# Automating Data-Throttling Analysis for Data-Intensive Workflow

**Ricardo J. Rodríguez**, Rafael Tolosana-Calasanz, Omer F. Rana {rjrodriguez, rafaelt}@unizar.es, o.f.rana@cs.cardiff.ac.uk



Universidad de Zaragoza Zaragoza, Spain



Cardiff University Cardiff, United Kingdom

May 15<sup>th</sup>, 2012

CCGrid'12: 12<sup>th</sup> IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing Ottawa, Canada

#### Motivation

### Outline

### 1 Motivation

- Knowing the problem
- Data-movement policies
- The goal & our approach

#### 2 Background

- Petri nets
- From DAG to PN
- Slack Concept

### 3 Automating Data-Throttling Analysis

#### Experiments and Results

- Impact on the Workflow Makespan
- Input Buffers and Network Bandwidth Usage

### Related Work

#### Conclusions and Future Work



### Motivation (I): data movement issue

### Scientific Workflows

- Control/data flow graph
  - Set of tasks
  - Dependencies



### Motivation (I): data movement issue

### Scientific Workflows

- Control/data flow graph
  - Set of tasks
  - Dependencies

### Data-Intensive Workflows

- Data workload ≥ computational workload
- How to deal with this?
  - Minimise data movement
  - Move data via higher capacity links as fast as possible



Motivation Data-movement policies

### Motivation (II): data-movement policies

### Pipeline workflow



• Transmit as fast as possible  $\rightarrow$  OK!



# Motivation (II): data-movement policies

### Pipeline workflow



• Transmit as fast as possible  $\rightarrow$  OK!

### Directed Acyclic Graph workflow



- Transmit as fast as possible  $\rightarrow$  WRONG!
- $\bullet \ \ Cluster \ tasks \rightarrow minimise \ data \\ movement$
- Problems may arise
  - Isolation view: networking + buffering (merge tasks)
  - In-the-large: harmful (finite buffer capacity → cannot be reused)

### The goal

- Balance workflow paths...
- ... eliminating unnecessary bandwidth usage



### The goal

- Balance workflow paths...
- ... eliminating unnecessary bandwidth usage
  - Efficient buffer usage
  - Efficient network bandwidth usage



### The goal

- Balance workflow paths...
- ... eliminating unnecessary bandwidth usage
  - Efficient buffer usage
  - Efficient network bandwidth usage

### Open question: WHAT strategy should I use?



### The goal

- Balance workflow paths...
- ... eliminating unnecessary bandwidth usage
  - Efficient buffer usage
  - Efficient network bandwidth usage

### Open question: WHAT strategy should I use?

### Our approach

• Automatically derive values for data-throttling

1542

( AKD

### The goal

- Balance workflow paths...
- ... eliminating unnecessary bandwidth usage
  - Efficient buffer usage
  - Efficient network bandwidth usage

### Open question: WHAT strategy should I use?

### Our approach

- Automatically derive values for data-throttling
- Directed Acyclic Graphs (DAGs)  $\rightarrow$  Petri nets (performance model)
- Analysis on Petri net model (explained later :))

1542

# Outline



#### • Knowing the problem

- Data-movement policies
- The goal & our approach

### 2 Background

- Petri nets
- From DAG to PN
- Slack Concept

3 Automating Data-Throttling Analysis

#### Experiments and Results

- Impact on the Workflow Makespan
- Input Buffers and Network Bandwidth Usage

### Related Work

#### Conclusions and Future Work



Background

Petri nets

# Background (I): Petri nets (PNs)





# Background (I): Petri nets (PNs)



- Mathematical formalism
- Places (circles, p<sub>X</sub>)
- Transitions (bars,  $t_X$ ). Associated delay
  - Immediate ( $\delta_{t_X} = 0$ )
  - Timed ( $\delta_{t_X}$  exponential distribution  $\rightarrow$  Stochastic Petri Nets)
- Tokens (black dots)
- Directed arcs
  - Place  $\rightarrow$  Transition
  - $\bullet \ \ {\sf Transition} \rightarrow {\sf Place}$
- Initial marking: no. tokens on places

