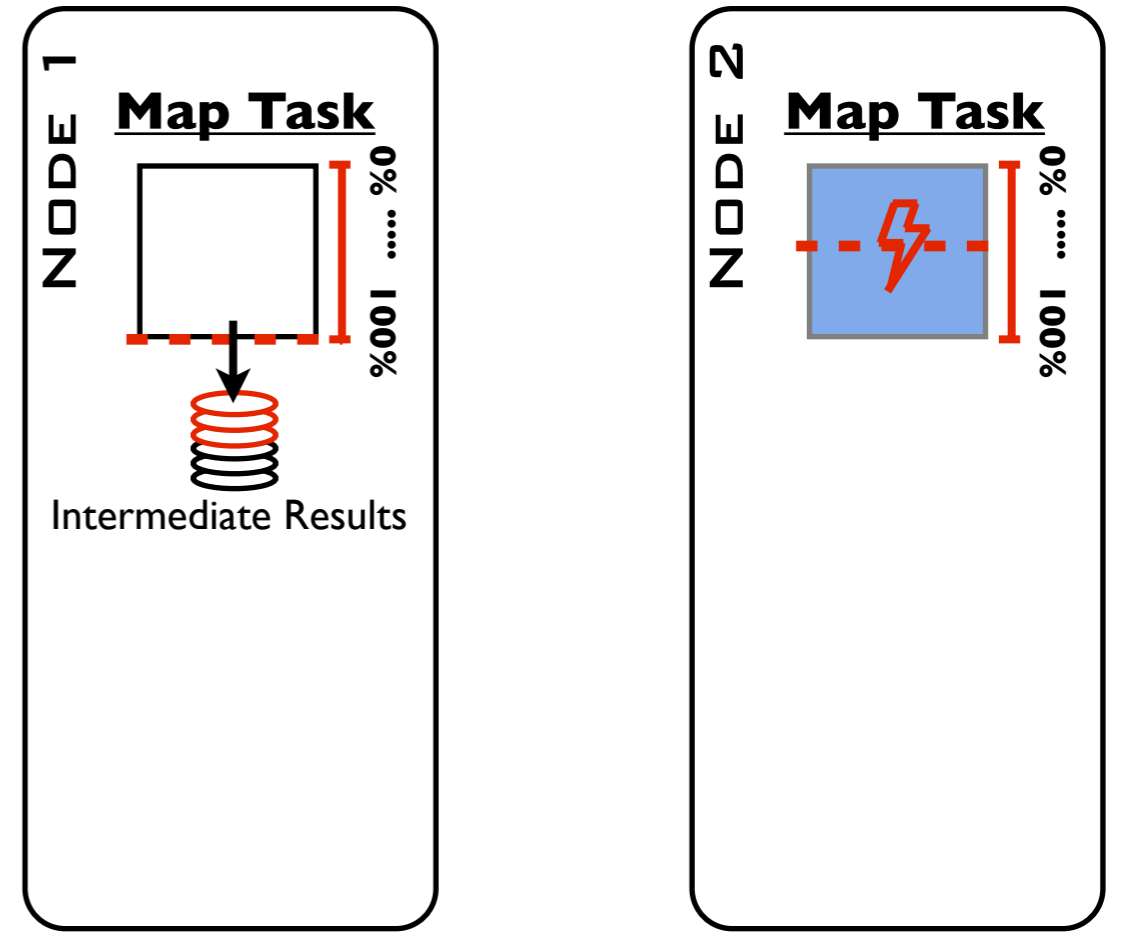
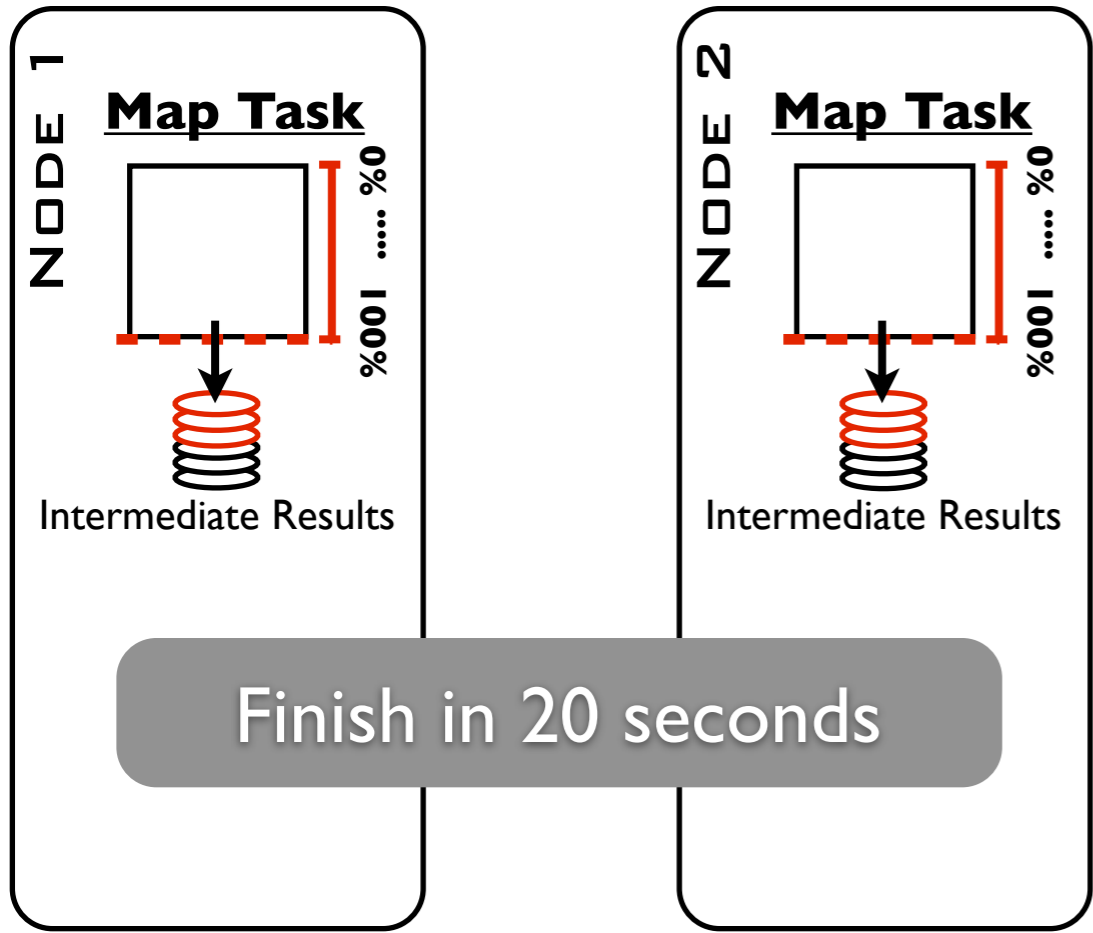
## Hadoop with Task Failures

Scheduler

Time:20s

Time:20s

Scheduler



# Current Approach

## Hadoop without Failures

Map Task: 20 seconds

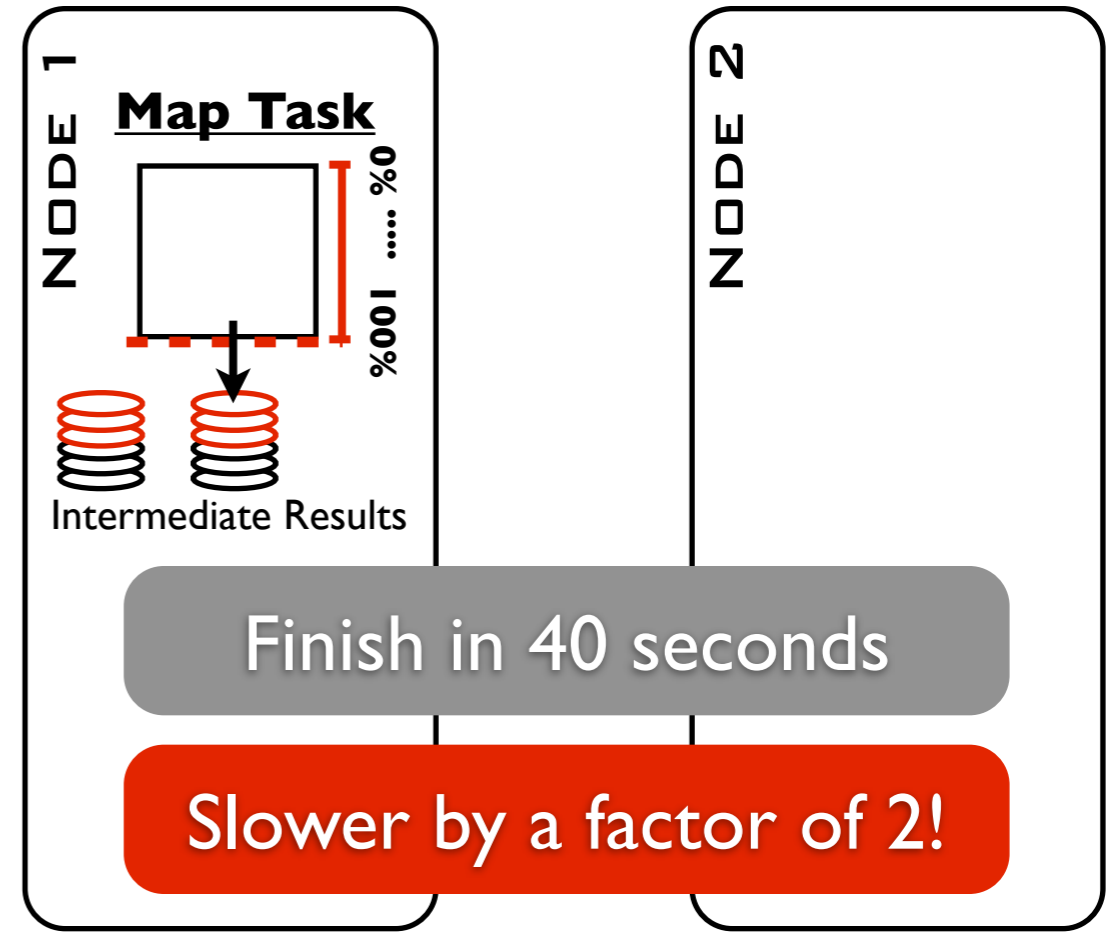
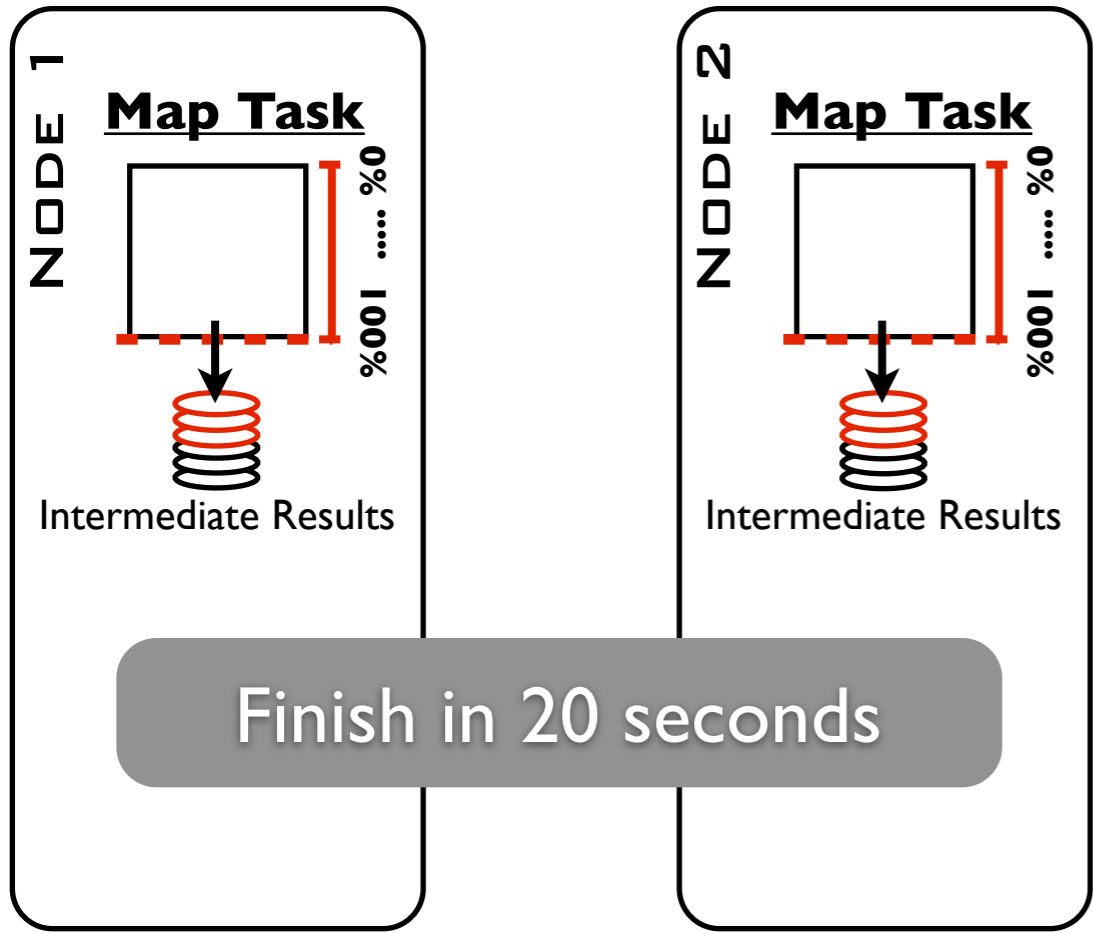
## Hadoop with Task Failures

Scheduler

Time:20s

Time:40s

Scheduler



# Solution: Local Checkpointing

Hadoop  
without Failures

Map Task: 20 seconds

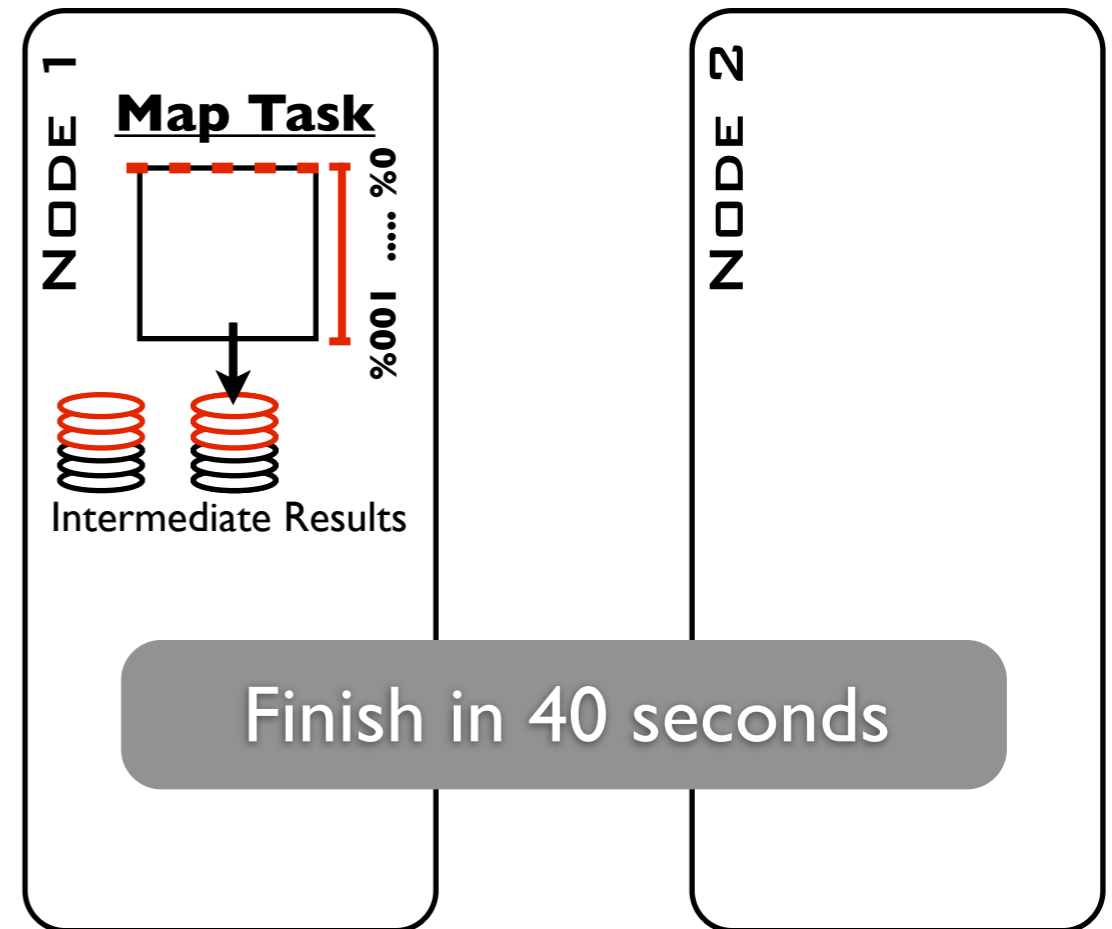
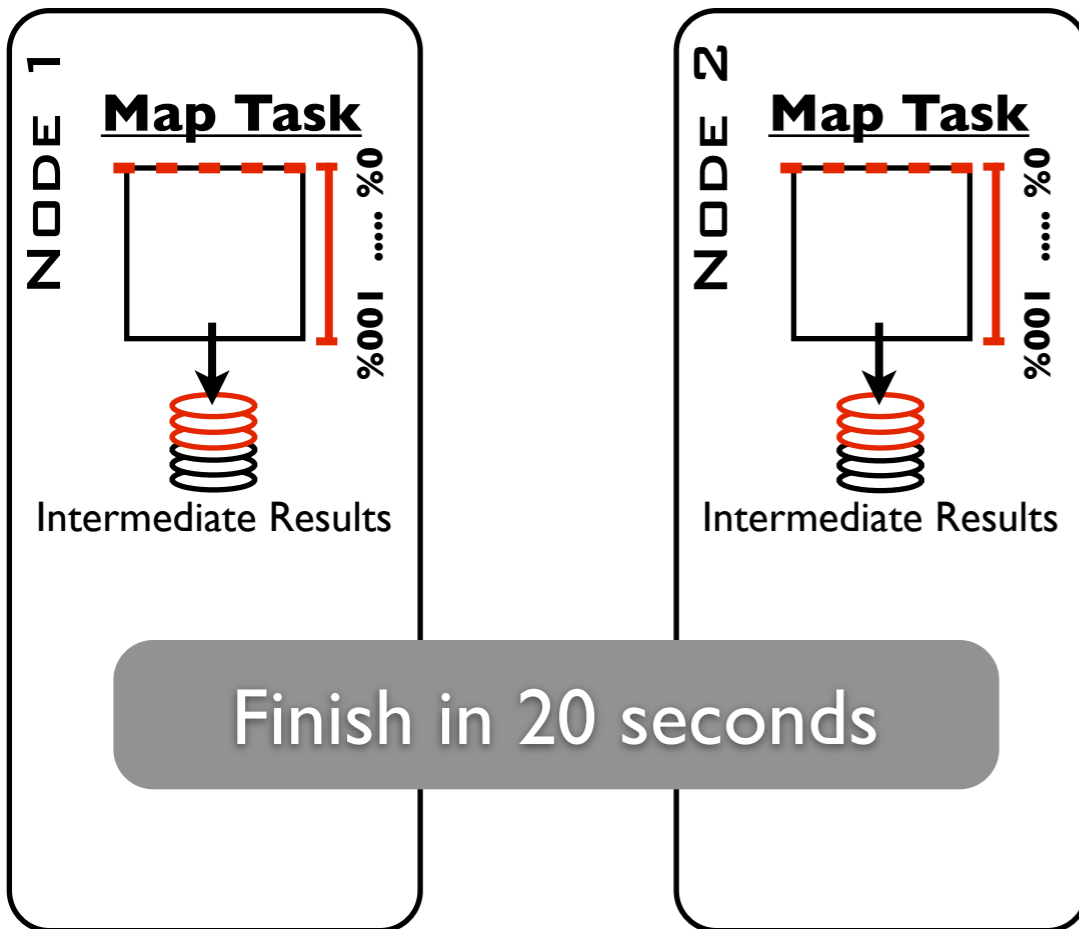
Hadoop  
with Task Failures

Scheduler

Time:20s

Time:40s

Scheduler



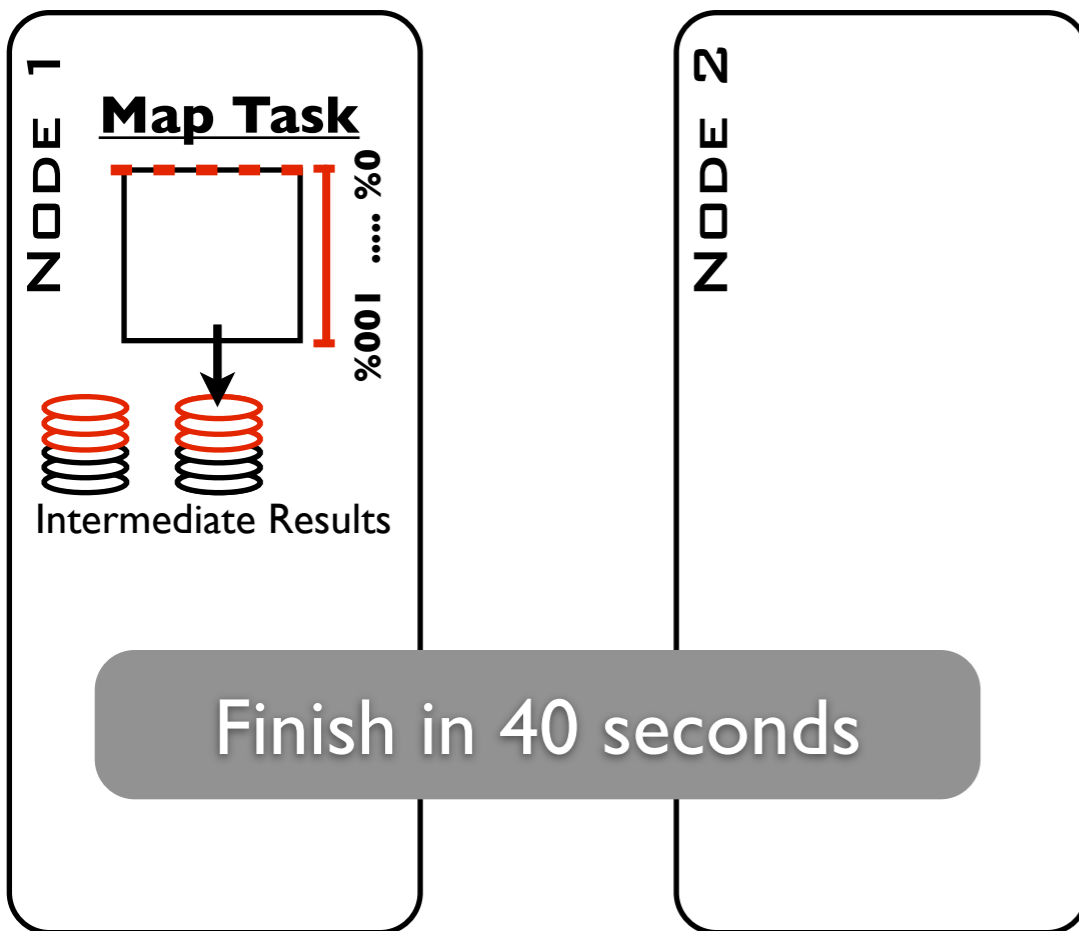
# Solution: Local Checkpointing

Hadoop

Map Task: 20 seconds

Time:40s

Scheduler



# Solution: Local Checkpointing

Hadoop

Map Task: 20 seconds

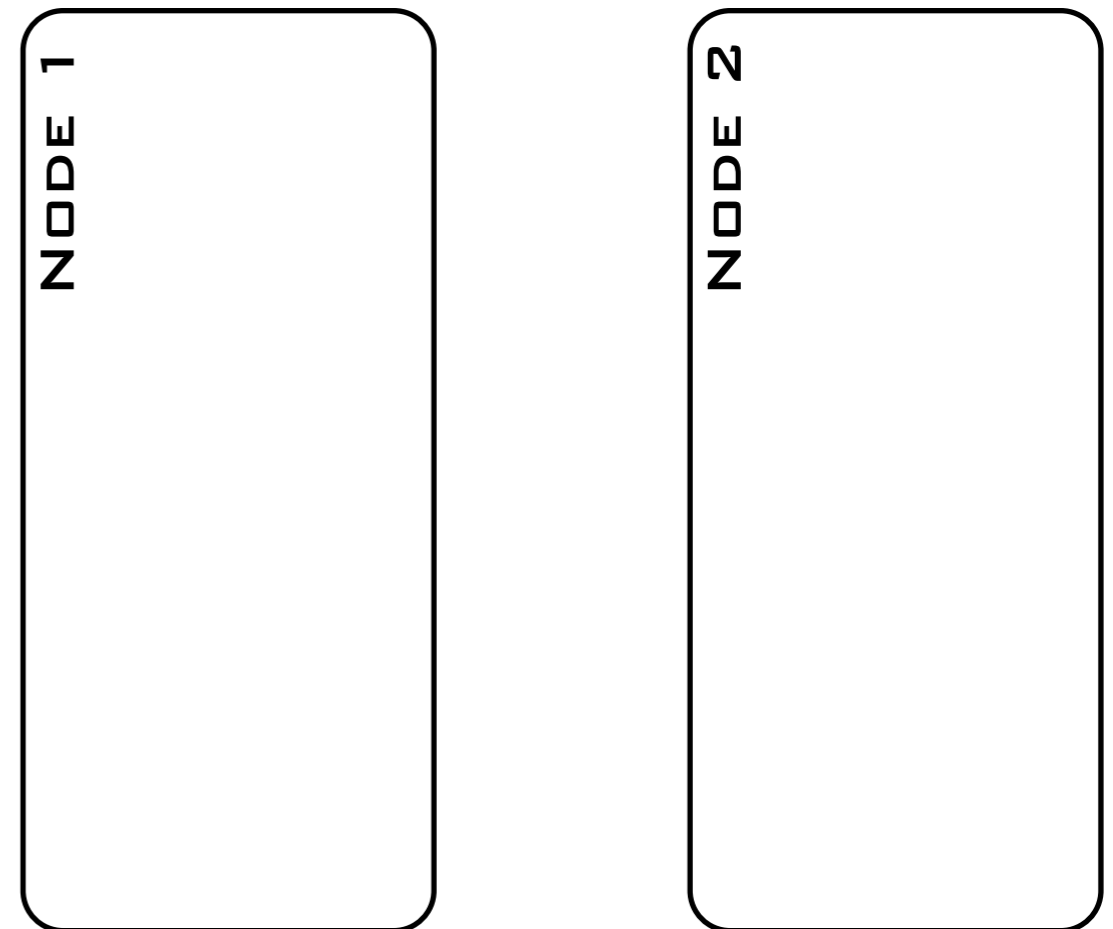
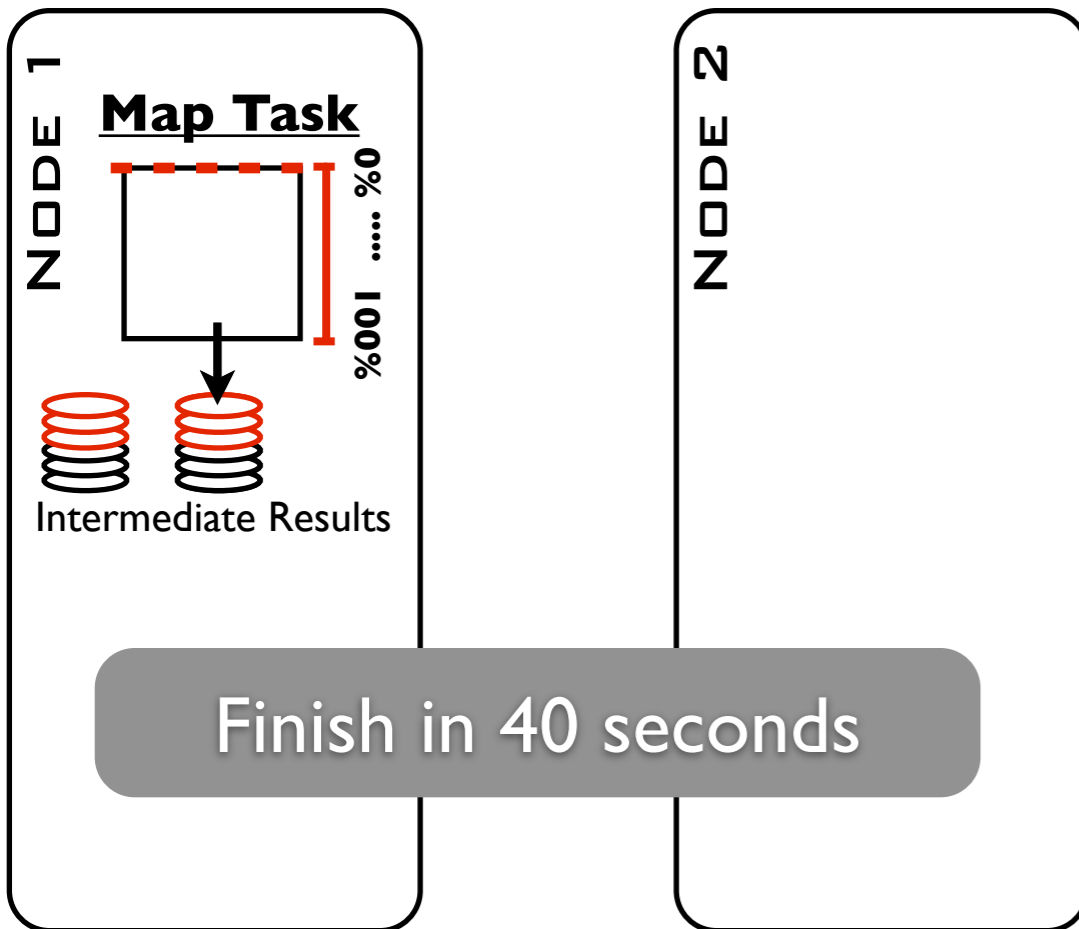
RAFT-LC

