

Combinatorial optimization done right:  
How to schedule networking tasks and deploy in-network  
processing services in an optimal fashion

Matthias Rost

Technische Universität Berlin

May 30th, 2014

ICSI, University of California, Berkeley

## Part I: How to schedule networking tasks

- How and when should data center synchronize their data?
- How and when should tasks be deployed *within* data centers?

## Part II: How to deploy in-network processing services

- How can we deploy a multicasting service, e.g. IPTV?
- How can we perform in-network analytics efficiently, e.g. distributed IDS?

## Part I: How to schedule networking tasks

- How and when should data center synchronize their data?
- How and when should tasks be deployed *within* data centers?

## Part II: How to deploy in-network processing services

- How can we deploy a multicasting service, e.g. IPTV?
- How can we perform in-network analytics efficiently, e.g. distributed IDS?

## Combinatorial optimization done right / *optimally*

- Mixed-Integer Programming (MIP) is used to tackle these problems
- We outline how to obtain *good* or *strong* formulations

## Part I: How to schedule networking tasks

- How and when should data center synchronize their data?
- How and when should tasks be deployed *within* data centers?

## Part II: How to deploy in-network processing services

- How can we deploy a multicasting service, e.g. IPTV?
- How can we perform in-network analytics efficiently, e.g. distributed IDS?

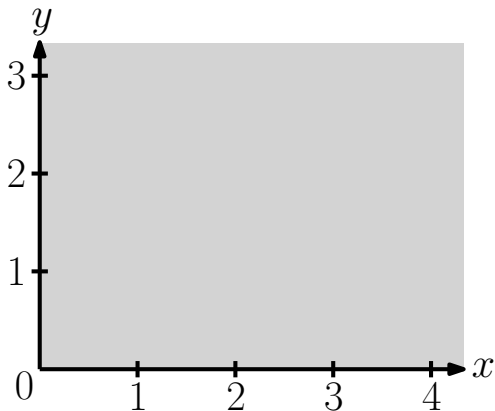
## Combinatorial optimization done right / *optimally*

- Mixed-Integer Programming (MIP) is used to tackle these problems
- We outline how to obtain *good* or *strong* formulations

Please, if you have any questions, ask right away!

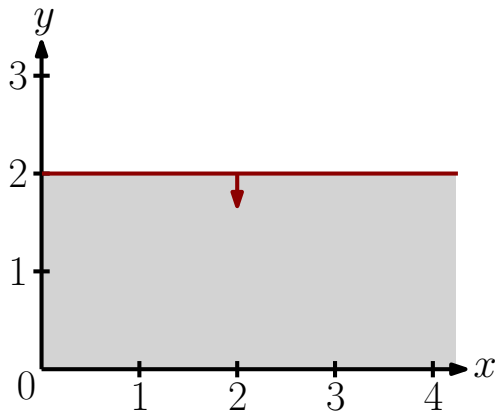
# Mixed-Integer Programming Crash Course

# Excursion: Linear Programming



$$x \geq 0$$
$$y \geq 0$$

# Excursion: Linear Programming

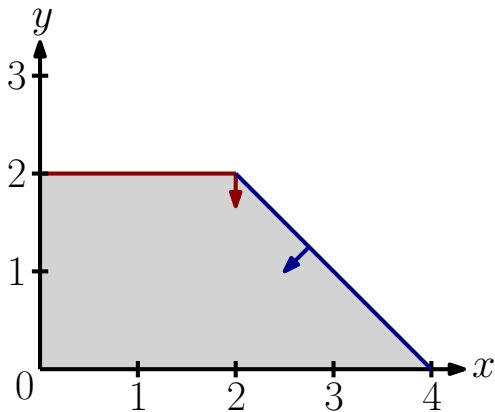


$$x \geq 0$$

$$y \geq 0$$

$$y \leq 2$$

# Excursion: Linear Programming



$$x \geq 0$$

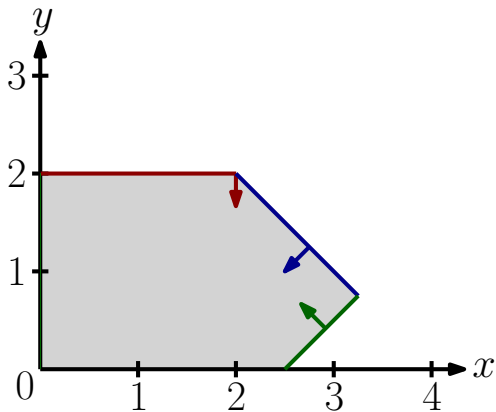
$$y \geq 0$$

$$y \leq 2$$

$$x + y \leq 4$$



# Excursion: Linear Programming



$$x \geq 0$$

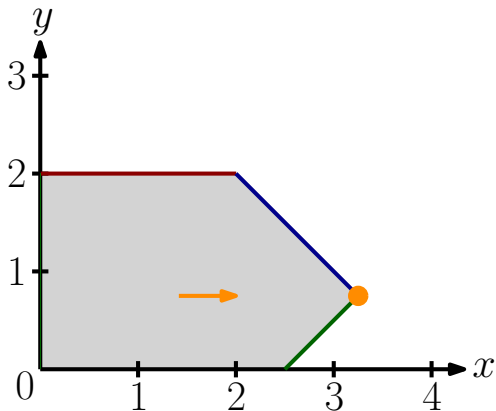
$$y \geq 0$$

$$y \leq 2$$

$$x + y \leq 4$$

$$x - y \geq 2.5$$

# Excursion: Linear Programming



$$x \geq 0$$

$$y \geq 0$$

$$y \leq 2$$

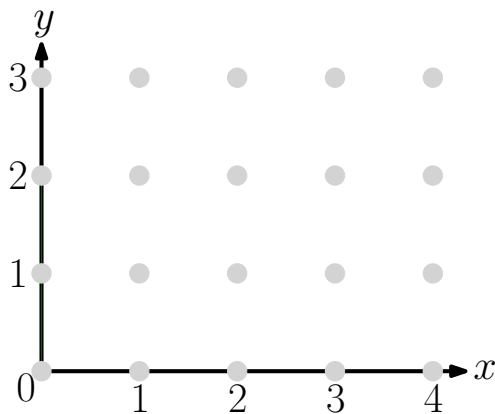
$$x + y \leq 4$$

$$x - y \geq 2.5$$

$$\max x$$

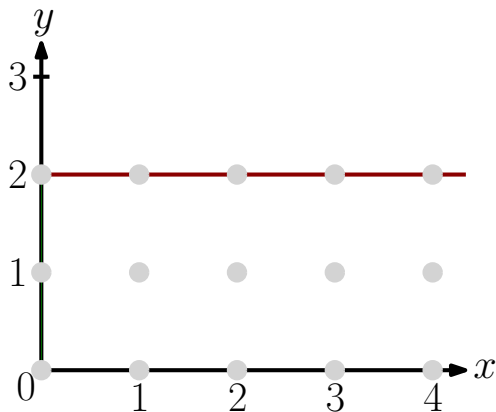
Can be solved in polynomial time!

# Excursion: Integer Programming



$$x \in \mathbb{N}$$
$$y \in \mathbb{N}$$

# Excursion: Integer Programming

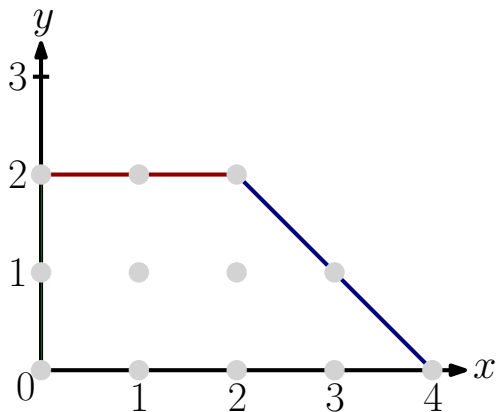


$$x \in \mathbb{N}$$

$$y \in \mathbb{N}$$

$$y \leq 2$$

# Excursion: Integer Programming



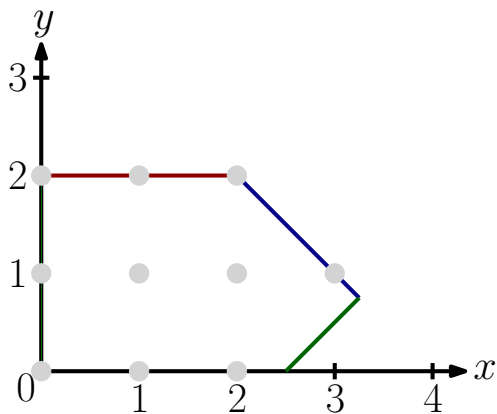
$$x \in \mathbb{N}$$

$$y \in \mathbb{N}$$

$$y \leq 2$$

$$x + y \leq 4$$

## Excursion: Integer Programming



$$x \in \mathbb{N}$$

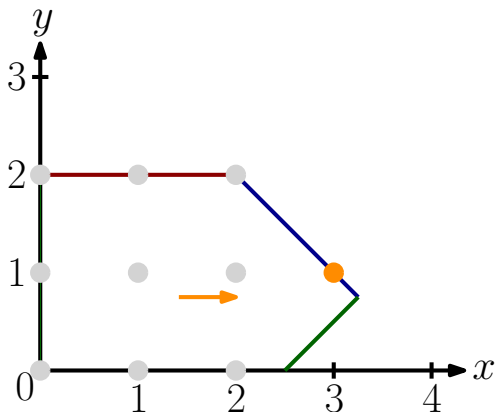
$$y \in \mathbb{N}$$

$$y \leq 2$$

$$x + y \leq 4$$

$$x - y \geq 2.5$$

## Excursion: Integer Programming



$$x \in \mathbb{N}$$

$$y \in \mathbb{N}$$

$$y \leq 2$$

$$x + y \leq 4$$

$$x - y \geq 2.5$$

$$\max x$$

In many cases: NP-hard!

