Introduction	Our approach 0	Tackling problem 1 000000000000000000000	Tackling problem 2 & 3	Conclusion

Toward an autonomic approach of workflows distribution on cloud

Hadrien Croubois

PhD Student at Avalon, Laboratoire de l'informatique du Parallélisme École Normale Supérieure de Lyon, France Supervised by Eddy Caron



Introduction	Our approach o	Tackling problem 1 00000000000000000000	Tackling problem 2 & 3	Conclusion
	Ŭ			

Introduction

Introduction	Our approach	Tackling problem 1	Tackling problem 2 & 3	Conclusion		
0						
The scheduling prob	The scheduling problem					



Introduction	Our approach	Tackling problem 1	Tackling problem 2 & 3	Conclusion		
0						
The scheduling prob	The scheduling problem					



• Jobs definition is changing,

Introduction	Our approach	Tackling problem 1	Tackling problem 2 & 3	Conclusion		
0						
The scheduling prob	The scheduling problem					



- Jobs definition is changing,
- Resources are changing.

Introduction	Our approach	Tackling problem 1	Tackling problem 2 & 3	Conclusion		
0						
The scheduling prob	The scheduling problem					



Challenge

The matching logic need to consider chose changes.



Challenge: Accounting for three factors

- Complex jobs (workflows);
- Multi-Tenant (collaborative);
- Dynamic platform (laaS cloud) with DaaS storage.



Figure 1: Rendering workflow in Natron



Challenge: Accounting for three factors

- Complex jobs (workflows);
- Multi-Tenant (collaborative);
- Dynamic platform (laaS cloud) with DaaS storage.



Figure 1: Rendering workflow in Natron

State of the art

Previous work considers at most 2 of those 3 factors.

Introduction	Our approach	Tackling problem 1	Tackling problem 2 & 3	Conclusion

Our approach



Figure 2: Framework architecture

Introduction	Our approach 0	Tackling problem 1	Tackling problem 2 & 3 00000	Conclusion

Tackling problem 1



Introduction	Our approach 0	Tackling problem 1 ●000000000000000000000000000000000000	Tackling problem 2 & 3 00000	Conclusion
Task clustering	of DAG scheduling - Th	e DCP algorithm		
Algor	ithm 1 DCP st	atic scheduling algo	rithm	
$\mathcal{C} \leftarrow$	- empty clusteri	ng	⊳ (one node per	task)
com	pute <i>BL</i> and 7	<i>L</i> for each task usin	ig ${\mathcal C}$	
whi	i le ∃ unmarked	dependency between	n tasks do	
	$(u, v) \leftarrow edge$	with the largest pa	ath length (most crit	tical).
Res	olve ties by edg	e size (select largest;	t).	
	$\mathcal{C}' \leftarrow \mathcal{C}.mergeC$	lusters(u, v)		
	compute <i>BL</i> ′ a	nd <i>TL</i> ' for each task	\mathfrak{c} using \mathcal{C}'	
	if DCPL(BL', 7	$TL') \leq DCPL(BL, TL)$	L) then	
	(\mathcal{C}, TL, BL)	$\leftarrow (\mathcal{C}', TL', BL')$		
	end if			
	mark(u, v)			
end	l while			

 $\textit{return} \ \mathcal{C}$

Y.-K. Kwok and I. Ahmad, "A static scheduling algorithm onto multiprocessors," in *Proceedings of the* 1994 International Conference on Parallel Processing.

Introduction	Our approach	Tackling problem 1	Tackling problem 2 & 3	Conclusion		
		000000000000000000000000000000000000000				
Task clustering of DAG scheduling - The DCP algorithm						

$$c(u, v) = \begin{cases} 0 & \text{if } \mathcal{C}(u) = \mathcal{C}(v) \\ \omega(u \to v) & \text{otherwise} \end{cases}$$
$$TL(v) = \begin{cases} 0 & \text{if } v \text{ has no predecessor} \\ \max(TL(u) + \omega(u) + c(u, v), \\ u \in pred(v) & avail_{TL}(\mathcal{C}, v)) \end{cases}$$
$$BL(u) = \begin{cases} \omega(u) & \text{if } u \text{ has no successor} \\ \omega(u) + \max(c(u, v) + BL(v), \\ v \in succ(u) & avail_{BL}(\mathcal{C}, u)) \end{cases}$$

