

Characterizing the Algorithmic Complexity of Reconfigurable Data Center Architectures

Klaus-T. Foerster (U. Vienna), Manya Ghobadi (Microsoft Research), Stefan Schmid (U. Vienna)

23 July 2018, IEEE/ACM ANCS 2018





Helios (core) Farrington et al., SIGCOMM '10



Rotornet (rotor switches) Mellette *et al.,* SIGCOMM '17



c-Through (HyPaC architecture) Wang et al., SIGCOMM '10



Solstice (architecture & scheduling) Liu et al., CoNEXT '15



REACTOR Liu *et al.,* NSDI '15



ProjecToR interconnect Ghobadi *et al.,* SIGCOMM '16



... and many more ...





- Results and conclusions often not portable
 - Between topologies/technologies



- Results and conclusions often not portable
 - Between topologies/technologies
- Assumption in routing takes away optimality



- Results and conclusions often not portable
 - Between topologies/technologies
- Assumption in routing takes away optimality
- We take a look from a theoretical perspective



- Results and conclusions often not portable
 - Between topologies/technologies
- Assumption in routing takes away optimality
- We take a look from a theoretical perspective
 - With average path length as an objective



- Results and conclusions often not portable
 - Between topologies/technologies
- Assumption in routing takes away optimality
- We take a look from a theoretical perspective
 - With average path length as an objective
 - For one switch (with/without this assumption)
 - Also briefly for multiple switches









Communication frequency: A→E: 10, A→G: 5





Communication frequency: A→E: 10, A→G: 5





Communication frequency: $A \rightarrow E: 10, A \rightarrow G: 5$

Weighted average path length: **4*10+6*5=70**





Communication frequency: $A \rightarrow E: 10, A \rightarrow G: 5$

Weighted average path length: 4*10+6*5=70

















Weighted average path length: **4*10+6*5=70** *static*









Weighted average path length: **1*10+6*5=40 1*10+(1+2)*5=25**

reconfig







reconfig **1*10+(1+2)*5=25** Weighted average path length: **1*10+6*5=40**



optimum



Beyond a Single Switch

• Especially important at scale: **multiple** reconfigurable switches



Rotornet Mellette *et al.,* SIGCOMM '17 A Tale of Two Topologies Xia et al., SIGCOMM '17





• Model: Either just 1 reconfig or just static



• Model: Either just 1 reconfig or just static



Communication frequency: $A \rightarrow E: 10, A \rightarrow G: 5$



• Model: Either just 1 reconfig or just static



Why this solution?

Communication frequency: $A \rightarrow E: 10, A \rightarrow G: 5$



• Model: Either just 1 reconfig or just static



Communication frequency: $A \rightarrow E: 10, A \rightarrow G: 5$

Why this solution?

Benefit of A→E: 10:

• Static-Reconfig: 40-10=30

Benefit of A→G: 5:

• Static-Reconfig: 30-5=25



• Model: Either just 1 reconfig or just static



- Model: Either just 1 reconfig or just static
- Optimal solution in polynomial time:



- Model: Either just 1 reconfig or just static
- Optimal solution in polynomial time:
 - 1. Compute & assign benefit to every matching edge



- Model: Either just 1 reconfig or just static
- Optimal solution in polynomial time:
 - 1. Compute & assign benefit to every matching edge
 - 2. Compute optimal weighted matching



- Model: Either just 1 reconfig or just static
- Optimal solution in polynomial time:
 - 1. Compute & assign benefit to every matching edge
 - 2. Compute optimal weighted matching
 - E.g., weighted Edmond's Blossom algorithm



- Model: Either just 1 reconfig or just static
- Optimal solution in polynomial time:
 - 1. Compute & assign benefit to every matching edge
 - 2. Compute optimal weighted matching
 - E.g., weighted Edmond's Blossom algorithm

• **Downside**: Only optimal under (artificially!?) segregated routing policy!