# Background (I): Petri nets (PNs)



- Mathematical formalism
- Places (circles,  $p_X$ )
- Transitions (bars,  $t_X$ ). Associated delay
  - Immediate ( $\delta_{t_X} = 0$ )
  - Timed ( $\delta_{t_X}$  exponential distribution  $\rightarrow$  Stochastic Petri Nets)
- Tokens (black dots)
- Directed arcs
  - Place  $\rightarrow$  Transition
  - $\bullet \ \ {\sf Transition} \rightarrow {\sf Place}$
- Initial marking: no. tokens on places
- Enables modelling
  - Concurrency
  - Synchronisation

### Background (II): From DAG to PN



### Transforming from a DAG to a PN

- Computation task  $\rightarrow$  place + timed transition
  - Delay equal to computation task time

zaragoza

(AERDYD)

## Background (II): From DAG to PN



### Transforming from a DAG to a PN

- Computation task  $\rightarrow$  place + timed transition
  - Delay equal to computation task time





### Background (II): From DAG to PN



### Transforming from a DAG to a PN

- Computation task  $\rightarrow$  place + timed transition
  - Delay equal to computation task time





### Background (II): From DAG to PN



### Transforming from a DAG to a PN

- Computation task  $\rightarrow$  place + timed transition
  - Delay equal to computation task time
- $\bullet$  Transmission link  $\rightarrow$  place + timed transition

• Delay equal to transmission time:  $\delta_t = \frac{data \ size \ transmitted}{link \ bandwidth}$ 

∠aragoza

1542

**AERD** 

#### Slack Concept

# Background (III): Slack concept (1)



Slack Concept

# Background (III): Slack concept (1)



Slack Concept

# Background (III): Slack concept (1)



Slack Concept

# Background (III): Slack concept (1)



Background

#### Slack Concept

# Background (III): Slack concept (1)

For human beings: an example



 $\delta$  execution time; tx transmission time



R.J. Rodríguez et al. Automating Data-Throttling Analysis for Data-Intensive Workflow CCGrid'12 9 / 25

#### Slack Concept

### Background (III): Slack concept (2) For tough guys/gals – Maths fans

- Upper performance bound vs. exact analysis
  - Exact analysis needs reachability graph  $\rightarrow$  NP-hard problem



### Background (III): Slack concept (2) For tough guys/gals – Maths fans

- Upper performance bound vs. exact analysis
  - $\bullet~\mbox{Exact}$  analysis needs reachability graph  $\rightarrow~\mbox{NP-hard}$  problem
- Use of Linear Programming (LP) techniques  $\rightarrow$  polynomial complexity
- Based on Petri net theory: tight marking (m̃) [p. 3 on the paper]
   Further reading: google "RJ-EPEW-10"

$$\tilde{\mathbf{m}}(p) \ge \delta(p^{\bullet}) \cdot \Theta \to \mathbf{m}(p) = \delta(p^{\bullet}) \cdot \Theta + \mu(p)$$
(1)

 $\delta(p^{\bullet})$  delay of transition after place p;  $\Theta$  inverse of execution time of slowest path;  $\mu(p)$  slack of place p



# Outline

### Motivation

- Knowing the problem
- Data-movement policies
- The goal & our approach

### 2) Background

- Petri nets
- From DAG to PN
- Slack Concept

### 3 Automating Data-Throttling Analysis

- Experiments and Results
  - Impact on the Workflow Makespan
  - Input Buffers and Network Bandwidth Usage

### Related Work

#### Conclusions and Future Work





- Inputs: Performance estimation (i.e., DAX annotations) + PN-based model
- Outputs: Data-throttling values + Analysis results



- Inputs: Performance estimation (i.e., DAX annotations) + PN-based model
- Outputs: Data-throttling values
   + Analysis results
- 4 steps
  - Compute slack values
  - Oluster slacks
    - Why? In the next slide!

