

# Reining in the **Outliers** in Map-Reduce Clusters using **Mantri**

## Motivation

**What are outliers?**

- An unexpected delay in completion of a task
- A task that takes significantly longer to complete than other tasks in the same job

**Why are outliers a problem?**

- They can cause a job to take significantly longer to complete than expected
- They can cause a job to consume significantly more resources than expected
- They can cause a job to be significantly more expensive than expected

**What causes outliers?**

**Hardware:** Network, Storage, CPU, Memory, etc.

**Software:** Application bugs, Configuration errors, etc.

**What has been done about the problem so far?**

- Manual intervention
- Heuristics
- Resource-aware scheduling

## The main idea

**Margin for improvement**

- Identifying outliers is not always a good idea
- It can be expensive
- It can be noisy
- It can be inaccurate
- It can be unreliable

**The Mantri approach**

- Don't recompute tasks that are running because of long tail of slow tasks
- Each recompute task is only triggered because of network congestion
- Recompute tasks are only triggered when network congestion is detected
- Place tasks in a network-aware way

**The estimate function**

We assume that task completion time can be approximated as a function of the following variables:

- Task size
- Network congestion
- Task placement

We will try to build an estimator of this function and we will use this estimator to decide which tasks to recompute.

## The setting

**Cosmos & Scope**

- Cosmos is a distributed system for managing large-scale data processing jobs
- Scope is a distributed system for managing large-scale data processing jobs

**The Big cluster**

- A cluster of 1000 nodes
- A cluster of 1000 nodes
- A cluster of 1000 nodes

**The input data**

- Log files from customer scheduler
- Log files from customer scheduler
- Log files from customer scheduler

## The scale of the problem

**Statistics**

**Recomputes**

**Task phases**

**Characteristics**

- High number of recomputes
- High number of recomputes
- High number of recomputes

**The impact of outliers**

## The implementation

**The concept**

**Network aware task placement**

- Network aware task placement
- Network aware task placement
- Network aware task placement

**Resource aware restart**

**Task result duplication**

- Task result duplication
- Task result duplication
- Task result duplication

## Evaluation

**Data clusters**

- Pre-production cluster & production cluster
- Pre-production cluster & production cluster
- Pre-production cluster & production cluster

**Results**

- Typical job speed-up: 25-28%
- 42% at 75-th percentile
- 3.1x better improvement as compared to the next best

Reining in the **Outliers** in  
Map-Reduce Clusters using

# Mantri

# Motivation

## What are outliers?

- A job - a single Map-Reduce execution scheme.
- A task - a chunk of computation run on a single machine
- A phase - a series of tasks run concurrently

An outlier - an abnormally long task

(for the purpose of this presentation: > 1.5 avg execution time)

## Why are outliers a problem?

- A single outlier can postpone completion of a whole Map-Reduce phase
- The bigger the cluster, the more serious the problem gets!
- If a task dies, the results used for its computation may not be available any more
- (cascade recomputation!)

Outliers inflate completion time of jobs by **34%** (at median)

## What causes outliers?

### Hardware

- Disk failure
- CPU congestion
- Memory congestion

### Network

- Transferring task data
- Moving data between racks

### Software

- Undivisible data chunks
- Suboptimal implementation
- Scheduling an additional task has a price
- More task cause additional network transfer

## What has been done about the problem so far?

- Mostly restarting outlier tasks
- Outlier detection at the end of each phase
- Restarting tasks without examining the root cause

# *What are outliers?*

- A job - a single Map-Reduce execution scheme.
- A task - a chunk of computation run on a single machine
- A phase - a series of tasks run concurrently

An outlier - an abnormally long task

(for the purpose of this presentation:  $> 1.5$  avg execution time)

# *What causes outliers?*

## Hardware

- Disk failure
- CPU congestion
- Memory congestion

## Network

- Transferring task data
- Moving data between racks

## Software

- Undivisible data chunks
- Suboptimal implementation
- Scheduling an additional task has a price
- More task cause additional network transfer

# *Why are outliers a problem?*

- A single outlier can postpone completion of a whole Map-Reduce phase
- The bigger the cluster, the more serious the problem gets!
- If a task dies, the results used for its computation may not be available any more
- (cascade recomputation!)

Outliers inflate completion time of jobs by **34%** (at median)

# *What has been done about the problem so far?*

- Mostly restarting outlier tasks
- Outlier detection at the end of each phase
- Restarting tasks without examining the root cause

# *The main idea*

## *Margin for improvement*

- Restarting outliers is not always a good choice
- Detecting outliers early
- Use progress reports!
- Network-aware task placement
- Duplicating resource-intensive results

## *The Mantri approach*

- Don't restart tasks that run long because of large amount of data to process
- Don't restart tasks that run long because of network congestion!\*
- \*Unless there is potential for better network transfer
- Place tasks in a network-aware way

## *The estimate function*

We assume that task completion time can be expressed as a function of the following arguments:

**f(data size, code, machine, network location)**

We will try to build an estimate of this function and we will use this estimate in our scheduler algorithm



# *Margin for improvement*

- Restarting outliers is not always a good choice
- Detecting outliers early
- Use progress reports!
- Network-aware task placement
- Duplicating resource-intensive results

# *The Mantri approach*

- Don't restart tasks that run long because of large amount of data to process
- Don't restart tasks that run long because of network congestion!\*

\*Unless there is potential for better network transfer

- Place tasks in a network-aware way

# *The estimate function*

We assume that task completion time can be expressed as a function of the following arguments:

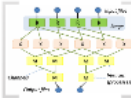
**f(data size, code, machine, network location)**

We will try to build an estimate of this function and we will use this estimate in our scheduler algorithm

# The setting

## Cosmos & Scope

- Cosmos - a commercial upgrade to Dryad
- Most of the jobs written in Scope language
  - Mash-up: SQL + user code (C#)
- Compiler transforms programs into DAGs of dependant tasks
- Compiler optimizes programs for maximal concurrency.