Introduction	Our approach	Tackling problem 1	Tackling problem 2 & 3	Conclusion		
		000000000000000000000000000000000000000				
Task clustering of DAG scheduling - The DCP algorithm						

Underlying network topology



Introduction	Our approach	Tackling problem 1	Tackling problem 2 & 3	Conclusion	
00		000000000000000	00000	0000	
Communication model					

Interference between concurrent communications



Figure 3: Transferring files between one node in sagittaire cluster (grid5000) and a DaaS (storage5K)

Introduction	Our approach	Tackling problem 1	Tackling problem 2 & 3	Conclusion
From dependencies t	o dataflows			

See dependencies from the Dataflow point of view



Figure 4: Representation of a fork-join DAG with n=5 independent jobs.



Figure 5: A generic model of DaaS-based network topology.





Figure 6: Preview of the communications between two tasks for a data-based workflow on a DaaS-based platform.

Worst case communications

15/34

$$c(u, v) = \sum_{\substack{d \in edges \\ u=d.src \\ v \notin d.dst}} \frac{d.size}{network_up}$$

$$+ \sum_{\substack{d \in edges \\ u=d.src \\ v \in d.dst}} \frac{d.size}{min (network_up, network_down)}$$

$$+ \max_{\substack{d \in edges \\ u=d.src \\ v \in d.dst}} \frac{d.size}{max (network_up, network_down)}$$

$$+ \sum_{\substack{d \in edges \\ u=d.src \\ v \in d.dst}} \frac{d.size}{network_down}$$

Introduction	Our approach	Tackling problem 1	Tackling problem 2 & 3	Conclusion
		00000000000000000		
Adapting the Critica	l Path computation			

Locality

$$islocal(d, v) = \begin{cases} 1 & \text{if } proc(d.src) = proc(v) \\ 0 & \text{otherwise} \end{cases}$$
$$islocal(d) = \prod_{v \in d.dst} islocal(d, v)$$



Worst case communications with locality

$$c_{loc}(u, v) = 0 \text{ if } proc(u) = proc(v)$$

$$c_{loc}(u, v) = \sum_{\substack{d \in dges \\ v \notin d. dst}} \frac{d.size}{network_up}$$

$$+ \sum_{\substack{d \in dges \\ v \notin d. dst}} \frac{d.size}{min (network_up, network_down)} + \max_{\substack{d \in dges \\ u = d.src}} \frac{d.size}{max (network_up, network_down)}$$

$$+ \sum_{\substack{d \in dges \\ v \notin d. dst}} \frac{d.size}{network_down}$$

$$+ \sum_{\substack{d \in dges \\ v \notin d. dst}} \frac{d.size}{network_down}$$



Worst case communications with locality

Different nodes w\ no locality, worst case communications $i \rightarrow j$



Same node, no communication i→j



Different nodes w\ locality, worst case communication $i \rightarrow j$





Figure 7: Preview of the critical path computation taking the machine network availability into account in DaaS-based platform.





Figure 8: Comparison of the different clustering policies (Gantt charts and their associated makespan) for multi-data fork-join DAG (n = 16).

20/34

Introduction	Our approach	Tackling problem 1	Tackling problem 2 & 3	Conclusion
		000000000000000000000000000000000000000		
Results				

DAC	Algorithm	#Nodos	Makespan	Cost
DAG	Algorithm	#Noues	(t)	(core×t)
	One task per node	18	22.024	67.204
Single Data Fork-join	Single node	1	18.000	18.000
	DCP	14	18.024	56.168
	DaaS aware DCP	2	13.012	20.012
	One task per node	18	37.024	82.204
Multiple Data Fark isin	Single node	1	18.000	18.000
Multiple Data Fork-join	DCP	14	33.803	70.156
	DaaS aware DCP	5	14.000	26.048

Figure 9: Cost and makespan details of the different clustering policies for single-data or multi-data fork-join DAG (n = 16).