# Why bother with (Mixed-)Integer Programming?

## In general:

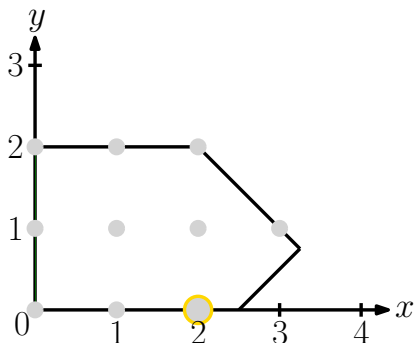
- Easy to formulate all kinds of optimization problems
- Highly optimized solvers which can even solve large problems

## Why we bother:

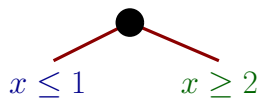
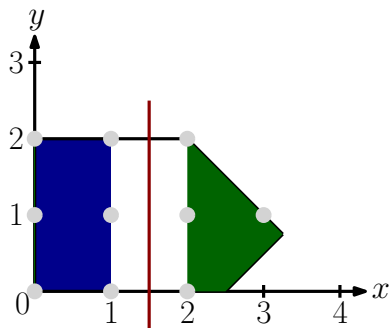
- Important as baseline for polynomial time heuristics
- Allows trading off runtime with solution quality



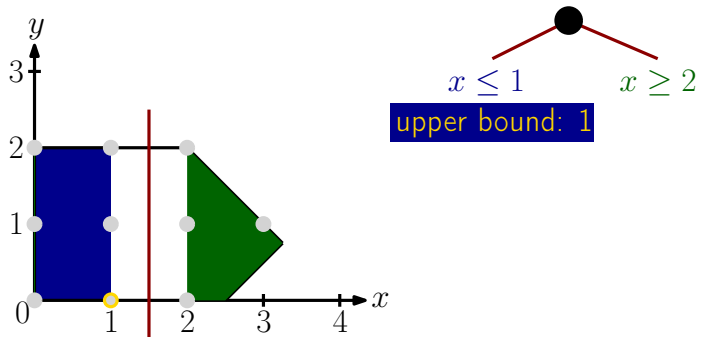
# Excursion: Solving MIPs via Branch-and-Bound



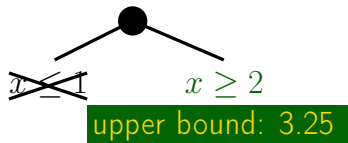
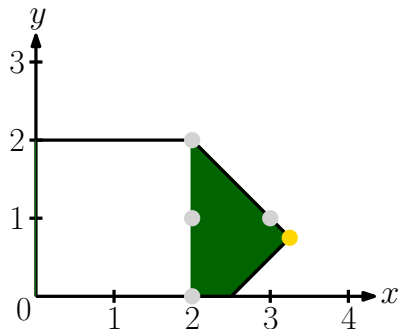
# Excursion: Solving MIPs via Branch-and-Bound



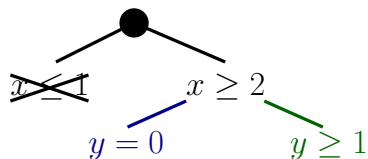
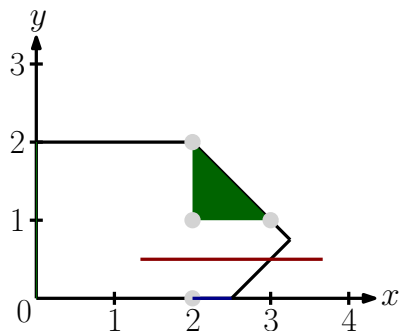
# Excursion: Solving MIPs via Branch-and-Bound



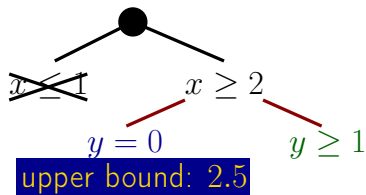
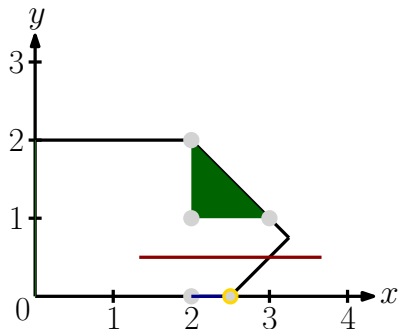
# Excursion: Solving MIPs via Branch-and-Bound



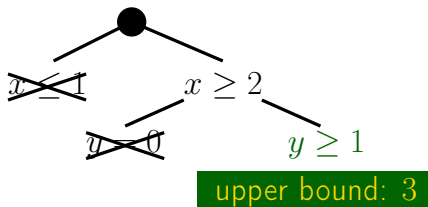
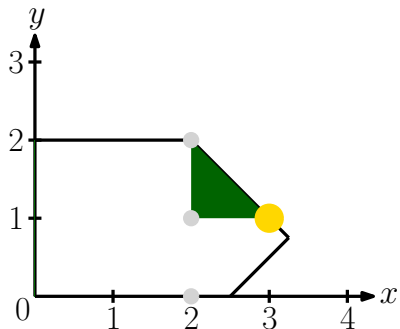
# Excursion: Solving MIPs via Branch-and-Bound



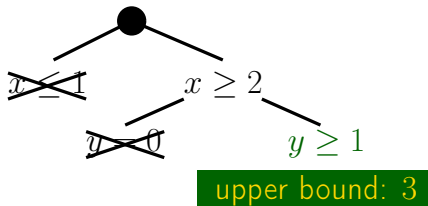
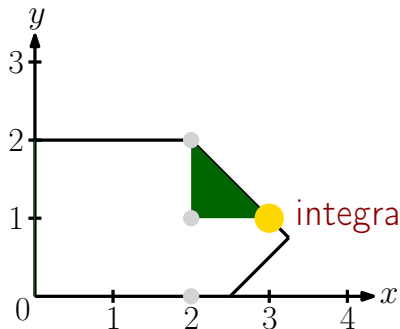
# Excursion: Solving MIPs via Branch-and-Bound



# Excursion: Solving MIPs via Branch-and-Bound



# Excursion: Solving MIPs via Branch-and-Bound

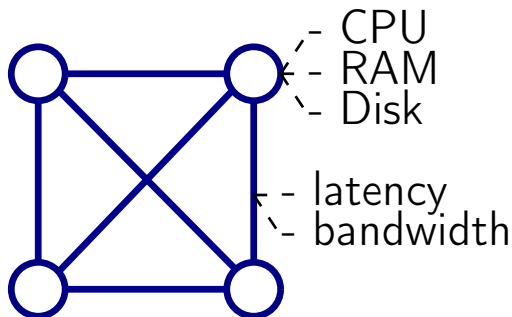




# General Setting

## Key Ingredient: Virtualization

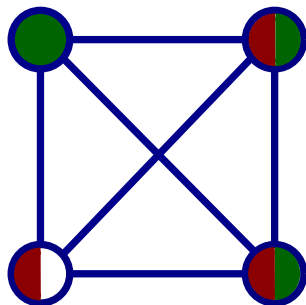
## Physical Network / Substrate



## Key Ingredient: Virtualization

### Node Virtualization

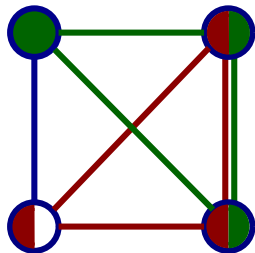
- data center: VMs
- wide-area: NFV



# Key Ingredient: Virtualization

## Link Virtualization

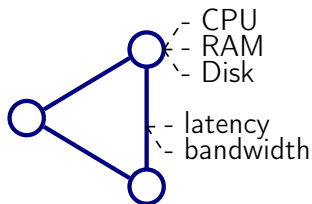
- data center: SDN
- wide-area: SDN / MPLS