- Model: Either just 1 reconfig or just static
- Optimal solution in polynomial time:
 - 1. Compute & assign benefit to every matching edge
 - 2. Compute optimal weighted matching
 - E.g., weighted Edmond's Blossom algorithm

- **Downside**: Only optimal under (artificially!?) segregated routing policy!
 - Not optimal under arbitrary routing policies







Can improve routing quality





Can improve routing quality



NP-hard to optimally compute





Can improve routing quality



NP-hard to optimally compute



Already for simple settings (sparse communication patterns, unit weights etc.)





Can improve routing quality



NP-hard to optimally compute



Already for simple settings (sparse communication patterns, unit weights etc.)



Approximation algorithms & restricted topologies





Can improve routing quality



NP-hard to optimally compute



Already for simple settings (sparse communication patterns, unit weights etc.)



Approximation algorithms & restricted topologies

Future Work





Can improve routing quality



NP-hard to optimally compute Already some work in different settings, e.g.:

- network forms a dynamic tree [Schmid et al., ToN '16]
- constant degree and sparse demands [Avin et al., DISC '17] •



Already for simple setting: (sparse communication patterns, unit weights etc.)

degree depends on node popularity [Avin et al., Inf. Pr. Let. '18] (these works assume all links are reconfigurable)



Approximation algorithms & restricted topologies

Future Work





• Makes the setting more scalable 🙂



• Makes the setting more scalable 😳

• But of course, still NP-hard 🛞

(already for one switch)



- Makes the setting more scalable 🙂
- But of course, still NP-hard 🙁

(already for one switch)

Let's make things simpler





• Can we optimize max. path length?



Can we optimize max. path length?
 o For 2 flows?



- Can we optimize max. path length?
 - For 2 flows?
 - -NP-hard again $\ensuremath{\mathfrak{S}}$













Consider weights





Consider weights





Consider weights





Consider weights





• Challenge:



- Challenge:
 - Proper matchings



- Challenge:
 - Proper matchings
 - Polynomial algorithm



- Challenge:
 - Proper matchings
 - Polynomial algorithm
- Idea: Use flow algorithms



- Challenge:
 - Proper matchings
 - Polynomial algorithm
- Idea: Use flow algorithms
 - Min-cost integral flow is polynomial



- Challenge:
 - Proper matchings
 - Polynomial algorithm
- Idea: Use flow algorithms
 - Min-cost integral flow is polynomial





- Challenge:
 - Proper matchings
 - Polynomial algorithm
- Idea: Use flow algorithms
 - Min-cost integral flow is polynomial





- Challenge:
 - Proper matchings
 - Polynomial algorithm
- Idea: Use flow algorithms
 - Min-cost integral flow is polynomial



Unidirectionality



- Challenge:
 - Proper matchings
 - Polynomial algorithm
- Idea: Use flow algorithms
 - Min-cost integral flow is polynomial



Unidirectionality



- Challenge:
 - Proper matchings
 - Polynomial algorithm
- Idea: Use flow algorithms
 - Min-cost integral flow is polynomial









- Challenge:
 - Proper matchings
 - Polynomial algorithm
- Idea: Use flow algorithms
 - Min-cost integral flow is polynomial



Unidirectionality





- Challenge:
 - Proper matchings
 - Polynomial algorithm
- Idea: Use flow algorithms
 - Min-cost integral flow is polynomial



Unidirectionality





- Challenge:
 - Proper matchings
 - Polynomial algorithm
- Idea: Use flow algorithms
 - Min-cost integral flow is polynomial



Unidirectionality



Same conceptual idea





*some small strings attached







Unidirectionality



• Challenge:

A

universität

wień

- Proper matchings
- Polynomial algorithm
- Idea: Use flow algorithms

capacity =1

Min-cost integral flow is polynomial





Summary and Outlook

- <u>one</u> reconfigurable switch

 segregated: Easy. Not optimal.
 not seg.: NP-hard. Improves solutions.
- <u>multiple</u> reconfigurable switches
 multiple flows: NP-hard
 - just one flow: Easy.
- <u>Next</u> steps
 - approximation algorithms
 - special topologies





Characterizing the Algorithmic Complexity of Reconfigurable Data Center Architectures

Klaus-T. Foerster (U. Vienna), Manya Ghobadi (Microsoft Research), Stefan Schmid (U. Vienna)

Thank you! 🙂

23 July 2018, IEEE/ACM ANCS 2018