Introduction	Our approach	Tackling problem 1	Tackling problem 2 & 3	Conclusion
		000000000000000000000000000000000000000		
Results				



Hadrien Croubois and Eddy Caron.

Communication-aware task placement for workflow scheduling on daas-based cloud.

In Submitted to PDCO2017 (IPDPS workshop), 2017.

Introduction	Our approach	Tackling problem 1	Tackling problem 2 & 3	Conclusion
		00000000000000000		
Future work				

Instance selection

Cloud platforms are heterogeneous ! We need to select which type of instance to run on.

Introduction	Our approach 0	Tackling problem 1 00000000000000000000	Tackling problem 2 & 3	Conclusion

Tackling problem 2 & 3



Introduction	Our approach	Tackling problem 1	Tackling problem 2 & 3	Conclusion	
			0000		
Platform management					

Multiple challenges

- Have the computing power needed to meet deadlines;
- Have the lowest deployment cost possible;
- React to sudden increase/decrease of the workload.

Introduction	Our approach 0	Tackling problem 1 0000000000000000000	Tackling problem 2 & 3 ●○○○○	Conclusion	
Platform management					

Multiple challenges

Challenge

Build a fully automatic solution.

• React to sudden increase/decrease of the workload.

Introduction	Our approach 0	Tackling problem 1 0000000000000000	Tackling problem 2 & 3 ○●○○○	Conclusion
Platform manageme	ıt			

Strategy

- Ranking of ready tasks based on priority;
- Available nodes execute tasks based on their priority (+ other parameter ?);
- Available nodes try to avoid wasting money by committing suicide when no work is available and remaining alive would cost money (end of hour);
- Nodes are deployed when current pool of workers is insufficient.

Introduction	Our approach 0	Tackling problem 1 0000000000000000	Tackling problem 2 & 3 ○●○○○	Conclusion	
Platform management					

Strategy

- Ranking of ready tasks based on priority;
- Available nodes execute tasks based on their priority (+ other parameter ?);

Challenge

Determine when to deploy new nodes : an autonomic loop ! cost money (end of nour);

• Nodes are deployed when current pool of workers is insufficient.





Figure 10: Preview of the scheduler and platform manager structure (threshold with autonomic controler).

Introduction	Our approach 0	Tackling problem 1 00000000000000000000	Tackling problem 2 & 3 ○○○●○	Conclusion
Autonomic loop				



Figure 11: Preview of the scheduler and platform manager structure (simulation).

Toward an autonomic approach of workflows distribution on cloud

Introduction	Our approach 0	Tackling problem 1	Tackling problem 2 & 3 ○○○○●	Conclusion
Autonomic loop				

Current & Future work



Introduction	Our approach	Tackling problem 1	Tackling problem 2 & 3	Conclusion
00	0	0000000000000000	00000	0000

Conclusion

Introduction	Our approach 0	Tackling problem 1 0000000000000000000	Tackling problem 2 & 3	Conclusion ●○○○
Current status				

Current status

- We have identified the problem;
- We have a framework for building a solution;
- We have submitted a paper (PDCO) dealing with the clustering;
- We have started a collaboration with Ctrl-A for building an autonomic loop.

Introduction	Our approach	Tackling problem 1	Tackling problem 2 & 3	Conclusion
				0000
What is next ?				

What is next ?

- Improving the static analysis step (instance selection),
- Autonomic ressource management,
 - $\rightarrow\,$ Co-design of the scheduling component with Aurélie (Roma)
- Implementing and validation the model.

Introduction	Our approach 0	Tackling problem 1 00000000000000000000	Tackling problem 2 & 3	Conclusion ○○●○
What is next ?				

Roadmap



Hadrien Croubois Toward an autonomic approach of workflows distribution on cloud

Introduction	Our approach	Tackling problem 1	Tackling problem 2 & 3	Conclusion
				0000
What is next ?				

Thank you for your attention.