Scheduler

Time:40s

Time:0s

Scheduler



# Solution: Local Checkpointing

Hadoop

Map Task: 20 seconds

Scheduler

Time:40s

NODE 1

Map Task

%0  
.....  
%100

Intermediate Results

NODE 2

Finish in 40 seconds

RAFT-LC

Time:0s

Scheduler

NODE 1

NODE 2



# Solution: Local Checkpointing

Hadoop

Map Task: 20 seconds

Scheduler

Time:40s

NODE 1

Map Task

%001 ..... %0

%001 ..... %0



Intermediate Results

NODE 2

Finish in 40 seconds

RAFT-LC

Time:8s

Scheduler

NODE 1

Map Task

%001 ..... %0

%001 ..... %0



NODE 2

Map Task

%001 ..... %0

%001 ..... %0



# Solution: Local Checkpointing

Hadoop

Map Task: 20 seconds

Scheduler

Time:40s

NODE 1

Map Task

%0  
...  
%100

Intermediate Results

NODE 2

Finish in 40 seconds

RAFT-LC

Time:8s

Scheduler

NODE 1

Map Task

N

Map Task

%0  
...  
%100

LC: Local Checkpoint

Offset	Spill Identifier
6501	Spill-1

# Solution: Local Checkpointing

Hadoop

Map Task: 20 seconds

Scheduler

Time: 40s

NODE 1

Map Task

%001 ..... %0

%001 ..... %0



Intermediate Results

NODE 2

Finish in 40 seconds

RAFT-LC

Time: 10s

Scheduler

NODE 1

Map Task

%001 ..... %0

%001 ..... %0



NODE 2

Map Task

%001 ..... %0

%001 ..... %0



# Solution: Local Checkpointing

Hadoop

Map Task: 20 seconds

Scheduler

Time: 40s

NODE 1

Map Task

%001 ..... %0

%001 ..... %0



Intermediate Results

NODE 2

Finish in 40 seconds

RAFT-LC

Time: 10s

Scheduler

NODE 1

Map Task

%001 ..... %0

%001 ..... %0



NODE 2

Map Task

%001 ..... %0

%001 ..... %0

# Solution: Local Checkpointing

Hadoop

Map Task: 20 seconds

Scheduler

Time: 40s

NODE 1

Map Task

%001 ..... %0

%100



Intermediate Results

NODE 2

Finish in 40 seconds

RAFT-LC

Time: 10s

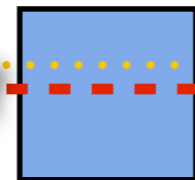
Scheduler

NODE 1

Map Task

%001 ..... %0

%100

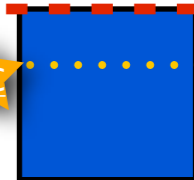


NODE 2

Map Task

%001 ..... %0

%100



# Solution: Local Checkpointing

Hadoop

Map Task: 20 seconds

Scheduler

Time: 40s

NODE 1

Map Task

%001 ..... %0

%001 ..... %0



Intermediate Results

NODE 2

Finish in 40 seconds

RAFT-LC

Time: 10s

Scheduler

NODE 1

Map Task

%001 ..... %0

%001 ..... %0



NODE 2

Map Task

%001 ..... %0

%001 ..... %0



# Solution: Local Checkpointing

Hadoop

Map Task: 20 seconds

Scheduler

Time: 40s

NODE 1

Map Task

%001 ..... %0

%001 ..... %0



Intermediate Results

NODE 2

Finish in 40 seconds

RAFT-LC

Time: 12s

Scheduler

NODE 1

Map Task

%001 ..... %0

%001 ..... %0



NODE 2

Map Task

%001 ..... %0

%001 ..... %0



# Solution: Local Checkpointing

Hadoop

Map Task: 20 seconds

Scheduler

Time: 40s

NODE 1

Map Task

%001 ..... %0

%001 ..... %100



Intermediate Results

NODE 2

Finish in 40 seconds

RAFT-LC

Time: 12s

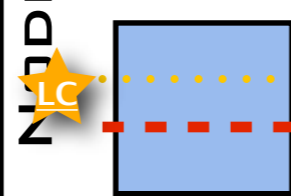
Scheduler

NODE 1

Map Task

%001 ..... %0

%001 ..... %100

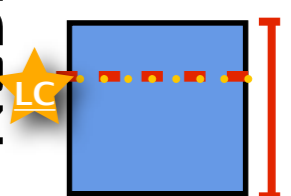


NODE 2

Map Task

%001 ..... %0

%001 ..... %100





# Solution: Local Checkpointing

Hadoop

Map Task: 20 seconds

Scheduler

Time: 40s

NODE 1

Map Task

%001 ..... %0

%001 ..... %100



Intermediate Results

NODE 2

Finish in 40 seconds

RAFT-LC

Time: 16s

Scheduler

NODE 1

Map Task

%001 ..... %0

%001 ..... %100



NODE 2

Map Task

%001 ..... %0

%001 ..... %100



# Solution: Local Checkpointing

Hadoop

Map Task: 20 seconds

Scheduler

Time:40s

NODE 1

Map Task

%001 ..... %0

%001 ..... %100



Intermediate Results

NODE 2

Finish in 40 seconds

RAFT-LC

Time:20s

Scheduler

NODE 1

Map Task

%001 ..... %0

%001 ..... %100



Intermediate Results

NODE 2

Map Task

%001 ..... %0

%001 ..... %100



# Solution: Local Checkpointing

Hadoop

Map Task: 20 seconds

Scheduler

Time:40s

NODE 1

Map Task

%001 ..... %0

%100



Intermediate Results

Finish in 40 seconds

NODE 2

RAFT-LC

Time:24s

Scheduler

NODE 1

Map Task

%001 ..... %0

%100



Intermediate Results

Finish in 24 seconds

Faster by a factor of ~1.7

NODE 2

Map Task

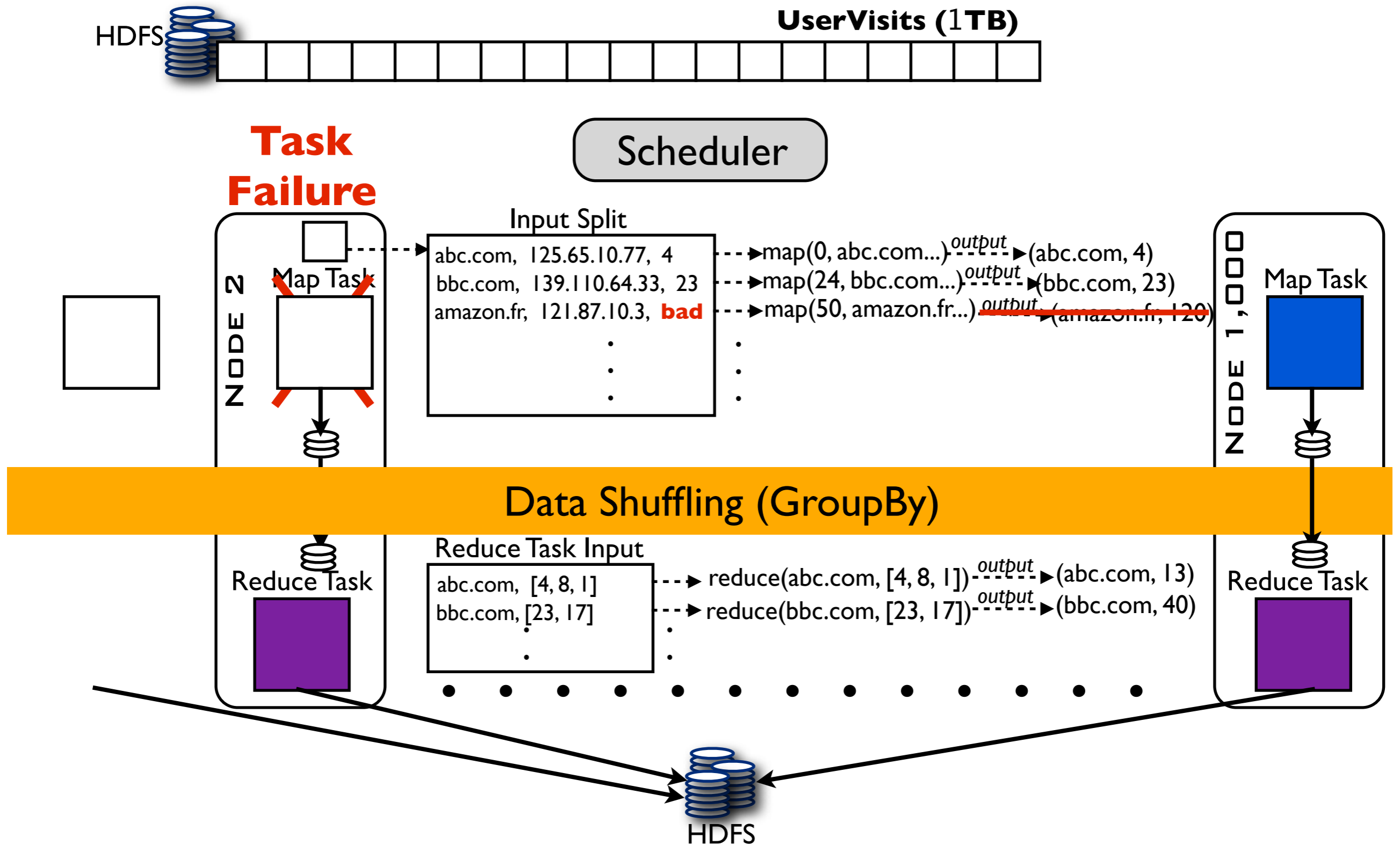
%001 ..... %0

%100

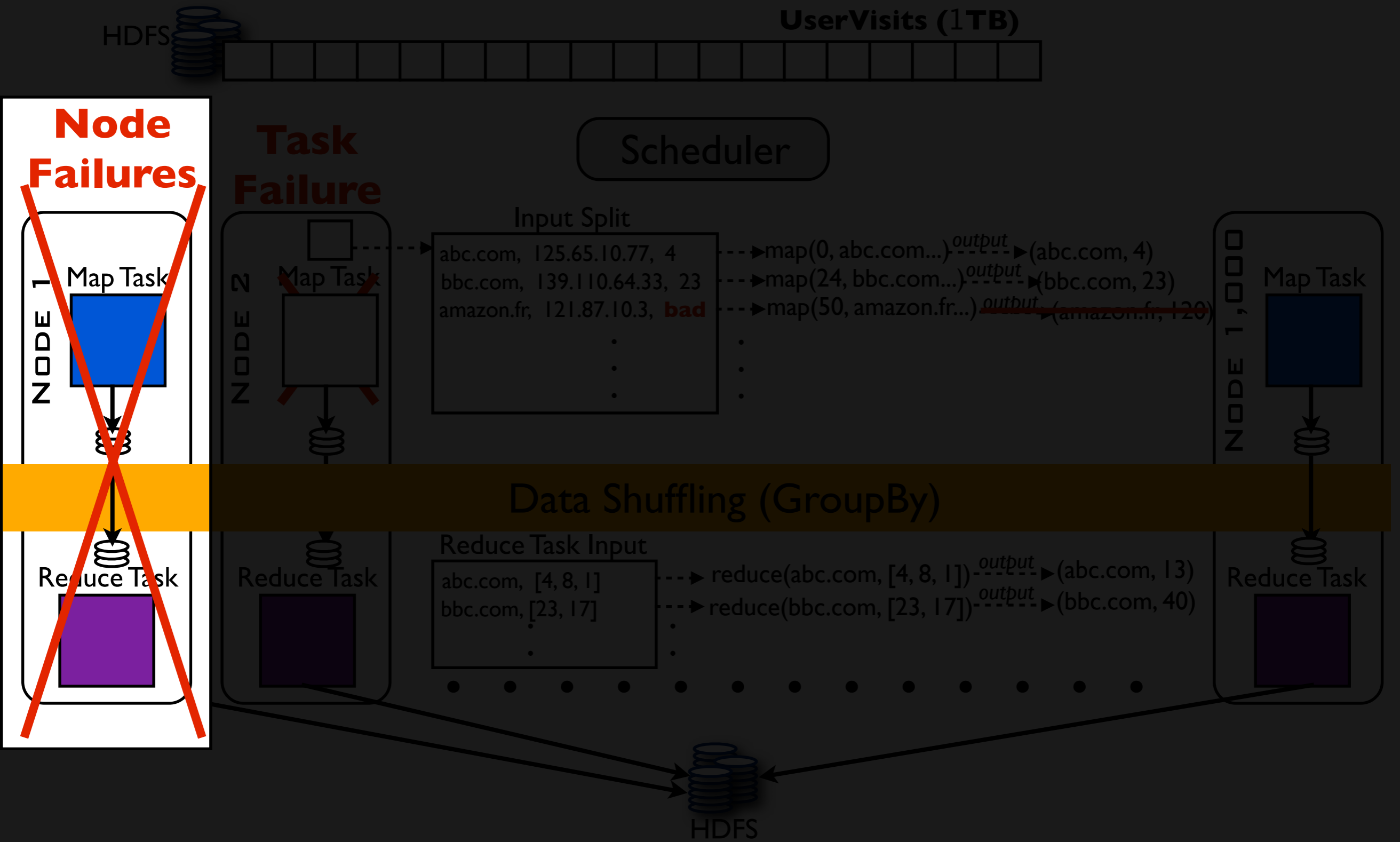


Intermediate Results

# Node Failures



# Node Failures



# Current Approach

## Hadoop without Failures

Map Task: 20 seconds  
Reduce Task: 30 seconds

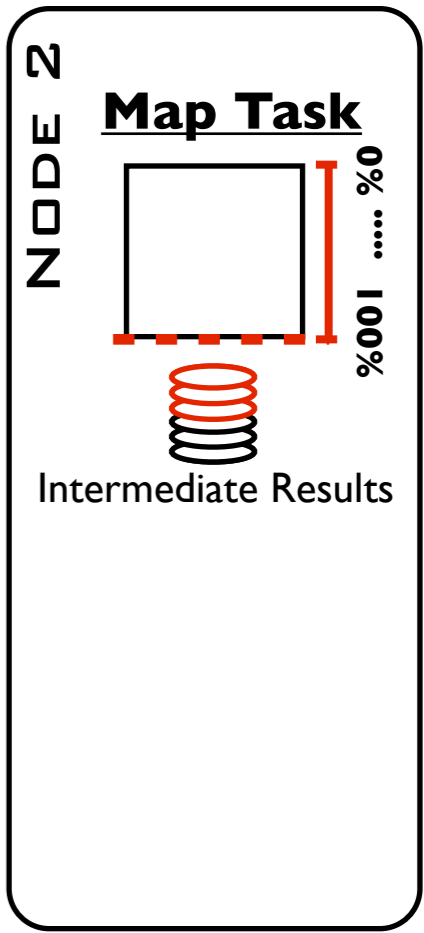
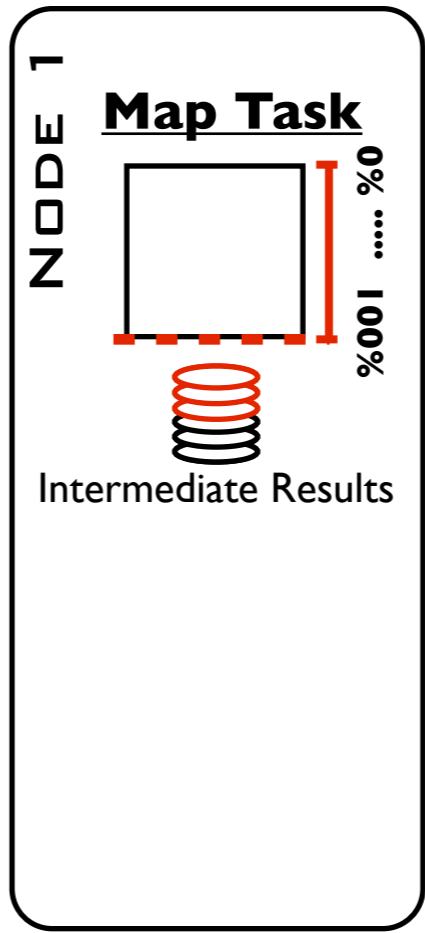
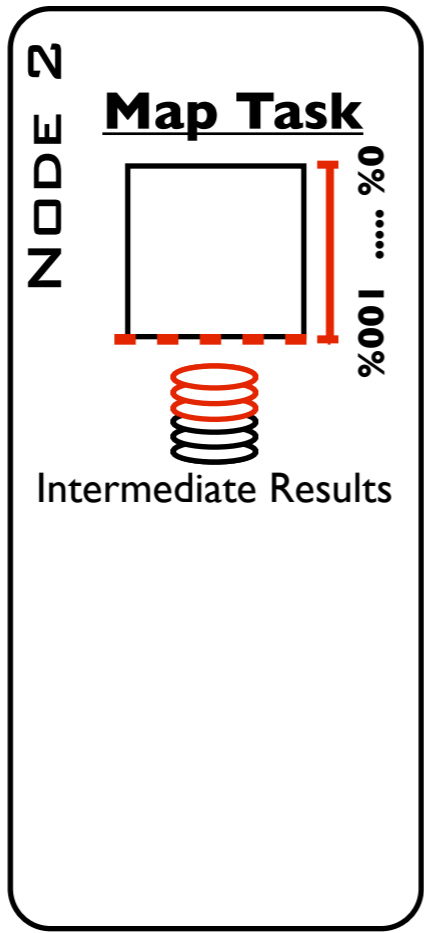
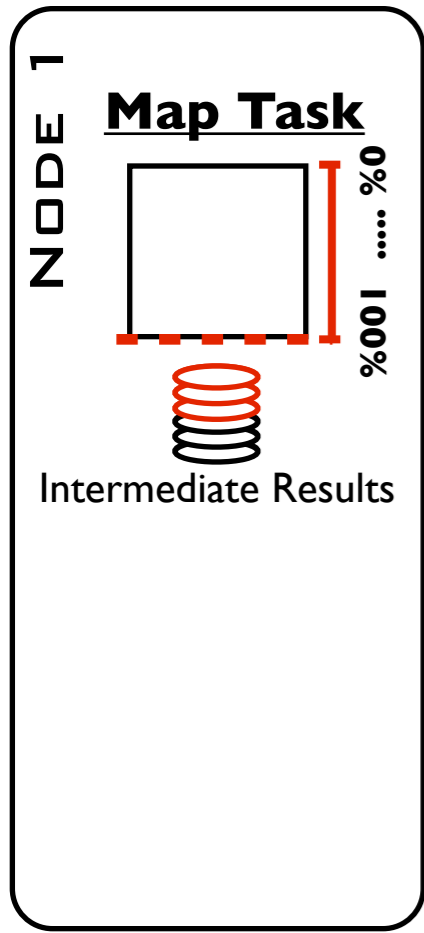
## Hadoop with Node Failures

Scheduler

Time:20s

Time:20s

Scheduler



# Current Approach

## Hadoop without Failures

Map Task: 20 seconds  
Reduce Task: 30 seconds

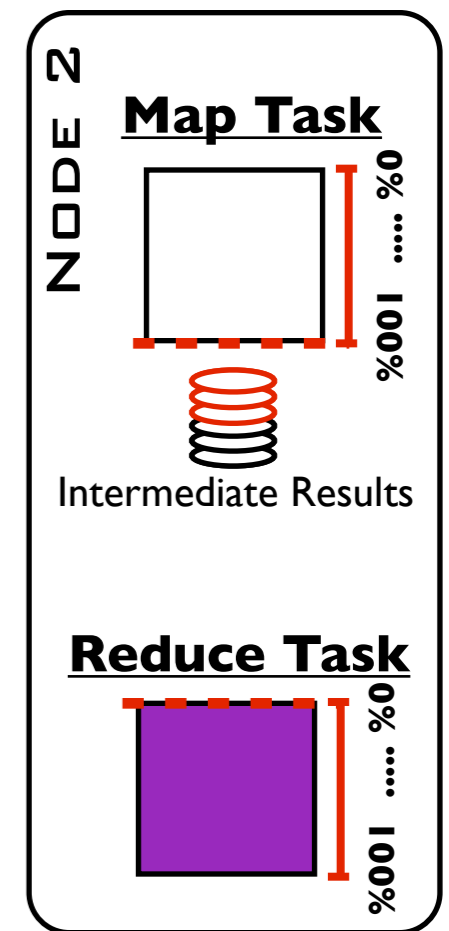
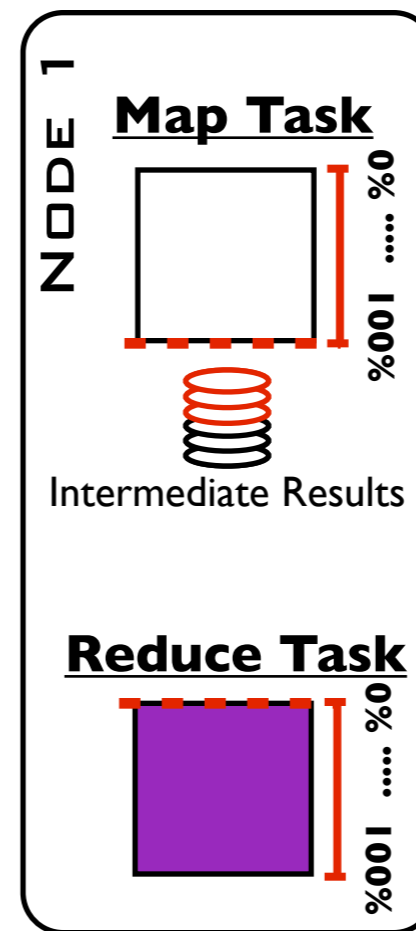
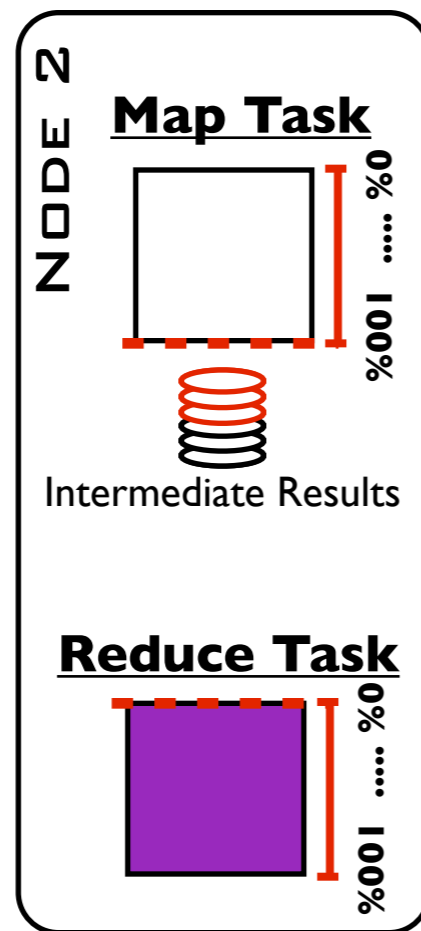
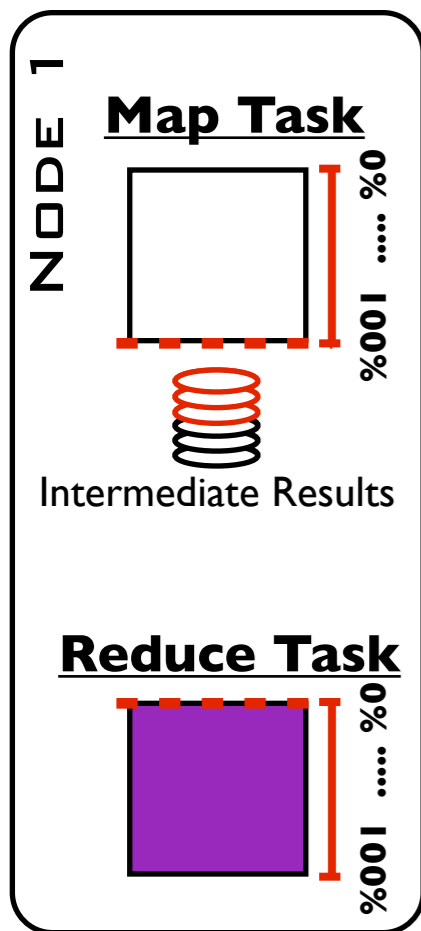
Scheduler

Time:20s

Time:20s

## Hadoop with Node Failures

Scheduler



# Current Approach

## Hadoop without Failures

Map Task: 20 seconds  
Reduce Task: 30 seconds

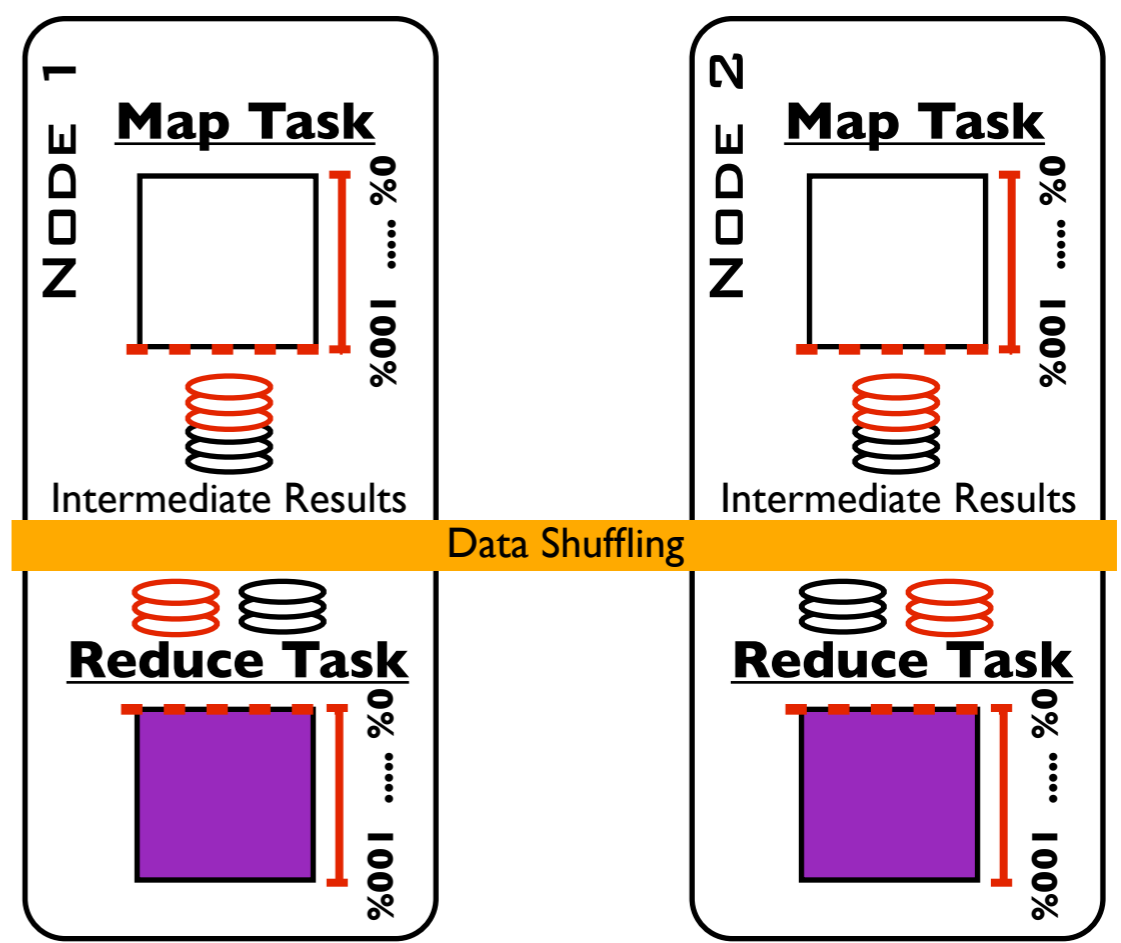
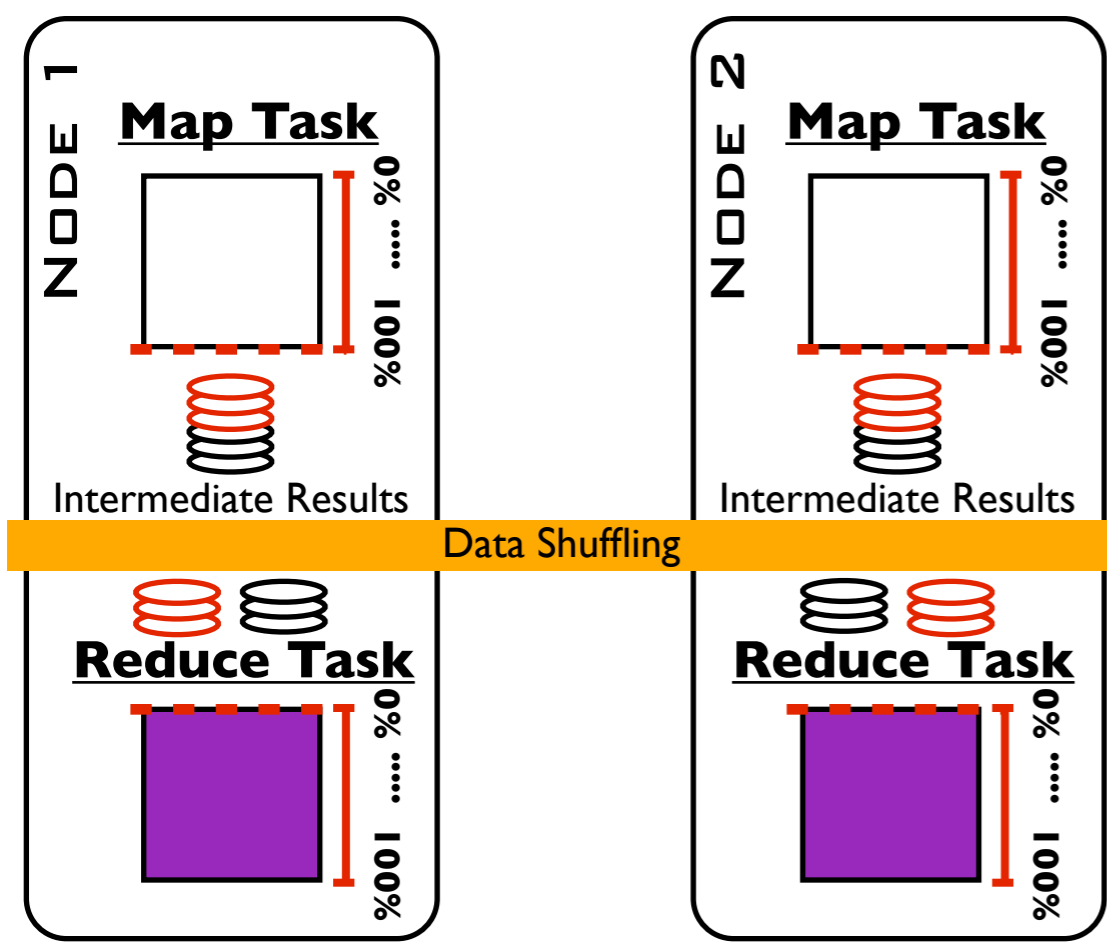
Scheduler