# Part I: Scheduling Network Tasks

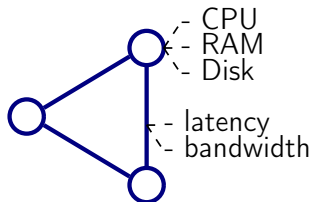
# The Virtual Network Embedding Problem (VNEP)

Physical Network

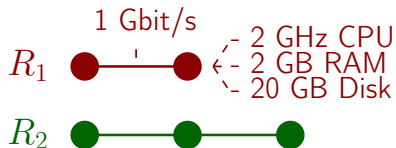


# The Virtual Network Embedding Problem (VNEP)

## Physical Network

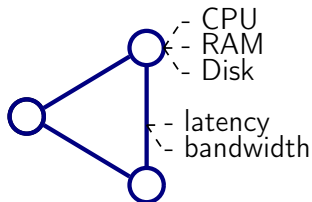


## Virtual Network Requests

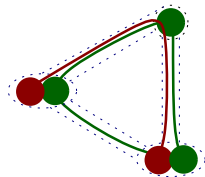
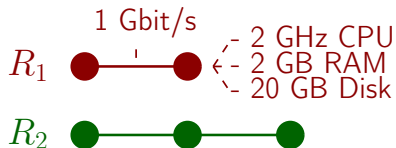


# The Virtual Network Embedding Problem (VNEP)

## Physical Network



## Virtual Network Requests



## Embedding

- map virtual onto substrate nodes
- map virtual links onto substrate paths
- obeying the substrate's capacities



# Facets of the VNEP

## Setting

(De)centralized

Multi-Provider

Reliability

Reconfigurations

## Objectives

Access Control

Load Balancing

Energy Savings

## Algorithms

Exact

Heuristic

# Related Work

TABLE III  
TAXONOMY OF CONSOLE VNE APPROACHES

| Category                           | Reference                           | Optimization                  | Coordination  | Contribution  |  |
|------------------------------------|-------------------------------------|-------------------------------|---|---|--|
| CNC                                | [35] Inführ and Kauf (2011)         | Exact                         | One-Stage   | Provides delay, location and routing constraints  |  |
|                                    | [37] Liu et al. (2011)              | Exact                         | One-Stage   | Exact VNE based on correspondence matrices  |  |
|                                    | [36] Trink et al. (2011)            | Exact                         | One-Stage   | Exact VNE problem with SLA QoS guarantees   |  |
|                                    | [32] Pagès et al. (2012)            | Exact/Metaheuristic           | One-Stage   | Introduces the VNE for optical networks   |  |
|                                    | [31] Liskhai and Kauf (2009)        | Heuristic                     | One-Stage   | Provides one-stage VNE. Based on RD   |  |
|                                    | [62] Di et al. (2010)               | Heuristic                     | One-Stage   | Improvement of the approach in [61]   |  |
|                                    | [63] Okunur and Sauman (2011)       | Heuristic                     | One-Stage   | Introduces hierarchical management of the SN  |  |
|                                    | [64] [63] Yu et al. (2011-2012)     | Heuristic                     | One-Stage   | First VNE approach in wireless multiplex networks. Introduces metrics and feasibility measures for wireless VNE     |  |
|                                    | [65] Chen et al. (2012)             | Heuristic                     | One-Stage   | Reduces resource fragmentation  |  |
|                                    | [66] Yu et al. (2012)               | Heuristic                     | One-Stage   | One-step VNE that increases coordination  |  |
|                                    | [67] Liu et al. (2011)              | Heuristic                     | Two-Stage   | Implements coordination based on nodes partition  |  |
|                                    | [21] [68] Sheng et al. (2011-2012)  | Heuristic                     | Two-Stage   | Opportunistic resource sharing to deal with load fluctuation  |  |
|                                    | [69] Li et al. (2012)               | Heuristic                     | Two-Stage   | Topology awareness to enforce VNE coordination  |  |
|                                    | [67] Lu and Tamer (2006)            | Heuristic                     | Uncoordinated   | Embedding in specific backbone-size VN topologies   |  |
|                                    | [62] Yu et al. (2008)               | Heuristic                     | Uncoordinated   | Utilizes the KSP algorithm [13] for VLM   |  |
|                                    | [33] Razaq and Siraj (2010)         | Heuristic                     | Uncoordinated   | Different K values in KSP based VLM   |  |
|                                    | [31] Razaq et al. (2011)            | Heuristic                     | Uncoordinated   | Investigates the VNE impact of bottlenecked nodes   |  |
|                                    | [32] Nagawa et al. (2011)           | Heuristic                     | Uncoordinated   | VNE considering SN resources heterogeneity  |  |
|                                    | [31] Leinfelder et al. (2011)       | Heuristic                     | Uncoordinated   | Introduces VNE for wireless network testbeds  |  |
|                                    | [32] [39] Botoni et al. (2011-2013) | Heuristic                     | Uncoordinated   | Introduces hidden hop constraints   |  |
| [34] Zhu and Ansumali (2006)       | Heuristic                           | Uncoordinated                 | Provides a balanced link and node stress in the SN  |   |  |
| [31] Fujitani et al. (2011)        | Metaheuristic                       | One-Stage                     | Max-Min Ant Colony metaheuristic VNE approach   |   |  |
| [32] Chang et al. (2012)           | Metaheuristic                       | One-Stage                     | Accelerates convergence of PSO VNE metaheuristic with topology aware node ranking [32]                |   |  |
| [32] Zhang et al. (2012)           | Heuristic                           | One-Stage                     | Maps one virtual node in several substrate nodes  |   |  |
| [32] Di et al. (2012)              | Heuristic                           | One-Stage                     | Coordinated VNE reducing the number of backtracks by carefully choosing the first virtual node to map |   |  |
| [32] Abdolrer and Baligh (2012)    | Heuristic                           | Uncoordinated                 | Introduces VNE in the optical domain trying to minimize the number of hop per link                    |   |  |
| [32] Aris Leivaditis et al. (2012) | Heuristic                           | Coordinated                   | Consistent importance of virtual nodes for embedding  |   |  |
| [32] Yao-Bin Lee et al. (2012)     | Heuristic                           | InterI#P                      | Clustering of virtual networks in multi-provider environment  |   |  |
| CVR                                | [22] Fujitani et al. (2011)         | Heuristic                     | One-Stage   | Migration of nodes with bottlenecked adjacent links   |  |
|                                    | [41] Buskowiak et al. (2010)        | Heuristic                     | Two-Stage   | Migration when service across position changes  |  |
|                                    | [34] Zhu and Ansumali (2006)        | Heuristic                     | Uncoordinated   | Reduce the cost of periodic reconfiguration   |  |
|                                    | [35] Fan and Ansumali (2006)        | Heuristic                     | Uncoordinated   | Reduces the cost of VNE re-configuration  |  |
|                                    | [41] Cai et al. (2010)              | Heuristic                     | Uncoordinated   | Reconfiguration based on SN evolution   |  |
|                                    | [41] Shan-li and Xue-song (2011)    | Heuristic                     | Uncoordinated   | Identifies mapped virtual nodes and links with not optimal mapping and migrate them to save SN resources            |  |
|                                    | [35] Sun et al. (2012)              | Heuristic                     | Uncoordinated   | Introduces the VNE problem for evolving VNRS  |  |
|                                    | DNC                                 | [36] [37] Houdi et al. (2010) | Heuristic   | Uncoordinated   | First distributed approach to solve VNE. Proposes a VNE process of using the communication among substrate nodes |
|                                    |                                     | [36] Xia et al. (2011)        | Heuristic   | InterI#P  | Introduces the InterI#P VNE for networked clouds   |
|                                    |                                     | [69] Lv et al. (2011)         | Heuristic   | InterI#P  | InterI#P VNE using hierarchical virtual resource allocation  |
| [42] Houdi et al. (2011)           |                                     | Exact/Metaheuristic           | InterI#P  | VNRS is optimized using each subVN in different h/P. Provides exact and heuristic splitting algorithms              |  |
| [36] Leivaditis et al. (2012)      |                                     | Heuristic                     | InterI#P  | Graph partitioning InterI#P VNE using a heuristic integrating a min-cost algorithm followed by subgraph isomorphism |  |
| D#VC                               | [61] Marqueton et al. (2010)        | Heuristic                     | Uncoordinated   | First distributed dynamic approach. Reorganizes the SN when VNs demands change                                      |  |