## The Bing cluster

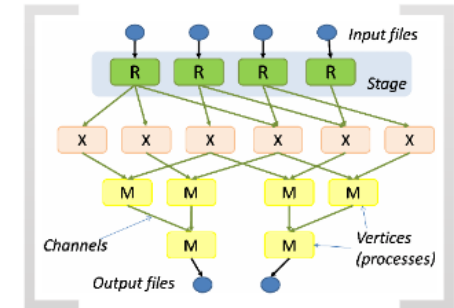
- Reading through network or from local disk storage
- Writing (always!) to local disk storage
- Data stored on the same machines that perform the computation
- By default Cosmos scheduler assigns tasks to machines where the data is available (data locality)
- Cross-rack traffic is costly
- Sum of data storage of the entire rack is smaller than the outgoing bandwidth

## The input data

- Log files from Cosmos scheduler
- Log files from Scope compiler
- Begin and end times for each task
- Input / output data sizes for each task
- Task workflow graph

# Cosmos & Scope

- Cosmos - a commercial upgrade to Dryad
- Most of the jobs written in Scope language
  - Mash-up: SQL + user code (C#)
- Compiler transforms programs into DAGs of dependant tasks
- Compiler optimizes programs for maximal concurrency.



```
SELECT query, COUNT(*) AS count
FROM "search.log" USING LogExtractor
GROUP BY query
HAVING count > 1000
ORDER BY count DESC;
OUTPUT TO "qcount.result";

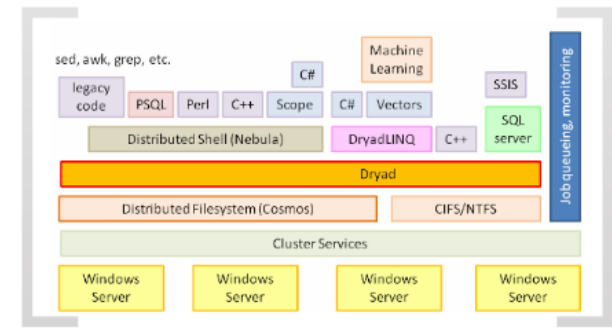
e = EXTRACT query
FROM "search.log"
USING LogExtractor;

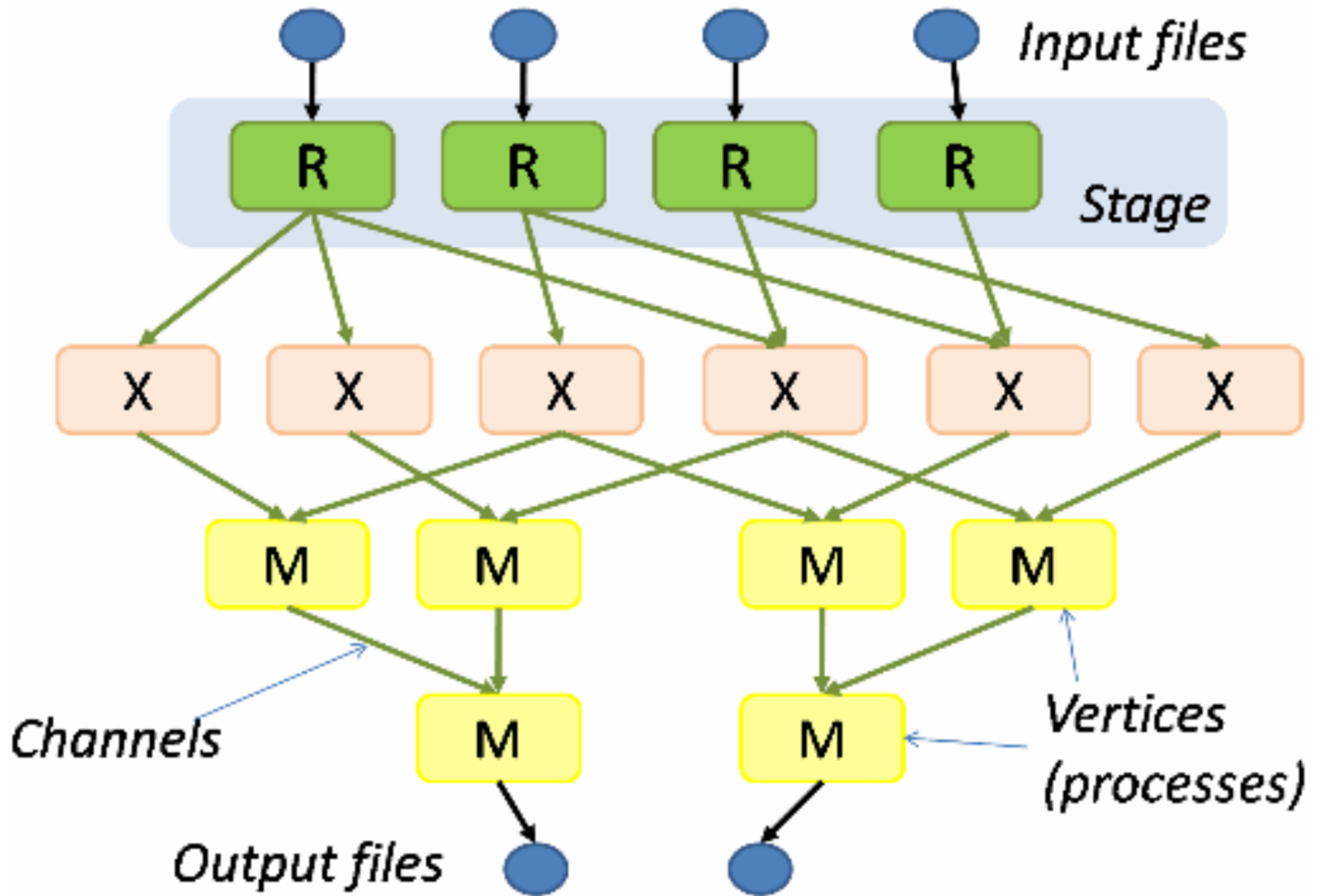
s1 = SELECT query, COUNT(*) as count
FROM e
GROUP BY query;

s2 = SELECT query, count
FROM s1
WHERE count > 1000;

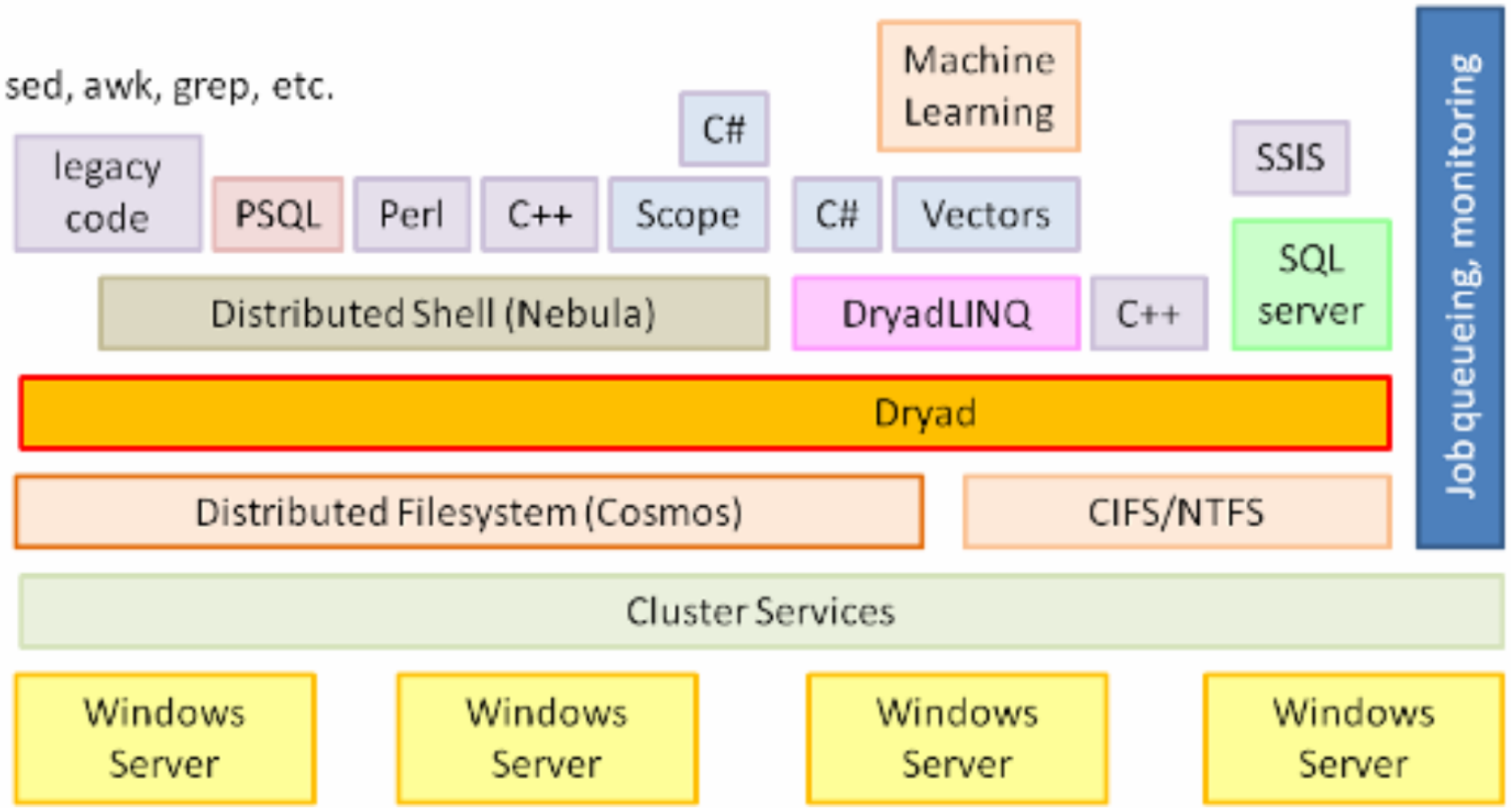
s3 = SELECT query, count
FROM s2
ORDER BY count DESC;

OUTPUT s3 TO "qcount.result"
```