Time:20s

Time:20s

## Hadoop with Node Failures

Scheduler





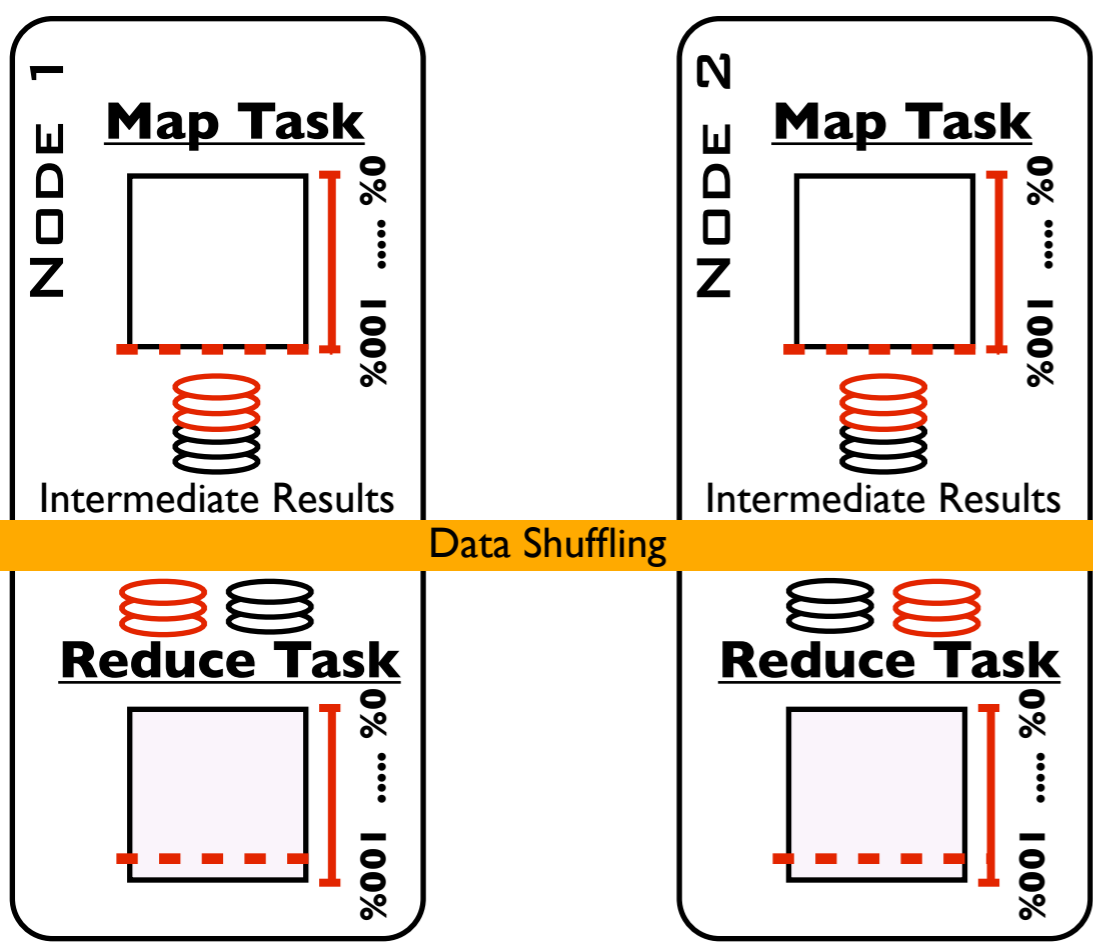
# Current Approach

## Hadoop without Failures

Map Task: 20 seconds  
Reduce Task: 30 seconds

Time:48s

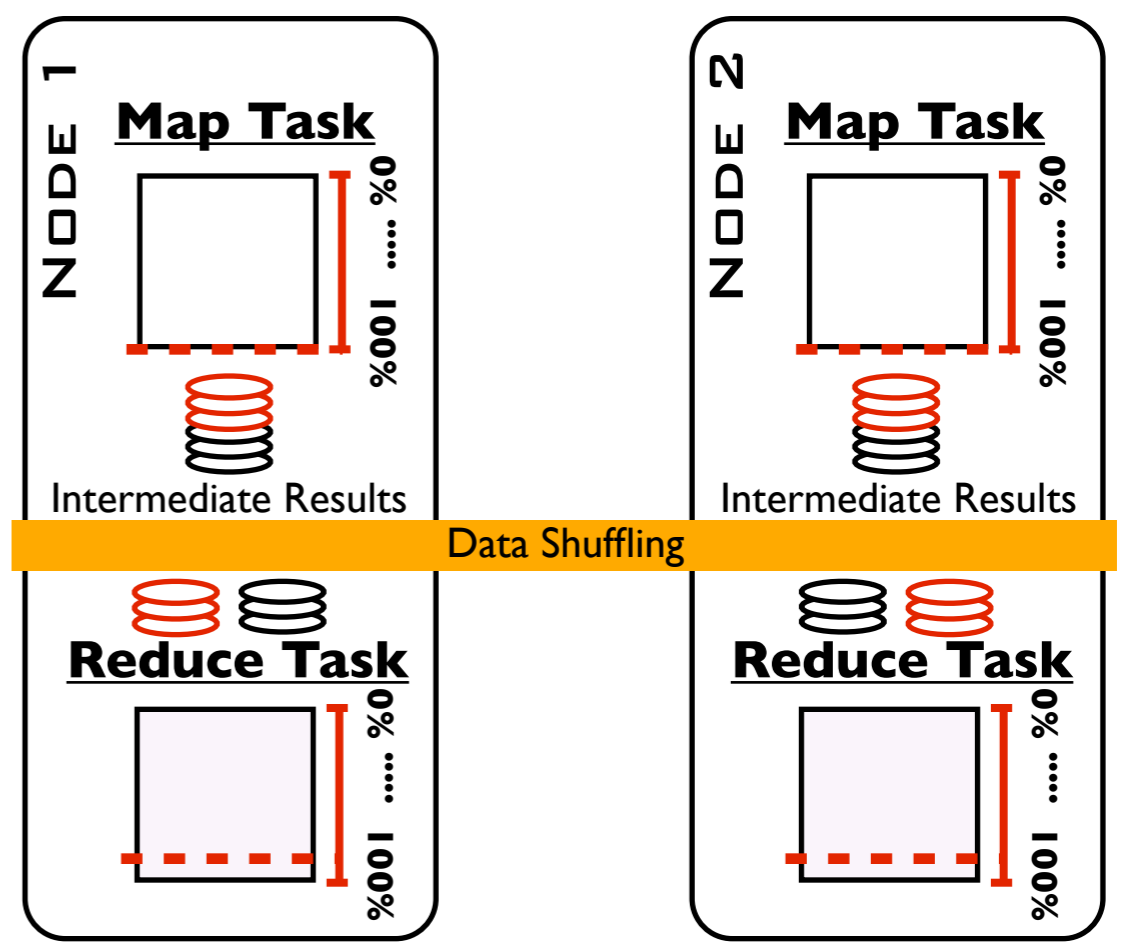
Scheduler



## Hadoop with Node Failures

Time:48s

Scheduler



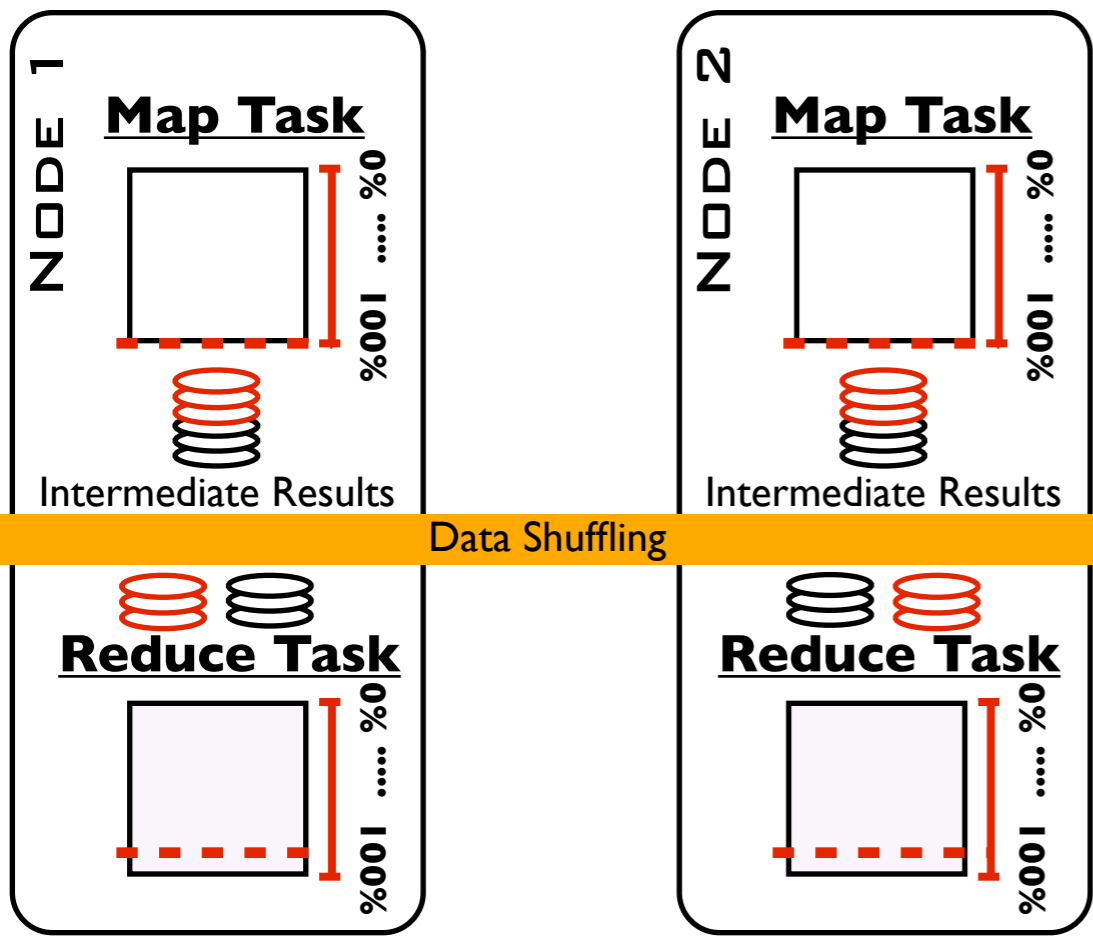
# Current Approach

## Hadoop without Failures

Map Task: 20 seconds  
Reduce Task: 30 seconds

Time:48s

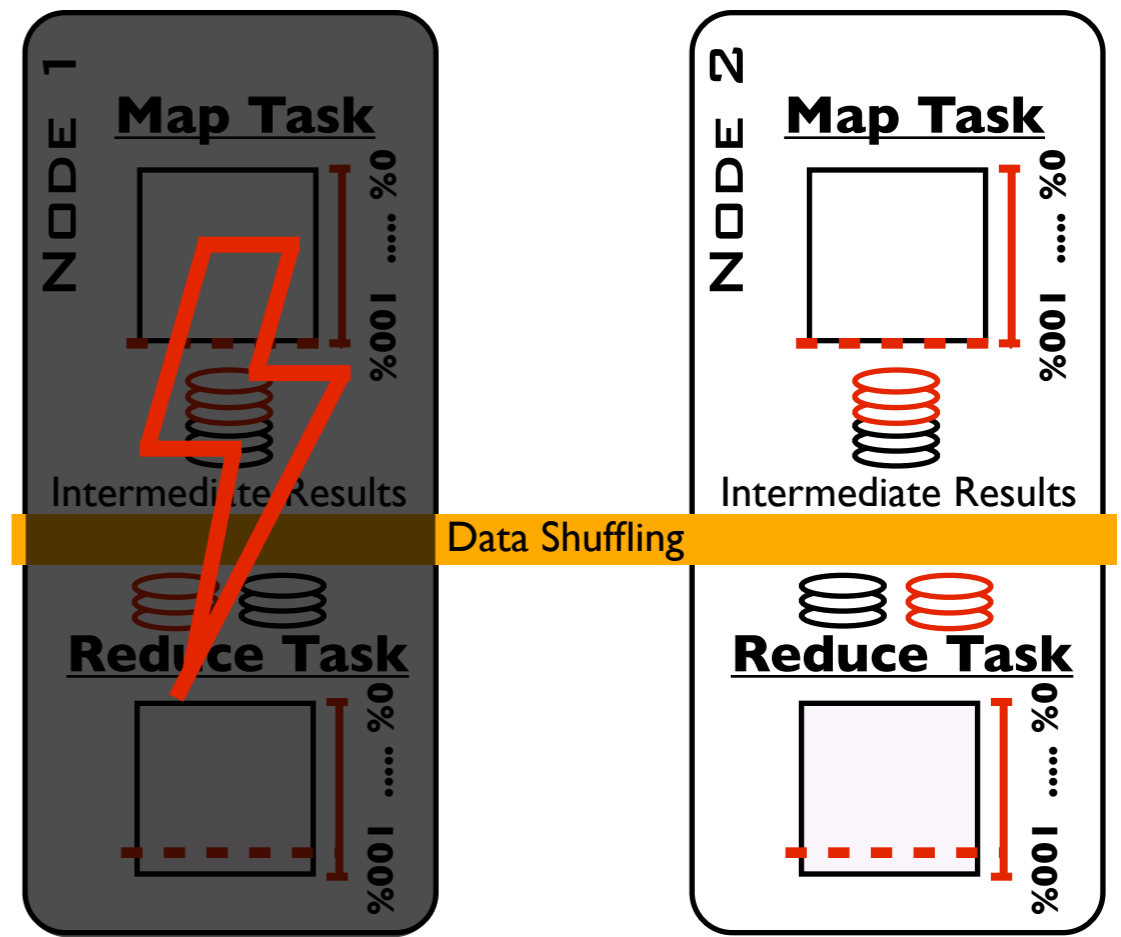
Scheduler



## Hadoop with Node Failures

Time:48s

Scheduler



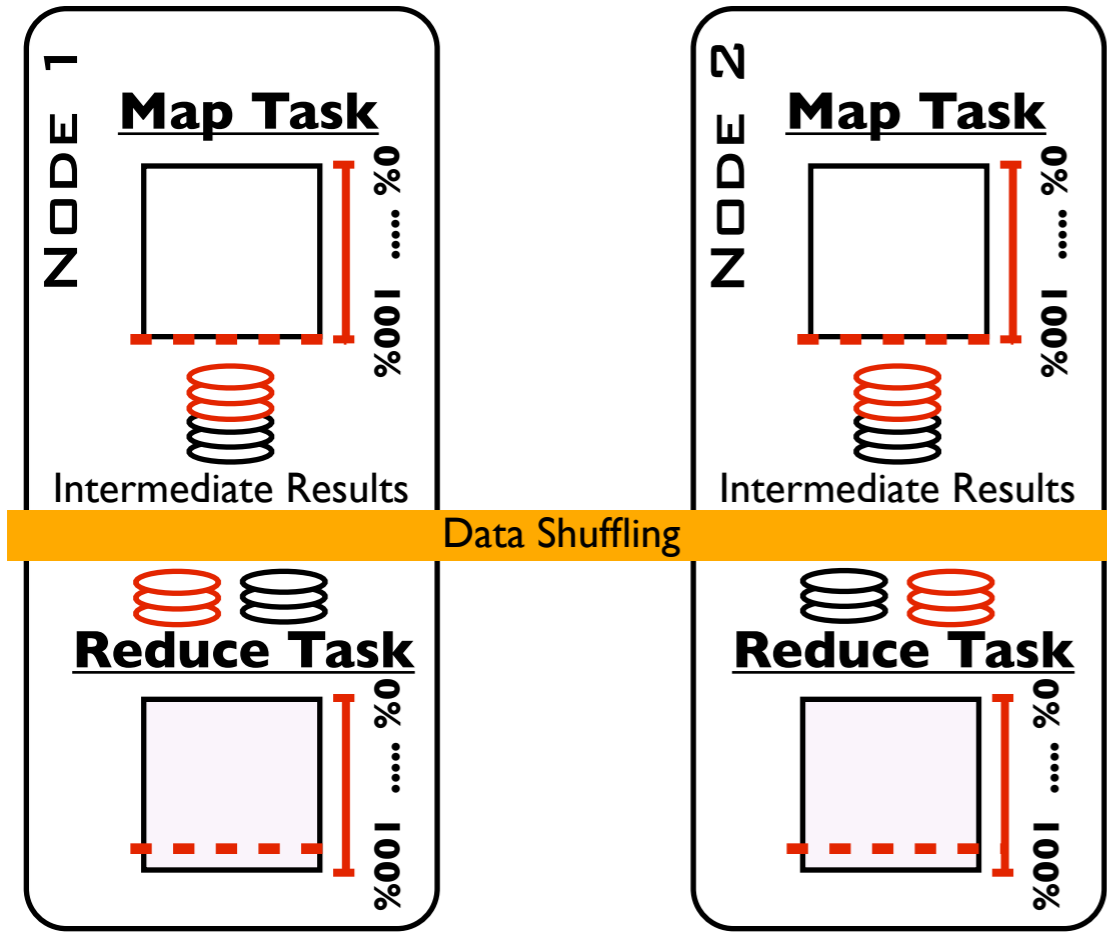
# Current Approach

## Hadoop without Failures

Map Task: 20 seconds  
Reduce Task: 30 seconds

Time:48s

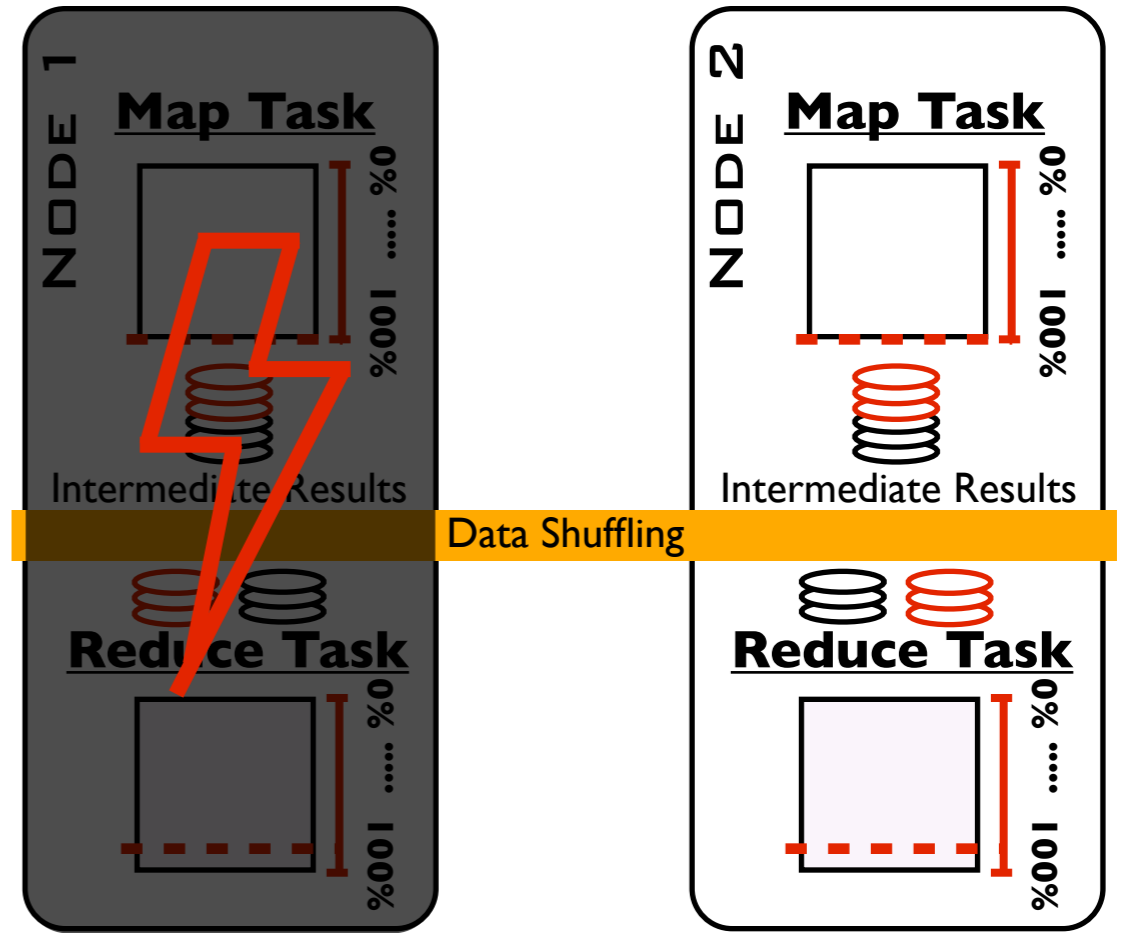
Scheduler



## Hadoop with Node Failures

Time:48s

Scheduler



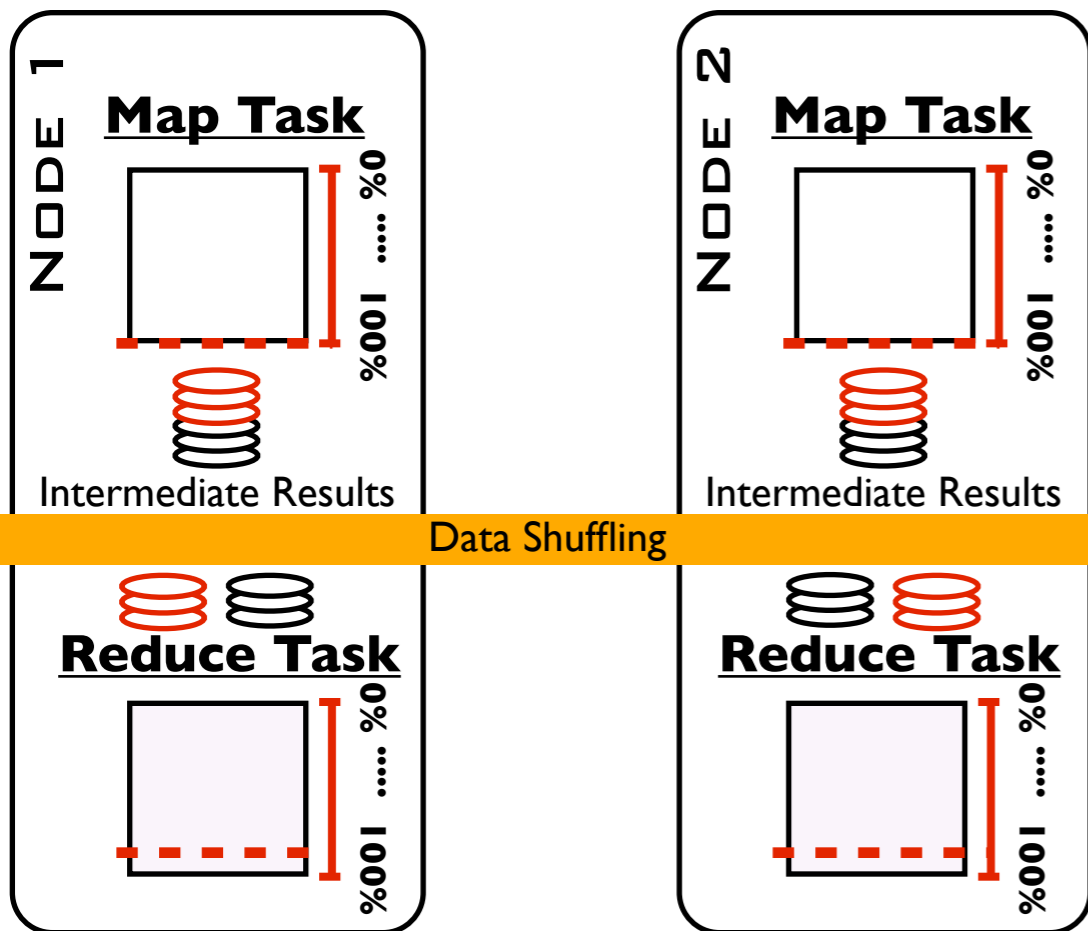
# Current Approach

## Hadoop without Failures

Map Task: 20 seconds  
Reduce Task: 30 seconds

Time:48s

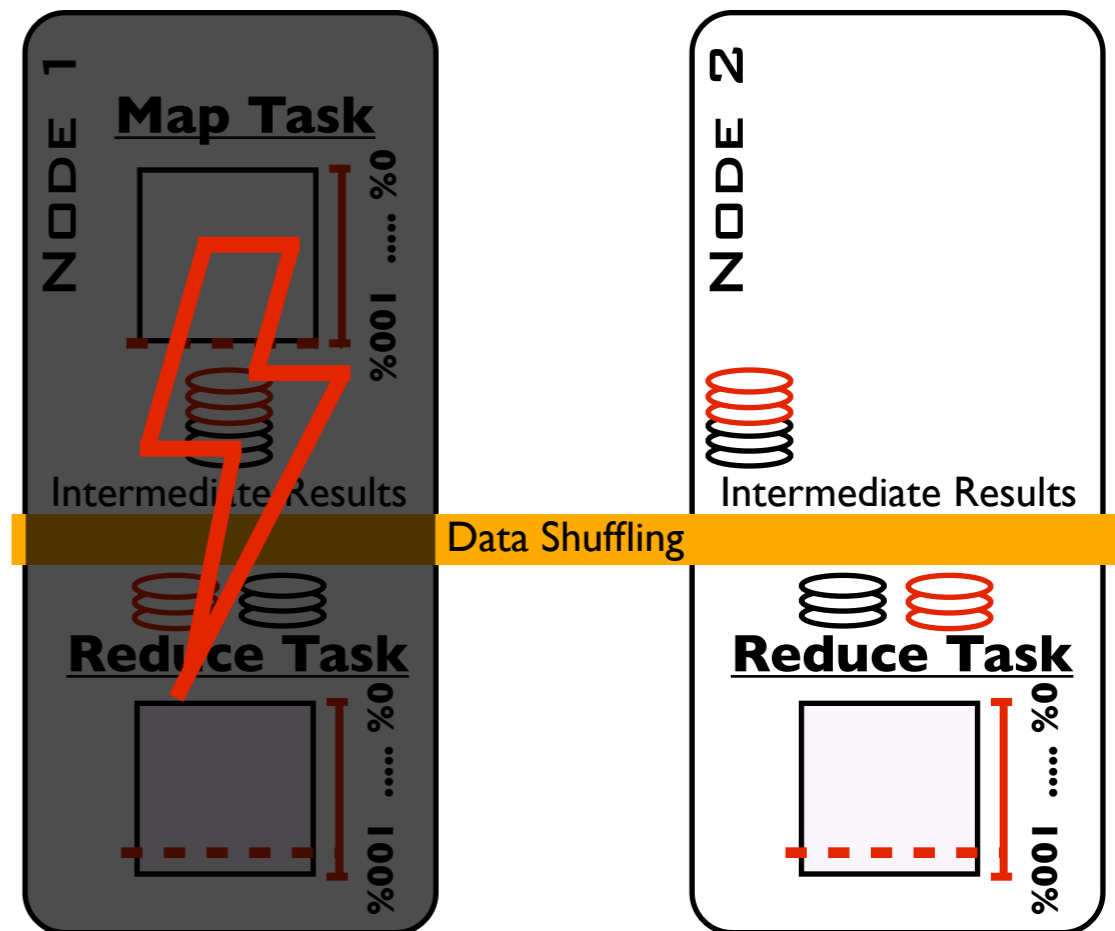
Scheduler



## Hadoop with Node Failures

Time:48s

Scheduler



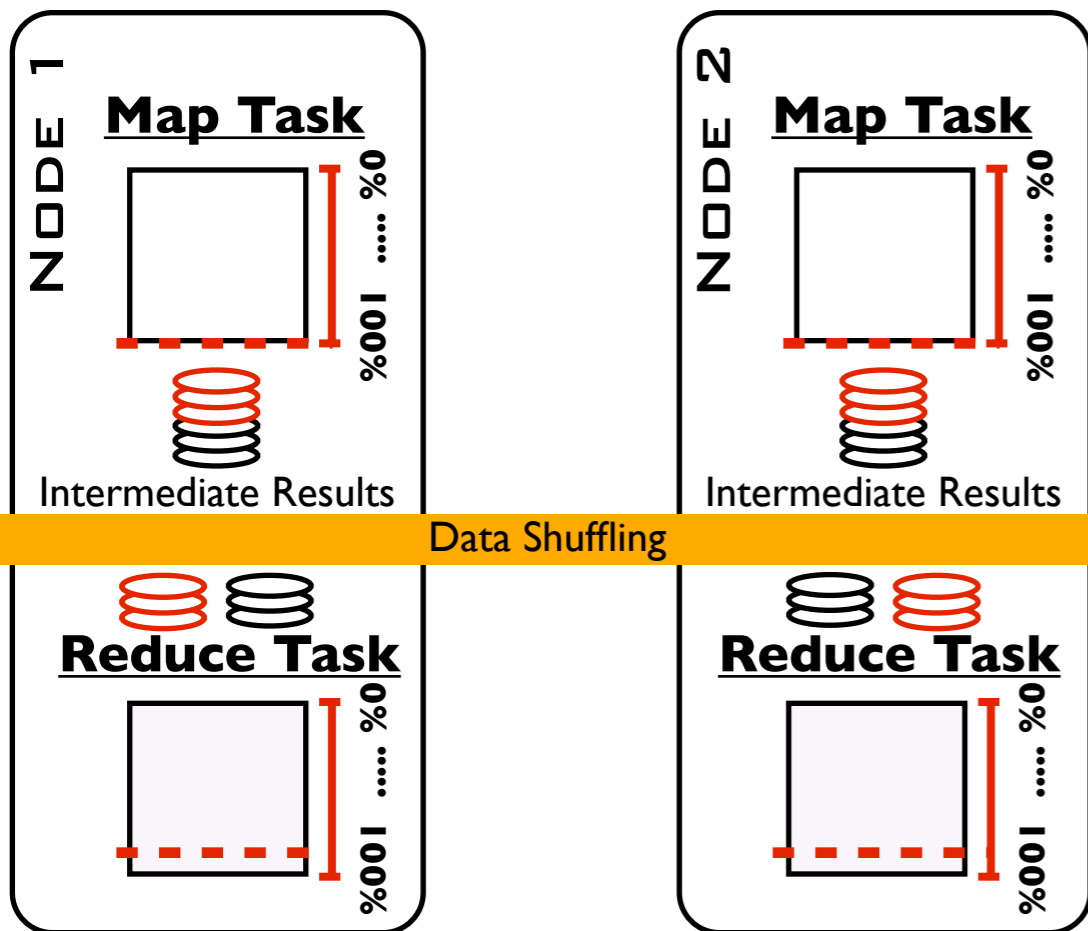
# Current Approach

## Hadoop without Failures

Map Task: 20 seconds  
Reduce Task: 30 seconds

Time:48s

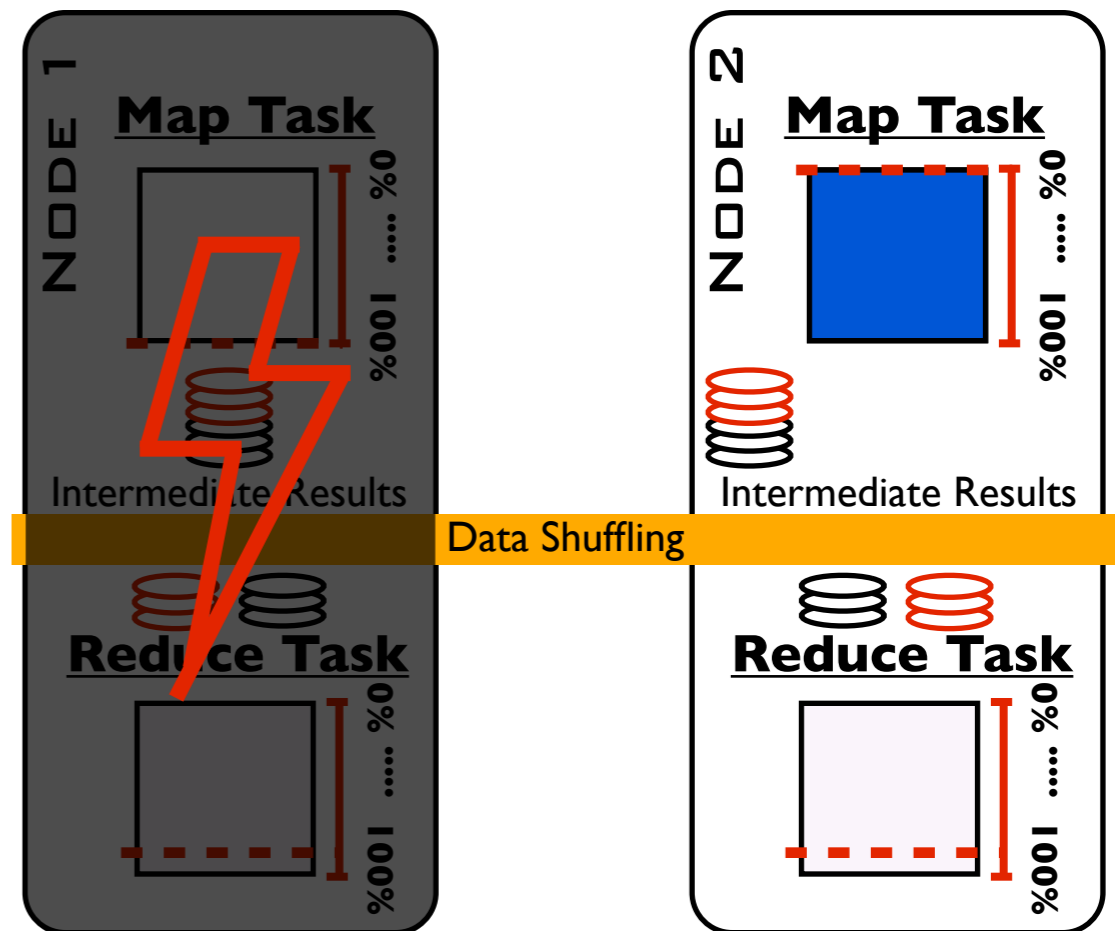
Scheduler



## Hadoop with Node Failures

Time:48s

Scheduler



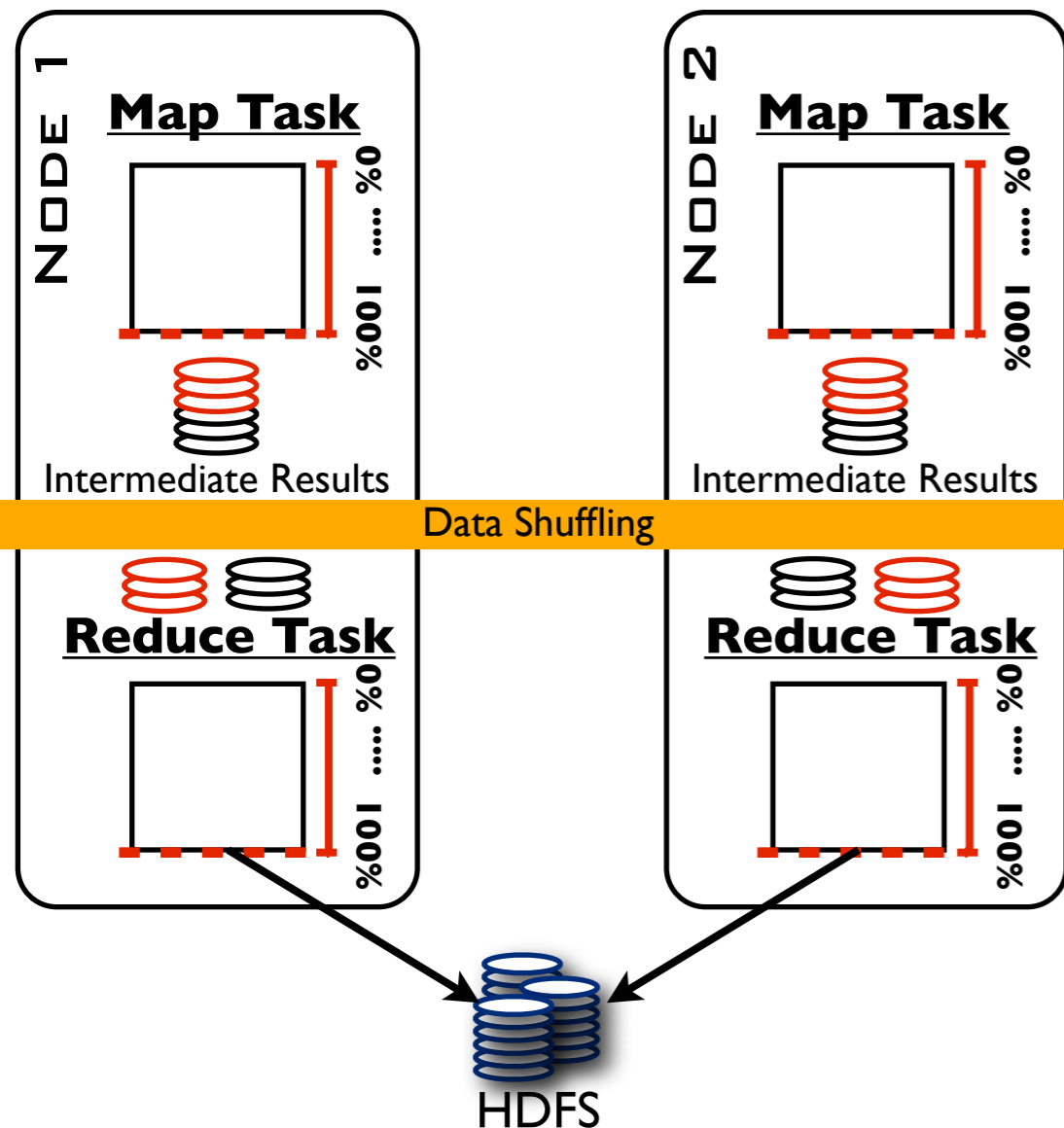
# Current Approach

## Hadoop without Failures

Map Task: 20 seconds  
Reduce Task: 30 seconds

Time: 50s

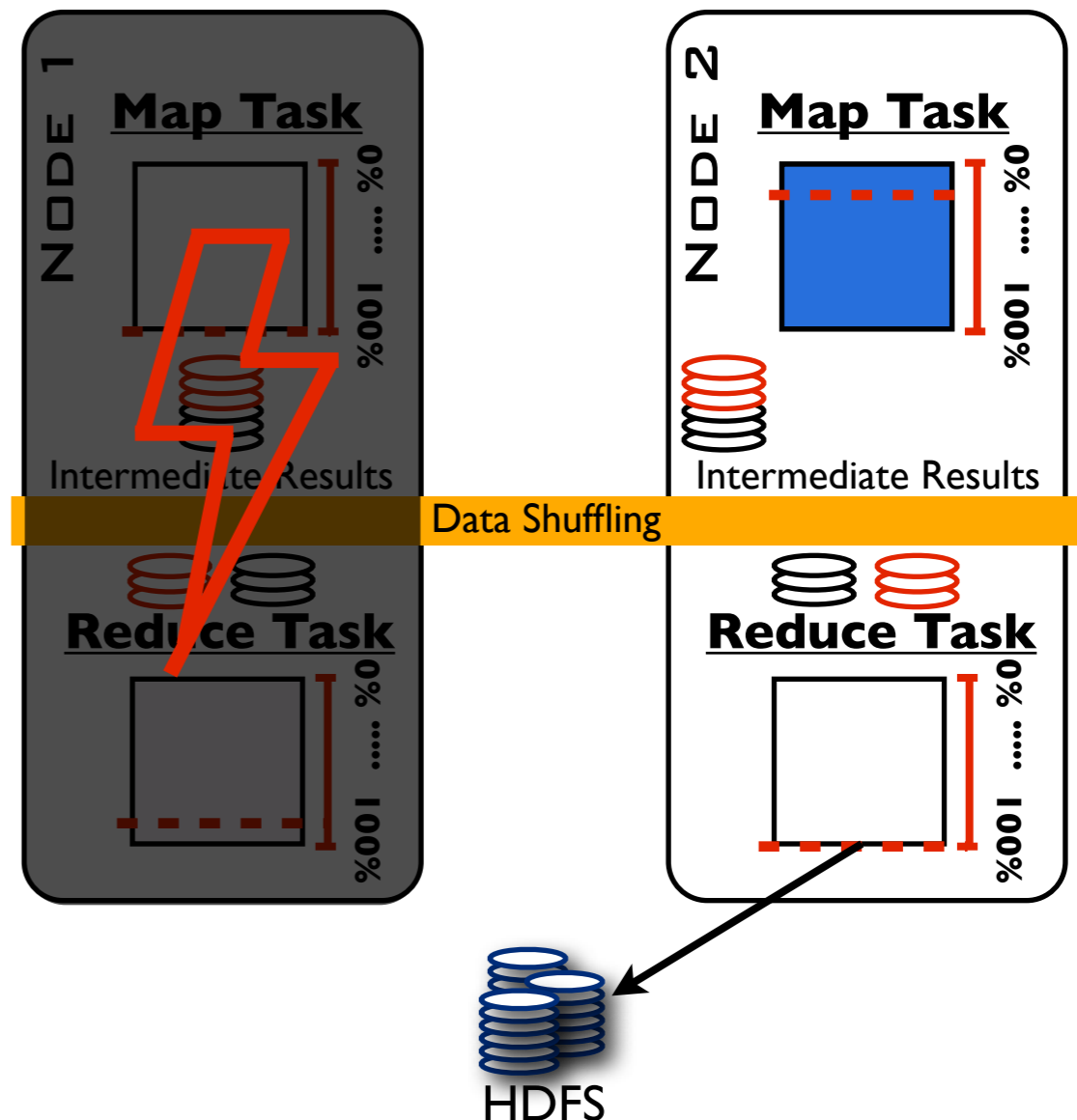
Scheduler