TABLE IV  
TAXONOMY OF REDUNDANT VNE APPROACHES

| Category                       | Reference                                       | Optimization  | Coordination  | Contribution  |
|--------------------------------|---|---------------|---|---|
| CNR                            | [43] Houdi et al. (2011)                        | Exact         | One-Stage   | First approach providing an ILP exact solution  |
|                                | [39] Zhang et al. (2011)                        | Exact         | One-Stage   | Optimal resilient solution attaining an enhanced QoS mapping. Provides diversified substrate back-up paths                    |
|                                | [44] Botoni et al. (2012)                       | Exact         | One-Stage   | Introduces the energy aware VNE   |
|                                | [39] Wang and Wolf (2011)                       | Exact         | One-Stage   | Reduces the VNR as a traffic matrix   |
|                                | [39] [65] [66] Shamsi and Rookmeyer (2007-2009) | Heuristic     | One-Stage   | Receiver link failures by providing backup paths with intermediate nodes  |
|                                | [66] Kozlovski et al. (2010)                    | Heuristic     | One-Stage   | Introduces reliability as a service offered by the I#P. Reliable VNEs based on subgraph isomorphism detection                 |
|                                | [66] Yu et al. (2010)                           | Heuristic     | One-Stage   | Introduces fault-dependent protection with a back-up solution for each regional failure                                       |
|                                | [69] Lv et al. (2012)                           | Heuristic     | One-Stage   | Introduces losses to multicast VNE in wireless mesh networks  |
|                                | [66] [37] Chowdhury et al. (2009-2011)          | Heuristic     | Two-Stage   | Coordination in VNE using multi-path for VLM  |
|                                | [45] Rahman et al. (2010)                       | Heuristic     | Two-Stage   | Uses of failure, the economic penalty is minimized by the pre-execution of a bandwidth quota for back-up in SN links          |
|                                | [45] Butt et al. (2010)                         | Heuristic     | Two-Stage   | VNE awareness of the SN bottlenecked resources  |
|                                | [32] Yoon et al. (2010)                         | Heuristic     | Two-Stage   | Introduces sharing among back-up resources. Reduces resources allocated for redundancy  |
|                                | [66] Sun et al. (2011)                          | Heuristic     | Two-Stage   | Resilient VNE optimizing the embedding cost and reducing computational complexity   |
|                                | [66] Yu et al. (2011)                           | Heuristic     | Two-Stage   | Resilient VNE analyzing failures in substrate nodes   |
|                                | [66] Yu et al. (2008)                           | Heuristic     | Uncoordinated   | Introduces the multi-path approach for VLM  |
|                                | [66] Gao et al. (2010)                          | Heuristic     | Uncoordinated   | Improvement of the approach in [45]   |
|                                | [66] Yang et al. (2010)                         | Heuristic     | Uncoordinated   | Divides the SN in regions to reduce VNE complexity  |
|                                | [66] Zhou et al. (2010)                         | Heuristic     | Uncoordinated   | Maps one virtual node to multiple substrate nodes   |
|                                | [66] Chen et al. (2010)                         | Heuristic     | Uncoordinated   | Reduces resiliency protection approach against failures during the online VNE process. Considers just-substrate link failures |
|                                | [66] Yu et al. (2011)                           | Heuristic     | Uncoordinated   | Proactive VNE approach offering protection against SN link failures for links with high stress                                |
| [66] Sun et al. (2011)         | Heuristic                                       | Uncoordinated | Introduces stochastic BW demand to the VNE  |   |
| [66] Lu et al. (2011)          | Heuristic                                       | Uncoordinated | Introduces load balancing in links  |   |
| [66] Guo et al. (2011)         | Heuristic                                       | Uncoordinated | Proactive resilient VLM approach sharing back-up paths                                      |   |
| [66] Cheng et al. (2011)       | Metaheuristic                                   | Two-Stage     | Introduces topology-awareness in VNE  |   |
| [66] Sheng et al. (2011)       | Metaheuristic                                   | Two-Stage     | Embedding time depends on VNR lifetime. Uses simulated annealing metaheuristic              |   |
| [66] Zhang et al. (2012)       | Metaheuristic                                   | Two-Stage     | Introduces particle swarm optimization (PSO) metaheuristic                                  |   |
| [66] Sun et al. (2012)         | Metaheuristic                                   | Two-Stage     | Introduces VNE in multi-datacenter environments   |   |
| [66] Lv et al. (2012)          | Metaheuristic                                   | Uncoordinated | Introduces VNE in wireless mesh networks  |   |
| [66] Leivaditis et al. (2012)  | Heuristic                                       | Two-Stage     | Uses the approach in [42] to solve the VNE for an arbitrary pool of heterogeneous resources |   |
| [66] Maiti and Raghavan (2012) | Heuristic                                       | Two-Stage     | VNE considering the residual capacity of the substrate links                                |   |
| [66] Zhang et al. (2012)       | Exact/Heuristic                                 | One-Stage     | Receiver link failures providing disjoint SN backup paths                                   |   |
| CVR#                           | [35] Butt et al. (2010)                         | Heuristic     | Two-Stage   | Receiver reconfiguration of virtual links and nodes causing reconfiguration to less critical SN regions                       |
|                                | [35] Yu et al. (2010)                           | Heuristic     | Uncoordinated   | Reconfigure the embedding by changing the splitting ratio in the multipath VLM solution                                       |
| DNR                            | [110] Schaffath et al. (2010)                   | Exact         | One-Stage   | ILP-based VNE. Dynamically reconfigures existing mappings   |
|                                | [111] Chen et al. (2011)                        | Heuristic     | Two-Stage   | Topology reconfiguration of SN nodes with high utilization  |
| DNR                            | [36] Chowdhury et al. (2010)                    | Heuristic     | InterI#P  | First InterI#P VNE proposal. Mediates between h/P and SP interactions. VNR is split across h/Ps and embedded to cloud         |
| D#VR                           | [112] Houdi et al. (2010)                       | Heuristic     | Two-Stage   | First-licent VNE that acts upon node and link failures  |