sed, awk, grep, etc.



```
SELECT query, COUNT(*) AS count
FROM "search.log" USING LogExtractor
GROUP BY query
HAVING count > 1000
ORDER BY count DESC;
OUTPUT TO "qcount.result";
```

```
e = EXTRACT query
FROM "search.log"
USING LogExtractor;
```

```
s1 = SELECT query, COUNT(*) as count
FROM e
GROUP BY query;
```

```
s2 = SELECT query, count
FROM s1
WHERE count > 1000;
```

```
s3 = SELECT query, count
FROM s2
ORDER BY count DESC;
```

```
OUTPUT s3 TO "qcount.result"
```



# *The input data*

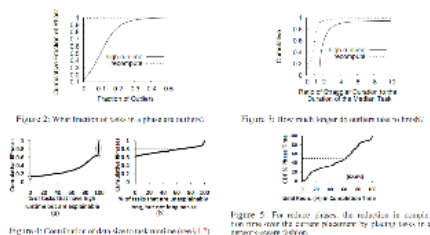
- Log files from Cosmos scheduler
- Log files from Scope compiler
- Begin and end times for each task
- Input / output data sizes for each task
- Task workflow graph

# *The Bing cluster*

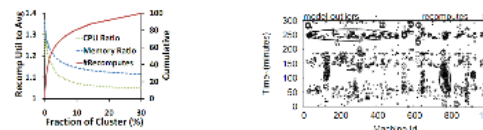
- Reading through network or from local disk storage
- Writing (always!) to local disk storage
- Data stored on the same machines that perform the computation
- By default Cosmos scheduler assigns tasks to machines where the data is available (data locality)
- Cross-rack traffic is costly
- Sum of data storage of the entire rack is smaller than the outgoing bandwidth

# The scale of the problem

## Statistics

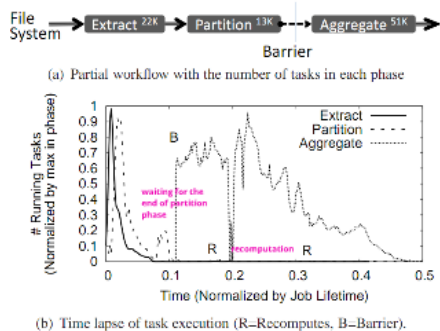


## Recomputes



- Recomputes appear locally
- Time locality
- Machine locality

## Task phases



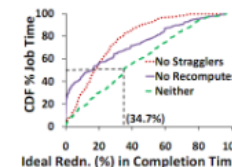
## Characteristics

- Begin of reduce phase is usually a border
- Reduce phase is not associative (no way to perform computation on functions)
- Reduce phase is not commutative
- Recomputation affects a small percentage of tasks.
- The consequences of recomputation are serious.
- 70% of data transfer is caused by reduce tasks

### 2 types of network outliers

- Non local tasks
- Right data block transfer to the bottleneck
- The reduce network congestion problem
- Big default tasks are assigned to the slower machines upon multiple runs