- Compute data-throttling values
- Performance analysis
  - With and w/out data-throttling

542

# Automating Data-Throttling Analysis (II)

The need of clustering slacks



- Influence between slacks
- Order of adjusting IS IMPORTANT





# Automating Data-Throttling Analysis (II)

The need of clustering slacks



- Influence between slacks
- Order of adjusting IS IMPORTANT
- From output to input
- Cluster slacks, then:

1542

- Compute data-throttling values
- Recompute slacks values

Zaragoza



# Automating Data-Throttling Analysis (II)

The need of clustering slacks



- Influence between slacks
- Order of adjusting IS IMPORTANT
- From output to input
- Cluster slacks, then:
  - Compute data-throttling values
  - 2 Recompute slacks values

Zaragoza

PRIFYSGOL

AERDY

In the example:

ÏĨĨ

1542

- $\{ \mu_{3,5}, \mu_{1,5} \}$
- **2** μ<sub>1,3</sub>

# Automating Data-Throttling Analysis (III)

Deriving data-throttling values from slack

μ(p) → some delay may be added on the path
tx<sub>new</sub> = tx<sub>old</sub> + α, α = μ(p)/Θ



# Automating Data-Throttling Analysis (III)

Deriving data-throttling values from slack

$$BW_{new} = rac{1}{rac{1}{BW_{old}} + rac{\mu(p)}{\Theta \cdot data \ size}}$$

 $\mu$  slack; *tx* transmission time;

 $\boldsymbol{\Theta}$  inverse of execution time of slowest path



(2)

# Automating Data-Throttling Analysis (IV)

Applying to an example



Recall: slacks on synchronisation points

Universidad

Zaragoza

......

1542

**^ARDIF** 

- BW=100Mbps, latency  $1e^{-4}s$
- Data-sets equal to 10MiB
- Dedicated network topology

# Automating Data-Throttling Analysis (IV)

Applying to an example



Recall: slacks on synchronisation points

- BW=100Mbps, latency  $1e^{-4}s$
- Data-sets equal to 10MiB
- Dedicated network topology

- Slowest path:  $2 \rightarrow 5 \rightarrow 6$
- Slacks: μ<sub>1,4</sub>, μ<sub>3,6</sub>, μ<sub>4,6</sub>

# Automating Data-Throttling Analysis (IV)

#### Applying to an example



### Recall: slacks on synchronisation points



Makespan: 5.6779 seconds

- Slowest path:  $2 \rightarrow 5 \rightarrow 6$
- Slacks: μ<sub>1,4</sub>, μ<sub>3,6</sub>, μ<sub>4,6</sub>

- BW=100Mbps, latency  $1e^{-4}s$
- Data-sets equal to 10MiB
- Dedicated network topology

# Automating Data-Throttling Analysis (IV)

#### Applying to an example



### Recall: slacks on synchronisation points



Makespan: 5.6779 seconds

### • Slowest path: $2 \rightarrow 5 \rightarrow 6$

- Slacks:  $\mu_{1,4}, \mu_{3,6}, \mu_{4,6}$ 
  - $1 \rightarrow 4$  adjust to 28.57%
  - $3 \rightarrow 6$  adjust to 35.15%
  - $4 \rightarrow 6$  adjust to 44.55%

- BW=100Mbps, latency  $1e^{-4}s$
- Data-sets equal to 10MiB
- Dedicated network topology

# Automating Data-Throttling Analysis (IV)

#### Applying to an example



### Assumptions

- BW=100Mbps, latency  $1e^{-4}s$
- Data-sets equal to 10MiB
- Dedicated network topology

### Recall: slacks on synchronisation points



- Slowest path:  $2 \rightarrow 5 \rightarrow 6$
- Slacks:  $\mu_{1,4}, \mu_{3,6}, \mu_{4,6}$ 
  - $1 \rightarrow 4$  adjust to 28.57%
  - $3 \rightarrow 6$  adjust to 35.15%
  - $4 \rightarrow 6$  adjust to 44.55%