# Related Work

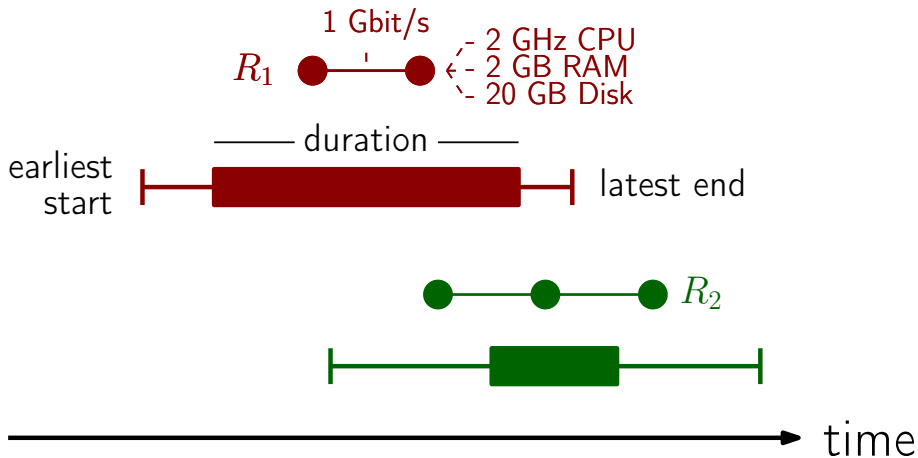
TABLE III  
TAXONOMY OF COEXIST VNE APPROACHES

| Category                           | Reference                           | Optimization                 | Coordination  | Contribution   |
|------------------------------------|-------------------------------------|------------------------------|---|--|
| C/NC                               | [26] Infilar and Radli (2011)       | Exact                        | One Stage   | Provides delay, location and routing constraints   |
|                                    | [27] Liu et al. (2011)              | Exact                        | One Stage   | Exact VNE based on correspondence matrices   |
|                                    | [28] Truh et al. (2011)             | Exact                        | One Stage   | Exact VNE problem with SLA QoS parameters  |
|                                    | [29] Pappas et al. (2012)           | Exact/Metaheuristic          | One Stage   | Introduces the VNE for optical networks  |
|                                    | [30] Lucchi and Kati (2009)         | Heuristic                    | One Stage   | Provides one stage VNE based on SD   |
|                                    | [31] Di et al. (2010)               | Heuristic                    | One Stage   | Improvement of the approach in [30]  |
|                                    | [32] Ghazizadeh and Saman (2011)    | Heuristic                    | One Stage   | Introduces hierarchical management of the SN   |
|                                    | [33], [34] Yoo et al. (2011-2012)   | Heuristic                    | One Stage   | First VNE approach in wireless multi-hop networks. Introduces metrics and feasibility measures for wireless VNE            |
|                                    | [35] Chen et al. (2012)             | Heuristic                    | One Stage   | Reduces resource fragmentation   |
|                                    | [36] Yu et al. (2012)               | Heuristic                    | One Stage   | One step VNE that increases coordination   |
|                                    | [37] Liu et al. (2011)              | Heuristic                    | Two Stages  | Improves coordination based on nodes proximity   |
|                                    | [38], [39] Sheng et al. (2011-2012) | Heuristic                    | Two Stages  | Opportunistic resource sharing to deal with load fluctuation   |
|                                    | [40] Li et al. (2012)               | Heuristic                    | Two Stages  | Topology awareness in wireless VNE coordination  |
|                                    | [41] Lu and Turner (2006)           | Heuristic                    | Uncoordinated   | Embedding in specific host-based star VN topologies  |
|                                    | [42] Yu et al. (2008)               | Heuristic                    | Uncoordinated   | Utilizes the KSP algorithm [43] for VLM  |
|                                    | [43] Razaq and Siraj (2009)         | Heuristic                    | Uncoordinated   | Different K values in KSP based VLM  |
|                                    | [44] Razaq et al. (2011)            | Heuristic                    | Uncoordinated   | Investigates the VNE impact of bottlenecked nodes  |
|                                    | [45] Nogueira et al. (2011)         | Heuristic                    | Uncoordinated   | VNE considering SN resources heterogeneity   |
|                                    | [46] Lefrançois et al. (2011)       | Heuristic                    | Uncoordinated   | Introduces VNE for wireless network topology   |
|                                    | [47] Dhillon et al. (2011-2013)     | Heuristic                    | Uncoordinated   | Introduces hybrid hop constraint based VNE   |
| [48] Zhang and Turner (2006)       | Heuristic                           | Uncoordinated                | Introduces hop constraint based VNE for SN  |  |
| [49] Pappas et al. (2012)          | Metaheuristic                       | One Stage                    | Multi-Modal Co-optimization of VNE and resource allocation  |  |
| [50] Ghazizadeh and Saman (2011)   | Metaheuristic                       | One Stage                    | Algorithm for VNE in multi-tenant networks  |  |
| [51] Zhang et al. (2012)           | Heuristic                           | Uncoordinated                | Maps one virtual node in several substrate nodes  |  |
| [52] Di et al. (2012)              | Heuristic                           | One Stage                    | Coordinated VNE reducing the number of backtracks by carefully choosing the first virtual node to map |  |
| [53] Abdelkar and Eshghi (2012)    | Heuristic                           | Uncoordinated                | Introduces VNE in the optical domain trying to minimize the number of hops per link                   |  |
| [54] Aris Lefrançois et al. (2012) | Heuristic                           | Coordinated                  | Considers importance of virtual nodes for embedding   |  |
| [55] Bao-Ho Lee et al. (2012)      | Heuristic                           | Interf@P                     | Clustering of virtual networks in multi-provider environment  |  |
| C/VC                               | [56] Fajant et al. (2011)           | Heuristic                    | One Stage   | Migration of nodes with bottlenecked adjacent links  |
|                                    | [57] Bonkowski et al. (2009)        | Heuristic                    | Two Stages  | Migrates when service across position changes  |
|                                    | [58] Zhu and Ansumali (2009)        | Heuristic                    | Uncoordinated   | Reduces the cost of periodic reconfiguration   |
|                                    | [59] Fan and Ansumali (2006)        | Heuristic                    | Uncoordinated   | Reduces the cost of VNE re-configuration   |
|                                    | [60] Cai et al. (2010)              | Heuristic                    | Uncoordinated   | Reconfiguration based on SN evolution  |
|                                    | [61] Shan-ji and Xie-song (2011)    | Heuristic                    | Uncoordinated   | Identifies mapped virtual nodes and links with not optimal mapping and migrate them to save SN resources                   |
|                                    | [62] Sun et al. (2012)              | Heuristic                    | Uncoordinated   | Introduces the VNE problem for existing VNMs   |
| D/NC                               | [63], [64] Hoadi et al. (2010)      | Heuristic                    | Uncoordinated   | First distributed approach to solve VNE. Proposes a VNE protocol to manage the communication among substrate nodes         |
|                                    | [65] Xia et al. (2011)              | Heuristic                    | Interf@P  | Introduces the Interf@P VNE for networked clouds   |
|                                    | [66] Li et al. (2011)               | Heuristic                    | Interf@P  | Interf@P VNE using hierarchical virtual resource organization  |
|                                    | [67] Hoadi et al. (2011)            | Exact/Metaheuristic          | Interf@P  | VNE is split among each sub-VN in different hFs. Provides exact and heuristic splitting approaches                         |
|                                    | [68] Lefrançois et al. (2012)       | Heuristic                    | Interf@P  | Graph partitioning Interf@P VNE using a heuristic integrating a min cut cost algorithm followed by sub-graph decomposition |
|                                    | D/VC                                | [69] Magagnoli et al. (2010) | Heuristic   | Uncoordinated  |

TABLE IV  
TAXONOMY OF REDUNDANT VNE APPROACHES

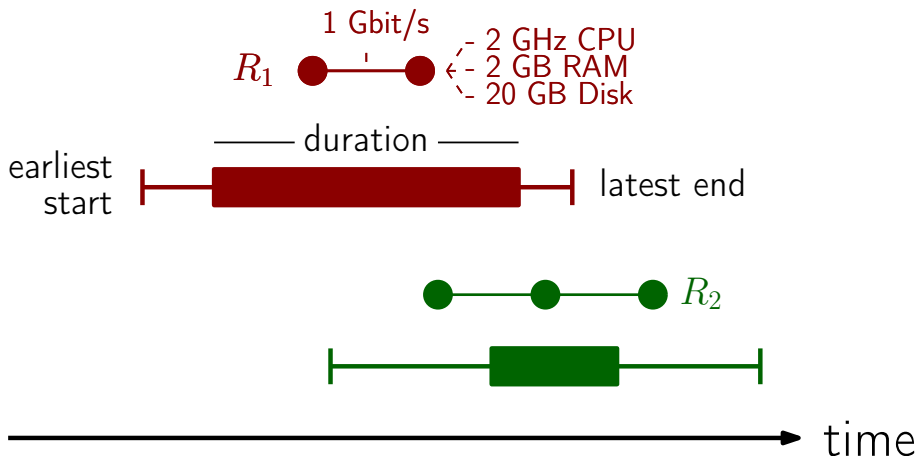
| Category                  | Reference  | Optimization    | Coordination   | Contribution   |
|---------------------------|--|-----------------|--|--|
| C/NC                      | [70] Hoadi et al. (2011)                         | Exact           | One Stage  | First approach providing an ILP-based solution   |
|                           | [72] Zhang et al. (2011)                         | Exact           | One Stage  | Optimal redundant solution attaining an enhanced QoS mapping. Provides distributed substrate back-up paths             |
|                           | [44] Botto et al. (2012)                         | Exact           | One Stage  | Introduces the energy aware VNE  |
|                           | [71] Wang and Wolf (2011)                        | Exact           | One Stage  | Reduces the VNR as a traffic matrix  |
|                           | [64], [65], [66] Shams and Bockmeier (2007-2009) | Heuristic       | One Stage  | Recover link failures by providing backup paths with intermediate nodes  |
|                           | [73] Koslowski et al. (2009)                     | Heuristic       | One Stage  | Introduces reliability as a service offered by the IAP. Reliable VNMs based on subgraph isomorphism detection          |
|                           | [68] Yu et al. (2010)                            | Heuristic       | One Stage  | Introduces failure-dependent protection with a backup solution for each regional failure                               |
|                           | [69] Li et al. (2012)                            | Heuristic       | One Stage  | Introduces losses to indicate VNE in wireless mesh networks  |
|                           | [38], [72] Choudhury et al. (2009-2011)          | Heuristic       | Two Stages   | Coordination in VNE using multi-paths for VLM  |
|                           | [74] Rahman et al. (2010)                        | Heuristic       | Two Stages   | Upon a failure, the economic priority is maintained by the pre-activation of a bandwidth quota for back-up in SN links |
|                           | [33] Butt et al. (2010)                          | Heuristic       | Two Stages   | VNE awareness of the SN bottlenecked resources   |
|                           | [75] Yoo et al. (2010)                           | Heuristic       | Two Stages   | Introduces sharing among back-up resources. Reduces resources allocated for redundancy                                 |
|                           | [109] Sun et al. (2011)                          | Heuristic       | Two Stages   | Resilient VNE optimizing the embedding cost and reducing computational complexity                                      |
|                           | [36] Yu et al. (2011)                            | Heuristic       | Two Stages   | Resilient VNE analyzing failures in substrate nodes  |
|                           | [37] Yu et al. (2011)                            | Heuristic       | Uncoordinated  | Introduces the multi-stage approach for VLM  |
|                           | [76] Guo et al. (2011)                           | Heuristic       | Uncoordinated  | Prevention of the multi-stage approach for VLM   |
|                           | [107] Yang et al. (2011)                         | Heuristic       | Uncoordinated  | Prevention of the multi-stage approach for VLM   |
|                           | [108] Zhu et al. (2011)                          | Heuristic       | Uncoordinated  | Prevention of the multi-stage approach for VLM   |
|                           | [110] Chen et al. (2010)                         | Heuristic       | Uncoordinated  | Prevention of the multi-stage approach for VLM   |
|                           | [111] Chen et al. (2010)                         | Heuristic       | Uncoordinated  | Prevention of the multi-stage approach for VLM   |
| [109] Yu et al. (2011)    | Heuristic  | Uncoordinated   | Proactive VNE approach offering protection against SN link failures for links with high status |  |
| [77] Sun et al. (2011)    | Heuristic  | Uncoordinated   | Introduces stochastic RW demand to the VNE   |  |
| [78] Lu et al. (2011)     | Heuristic  | Uncoordinated   | Introduces load balancing in links   |  |
| [79] Guo et al. (2011)    | Heuristic  | Uncoordinated   | Proactive resilient VLM approach sharing back-up paths   |  |
| [80] Cheng et al. (2011)  | Metaheuristic                                    | Two Stages      | Introduces topology-awareness in VNE   |  |
| [109] Sheng et al. (2011) | Metaheuristic                                    | Two Stages      | Embedding time depends on VNR lifetime. Uses simulated annealing metaheuristic                 |  |
| C/VC                      | [81] Zhang et al. (2012)                         | Metaheuristic   | Two Stages   | Introduces particle swarm optimization (PSO) metaheuristic   |
|                           | [82] Sun et al. (2012)                           | Metaheuristic   | Two Stages   | Introduces VNE in multi-tenant environments  |
|                           | [83] Lu et al. (2012)                            | Metaheuristic   | Uncoordinated  | Introduces VNE in wireless mesh networks   |
|                           | [84] Lefrançois et al. (2012)                    | Heuristic       | Two Stages   | Uses the approach in [82] to solve the VNE for an arbitrary pool of heterogeneous resources                            |
|                           | [85] Masi and Righiani (2012)                    | Heuristic       | Two Stages   | VNE considering the residual capacity of the substrate links   |
|                           | [109] Zhang et al. (2012)                        | Exact/Heuristic | One Stage  | Recover link failures providing disjoint SN backup paths   |
| C/VR                      | [86] Butt et al. (2010)                          | Heuristic       | Two Stages   | Reactive reconfiguration of virtual links and nodes causing reaction to less critical SN regions                       |
|                           | [87] Yu et al. (2010)                            | Heuristic       | Uncoordinated  | Reconfigure the embedding by changing the splitting ratio in the multiple VLM solution                                 |
|                           | [110] Schaffrath et al. (2010)                   | Exact           | One Stage  | ILP-based VNE. Dynamically reconfigure existing mappings   |
| D/NC                      | [88] Chen et al. (2011)                          | Heuristic       | Two Stages   | Periodic reconfiguration of SN nodes with high utilization   |
|                           | [89] Choudhury et al. (2010)                     | Heuristic       | Interf@P   | First Interf@P VNE proposal. Modulates Interf@P and SP metrics. VNOR is split across hFs and methods locally           |
| D/VR                      | [112] Hoadi et al. (2010)                        | Heuristic       | Two Stages   | Fast-redundant VNE that acts upon node and link failures   |