## The impact of outliers

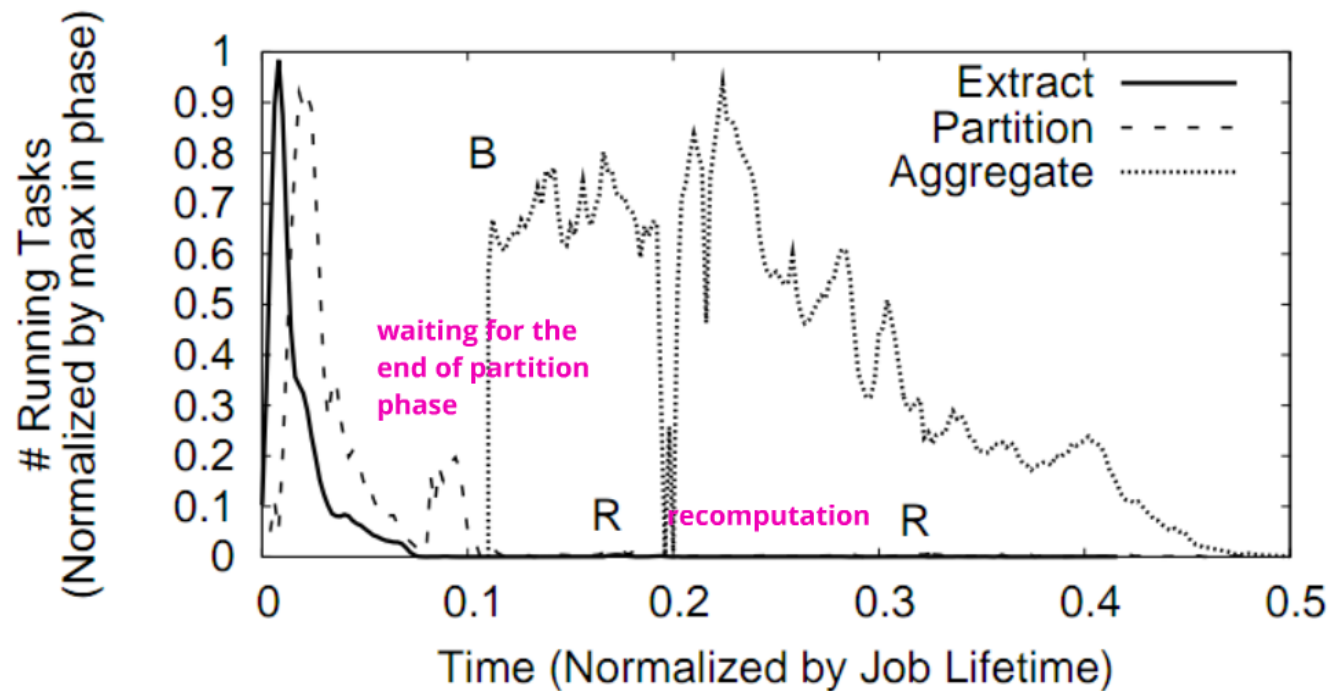


- Data from a simulator (not a real life scenario)
- The simulation took 2 weeks to compute

# Task phases



(a) Partial workflow with the number of tasks in each phase



(b) Time lapse of task execution (R=Recomputes, B=Barrier).

# Characteristics

- Begin of reduce phase is usually a border
- Reduce phase is not associative (no way to perform computation on functions)
- Reduce phase is not commutative
- Recomputation affects a small percentage of tasks.
- The consequences of recomputation are serious.
- 70% of data transfer is caused by reduce tasks

## *2 types of network outliers*

- Non-local tasks
- Again: cross-rack transfer is the bottleneck
- The reduce network congestion problem
- By default tasks are assigned to whichever rack has spare machines

## *2 types of network outliers*

- Non-local tasks
- Again: cross-rack transfer is the bottleneck
- The reduce network congestion problem
- By default tasks are assigned to whichever rack has spare machines

# Statistics

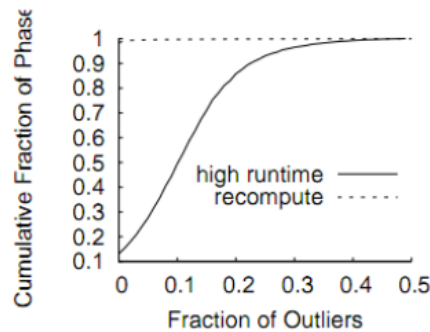


Figure 2: What fraction of tasks in a phase are outliers?

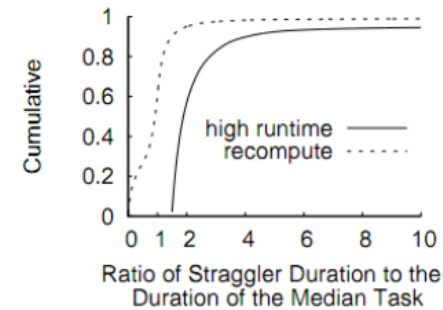


Figure 3: How much longer do outliers take to finish?

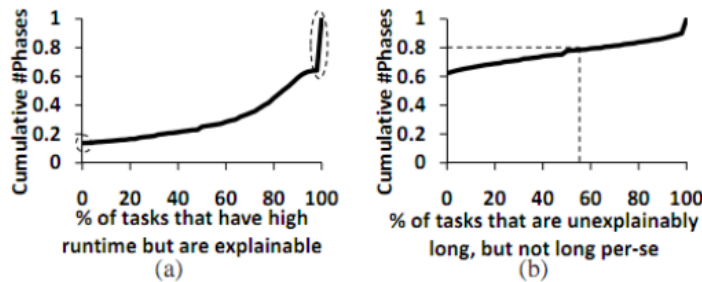


Figure 4: Contribution of data size to task runtime (see §4.2)

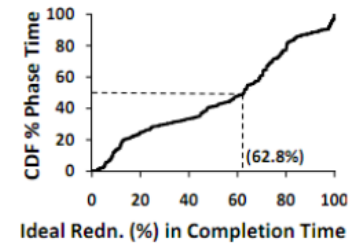
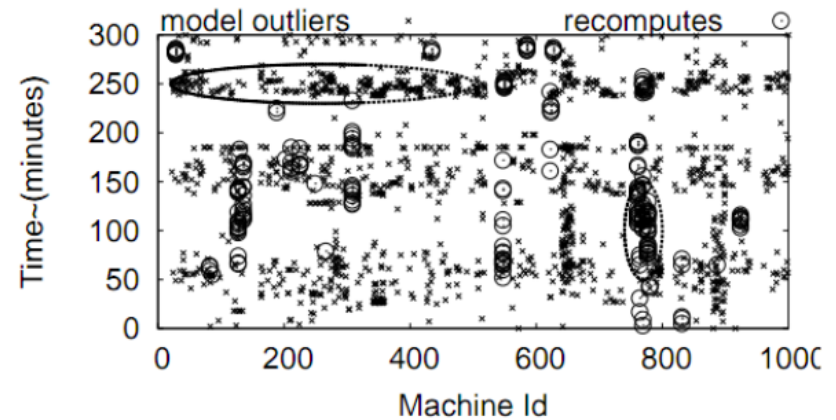
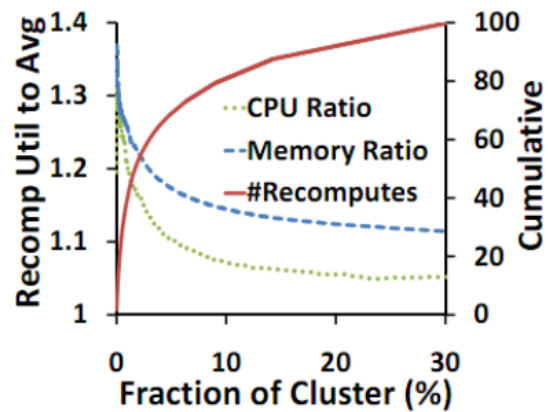


Figure 5: For reduce phases, the reduction in completion time over the current placement by placing tasks in a network-aware fashion.