## Hadoop with Node Failures

Time: 50s

Scheduler



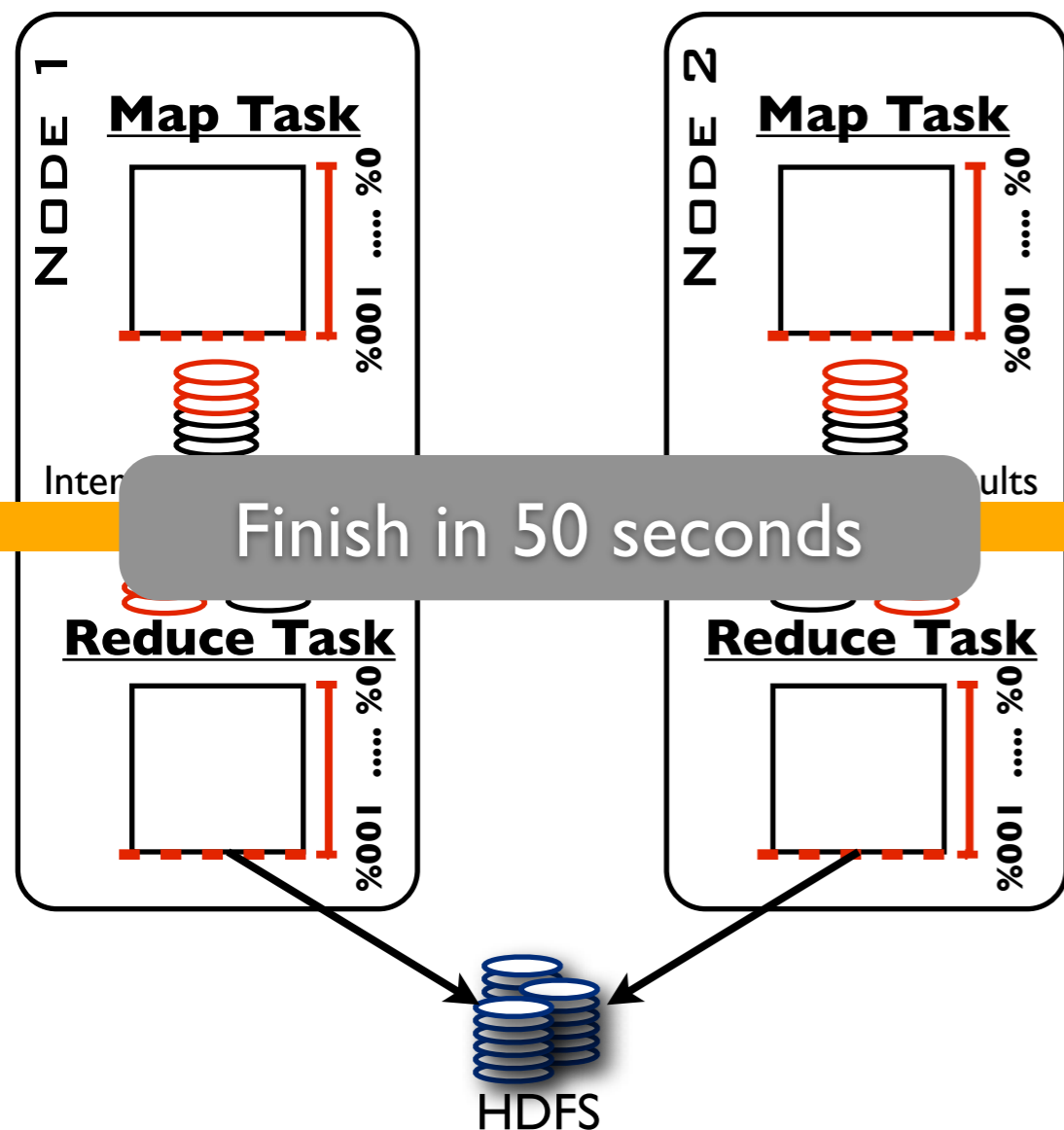
# Current Approach

## Hadoop without Failures

Map Task: 20 seconds  
Reduce Task: 30 seconds

Time: 50s

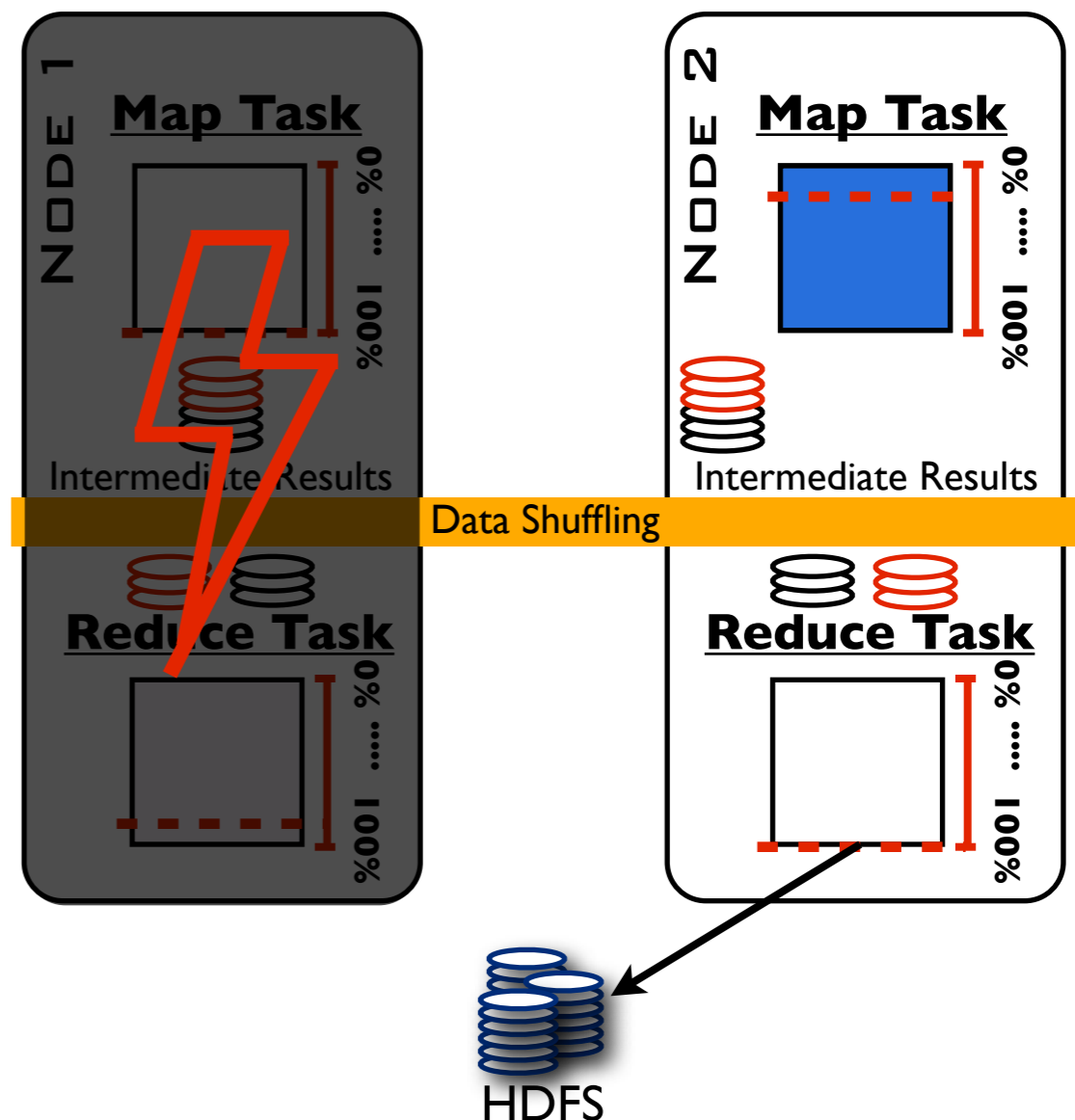
Scheduler



## Hadoop with Node Failures

Time: 50s

Scheduler



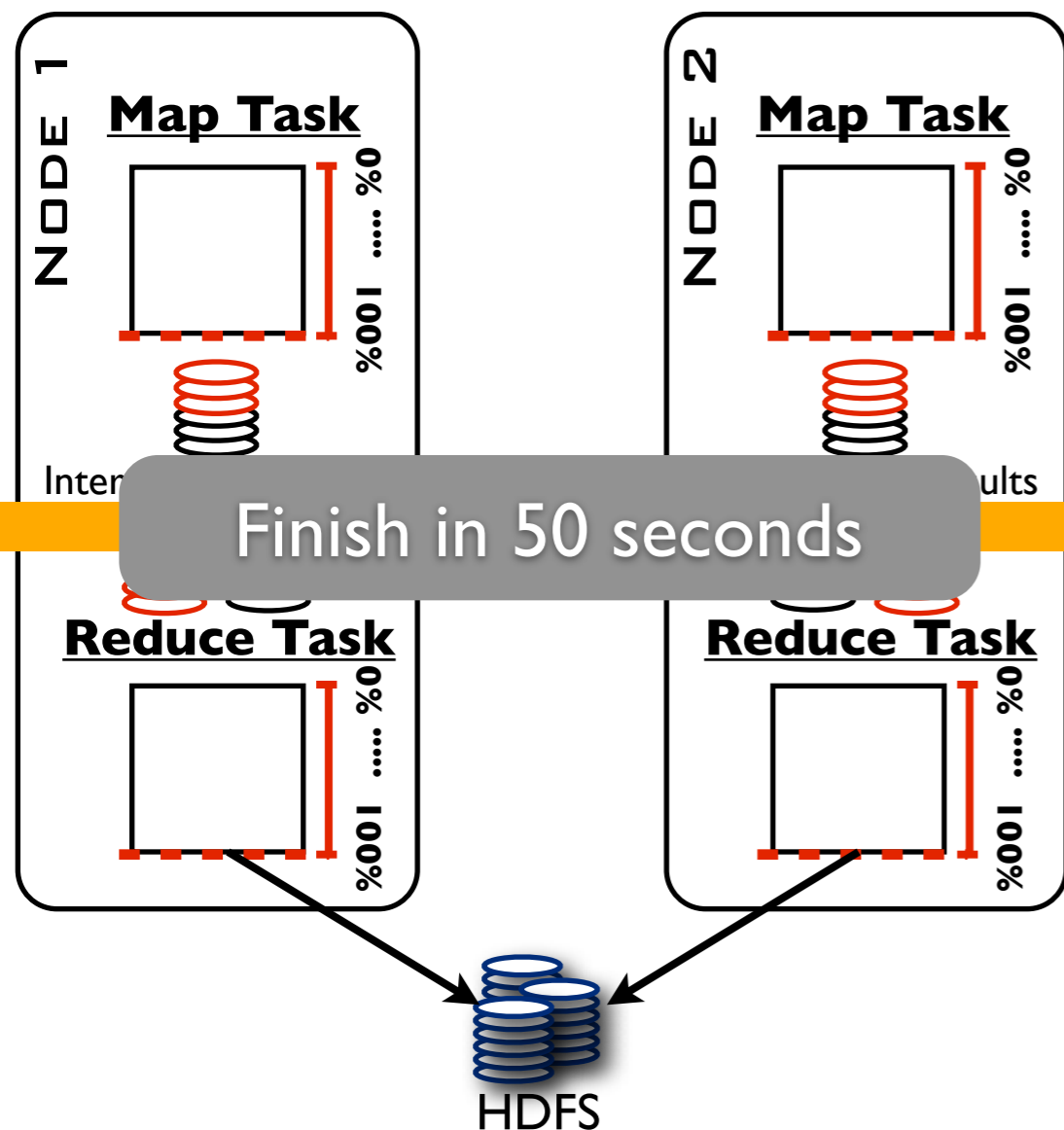
# Current Approach

## Hadoop without Failures

Map Task: 20 seconds  
Reduce Task: 30 seconds

Time:50s

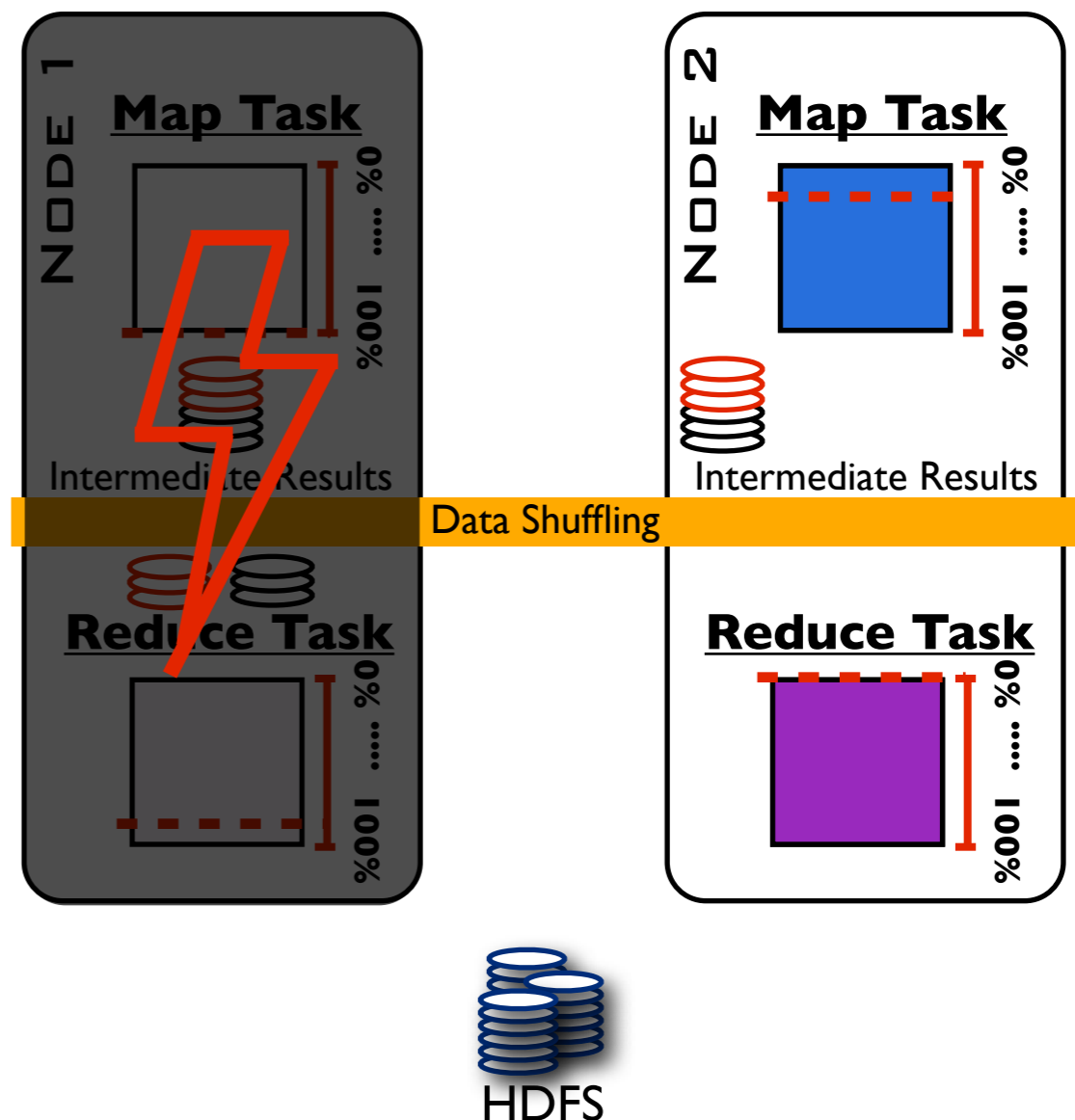
Scheduler



## Hadoop with Node Failures

Time:50s

Scheduler





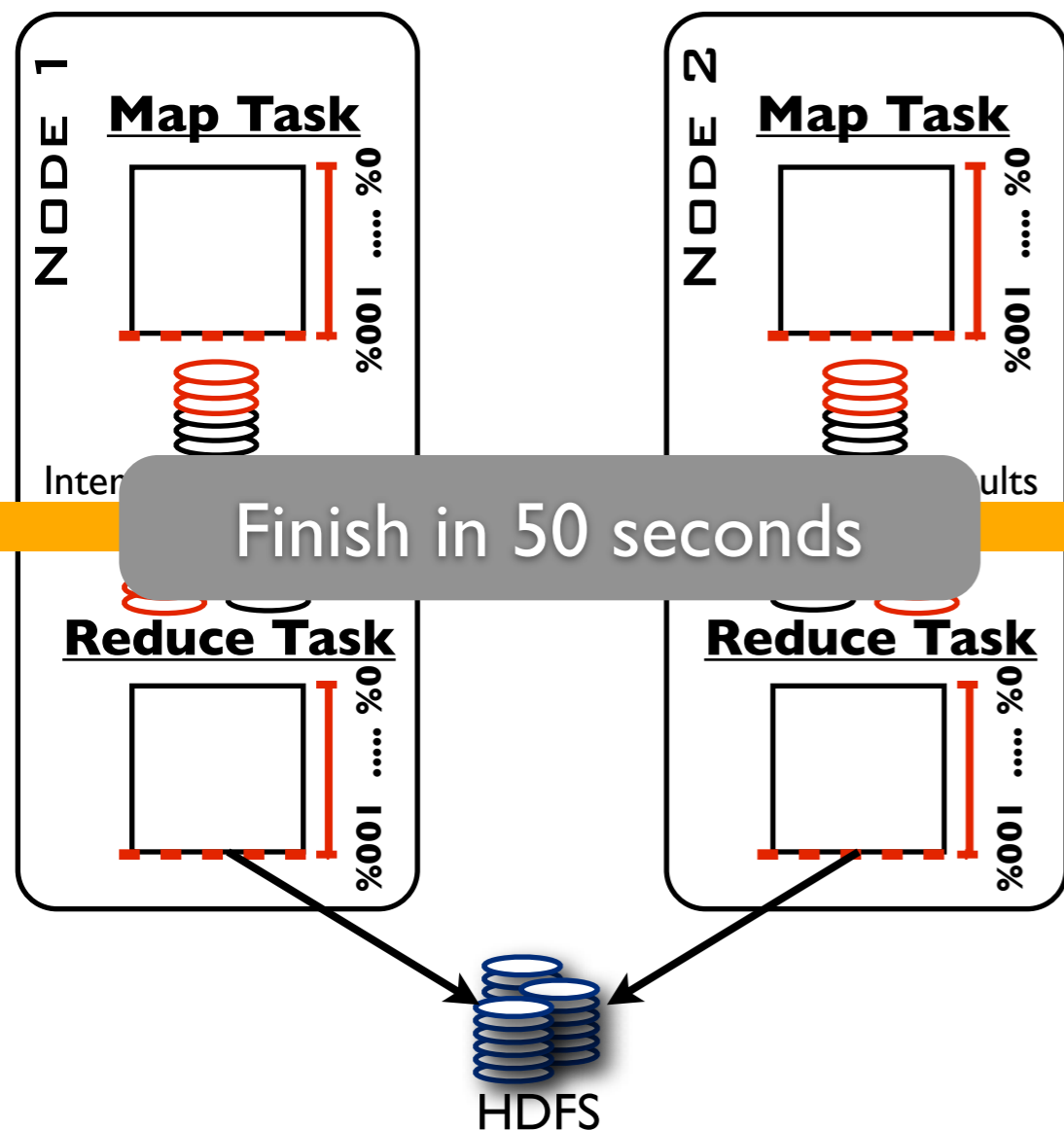
# Current Approach

## Hadoop without Failures

Map Task: 20 seconds  
Reduce Task: 30 seconds

Time: 50s

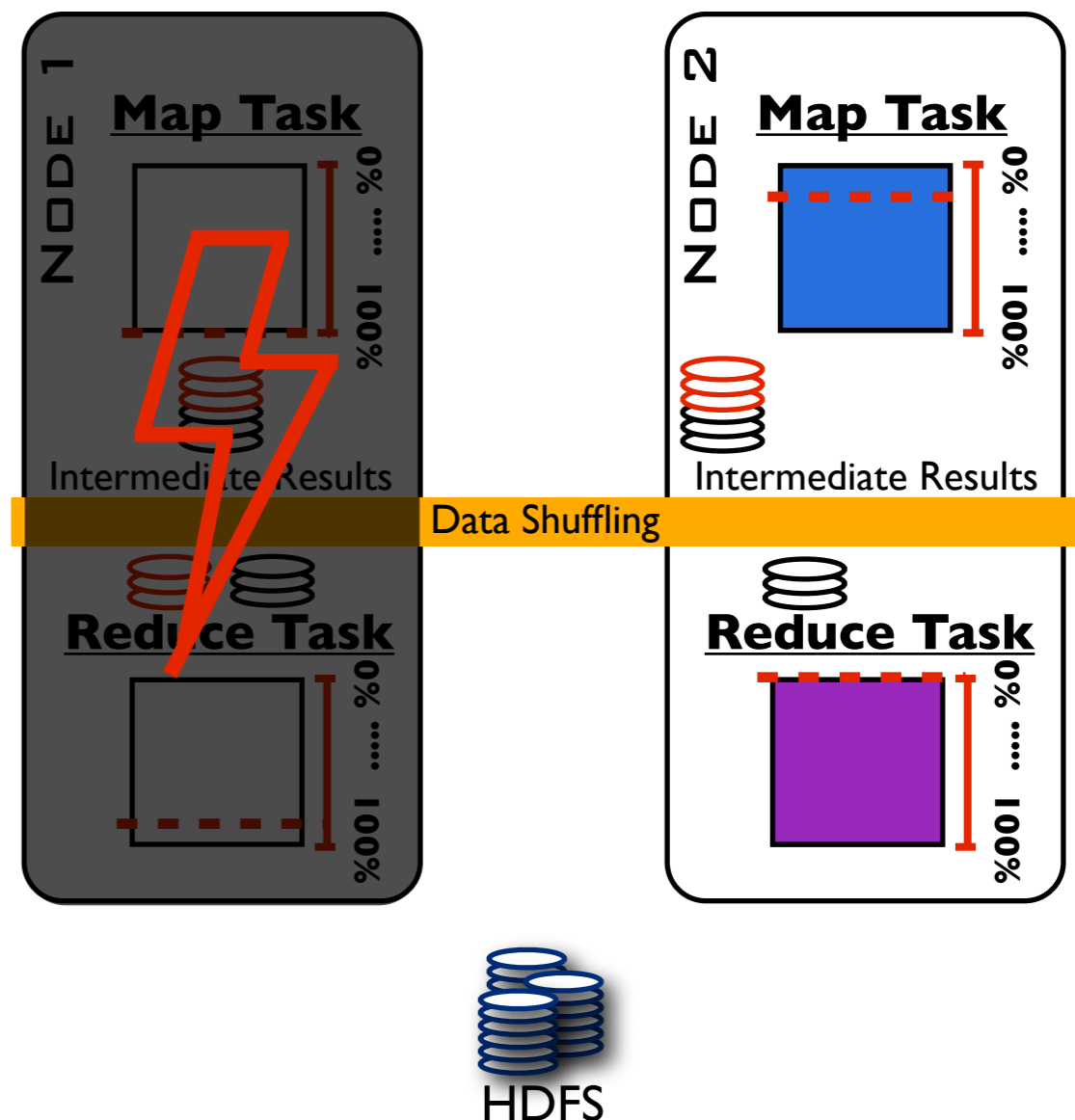
Scheduler



## Hadoop with Node Failures

Time: 50s

Scheduler



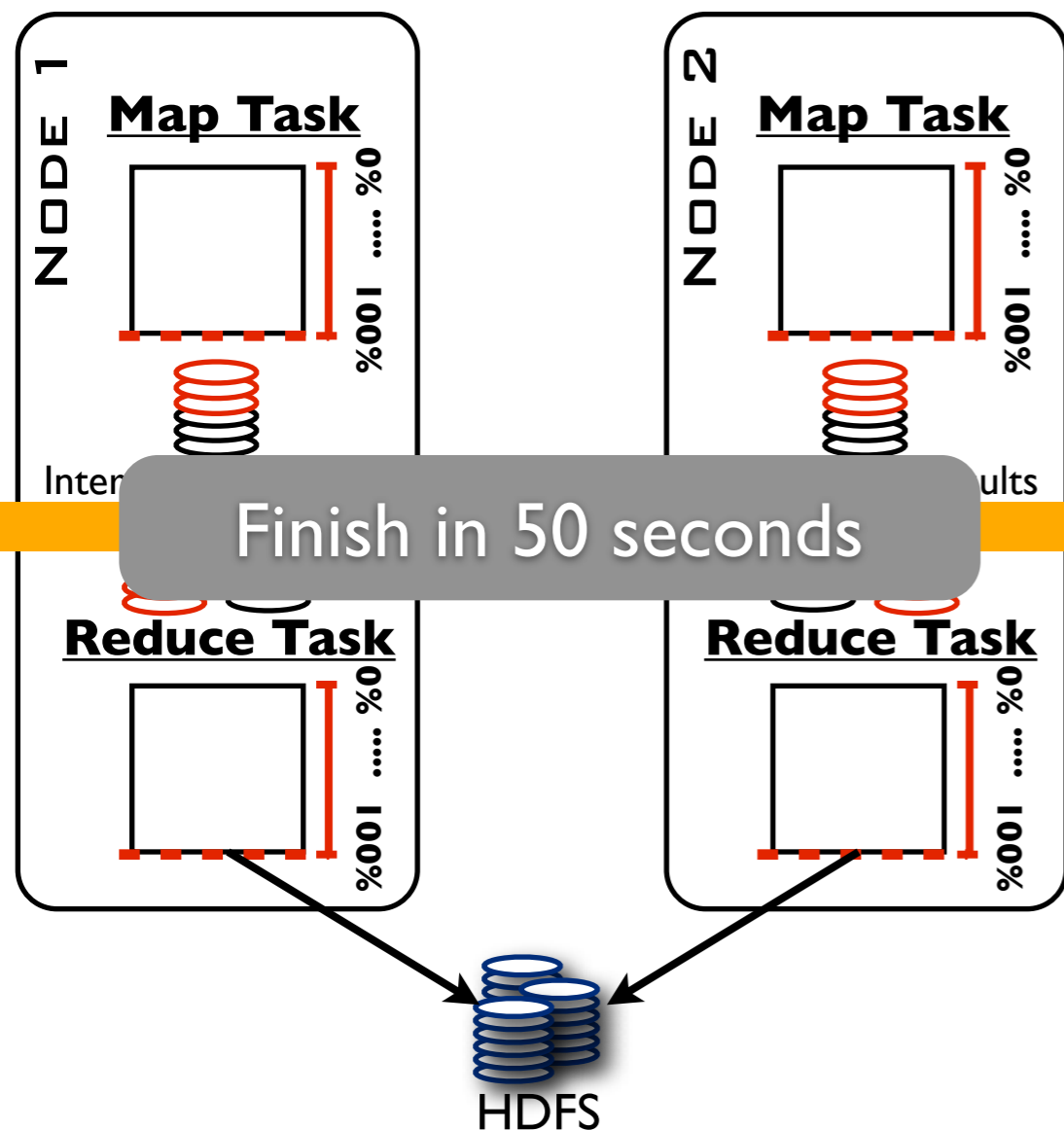
# Current Approach

## Hadoop without Failures

Map Task: 20 seconds  
Reduce Task: 30 seconds

Time: 50s

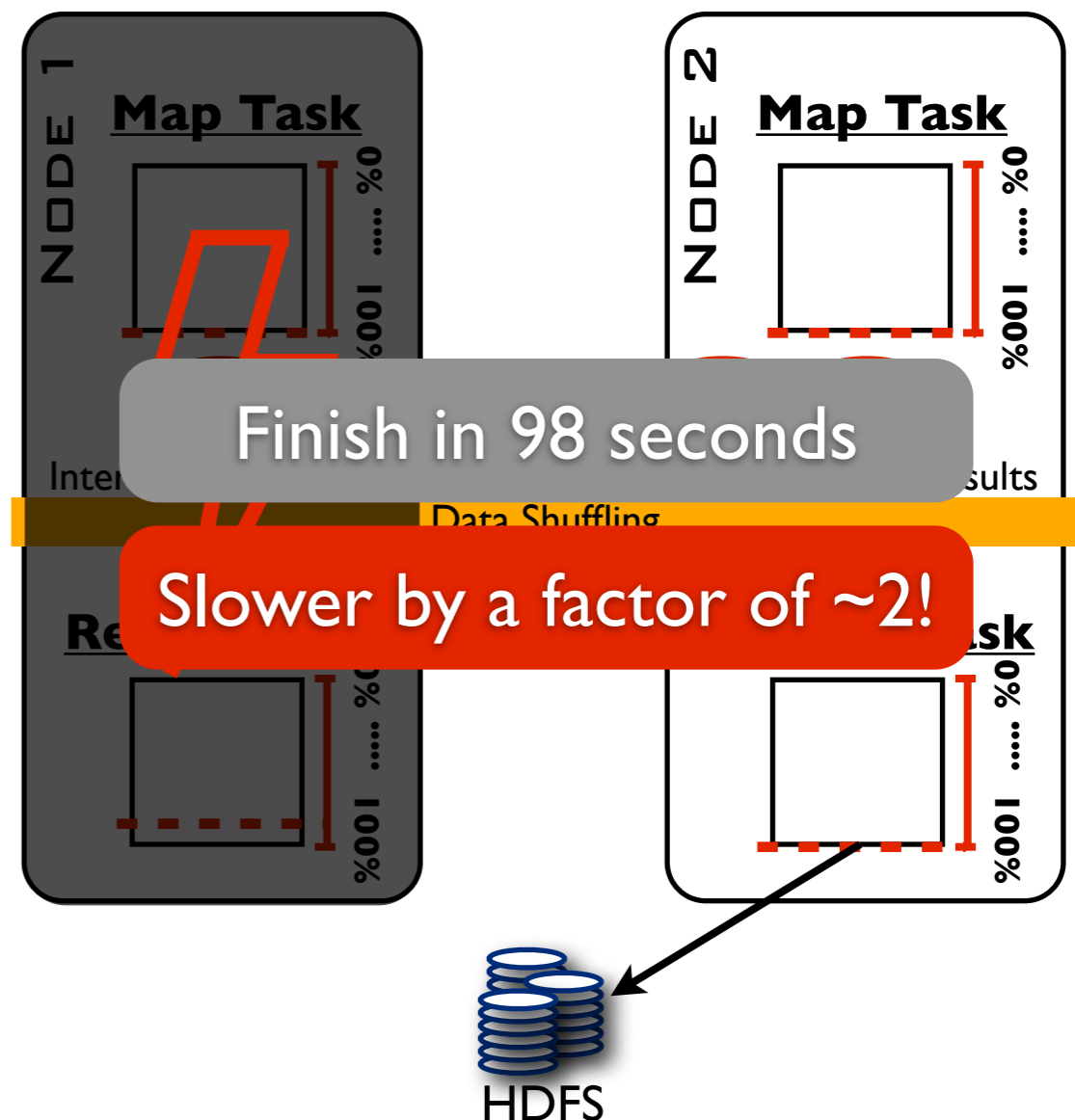
Scheduler



## Hadoop with Node Failures

Time: 98s

Scheduler



# Solution: Query Metadata Checkpointing

---

# Solution: Query Metadata Checkpointing

---

**Ideas:** (1) forward tracing for each key-value pair,

# Solution: Query Metadata Checkpointing

---

- Ideas:** (1) forward tracing for each key-value pair,  
(2) push intermediate results to all reducers

# Solution: Query Metadata Checkpointing

- Ideas:** (1) forward tracing for each key-value pair,  
(2) push intermediate results to all reducers