# Our Model



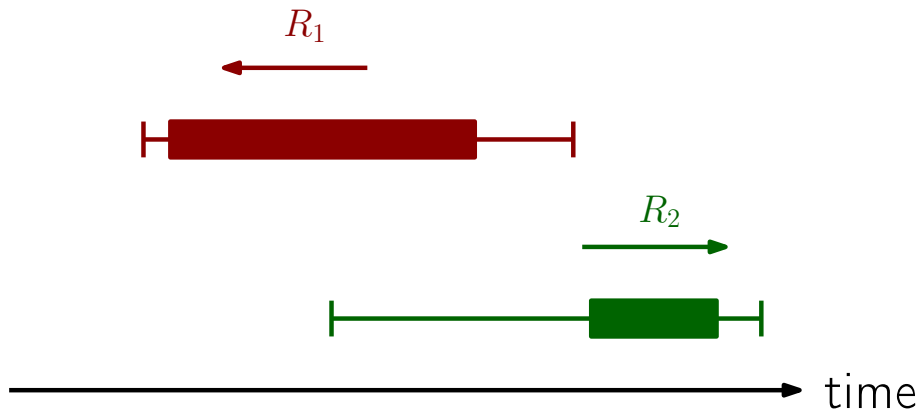
## Our Model

## Offline scenario

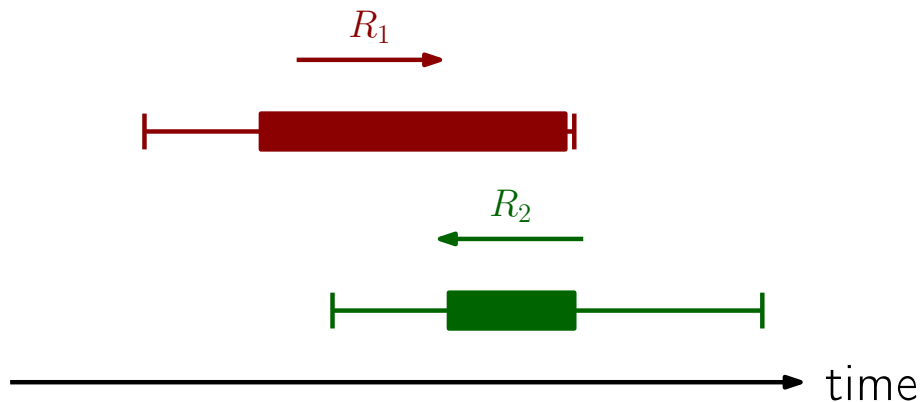


Motivation #1: Business

## Provider Incentives: Minimizing Load



# Provider Incentives: Maximizing Utilization by Collocation



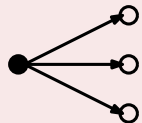


## Motivation #2: Modeling Opportunities

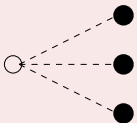
# Modeling Opportunities: Evolution of VNets

## MapReduce [6]

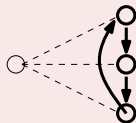
Distribution



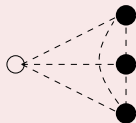
Mapping



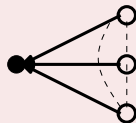
Shuffle



Reduce



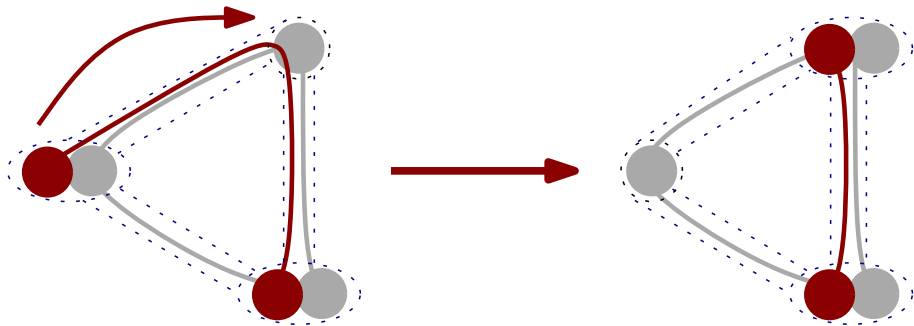
Storing



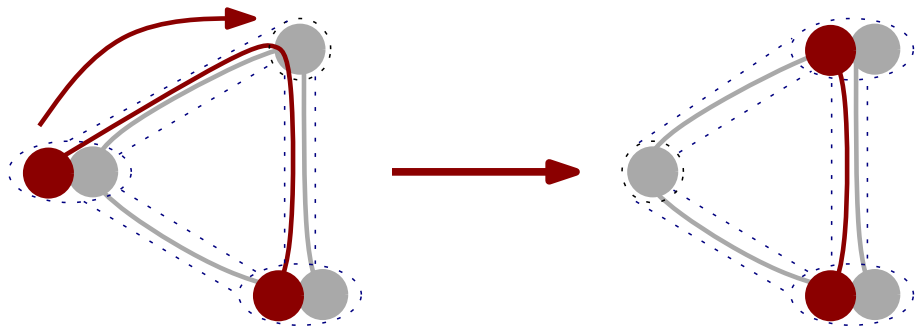
time →

Reservation of maximal allocations over the whole time?