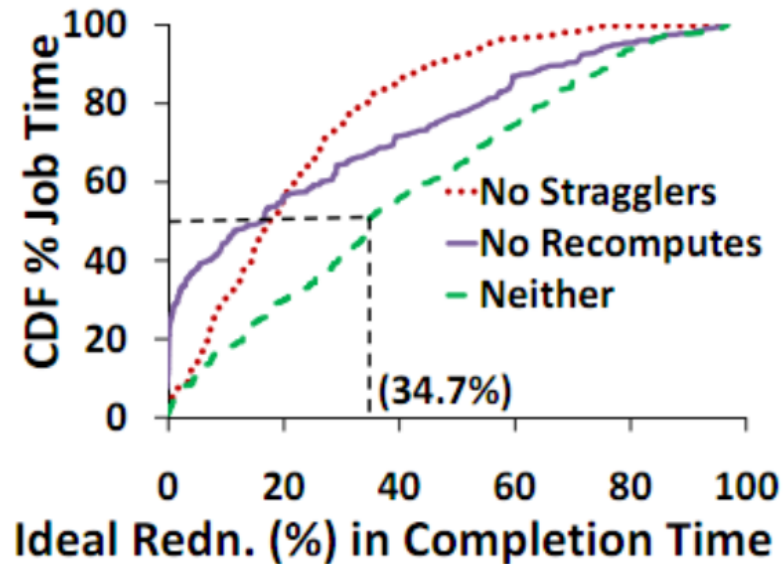
# Recomputes



- Recomputes appear locally
- Time locality
- Machine locality



# *The impact of outliers*



- Data from a simulator (not a real life scenario)
- The simulation took 2 weeks to compute

# The implementation

## The concept

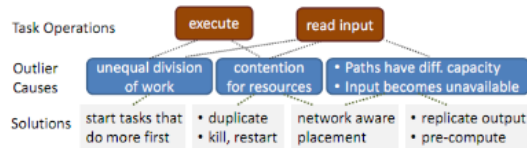


Figure 9: The Outlier Problem: Causes and Solutions

## Network aware task placement

- Attempting to overcome the reduce network congestion problem
- Analyzing every assignment of reduce tasks to racks
- Required transfer bandwidths are known ahead of the time
- Restarts are network aware as well

## Resource-aware restart

```

1: let  $\Delta$  = period of progress reports
2: let  $r$  = number of copies of a task
3: periodically, for each running task, kill all but the fastest  $\alpha$  copies after  $\Delta$  time has passed since begin
4: while slots are available do
5:   if tasks are waiting for slots then
6:     kill, restart task if  $t_{exec} > E(t_{exec}) + \Delta$ , stop at  $\tau$  restarts
7:     duplicate if  $P(t_{exec} > t_{exec} \frac{r-1}{r}) > \delta$ 
8:     start the waiting task that has the largest data to read
9:   else
10:    duplicate if  $E(t_{exec} - t_{exec}) > \rho\Delta$ 
11:  end if
12: end while
  
```

Pseudocode 1: Algorithm for Resource-aware restarts (simplified).

- Consider kill & restart or duplication
- Keep number of running copies constant
- Duplicate if the expected amount of resources needed to complete the task would decrease
- Kill tasks when duplicates exceed budget cap

## Task result duplication

- Replicate results if they are small
- Replicate when the result is costly to recompute
- Replicate the tasks that take the longest to recompute
- Budget cap on the amount of replicated data
- Replication throttled through a system of tokens