## QMC: Query Metadata Checkpoint

Offset	Reduce Task ID



# Solution: Query Metadata Checkpointing

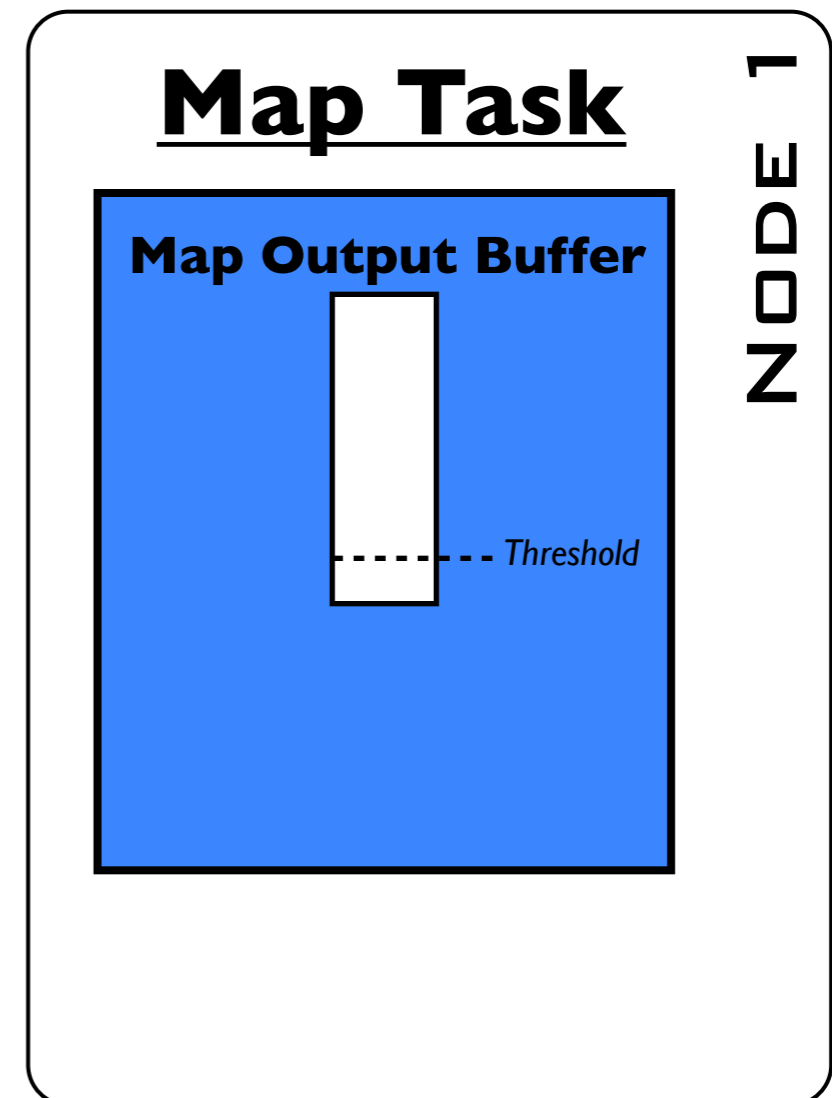
- Ideas:** (1) forward tracing for each key-value pair,  
(2) push intermediate results to all reducers

## Input Split

Offset	Contents
⋮	⋮
6500	www.bbc.com, 125.65.10.77, ...
6501	skysports.com, 125.65.10.77, ...
⋮	⋮
10000	www.msn.com, 125.65.10.77, ...
10001	www.abc.co.uk, 125.65.10.77, ...
⋮	⋮

## QMC: Query Metadata Checkpoint

Offset	Reduce Task ID



# Solution: Query Metadata Checkpointing

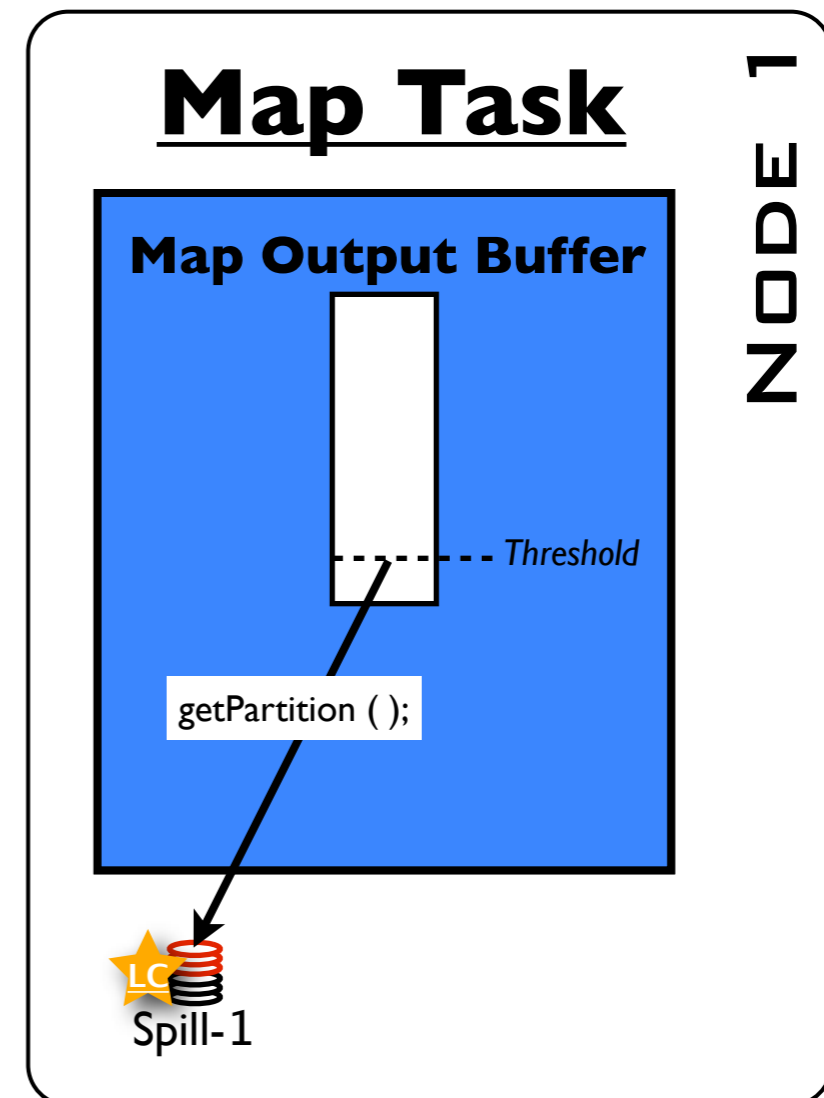
- Ideas:** (1) forward tracing for each key-value pair,  
(2) push intermediate results to all reducers

## Input Split

Offset	Contents
⋮	⋮
6500	www.bbc.com, 125.65.10.77, ...
6501	skysports.com, 125.65.10.77, ...
⋮	⋮
10000	www.msn.com, 125.65.10.77, ...
10001	www.abc.co.uk, 125.65.10.77, ...
⋮	⋮

## QMC: Query Metadata Checkpoint

Offset	Reduce Task ID





# Solution: Query Metadata Checkpointing

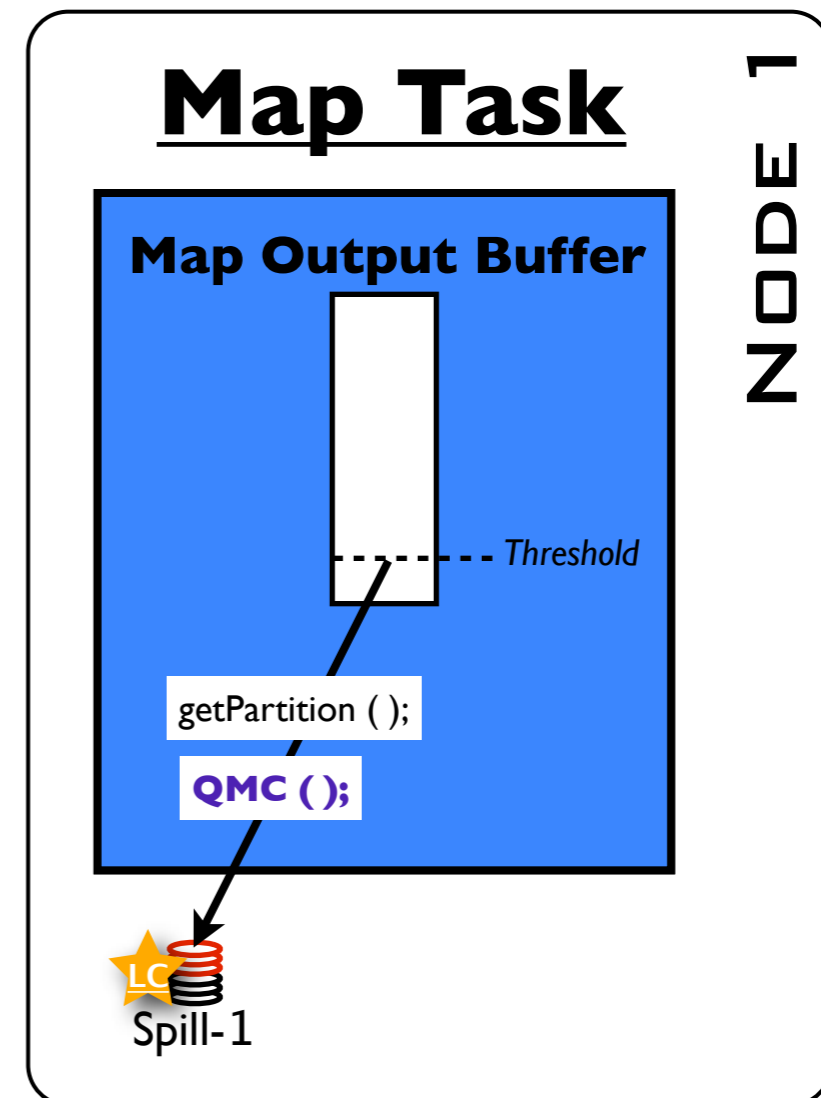
- Ideas:** (1) forward tracing for each key-value pair,  
(2) push intermediate results to all reducers

## Input Split

Offset	Contents
⋮	⋮
6500	www.bbc.com, 125.65.10.77, ...
6501	skysports.com, 125.65.10.77, ...
⋮	⋮
10000	www.msn.com, 125.65.10.77, ...
10001	www.abc.co.uk, 125.65.10.77, ...
⋮	⋮

## QMC: Query Metadata Checkpoint

Offset	Reduce Task ID
⋮	⋮
6500	Red-1
6501	Red-2



# Solution: Query Metadata Checkpointing

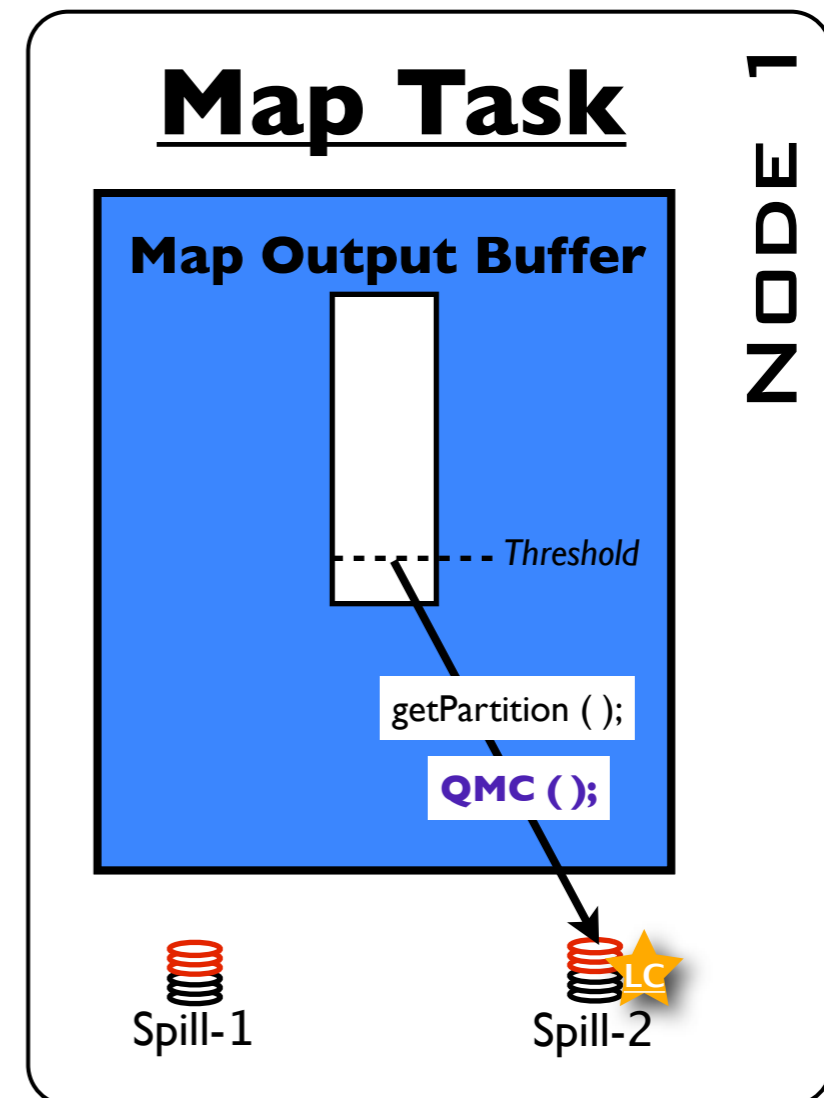
- Ideas:** (1) forward tracing for each key-value pair,  
(2) push intermediate results to all reducers

## Input Split

Offset	Contents
⋮	⋮
6500	www.bbc.com, 125.65.10.77, ...
6501	skysports.com, 125.65.10.77, ...
⋮	⋮
10000	www.msn.com, 125.65.10.77, ...
10001	www.abc.co.uk, 125.65.10.77, ...
⋮	⋮

## QMC: Query Metadata Checkpoint

Offset	Reduce Task ID
⋮	⋮
6500	Red-1
6501	Red-2
⋮	⋮
10000	Red-2
10001	Red-1



# Solution: Query Metadata Checkpointing

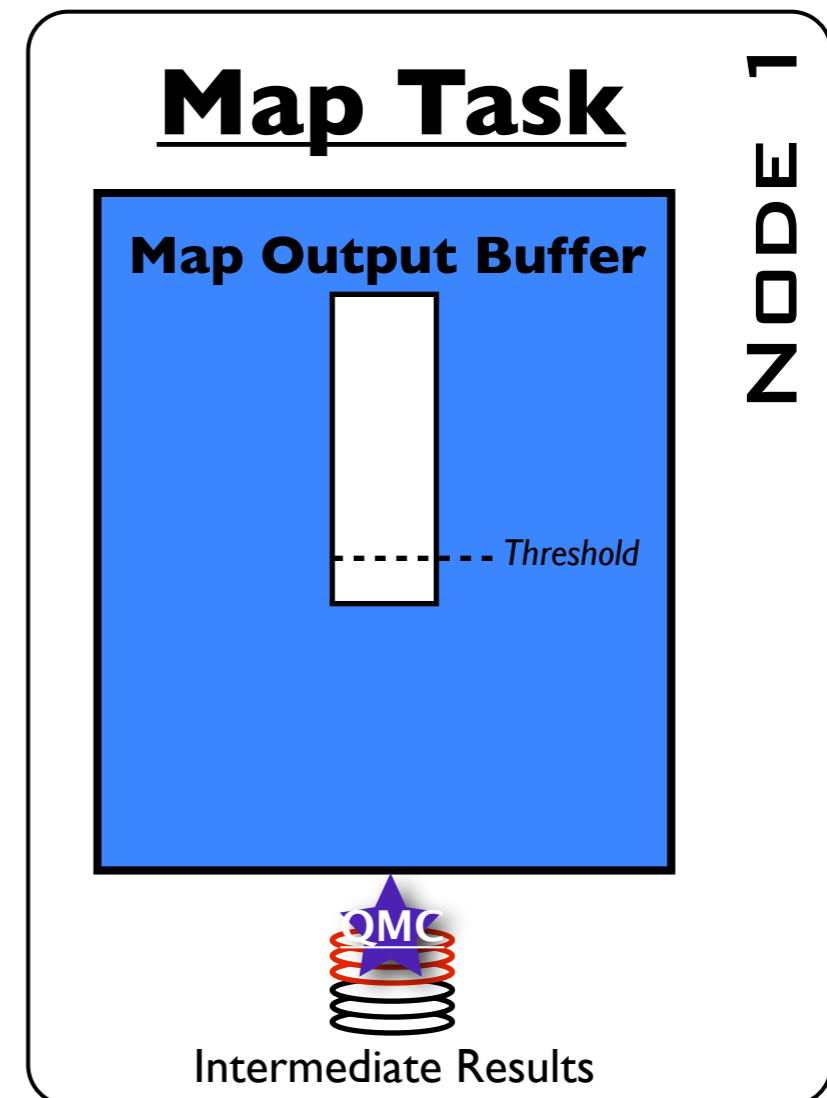
- Ideas:** (1) forward tracing for each key-value pair,  
(2) push intermediate results to all reducers

## Input Split

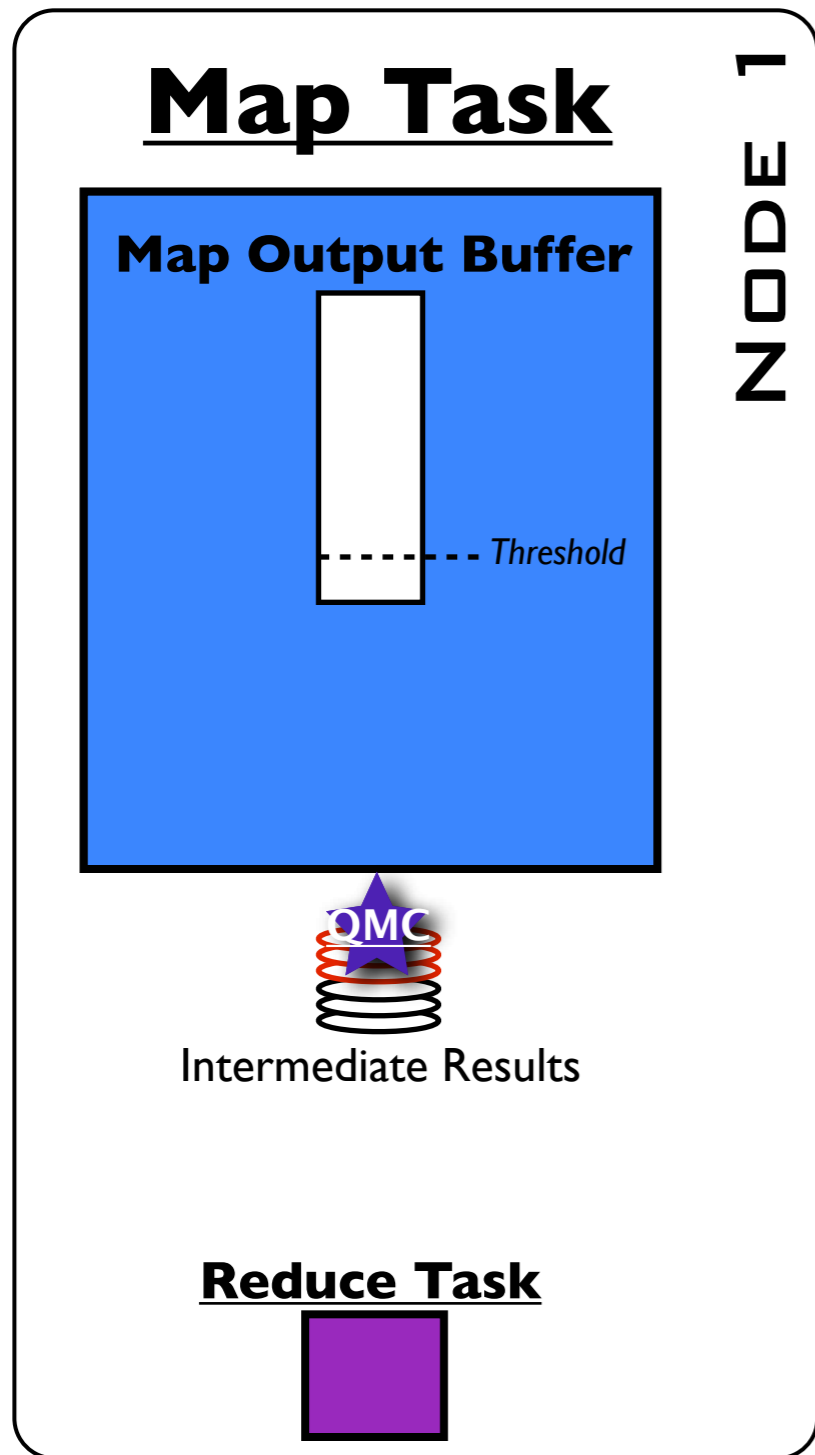
Offset	Contents
⋮	⋮
6500	www.bbc.com, 125.65.10.77, ...
6501	skysports.com, 125.65.10.77, ...
⋮	⋮
10000	www.msn.com, 125.65.10.77, ...
10001	www.abc.co.uk, 125.65.10.77, ...
⋮	⋮

## QMC: Query Metadata Checkpoint

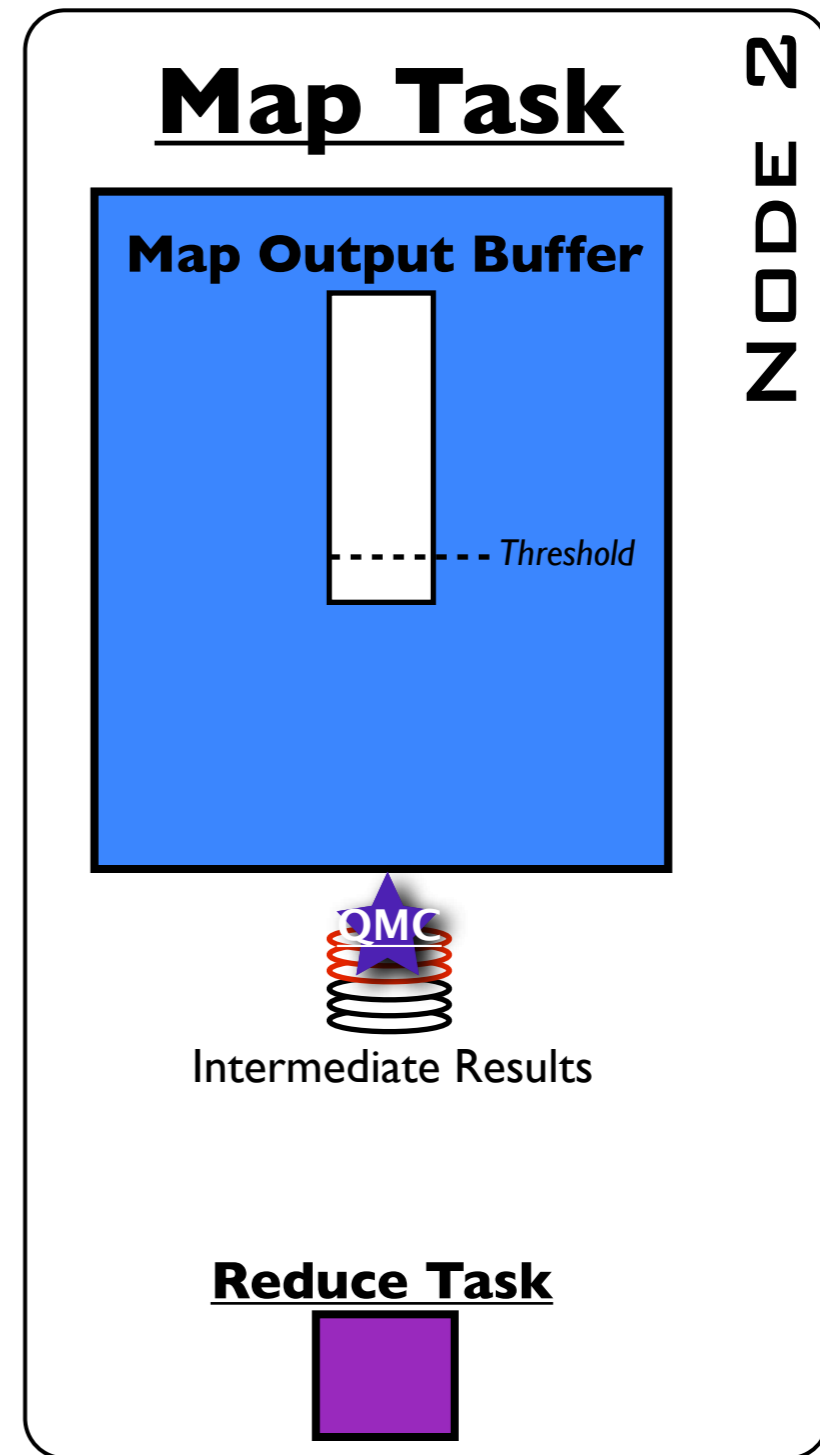
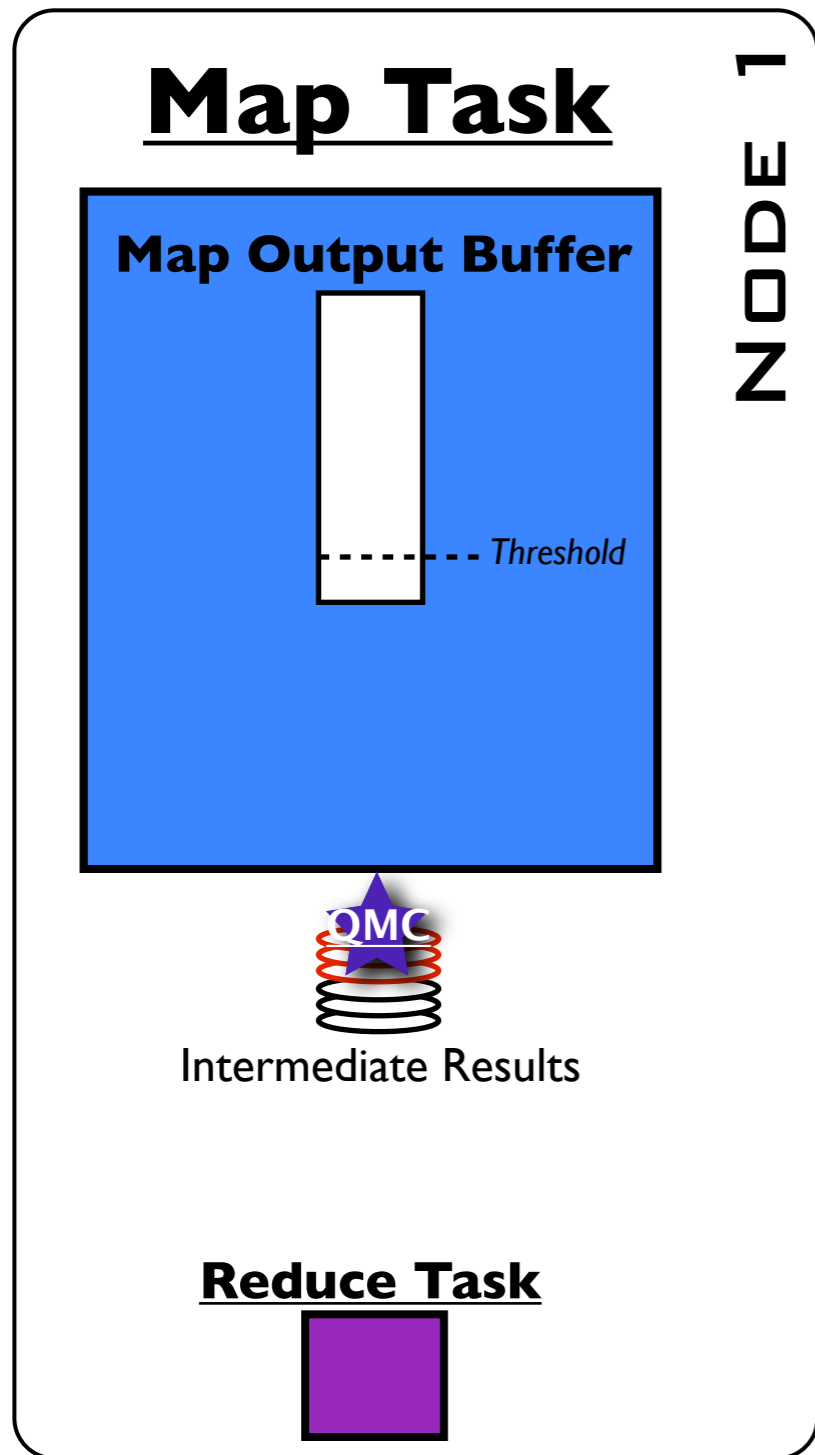
Offset	Reduce Task ID
⋮	⋮
6500	Red-1
6501	Red-2
⋮	⋮
10000	Red-2
10001	Red-1
⋮	⋮



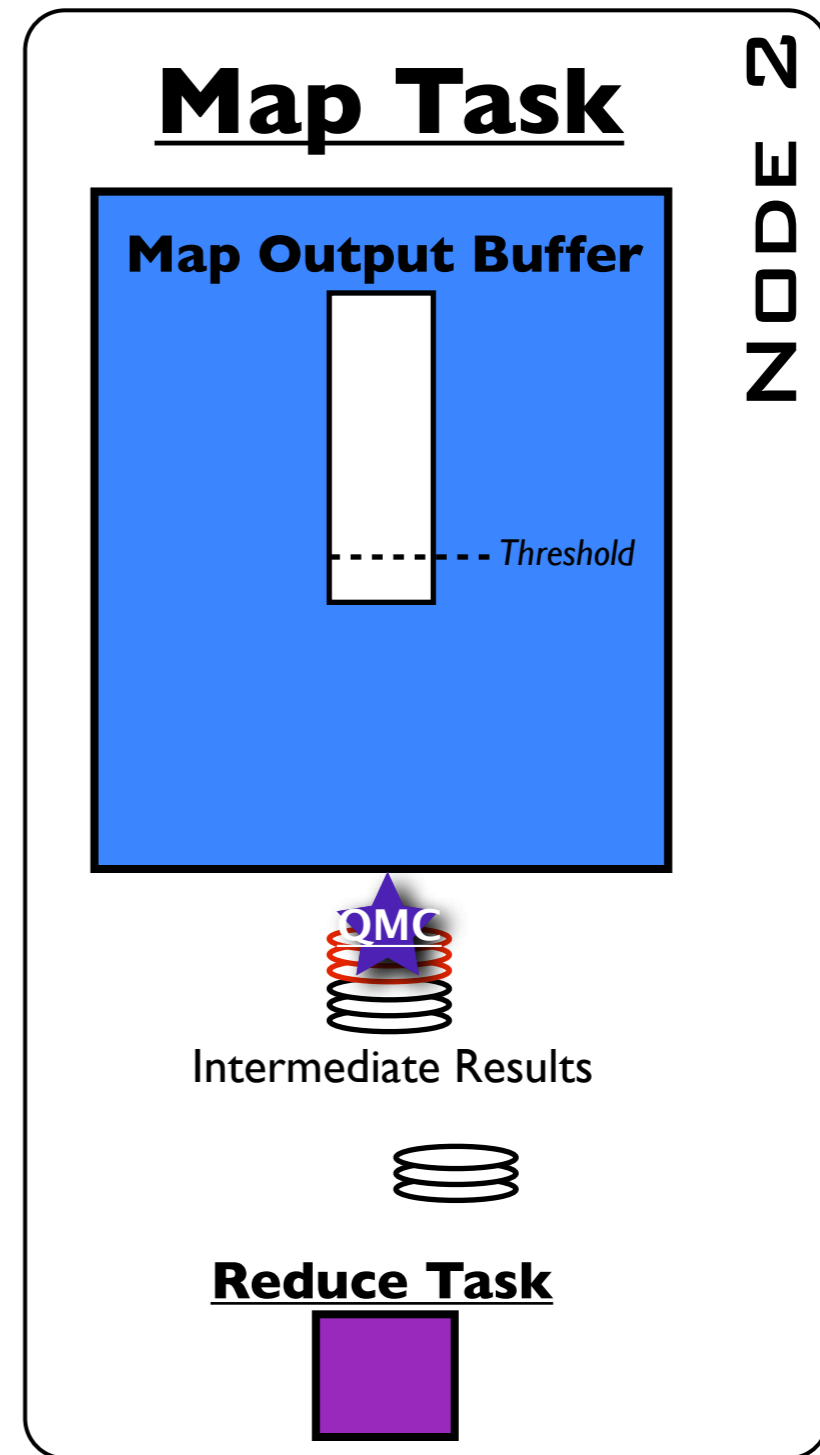
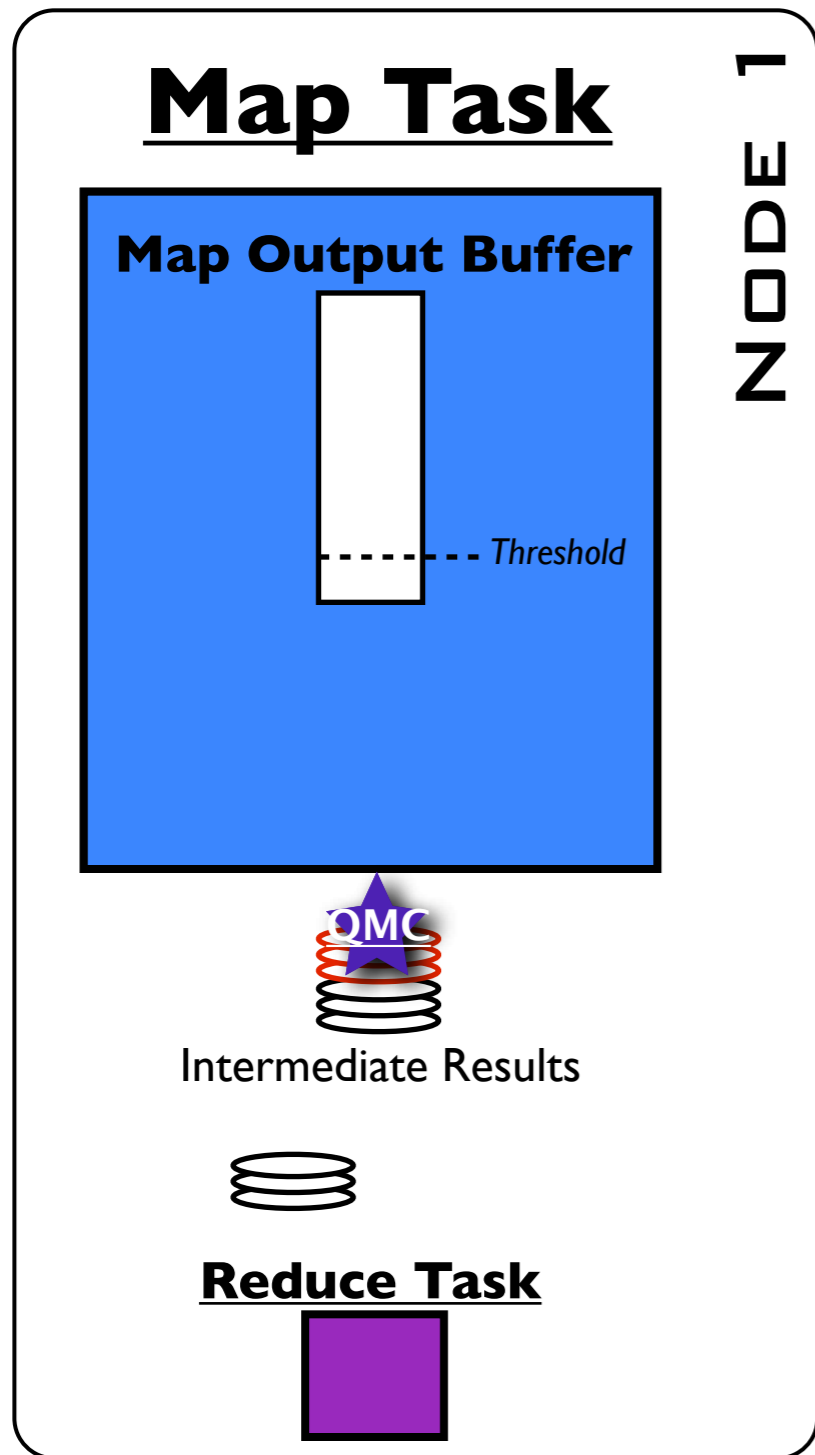
# Replicating QMC vs Local Partition