# Outline

### Motivation

- Knowing the problem
- Data-movement policies
- The goal & our approach

### 2) Background

- Petri nets
- From DAG to PN
- Slack Concept

### Automating Data-Throttling Analysis

- Experiments and Results
  - Impact on the Workflow Makespan
- Input Buffers and Network Bandwidth Usage

#### Related Work

#### Conclusions and Future Work



### Experiments and Results (I): Description





# Experiments and Results (I): Description



Montage Workflow – 5inputs

 Performance estimation from DAX (Pegasus Wf System)



17 / 25

# Experiments and Results (I): Description



- Montage Workflow 5inputs
- Performance estimation from DAX (Pegasus Wf System)
- Network topologies assumed:
  - Single vs. dedicated data-link (throttling in dedicated)



17 / 25

# Experiments and Results (I): Description



- Montage Workflow 5inputs
- Performance estimation from DAX (Pegasus Wf System)
- Network topologies assumed:
  - Single vs. dedicated data-link (throttling in dedicated)

#### Tools used

- Slack computation: MATLAB
- Performance analysis: SimGrid



# Experiments and Results (I): Description



- Montage Workflow 5inputs
- Performance estimation from DAX (Pegasus Wf System)
- Network topologies assumed:
  - Single vs. dedicated data-link (throttling in dedicated)

### Tools used

- Slack computation: MATLAB
- Performance analysis: SimGrid

### Experiments performed

- Impact on makespan
- Input buffers & net BW usage

Lalayuza

1542

CAERDY

# Experiments and Results (II): Experiments (1)

Impact on the Workflow Makespan

Network topology	Network bandwidth		
	10Mbps	100Mbps	1Gbps
Single output	193.20 <i>s</i>	61.18 <i>s</i>	47.98 <i>s</i>
PP without BW throttling	153.15 <i>s</i>	57.17 <i>s</i>	47.58 <i>s</i>
PP with BW throttling	153.32 <i>s</i>	57.21 <i>s</i>	47.58 <i>s</i>

- Data-throttling looses (insignificantly)
- Correlation <u>data transmission</u> (as indicated by Park & Humphrey) <u>computation</u>
  - Verified by our results

R.J. Rodríguez et al. Automating Data-Throttling Analysis for Data-Intensive Workflow CCGrid'12 18 / 25

Input Buffers and Network Bandwidth Usage

# Experiments and Results (II): Experiments (2)

Input Buffers and Network Bandwidth Usage - some plots



• Data-throttling has great impact on input buffers • Outperforms both other topologies

R.J. Rodríguez et al. Automating Data-Throttling Analysis for Data-Intensive Workflow CCGrid'12 19 / 25

#### Related Work

### Outline

### Motivation

- Knowing the problem
- Data-movement policies
- The goal & our approach

#### 2 Background

- Petri nets
- From DAG to PN
- Slack Concept

### 3 Automating Data-Throttling Analysis

#### Experiments and Results

- Impact on the Workflow Makespan
- Input Buffers and Network Bandwidth Usage

### Related Work

#### Conclusions and Future Work



### Related Work

- PNs already used in scientific workflow community
  - GWorkflowDL, GridFlow, FlowManager
- Overhead analysis (Nerieri et al.)
  - Load imbalance + data movement
- Data throttling issue (Park & Humphrey)

1542



### Related Work

- PNs already used in scientific workflow community
  - GWorkflowDL, GridFlow, FlowManager
- Overhead analysis (Nerieri et al.)
  - Load imbalance + data movement
- Data throttling issue (Park & Humphrey)
- Performance analysis
  - Hybrid Bayesian-neural network (Duan et al., predicts execution time of tasks)
  - Parametrised PN-based model (Tolosana-Calasanz et al.)
- Structural analysis (Van der Alst & Van Hee)
  - WF-nets
    - Operations: sequence, choice, synchroniser, fork, merge
    - Analysis: correctness, deadlock analysis, liveness

(AERDY

Nice work, but...