## The concept

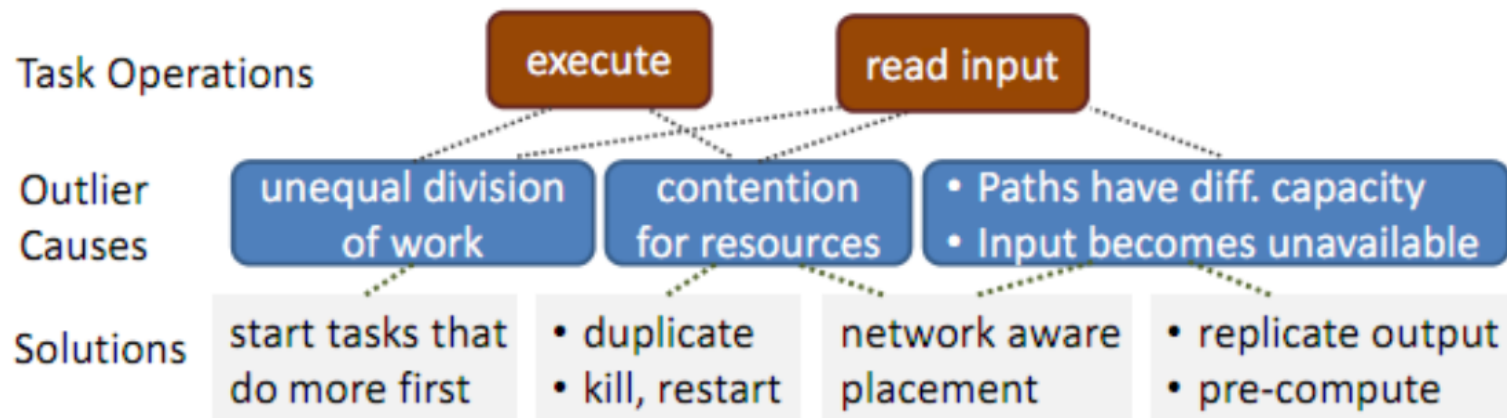


Figure 9: The Outlier Problem: Causes and Solutions

# Resource-aware restart

```
1: let  $\Delta$  = period of progress reports
2: let  $c$  = number of copies of a task
3: periodically, for each running task, kill all but the fastest  $\alpha$  copies after  $\Delta$  time has passed since begin
4: while slots are available do
5:   if tasks are waiting for slots then
6:     kill, restart task if  $t_{rem} > \mathbb{E}(t_{new}) + \Delta$ , stop at  $\gamma$  restarts
7:     duplicate if  $\mathbb{P}(t_{rem} > t_{new} \frac{c+1}{c}) > \delta$ 
8:     start the waiting task that has the largest data to read
9:   else
10:    duplicate iff  $\mathbb{E}(t_{new} - t_{rem}) > \rho\Delta$ 
11:  end if
12: end while
```

- Typically Delta equals 10 sec
- Reschedule requires additional time
- Don't duplicate recently duplicated tasks ▷ all tasks have begun
- No response, assume no progress

**Pseudocode 1:** Algorithm for Resource-aware restarts (simplified).

- Consider kill & restart or duplication
- Keep number of running copies constant
- Duplicate if the expected amount of resources needed to complete the task would decrease
- Kill tasks when duplicates exceed budget cap

# *Network aware task placement*

- Attempting to overcome the reduce network congestion problem
- Analyzing every assignment of reduce tasks to racks
- Required transfer bandwidths are known ahead of the time
- Restarts are network aware as well

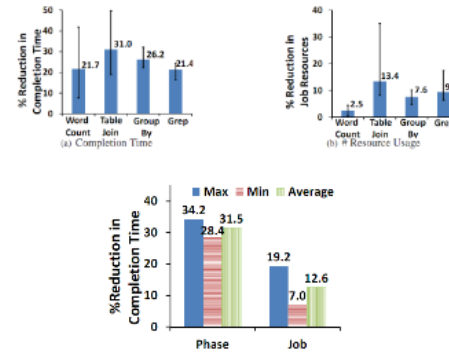
# *Task result duplication*

- Replicate results if they are small
- Replicate when the result is costly to recompute
- Replicate the tasks that take the longest to recompute
- Budget cap on the amount of replicated data
- Replication throttled through a system of tokens

# Evaluation

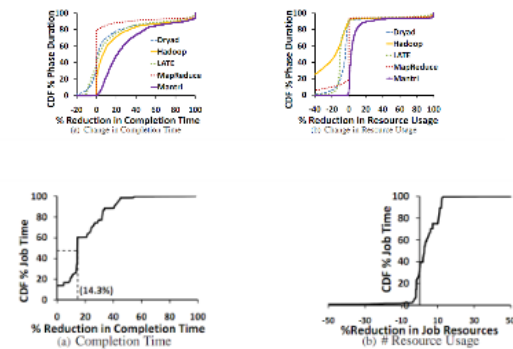
## Data clusters

- Preproduction cluster & production cluster
- Baseline: an unmodified Cosmos scheduler
- Comparison to Hadoop, Dryad, MapReduce, LATE & modified LATE



## Results

- Typical job speed-up: 25-28%
- 42% at 75-th percentile
- 3.1x better improvement as compared to the next best