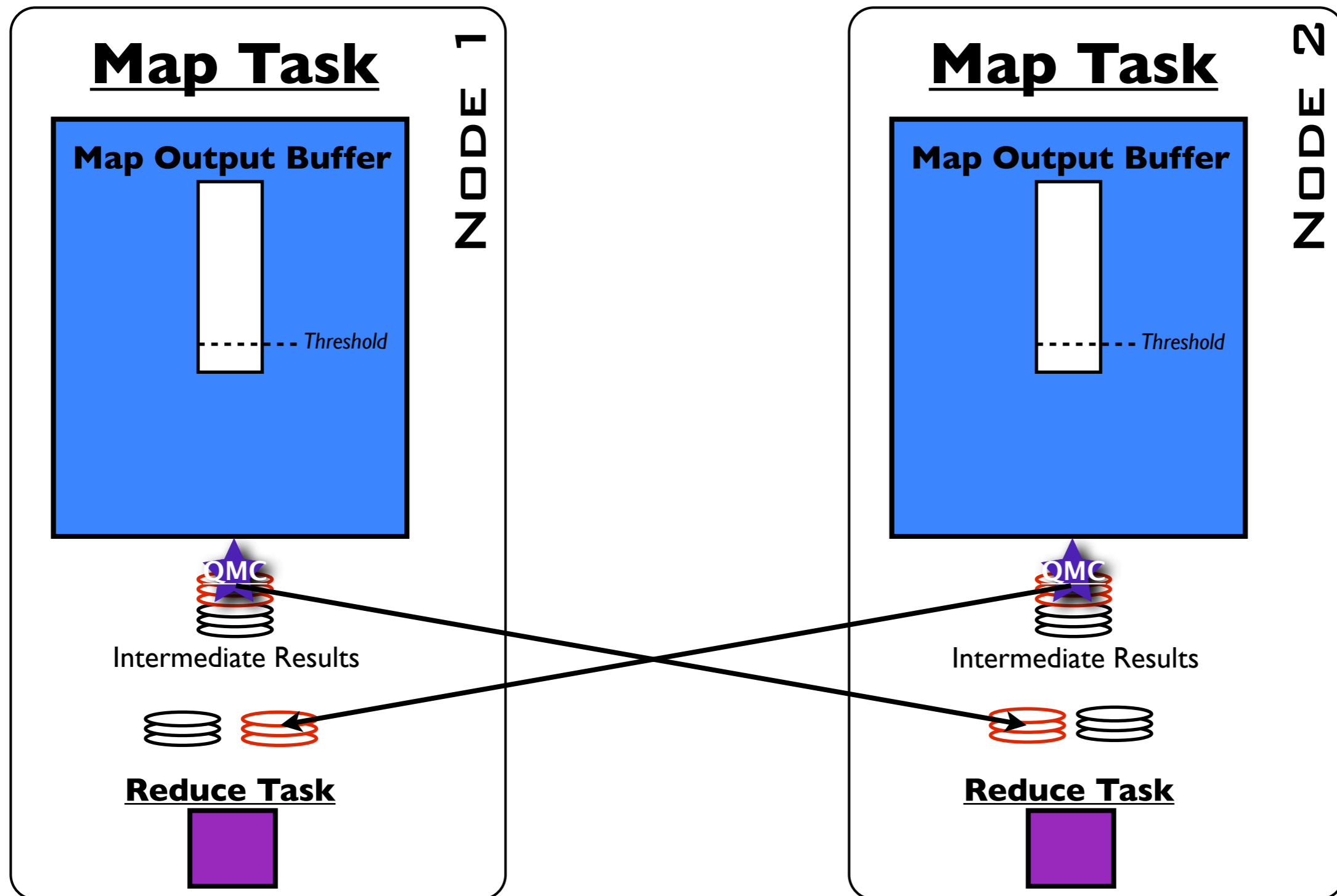
# Replicating QMC vs Local Partition



# Replicating QMC vs Local Partition

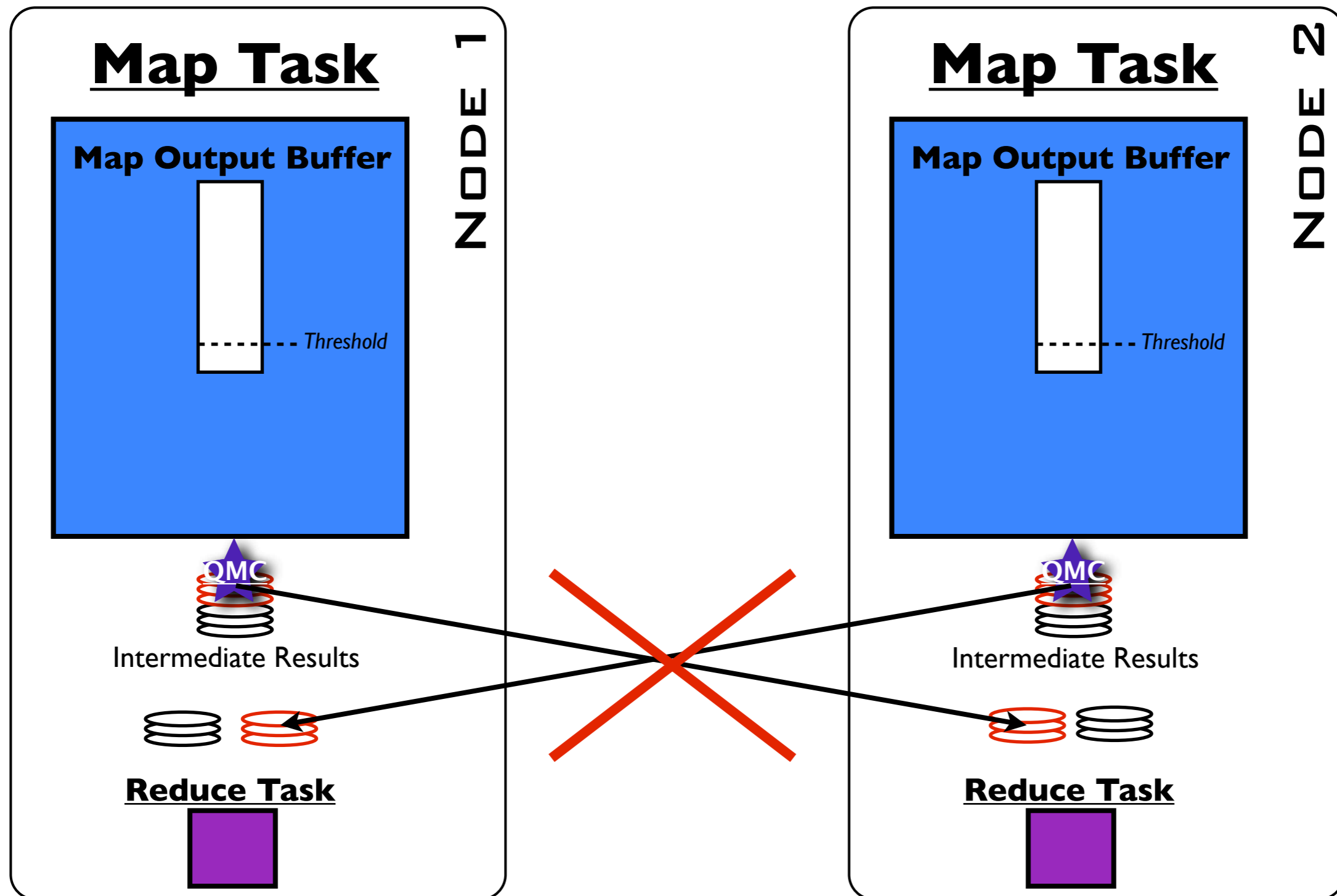


# Replicating QMC vs Local Partition



[MapReduce Online, NSDI'10]

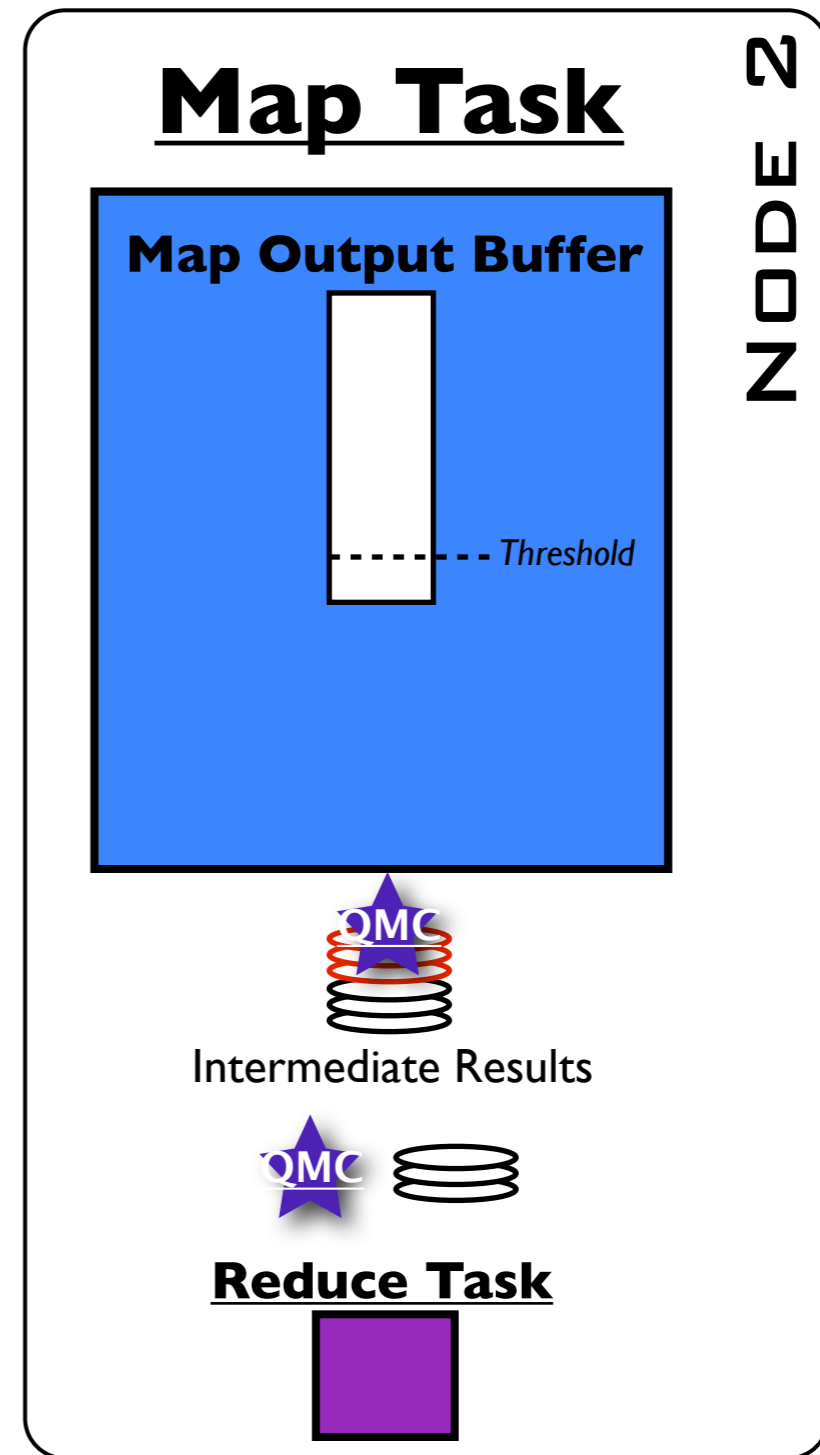
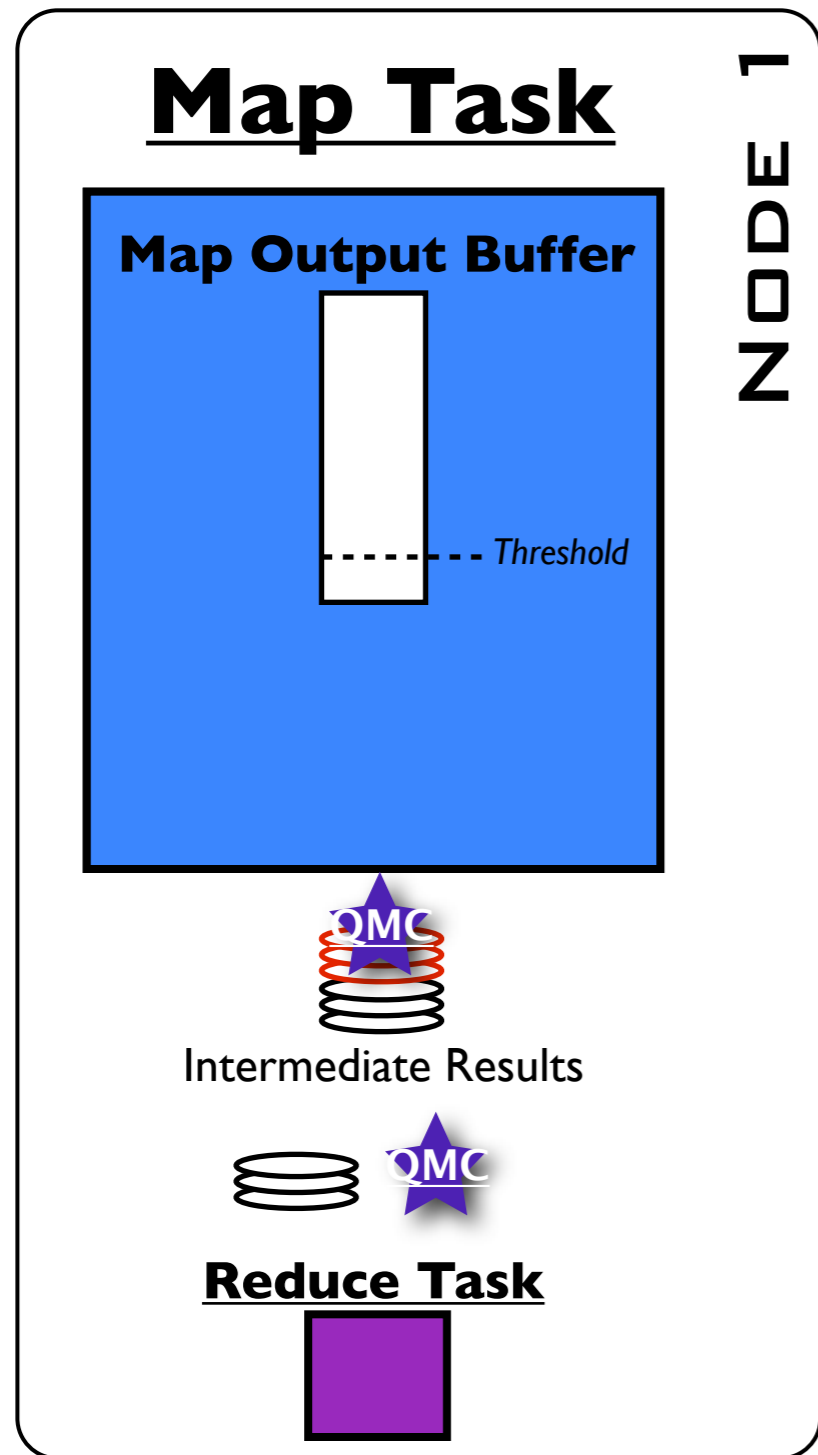
# Replicating QMC vs Local Partition



[MapReduce Online, NSDI'10]



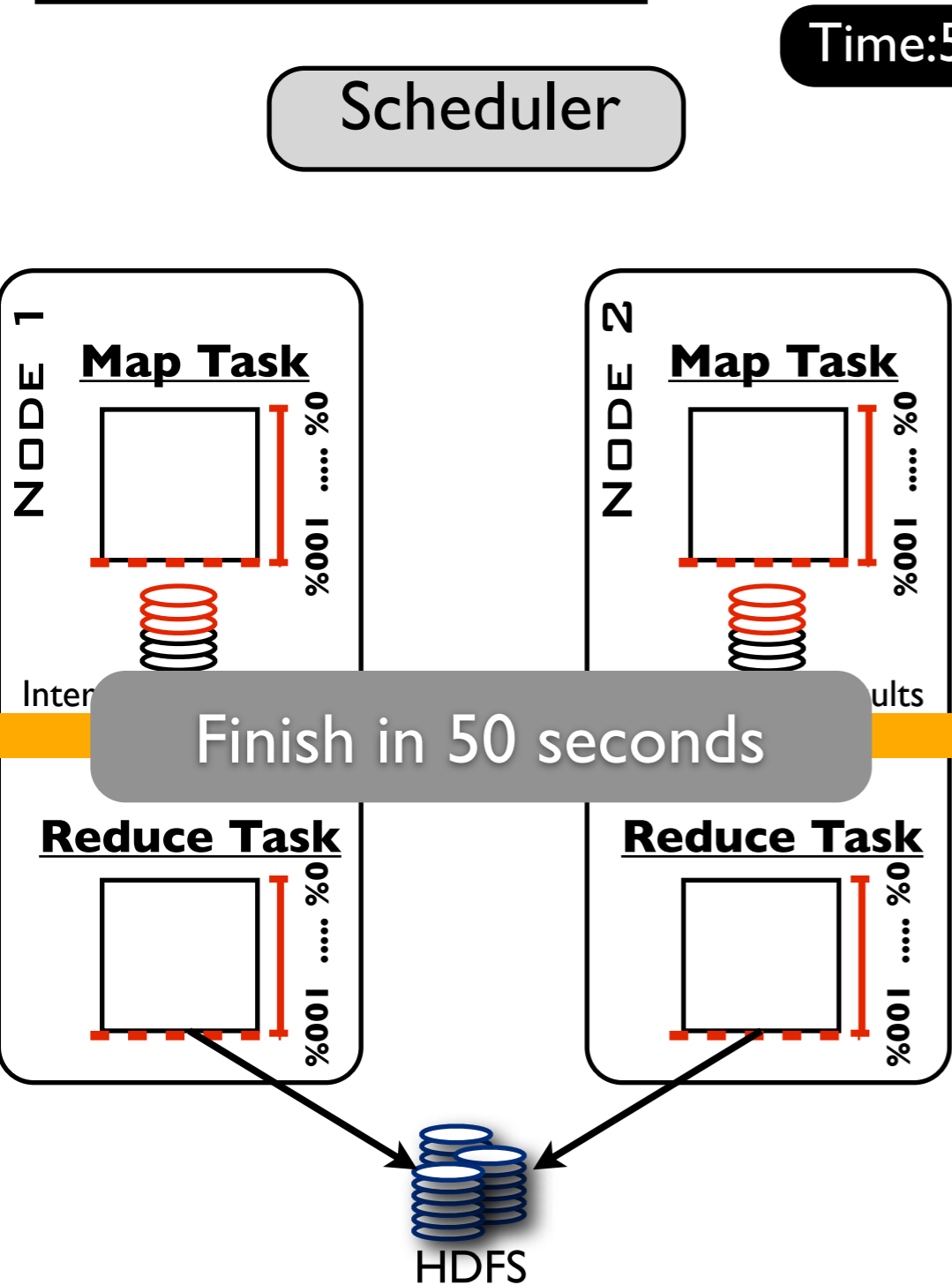
# Replicating QMC vs Local Partition



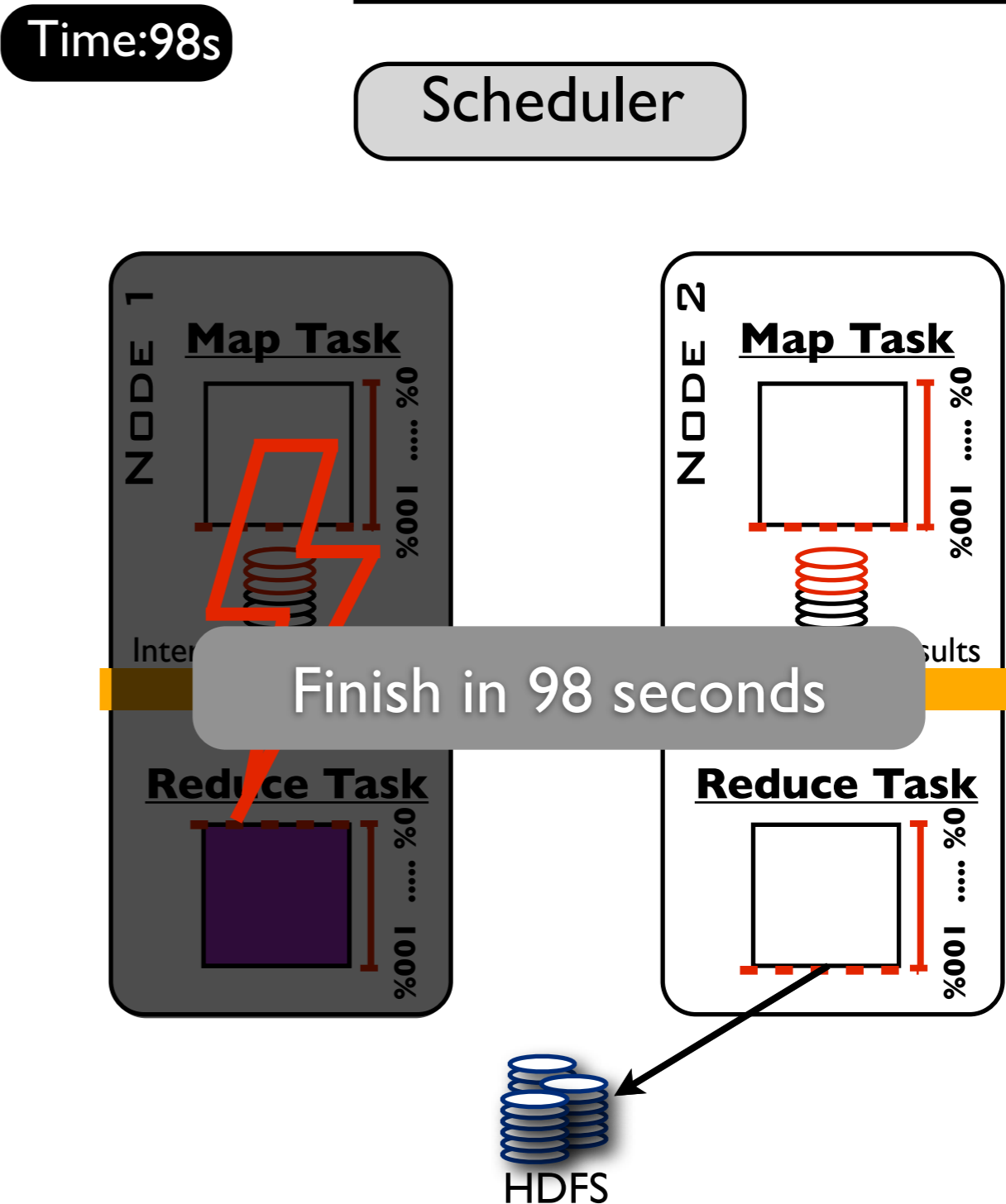
[MapReduce Online, NSDI'10]