# Modeling Opportunities: Migrations

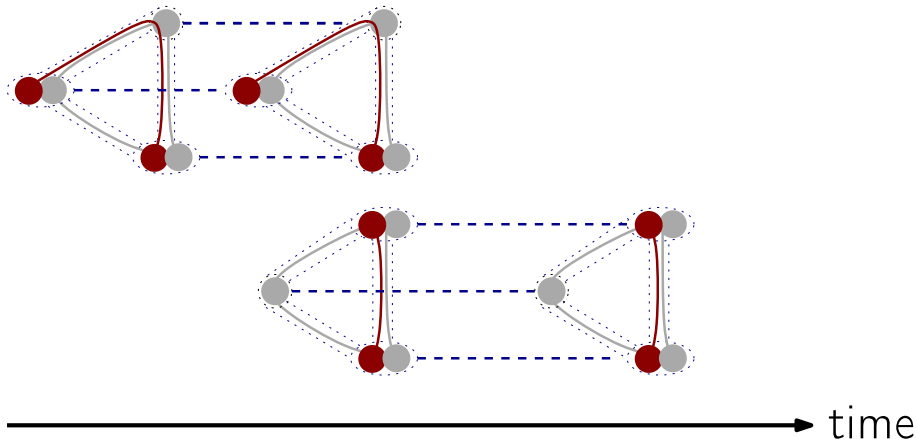


# Modeling Opportunities: Migrations

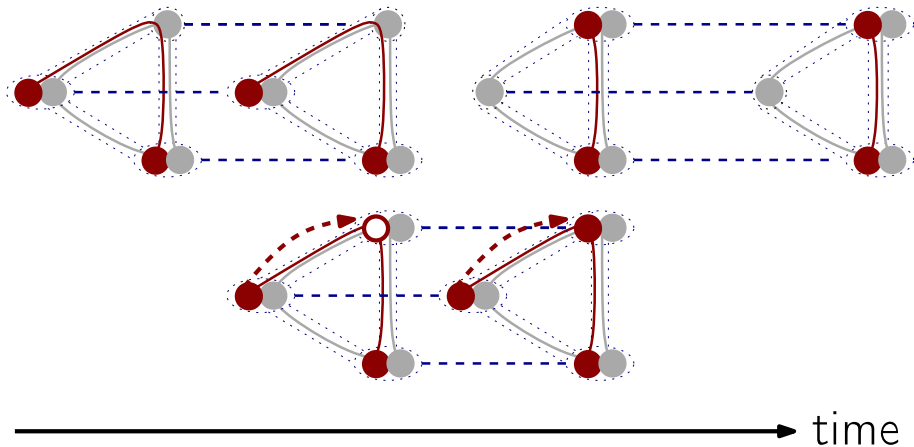


In previous work instantaneous operation!

# Modeling Opportunities: Migrations

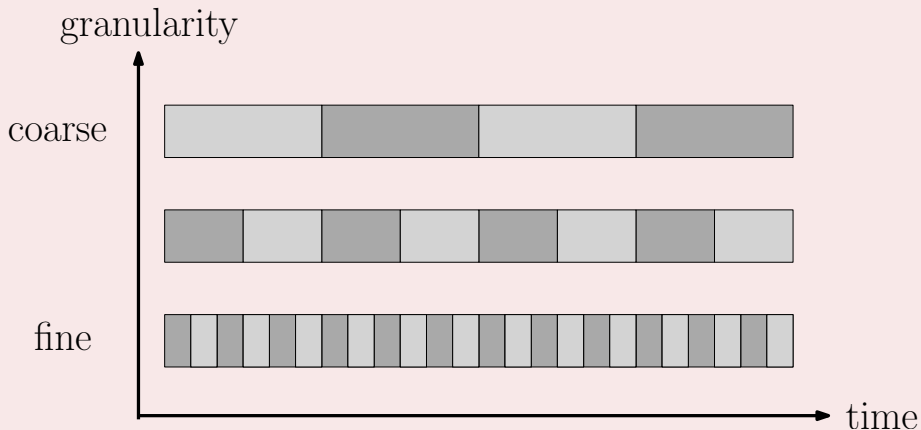


# Modeling Opportunities: Fine-grained Migrations



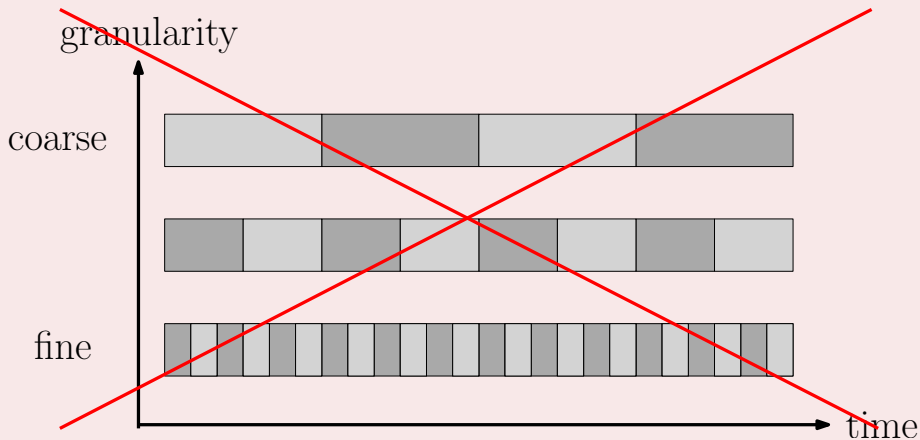
# Important Decision: Continuous-Time Model!

## Discretization



# Important Decision: Continuous-Time Model!

No Discretization!





# Problem Statement

# Temporal Virtual Network Embedding Problem (TVNEP)

- Access Control      Decide which of the requests to embed.
- Resource Mapping   Find embeddings for requests.
- Scheduling        Find appropriate start and end times for requests.
- Feasibility        For *each* point in time the substrate's capacity is not exceeded.

# Local Embedding

## Classic VNEP Task

**Access Control**      Decide which of the requests to embed.

**Resource Mapping**   Find embeddings for requests.

Mapping can be easily done with Mixed-Integer Programming.  
Not explained here.

# Overview

# Overview

## Contributions

- 1 Continuous-time Mixed-Integer Programming formulations for TVNEP
- 2  $c\Sigma$ -Model utilizes state-space and symmetry reductions to render solving TVNEP (computationally) feasible
- 3 Greedy polynomial time heuristic which is based on  $c\Sigma$ -Model
- 4 Initial computational evaluation

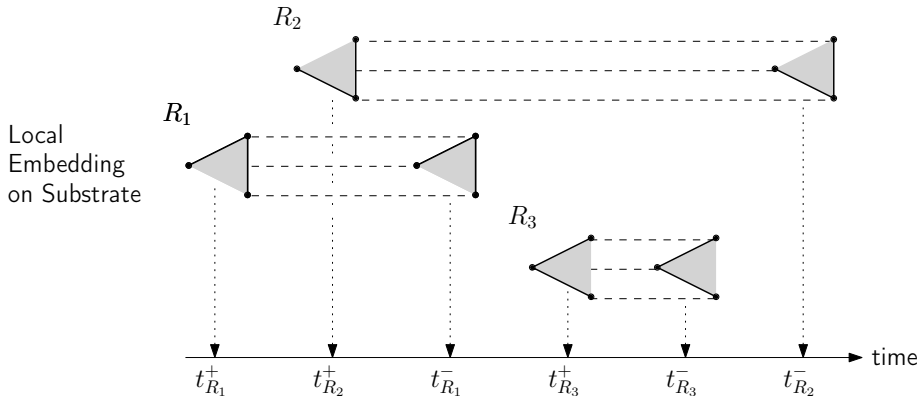
Joint work with Stefan Schmid, Anja Feldmann: IEEE IPDPS 2014.

## Modeling Continuous-Time: Checking Feasibility

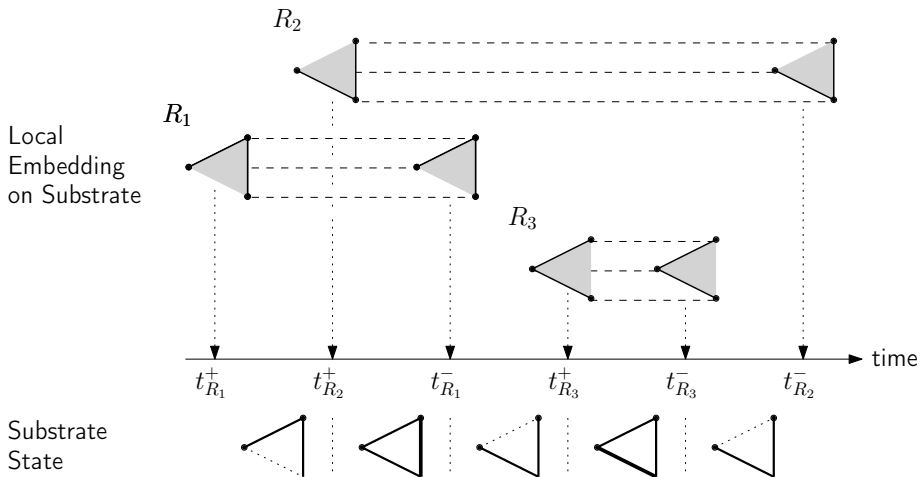
Assume for now:

Local embeddings and start / end times are fixed.