## How can I use this \*fancy\* approach?



Nice work, but...

# How can I use this \*fancy\* approach?

### Workflows characteristics

- Synchronisation points/merge tasks
- DAX (DAG in XML format)  $\rightarrow$  PNML (PN in XML format)
  - Automatic transformation



Nice work, but...

# How can I use this \*fancy\* approach?

### Workflows characteristics

- Synchronisation points/merge tasks
- DAX (DAG in XML format)  $\rightarrow$  PNML (PN in XML format)
  - Automatic transformation

<pre>k?xml version="1.0" encoding="UTF-8"?&gt; <l 2008-09-24t14:28:09-07:00="" generated:=""> <l [7?]="" by:="" generated="" shishir=""> <l [7?]="" by:="" generated="" shishir=""> <adog 1.0"="" ?="" childcount="; &lt;l part 1: list of all referenced files (may be empty)&gt; &lt;l-&gt; part 2: definition of all jobs (at least one)&gt; &lt;li&gt;&lt;job id=" dax"="" encoding="iso-8859-1" file="cellsa-toldex=00000e-citb20000e-citb20000e-titb2000e-titb200e-titb2000e-titb2000e-titb2&lt;/th&gt;&lt;th&gt;&lt;pre&gt;k?xml version=" filecount="0" http:="" id08000"="" index="0" jobcount="25" link="input" name="http://powers" nomespace="Montage" pegasus.isi.edu="" pi="" register="true" schema="" thtp:="" true"="" trues="" uses="" xmln="http://pegasus.isi.edu/schema/DAX" xmlns:xsi="http xsi:schemaLocation="> <pre>cpnml&gt; <net id="Net-One" type="P/T net"> <pre>cplace id="P0"&gt; <pre>cgraphics&gt; <position x="180.0" y="50.0"></position>  <name> <sulue>P0 <graphics> <offset x="-15.0" y="15.0"></offset> </graphics> <place-transitiontype> <value>D</value> </place-transitiontype></sulue></name></pre></pre></net></pre></adog></l></l></l></pre>	
<pre></pre>	<pre><place-transitiontype> <value>D</value> </place-transitiontype> <showlabel> <position1>1</position1> <position2>@</position2> <position3>@</position3></showlabel></pre>
R.J. Rodríguez et al. Automating Data-Throttling Anal	lysis for Data-Intensive Workflow CCGrid'12 22 / 25

### Outline

### Motivation

- Knowing the problem
- Data-movement policies
- The goal & our approach

#### 2 Background

- Petri nets
- From DAG to PN
- Slack Concept

### 3 Automating Data-Throttling Analysis

#### Experiments and Results

- Impact on the Workflow Makespan
- Input Buffers and Network Bandwidth Usage

### Related Work

### Conclusions and Future Work



Conclusions and Future Work

# Conclusions and Future Work

### Conclusions

- Transfer as-fast-as-possible is not always the best option
  - Network bandwidth misuse
  - Input buffer misuse



Conclusions and Future Work

# Conclusions and Future Work

### Conclusions

- Transfer as-fast-as-possible is not always the best option
  - Network bandwidth misuse
  - Input buffer misuse
- Strategy for computing data-throttling values proposed
  - Main drawbacks
    - Needs previous performance information
    - Scalability

### Future Work

- Test in more realistic environments
- Extend to other workflows
- Improve strategy computation (reduce its complexity)

1542

( A<sup>E</sup>RDY

# Automating Data-Throttling Analysis for Data-Intensive Workflow

**Ricardo J. Rodríguez**, Rafael Tolosana-Calasanz, Omer F. Rana {rjrodriguez, rafaelt}@unizar.es, o.f.rana@cs.cardiff.ac.uk



Universidad de Zaragoza Zaragoza, Spain



Cardiff University Cardiff, United Kingdom

May 15<sup>th</sup>, 2012

CCGrid'12: 12<sup>th</sup> IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing Ottawa, Canada