# Recovery from Node Failures

## Hadoop without Failures



## Hadoop with Node Failures



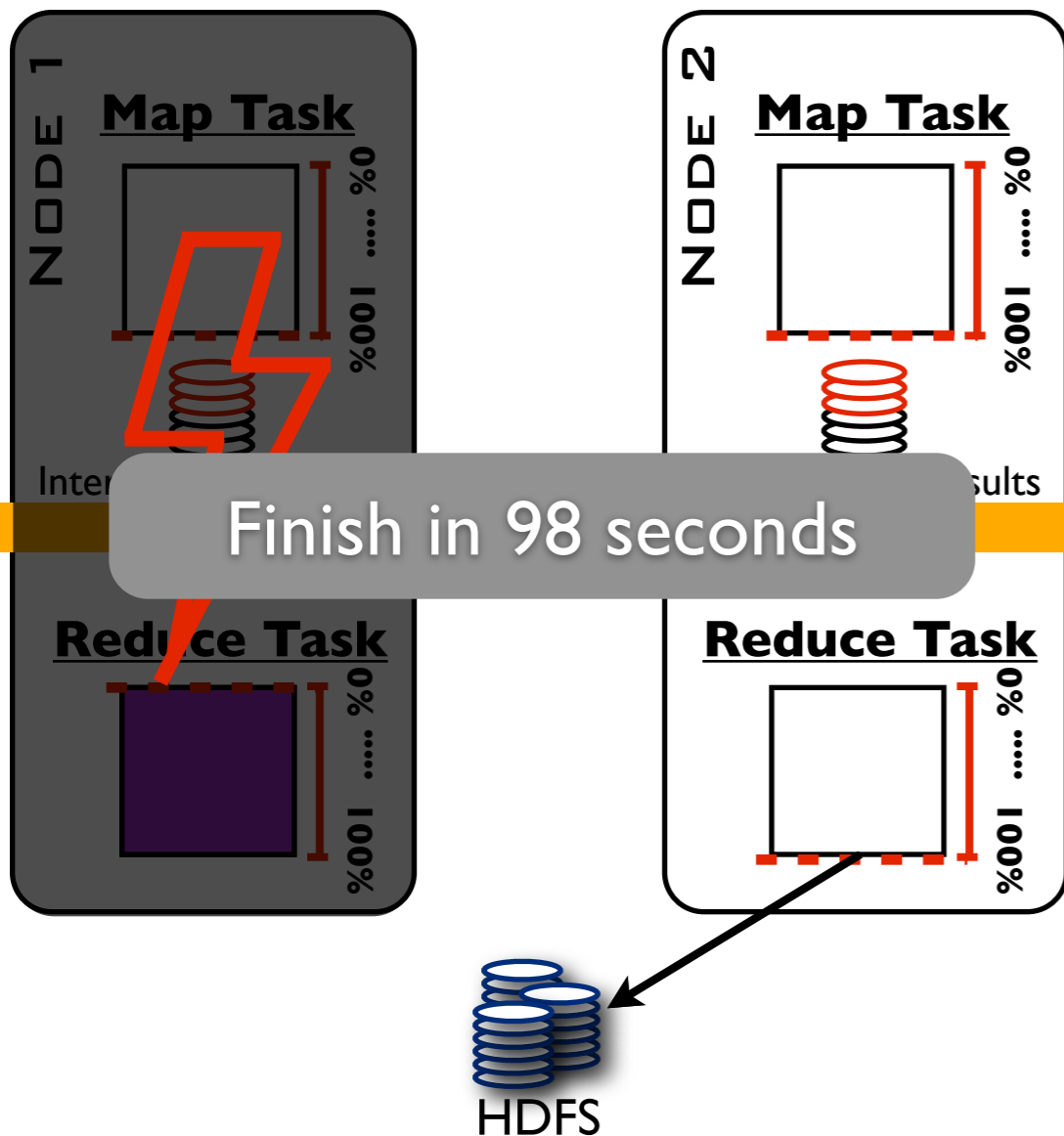
# Recovery from Node Failures

Hadoop

Map Task: 20 seconds  
Reduce Task: 30 seconds

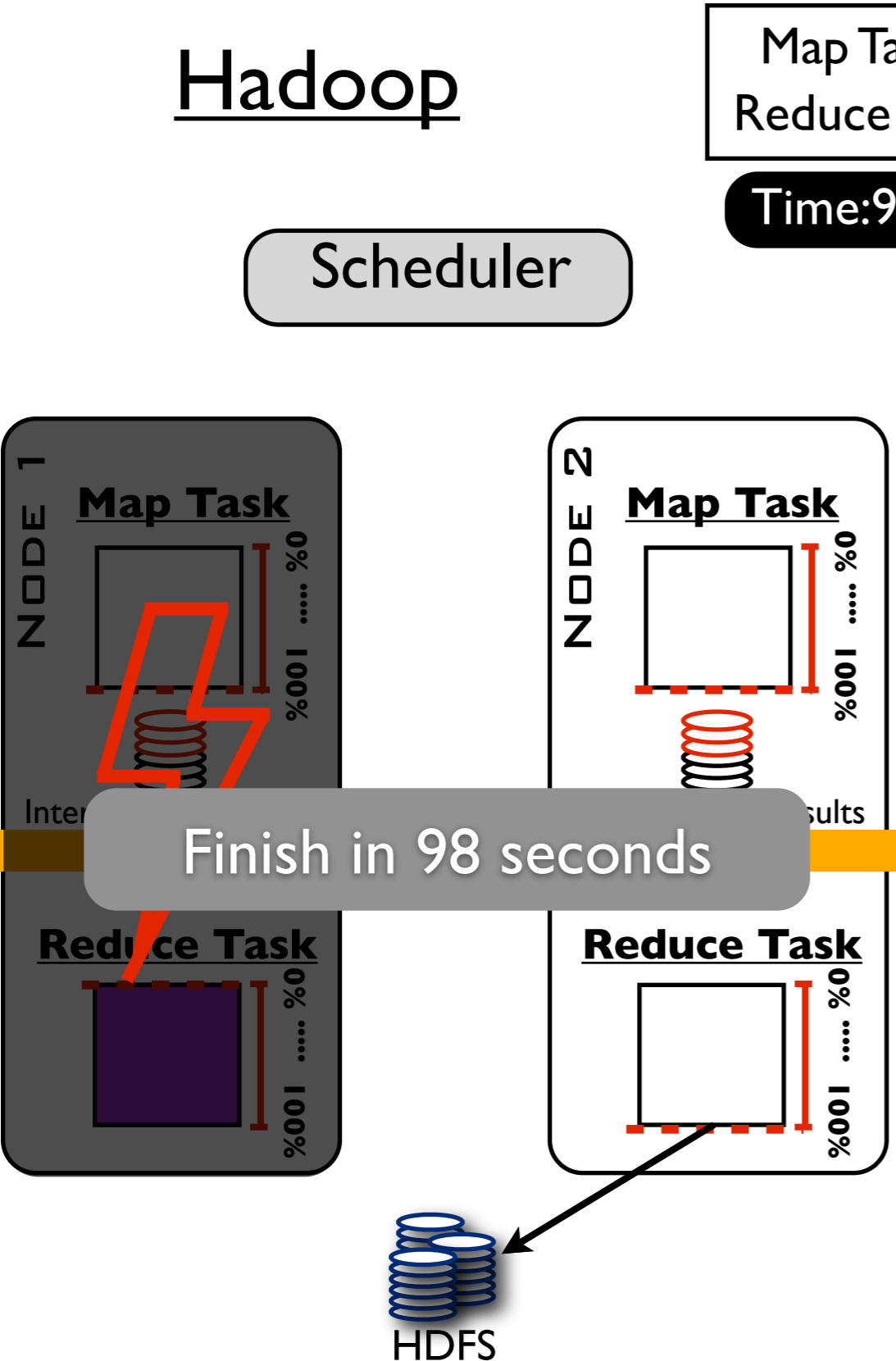
Time: 98s

Scheduler

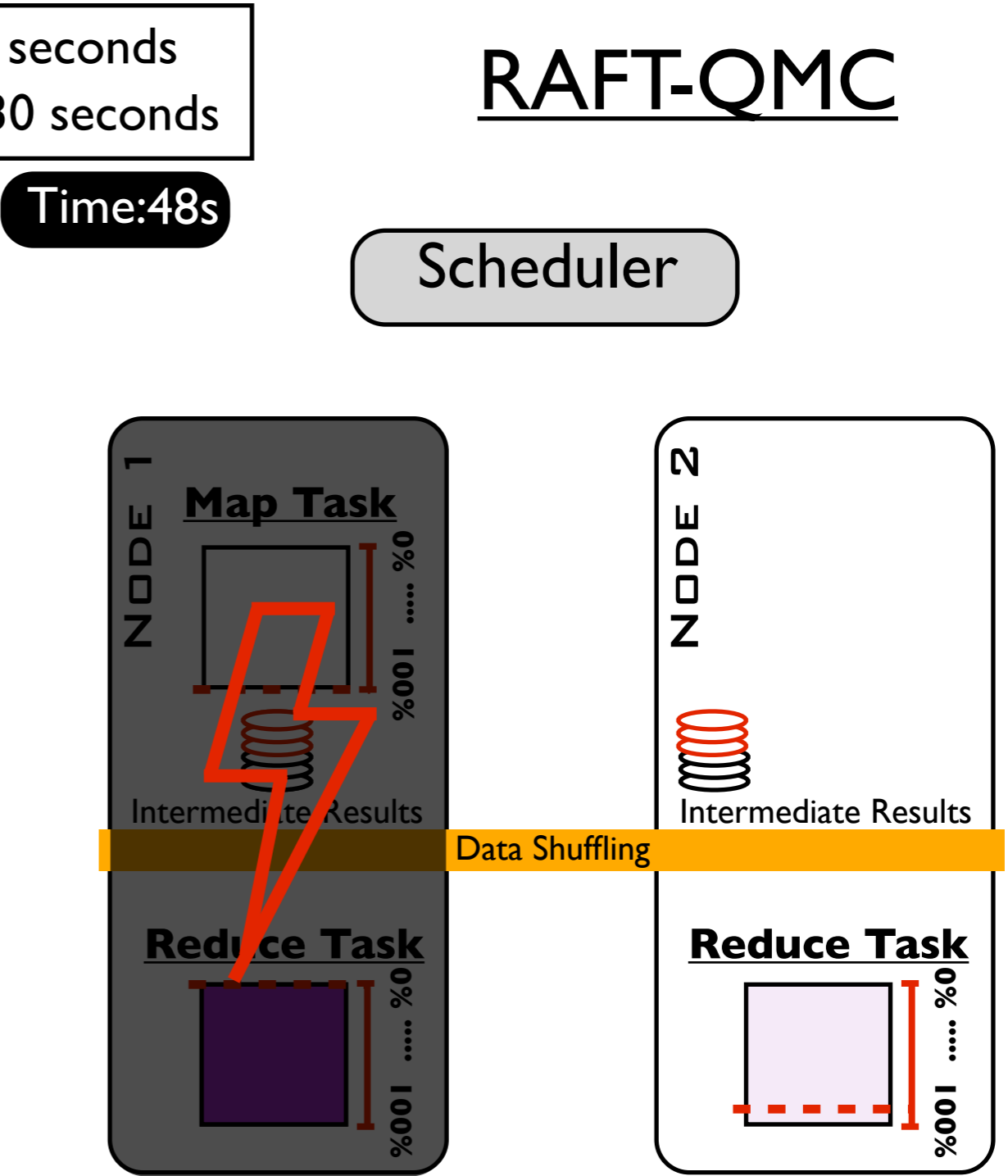


# Recovery from Node Failures

## Hadoop



## RAFT-QMC



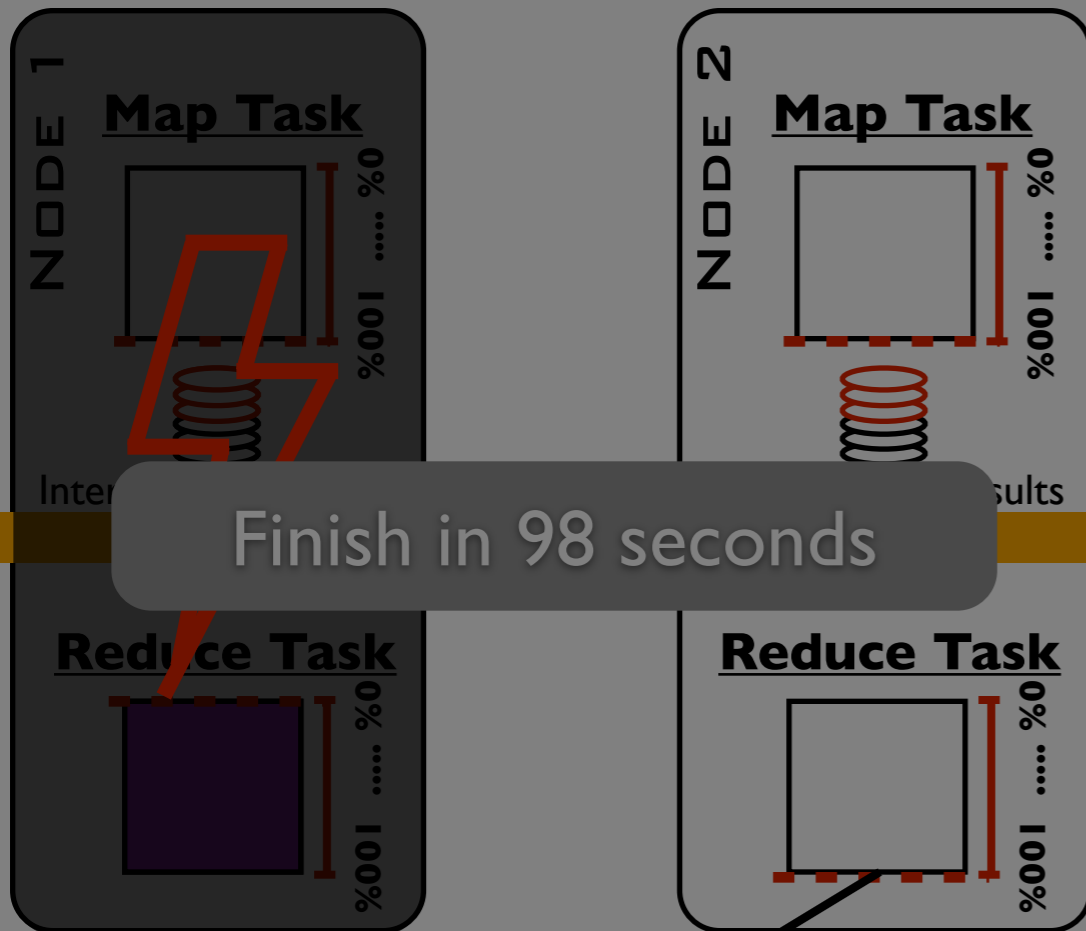
# Recovery from Node Failures

## Hadoop

Map Task: 20 seconds  
Reduce Task: 30 seconds

Scheduler

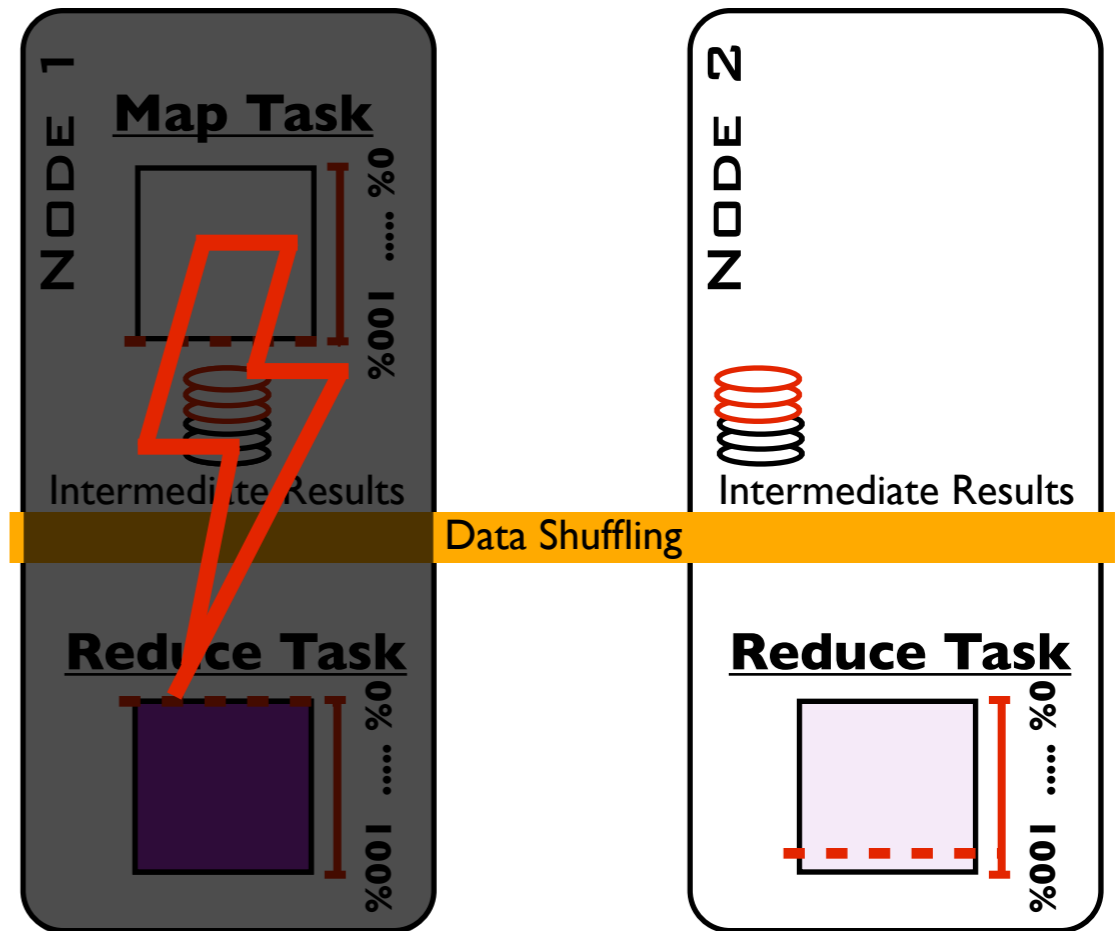
Time:98s



## RAFT-QMC

Time:48s

Scheduler



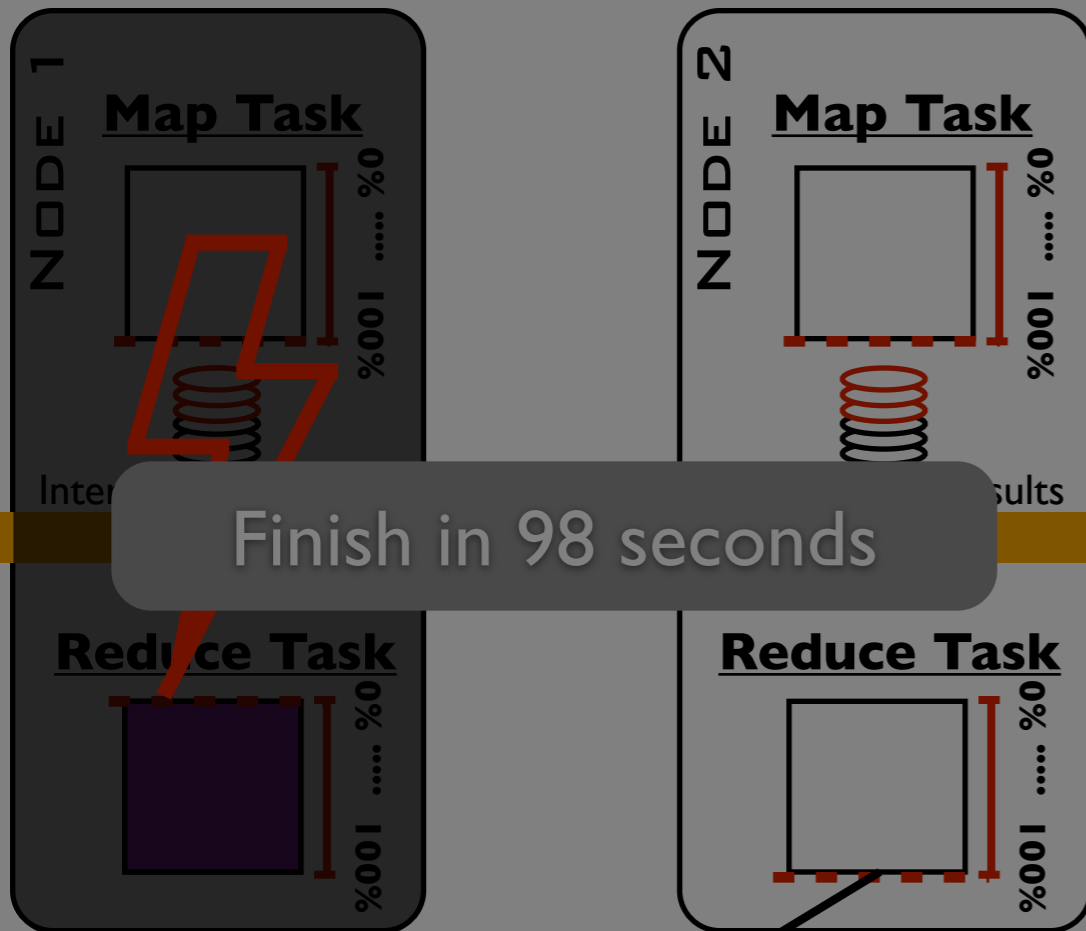
# Recovery from Node Failures

## Hadoop

Map Task: 20 seconds  
Reduce Task: 30 seconds

Scheduler

Time:98s

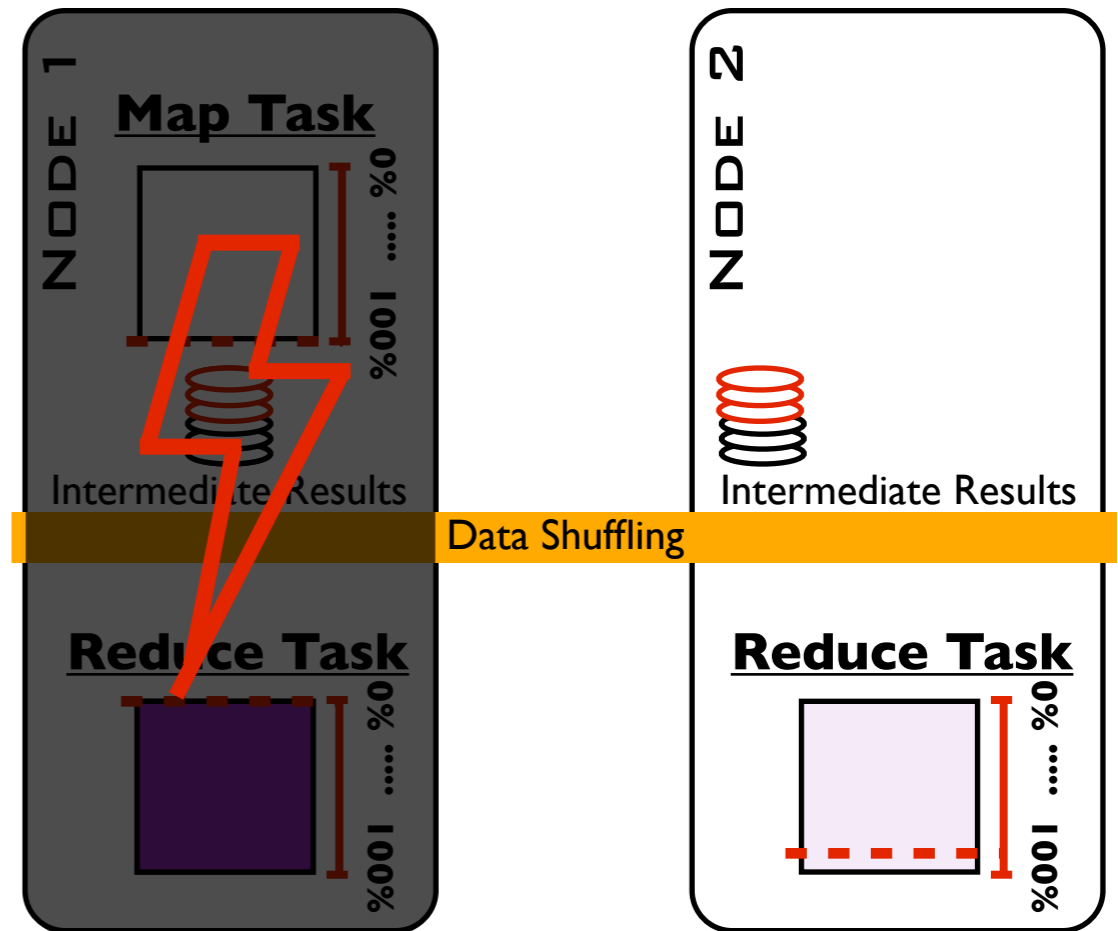


HDFS

## RAFT-QMC

Scheduler

Time:48s



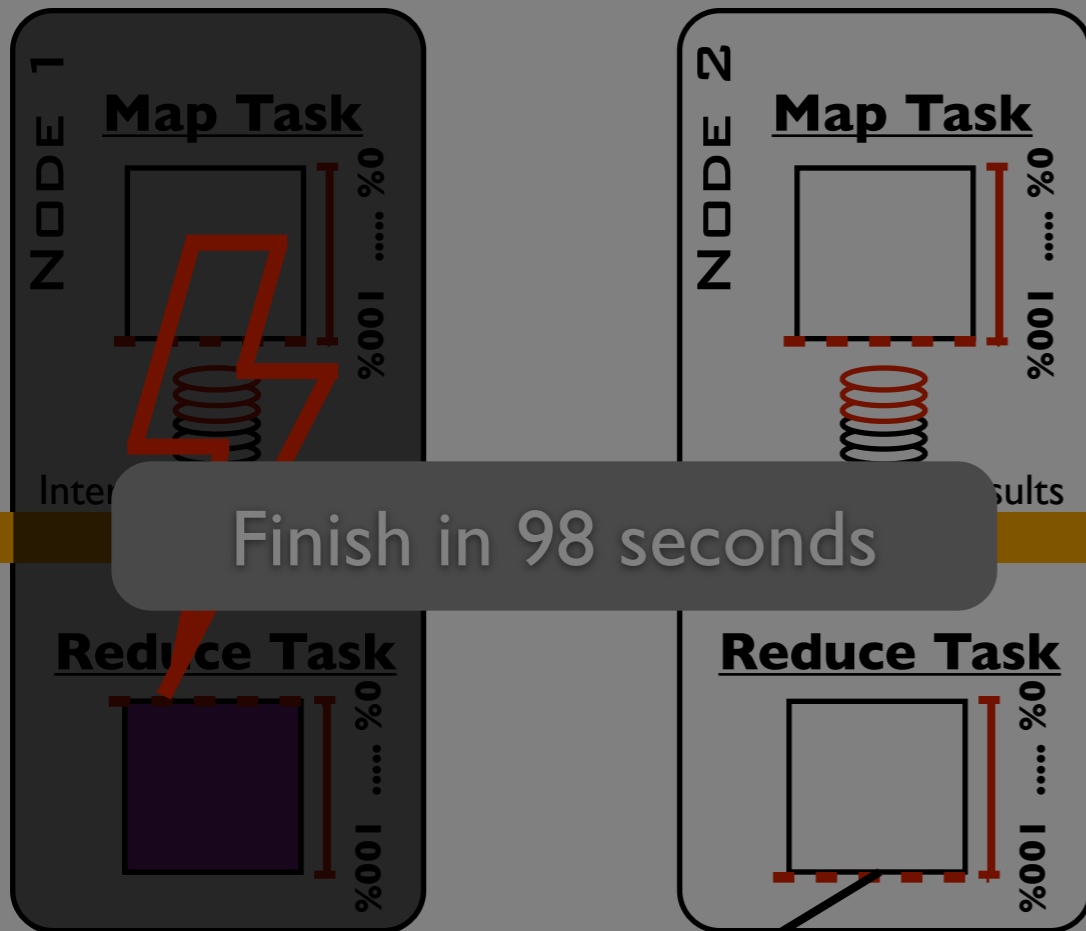
# Recovery from Node Failures

## Hadoop

Map Task: 20 seconds  
Reduce Task: 30 seconds

Scheduler

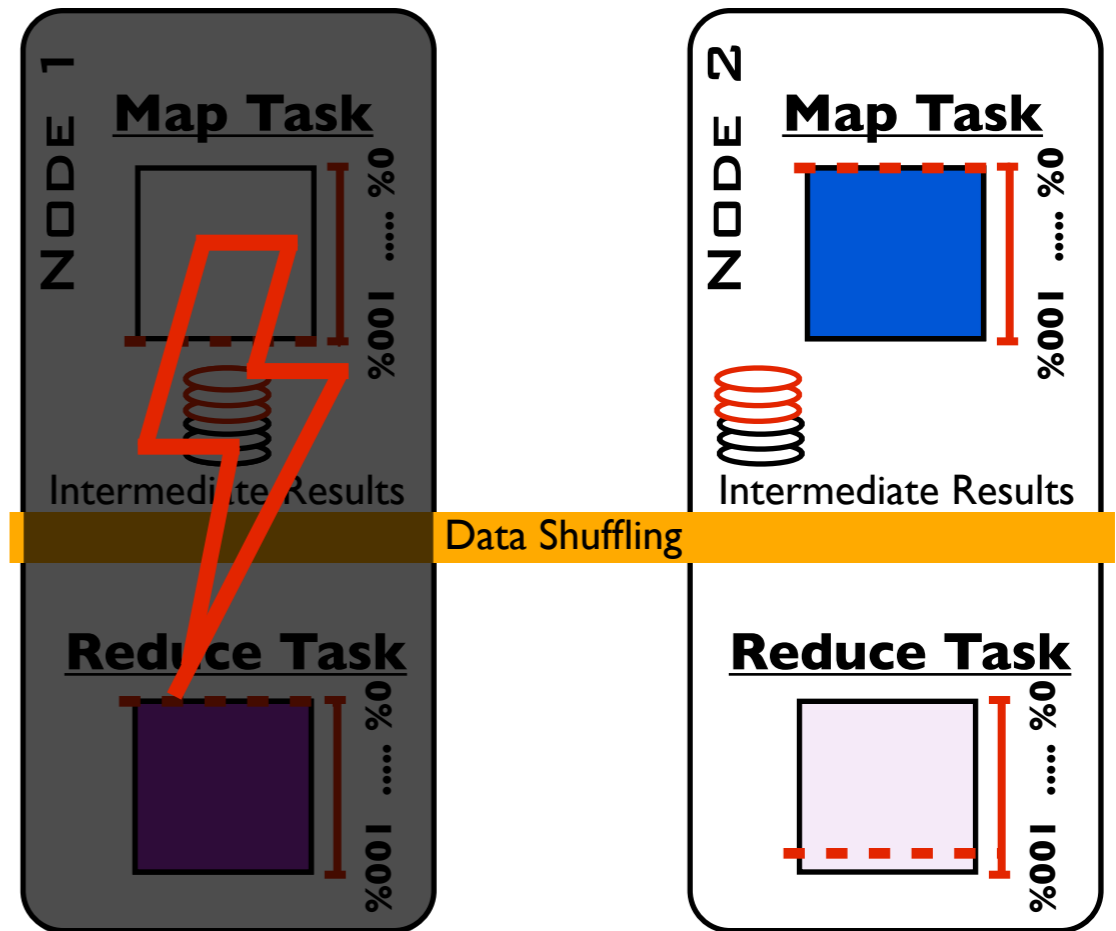
Time:98s



## RAFT-QMC

Time:48s

Scheduler



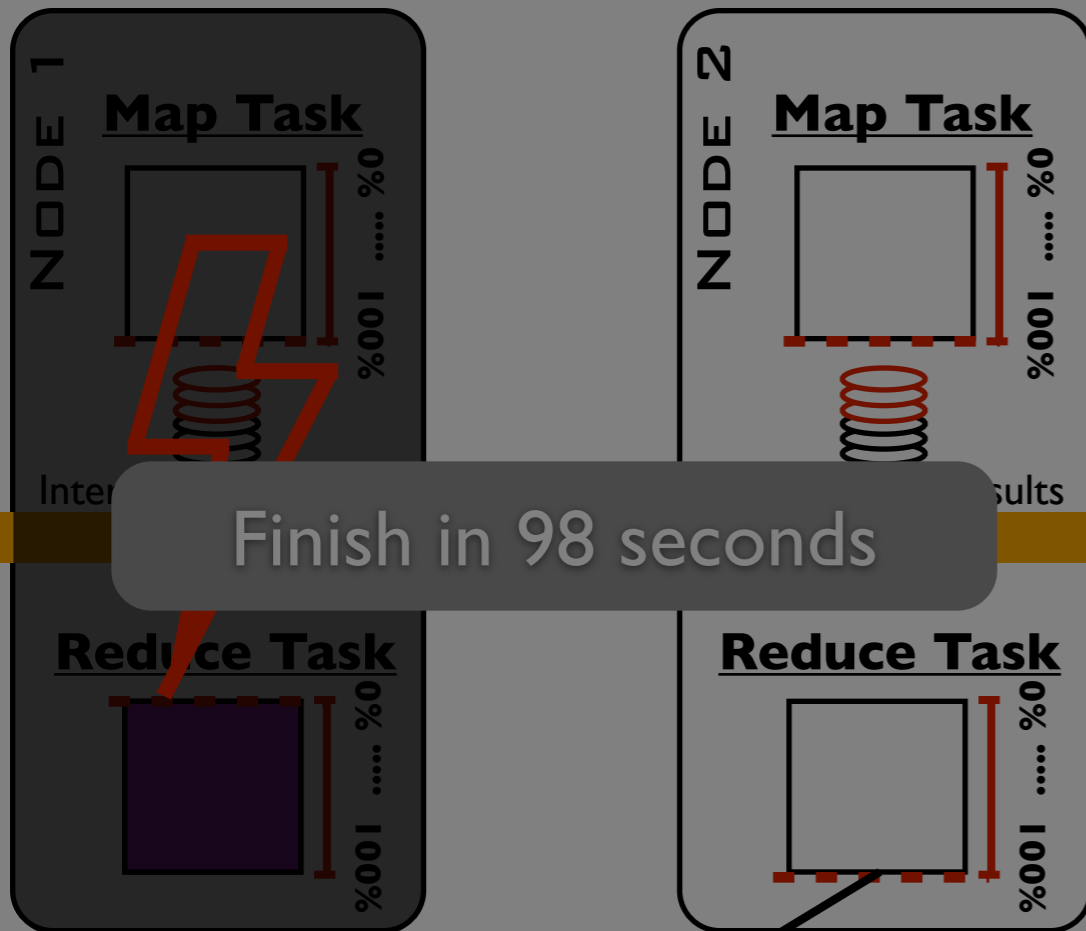
# Recovery from Node Failures

## Hadoop

Map Task: 20 seconds  
Reduce Task: 30 seconds

Scheduler

Time:98s

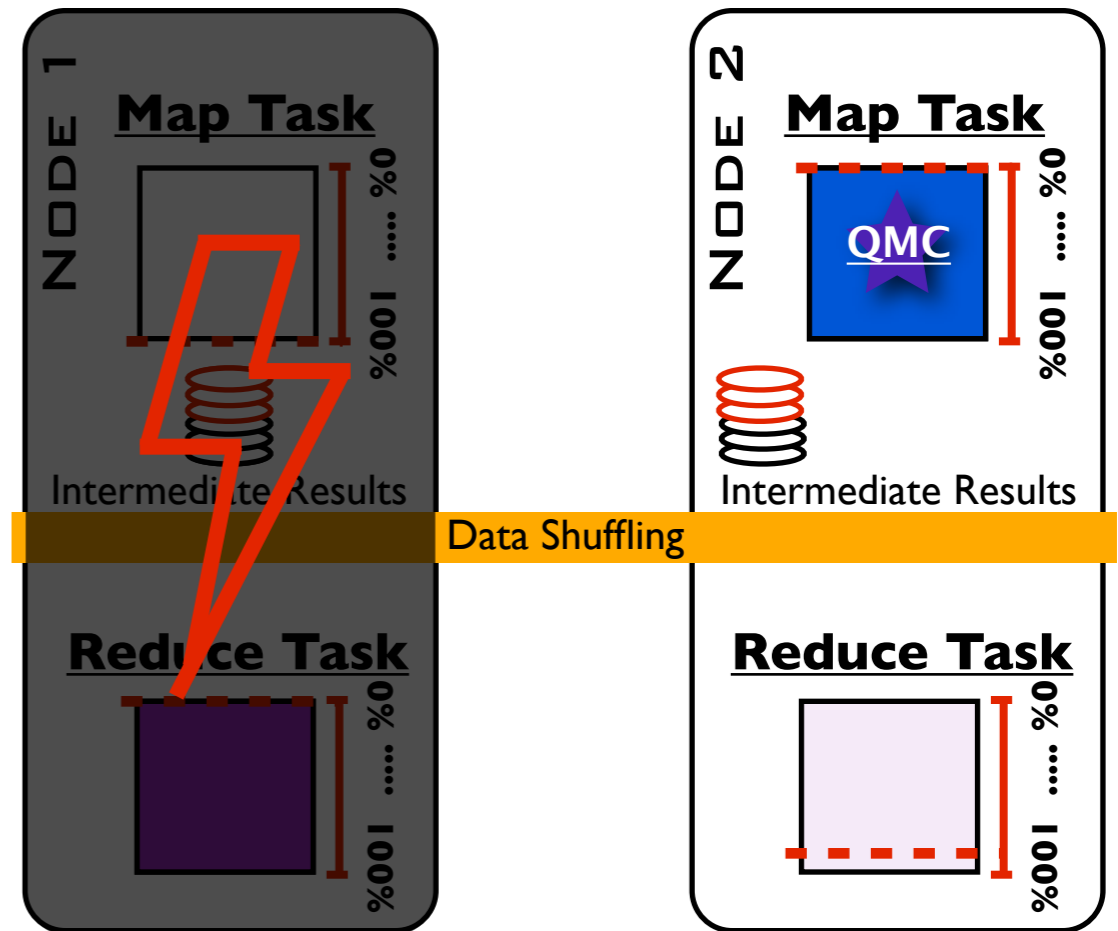


HDFS

## RAFT-QMC

Time:48s

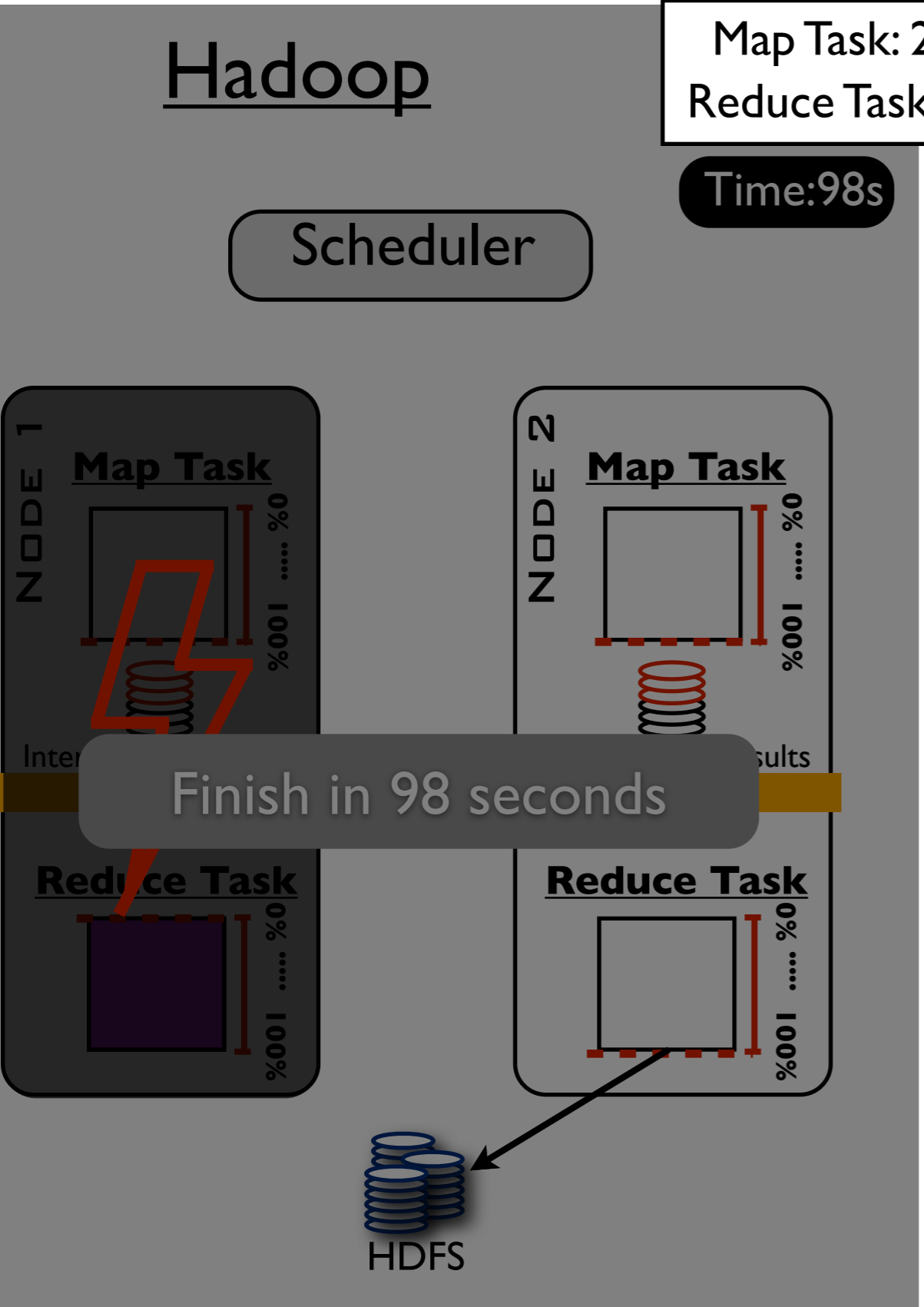
Scheduler



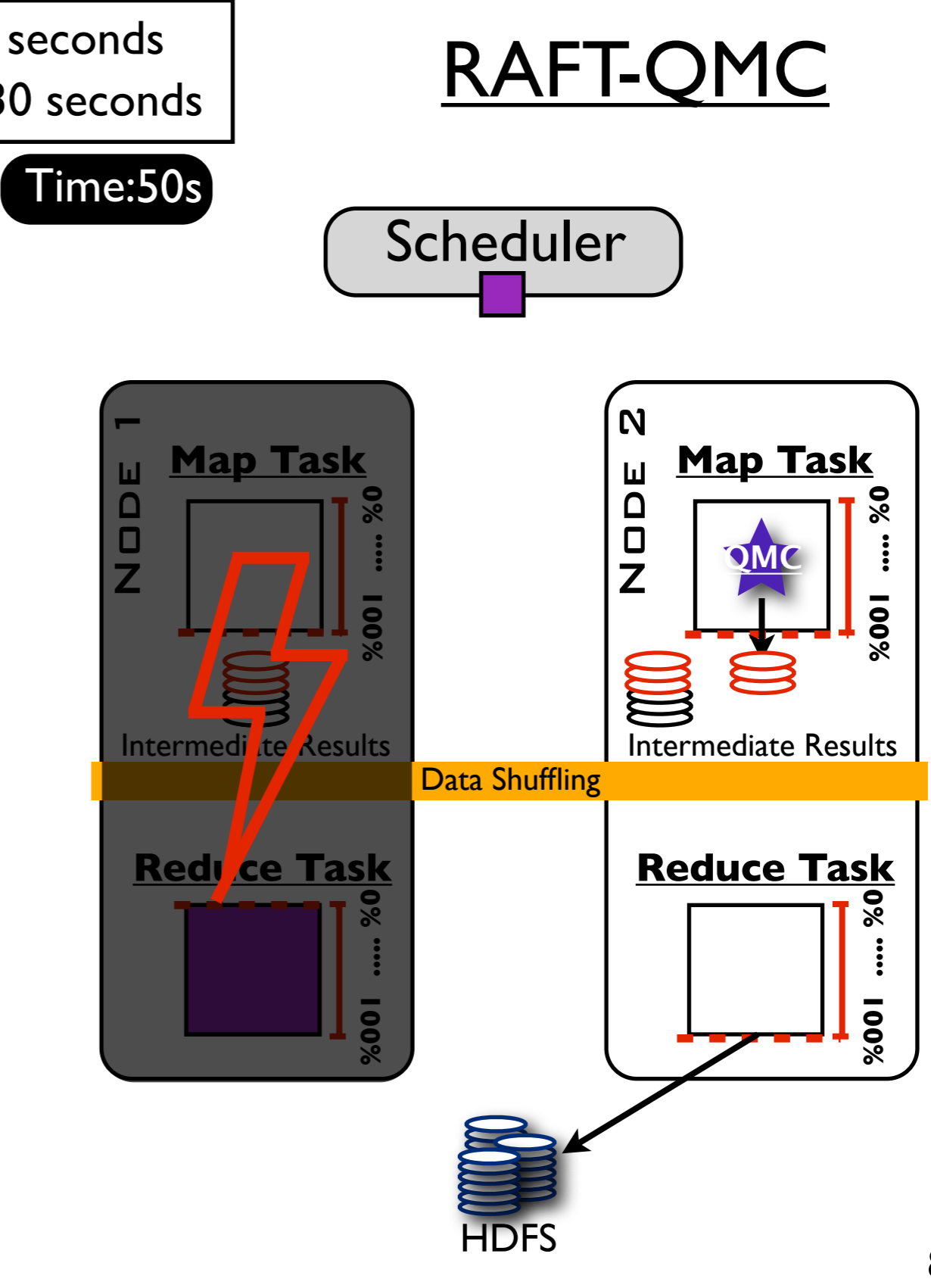


# Recovery from Node Failures

## Hadoop



## RAFT-QMC



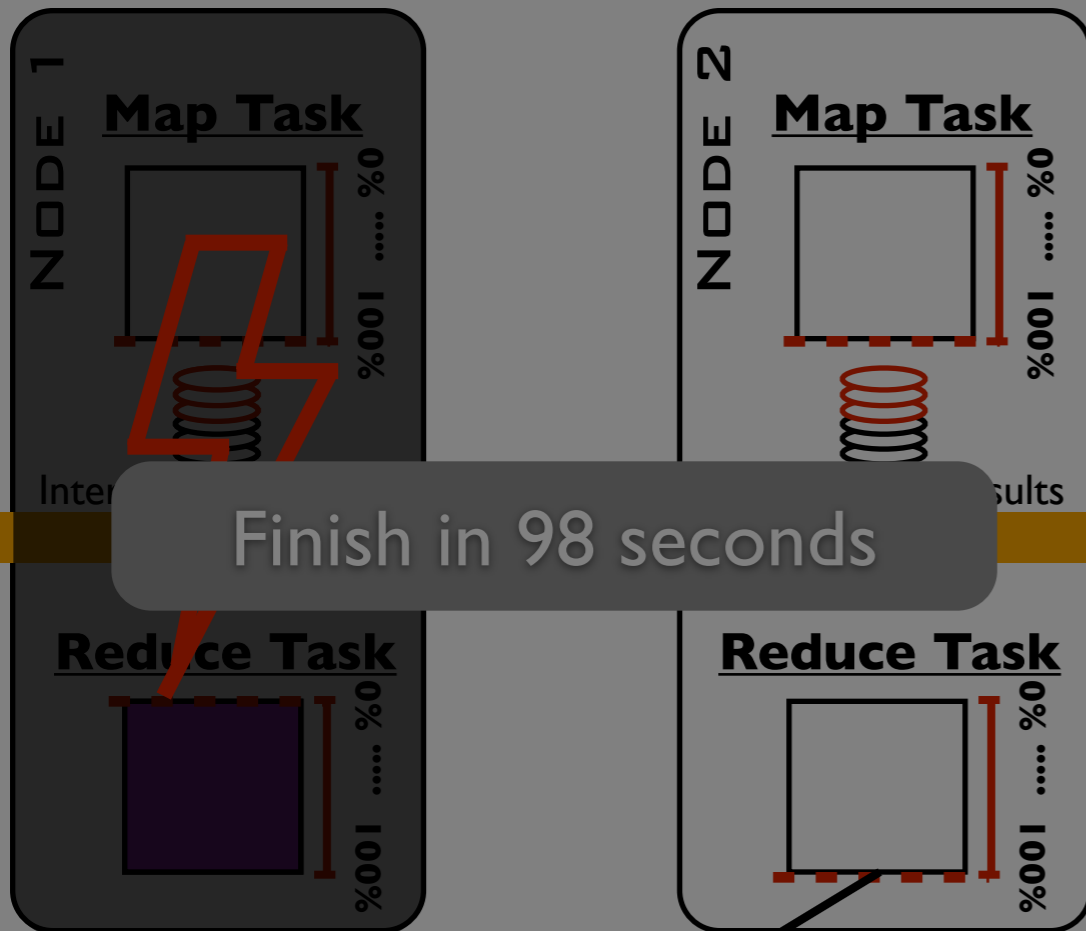
# Recovery from Node Failures

## Hadoop

Map Task: 20 seconds  
Reduce Task: 30 seconds

Scheduler

Time:98s

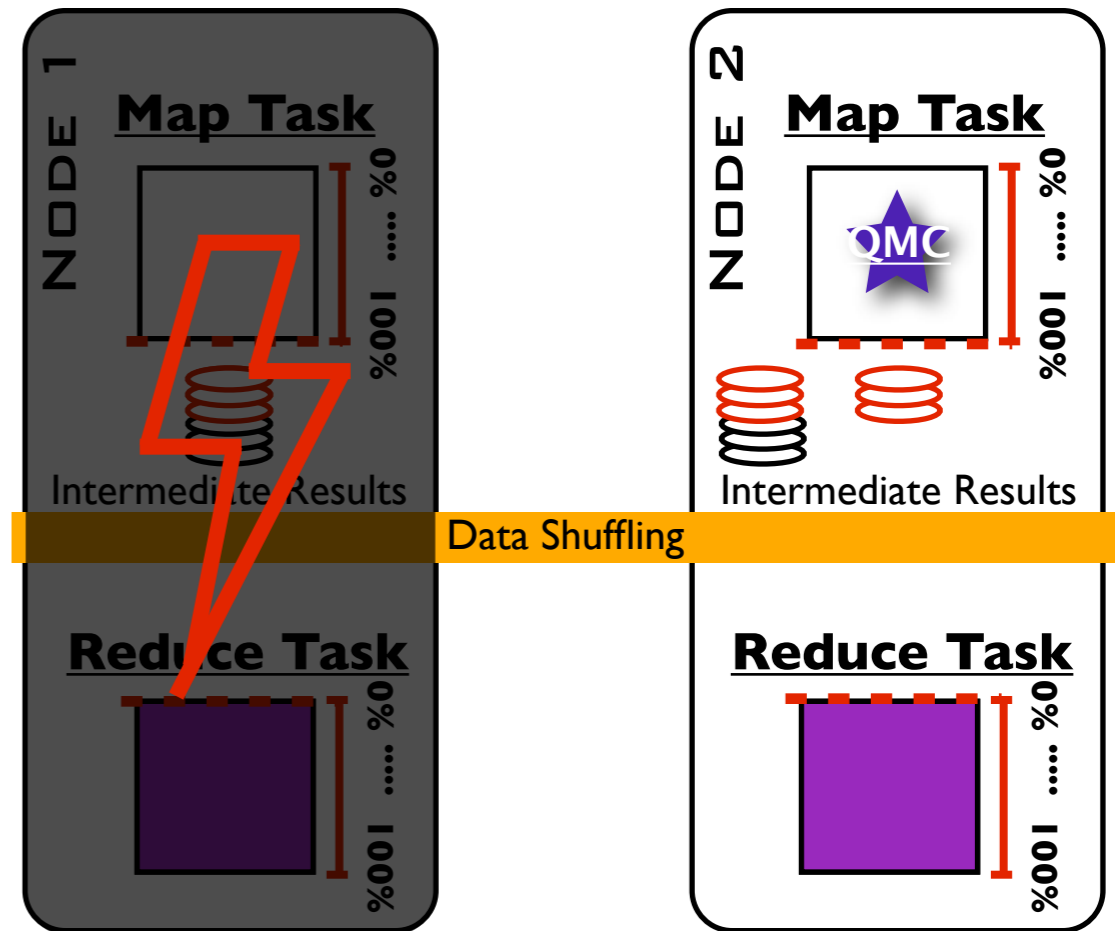


HDFS

## RAFT-QMC

Time:50s

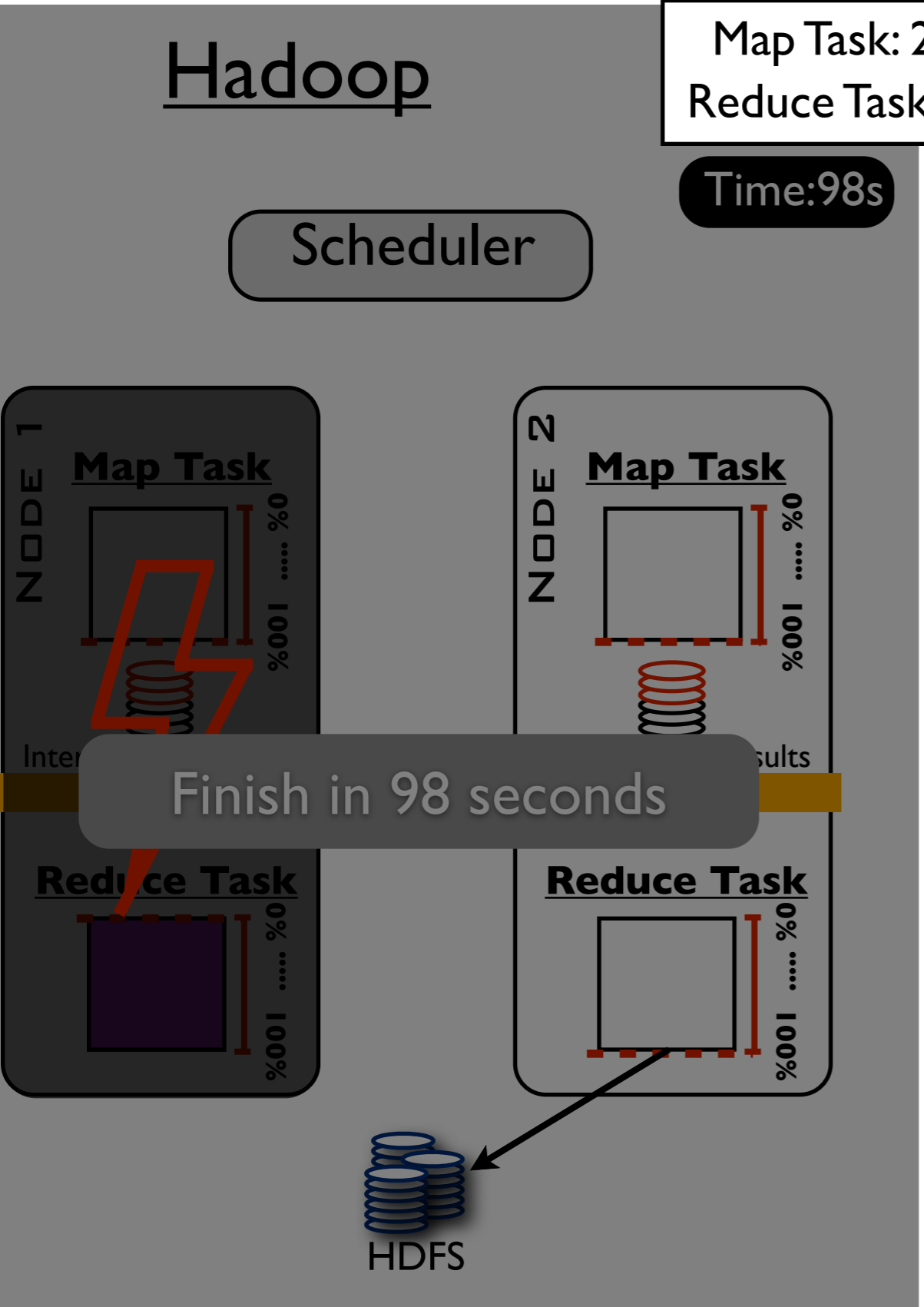
Scheduler



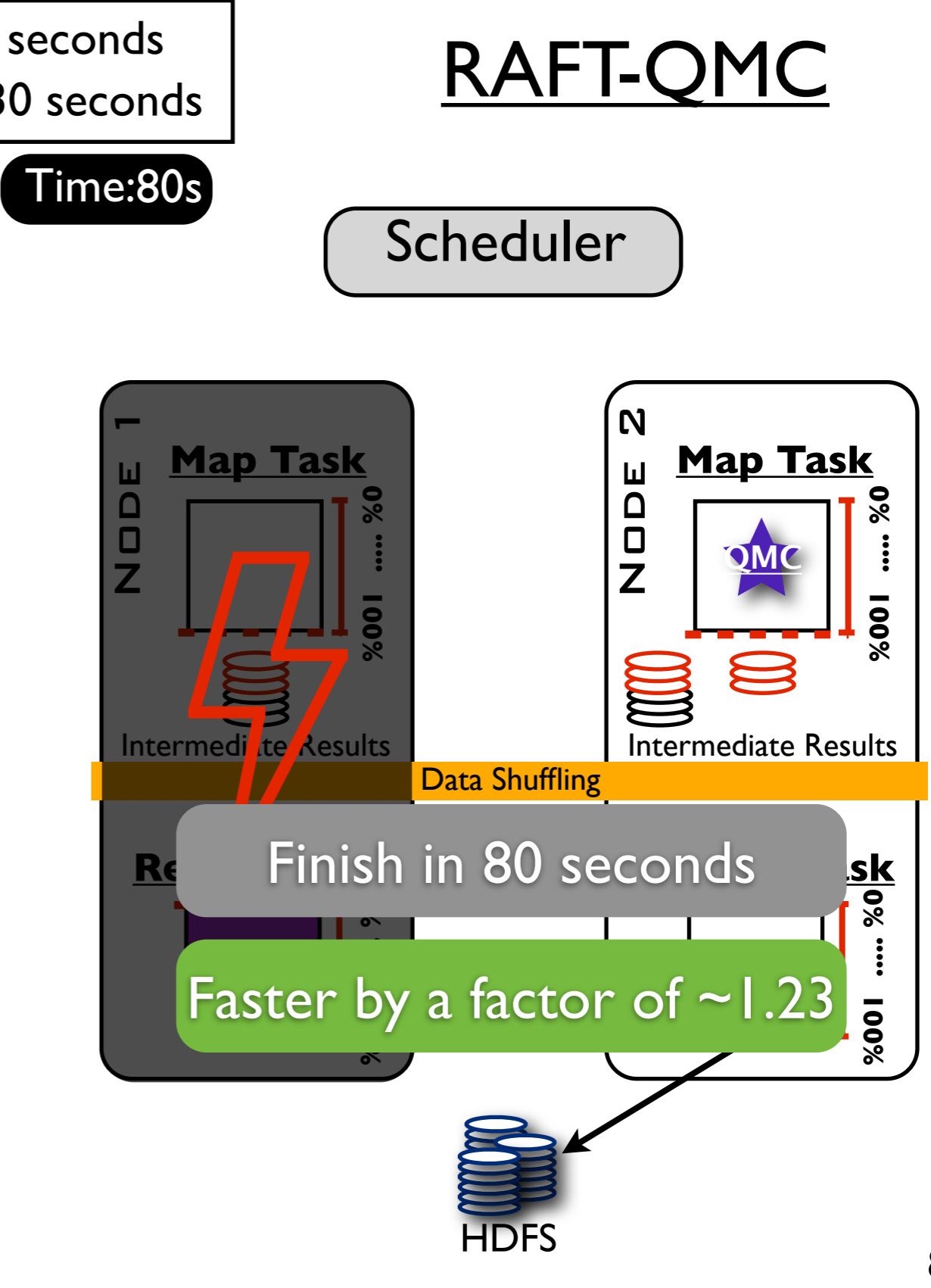
HDFS

# Recovery from Node Failures

## Hadoop



## RAFT-QMC



# Experiments Setup

---

# Experiments Setup

---

- Tested methods:

- RAFT (LC + QMC)
- RAFT-QMC
- RAFT-RC
- RAFT-LC

In this talk

# Experiments Setup

- Tested methods:
    - RAFT (LC + QMC)
    - RAFT-QMC
    - RAFT-RC
    - RAFT-LC
  - Benchmark:
    - Q1: Selection Task
    - Q2: Aggregation Task
    - Q3: Selective Aggregation Task
- In this talk

# Experiments Setup

- Tested methods:
  - RAFT (LC + QMC)
  - RAFT-QMC
  - RAFT-RC
  - RAFT-LC
- Benchmark:
  - **Q1**: Selection Task
  - **Q2**: Aggregation Task
  - **Q3**: Selective Aggregation Task
- **1** TB data set (web logs data), [Pavlo et al., SIGMOD 2009]

In this talk

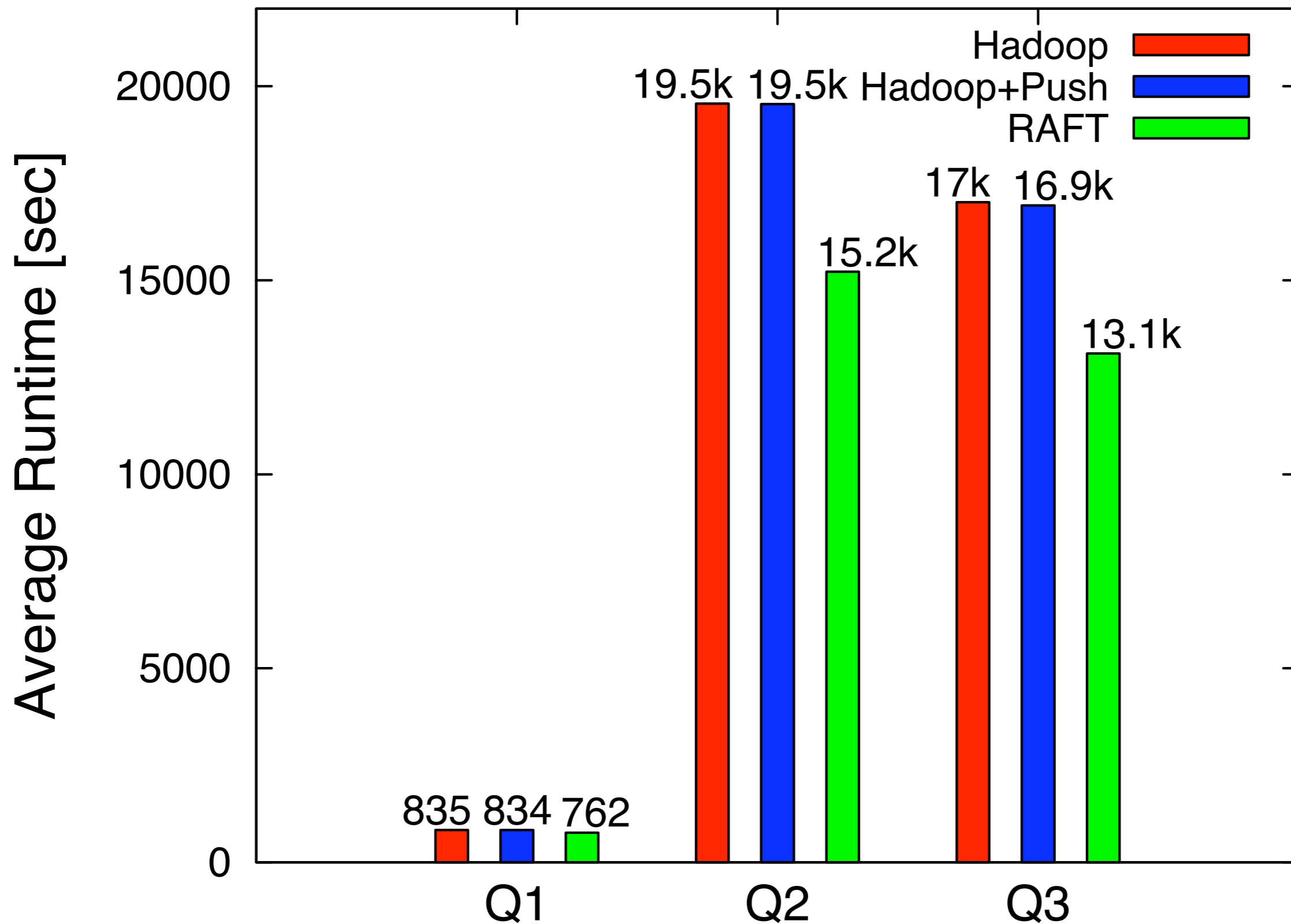
# Experiments Setup

- Tested methods:
  - RAFT (LC + QMC)
  - RAFT-QMC
  - RAFT-RC
  - RAFT-LC
- Benchmark:
  - **Q1**: Selection Task
  - **Q2**: Aggregation Task
  - **Q3**: Selective Aggregation Task
- **1** TB data set (web logs data), [Pavlo et al., SIGMOD 2009]
- **45**-virtual node cluster

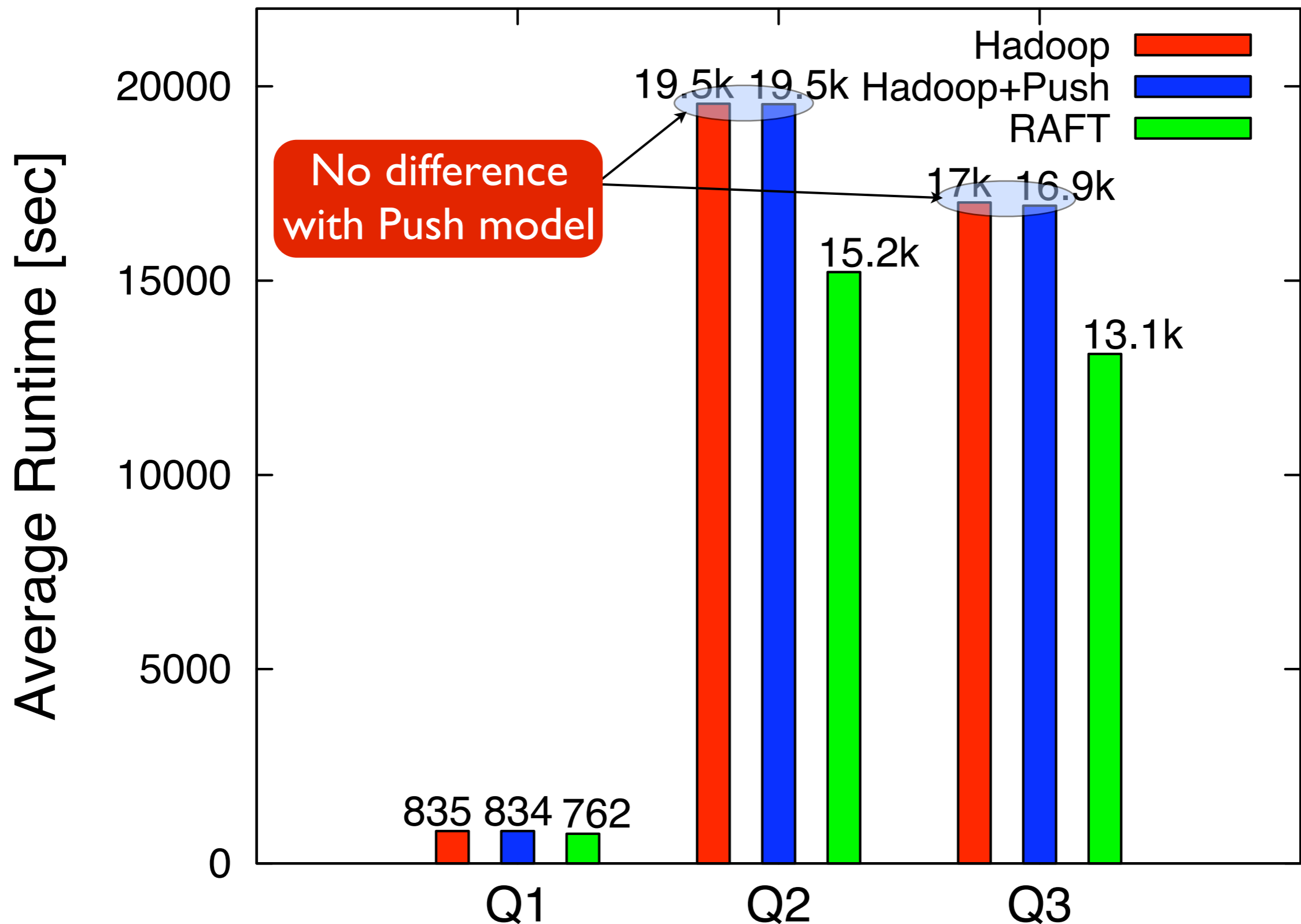
In this talk



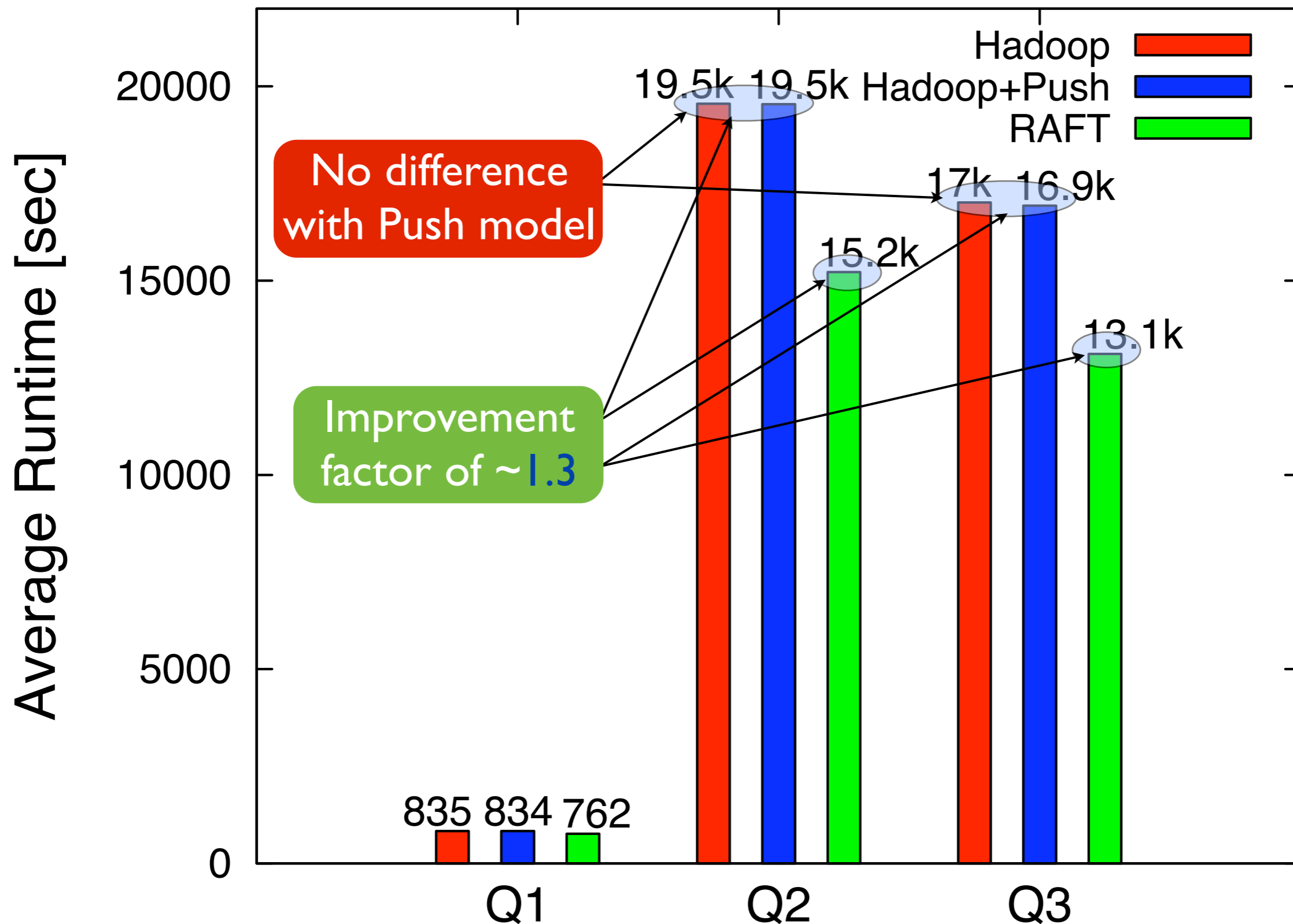
# RAFT Performance



# RAFT Performance



# RAFT Performance



# RAFT Overhead

---

	<b>RAFT</b>
<b>Q1</b>	5.7%
<b>Q2</b>	6.6%
<b>Q3</b>	8.7%

# RAFT Overhead

---

	<b>RAFT</b>
<b>Q1</b>	5.7%