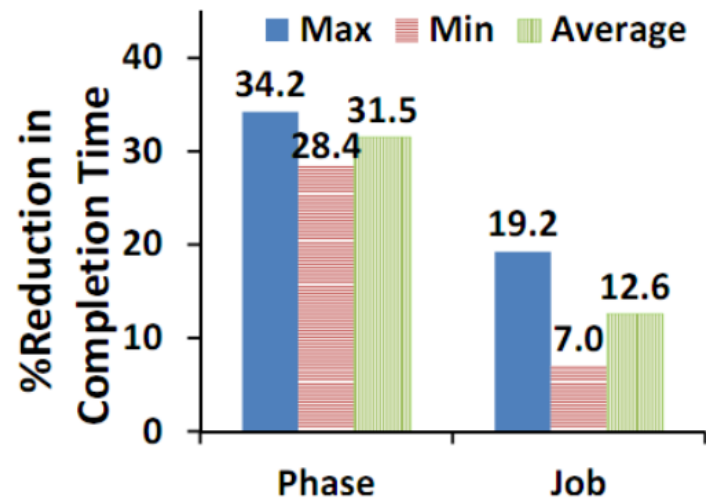
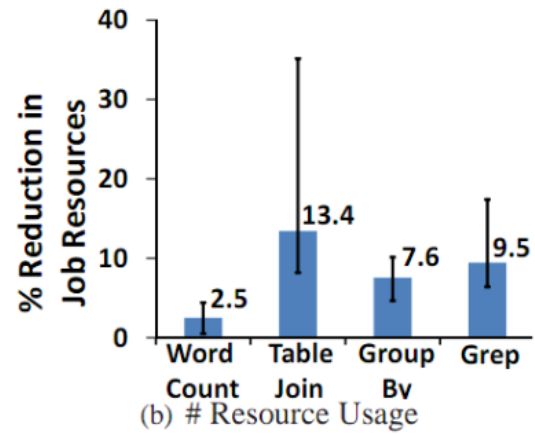
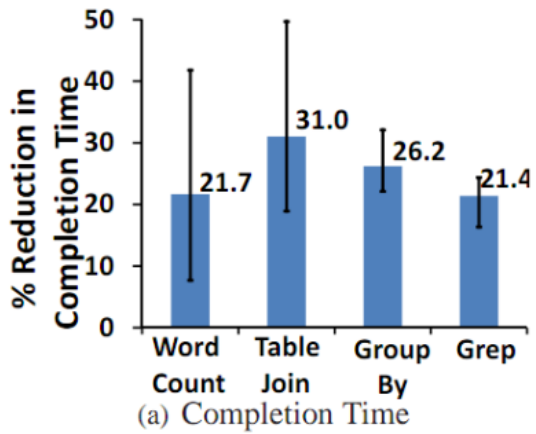
# *Data clusters*

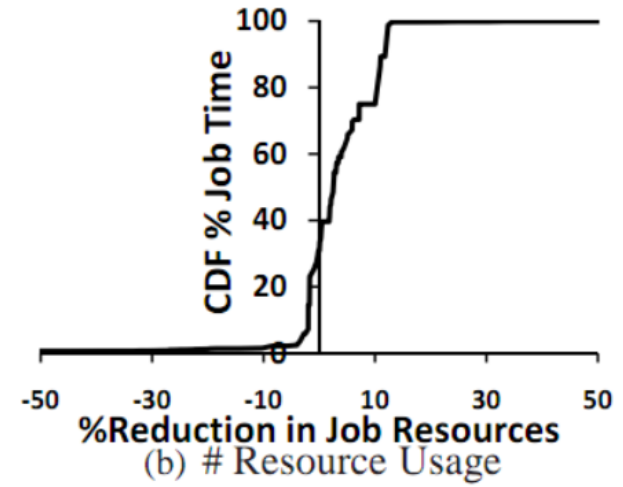
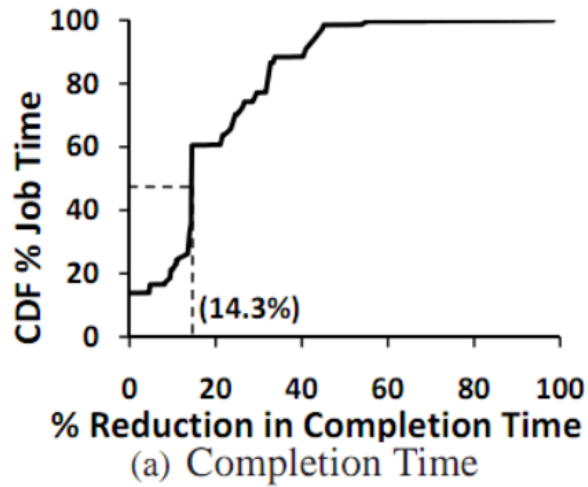
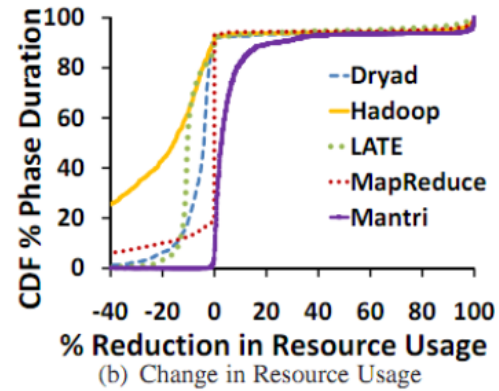
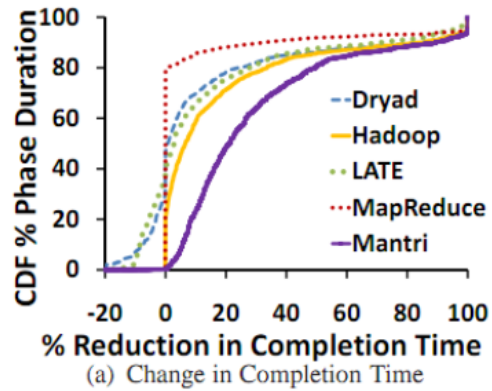
- Preproduction cluster & production cluster
- Baseline: an unmodified Cosmos scheduler
- Comparison to Hadoop, Dryad, MapReduce, LATE & modified LATE



# *Results*

- Typical job speed-up: 25-28%
- 42% at 75-th percentile
- 3.1x better improvement as compared to the next best





Reining in the **Outliers** in  
Map-Reduce Clusters using

# Mantri