<b>Q2</b>	6.6%
<b>Q3</b>	8.7%

# RAFT Overhead

---

	<b>RAFT</b>	<b>RAFT-QMC</b>	<b>RAFT-RC</b>
<b>Q1</b>	5.7%	3.3%	13.7%
<b>Q2</b>	6.6%	4.3%	9.5%
<b>Q3</b>	8.7%	4.7%	9.8%

# RAFT Overhead

	<b>RAFT</b>	<b>RAFT-QMC</b>	<b>RAFT-RC</b>
<b>Q1</b>	5.7%	3.3%	13.7%
<b>Q2</b>	6.6%	4.3%	9.5%
<b>Q3</b>	8.7%	4.7%	9.8%

# RAFT Overhead

	<b>RAFT</b>	<b>RAFT-QMC</b>	<b>RAFT-RC</b>
<b>Q1</b>	5.7%	3.3%	13.7%
<b>Q2</b>	6.6%	4.3%	9.5%
<b>Q3</b>	8.7%	4.7%	9.8%



# RAFT Overhead

	<b>RAFT</b>	<b>RAFT-QMC</b>	<b>RAFT-RC</b>
<b>Q1</b>	5.7%	3.3%	13.7%
<b>Q2</b>	6.6%	4.3%	9.5%
<b>Q3</b>	8.7%	4.7%	9.8%

RAFT performance  
come almost for free!

# Conclusion

---

**Issue:** failures decrease performance significantly!

# Conclusion

---

**Issue:** failures decrease performance significantly!

**Proposal:** a family of Recovery Algorithms for Fast-Tracking (RAFT) MapReduce:

- \* Local Checkpointing,
- \* Query Metadata Checkpointing,
- \* Remote Checkpointing

# Conclusion

---

**Issue:** failures decrease performance significantly!

**Proposal:** a family of Recovery Algorithms for Fast-Tracking (RAFT) MapReduce:

- \* Local Checkpointing,
- \* Query Metadata Checkpointing,
- \* Remote Checkpointing

**In this talk**

# Conclusion

---

**Issue:** failures decrease performance significantly!

**Proposal:** a family of Recovery Algorithms for Fast-Tracking (RAFT) MapReduce:

- \* Local Checkpointing,
- \* Query Metadata Checkpointing,
- \* Remote Checkpointing

**In this talk**

RAFT algorithms are **fast** and **inexpensive**

# Conclusion

**Issue:** failures decrease performance significantly!

**Proposal:** a family of Recovery Algorithms for Fast-Tracking (RAFT) MapReduce:

- \* Local Checkpointing,
- \* Query Metadata Checkpointing,
- \* Remote Checkpointing

**In this talk**

RAFT algorithms are **fast** and **inexpensive**

**Results:** RAFT outperforms Hadoop by a factor of  $\sim 1.3$  on average, still a negligible runtime overhead!

# Conclusion

**Issue:** failures decrease performance significantly!

**Proposal:** a family of Recovery Algorithms for Fast-Tracking (RAFT) MapReduce:

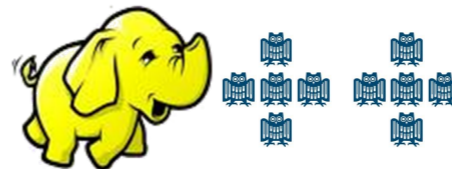
- \* Local Checkpointing,
- \* Query Metadata Checkpointing,
- \* Remote Checkpointing

In this talk

RAFT algorithms are **fast** and **inexpensive**

**Results:** RAFT outperforms Hadoop by a factor of  $\sim 1.3$  on average, still a negligible runtime overhead!

PART OF A BIGGER PROJECT:



Hadoop++, <http://infosys.cs.uni-saarland.de/hadoop++.php>

# Conclusion

**Issue:** failures decrease performance significantly!

**Proposal:** a family of Recovery Algorithms for Fast-Tracking (RAFT) MapReduce:

- \* Local Checkpointing,
- \* Query Metadata Checkpointing,
- \* Remote Checkpointing

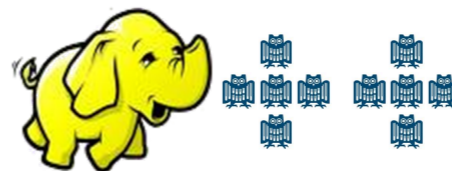
In this talk

RAFT algorithms are **fast** and **inexpensive**

**Results:** RAFT outperforms Hadoop by a factor of  $\sim 1.3$  on average, still a negligible runtime overhead!

**SEE YOU AT SIGMOD' 2011  
FOR A RAFT DEMONSTRATION**

PART OF A BIGGER PROJECT:



Hadoop++, <http://infosys.cs.uni-saarland.de/hadoop++.php>