# Modeling Continuous-Time: Checking Feasibility



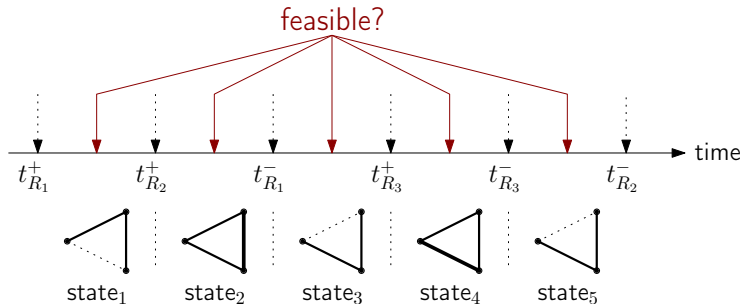
# Modeling Continuous-Time: Checking Feasibility





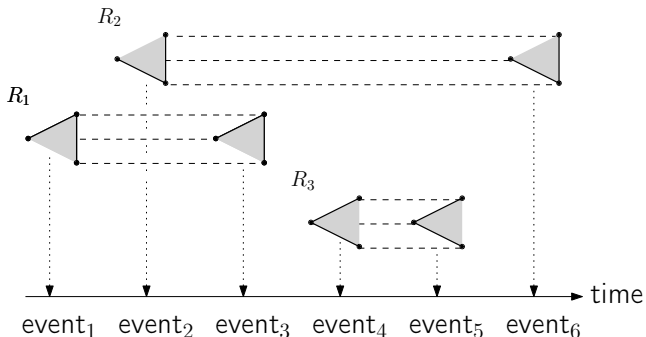
# Modeling Continuous-Time: Checking Feasibility

Check the feasibility of the  $2 \cdot |\mathcal{R}| - 1$  states.



# Abstract Event Model

# Modeling Continuous-Time: Abstract Event Model



## Mapping Variables

$$\mathcal{E} = \{\mathbf{e}_1, \dots, \mathbf{e}_{2 \cdot |\mathcal{R}|}\}$$

$$\forall R \in \mathcal{R}. \chi_R^+ : \mathcal{E} \rightarrow \mathbb{B}$$

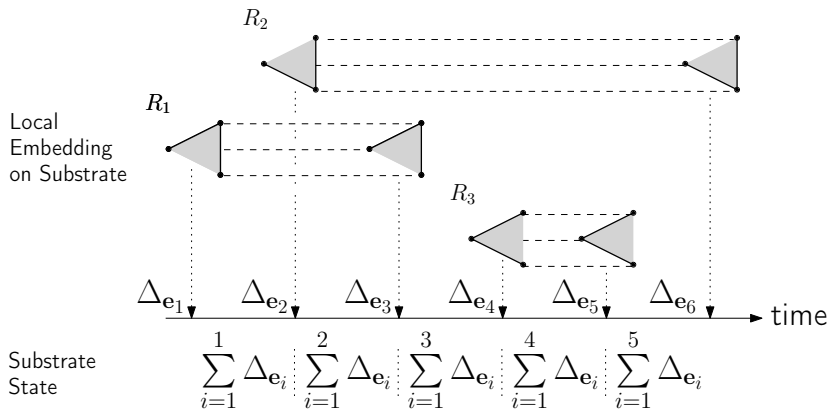
$$\forall R \in \mathcal{R}. \chi_R^- : \mathcal{E} \rightarrow \mathbb{B}$$

## Bijjective Mapping

$$\forall R \in \mathcal{R}. \sum_{\mathbf{e}_i \in \mathcal{E}} \chi_R^+(\mathbf{e}_i) = 1 \wedge \sum_{\mathbf{e}_i \in \mathcal{E}} \chi_R^-(\mathbf{e}_i) = 1$$

$$\forall \mathbf{e}_i \in \mathcal{E}. \sum_{R \in \mathcal{R}} \chi_R^+(\mathbf{e}_i) + \chi_R^-(\mathbf{e}_i) = 1$$

$\Delta$ -Model

Reconstructing States:  $\Delta$ -Model

## Idea

- 1 Compute state *changes*:  $\Delta_{e_i} : \mathbf{V}_S \cup \mathbf{E}_S \rightarrow \mathbb{R}$  via  $\chi_R^+(e_i), \chi_R^-(e_i)$
- 2 Enforce  $\sum_{j=1}^i \Delta_{e_j} \leq \mathbf{c}_S$  for each state

$\Delta$ -Model: Computing State Changes

## Conditional Assignment

 $\forall \mathbf{e}_i \in \mathcal{E}$ .

$$\Delta_{\mathbf{e}_i} = \begin{cases} +alloc(R_1) & , \text{ if } \chi_{R_1}^+(\mathbf{e}_i) = 1 \\ -alloc(R_1) & , \text{ if } \chi_{R_1}^-(\mathbf{e}_i) = 1 \\ \vdots & \\ +alloc(R_k) & , \text{ if } \chi_{R_k}^+(\mathbf{e}_i) = 1 \\ -alloc(R_k) & , \text{ if } \chi_{R_k}^-(\mathbf{e}_i) = 1 \end{cases}$$

## $\Delta$ -Model: Computing State Changes

### Conditional Assignment via Big-M Constraints

$\forall R \in \mathcal{R}. \forall \mathbf{e}_i \in \mathcal{E}.$

$$\Delta_{\mathbf{e}_i} \leq + \text{alloc}(R) + \mathbf{cs}(1 - \chi_R^+(\mathbf{e}_i))$$

$$\Delta_{\mathbf{e}_i} \geq + \text{alloc}(R) - \mathbf{cs}(1 - \chi_R^+(\mathbf{e}_i)) \cdot 2$$

$$\Delta_{\mathbf{e}_i} \leq - \text{alloc}_V(R) + \mathbf{cs}(1 - \chi_R^-(\mathbf{e}_i)) \cdot 2$$

$$\Delta_{\mathbf{e}_i} \geq - \text{alloc}_V(R) - \mathbf{cs}(1 - \chi_R^-(\mathbf{e}_i))$$

$\Delta$ -Model: Computing State Changes

## Big-M Assignment of Start

 $\forall R \in \mathcal{R}. \forall e_i \in \mathcal{E}.$ 

$$\Delta_{e_i} \leq + \text{alloc}(R) + c_S(1 - \chi_{R_1}^+(e_i))$$

$$\Delta_{e_i} \geq + \text{alloc}(R) - c_S(1 - \chi_{R_1}^+(e_i)) \cdot 2$$

$$\chi_R^+(e_i) = 0$$

$$\Delta_{e_i} \leq + \text{alloc}(R) + c_S$$

$$\Delta_{e_i} \geq + \text{alloc}(R) - 2 \cdot c_S$$



unbounded

$$\Delta_{e_i}(N_S) \leq c_S$$

$$\Delta_{e_i}(N_S) \geq -c_S$$



$\Delta$ -Model: Computing State Changes

## Big-M Assignment of Start

$$\forall R \in \mathcal{R}. \forall \mathbf{e}_i \in \mathcal{E}.$$

$$\Delta_{\mathbf{e}_i} \leq + \text{alloc}(R) + \mathbf{c}_S(1 - \chi_{R_1}^+(\mathbf{e}_i))$$

$$\Delta_{\mathbf{e}_i} \geq + \text{alloc}(R) - \mathbf{c}_S(1 - \chi_{R_1}^+(\mathbf{e}_i)) \cdot 2$$

$$\chi_{R_1}^+(\mathbf{e}_i) = 1$$

$$\Delta_{\mathbf{e}_i} \leq + \text{alloc}(R)$$

$$\Delta_{\mathbf{e}_i} \geq + \text{alloc}(R)$$

$$\Rightarrow$$

equal

$$\Delta_{\mathbf{e}_i} = \text{alloc}(R)$$

## $\Delta$ -Model Issue: LP Smearings!

### LP Relaxation Example

$\chi_{R_j}^+(\mathbf{e}_j) = 0.5$  for  $j \in \{1, 2\}$ :

$$-c_s + \text{alloc}(R_j) \leq \Delta_{\mathbf{e}_j} \leq \text{alloc}(R_j) + 0.5 \cdot c_s$$

## $\Delta$ -Model Issue: LP Smearings!

### LP Relaxation Example

$\chi_{R_j}^+(\mathbf{e}_j) = 0.5$  for  $j \in \{1, 2\}$ :

$$-c_s + \text{alloc}(R_j) \leq \Delta_{\mathbf{e}_j} \leq \text{alloc}(R_j) + 0.5 \cdot c_s$$

### Implications

- 1  $\Delta_{\mathbf{e}_j} \leq 0$  is always feasible (when  $\chi_{R_j}^+(\mathbf{e}_j) = 0.5$  for  $j \in \{1, 2\}$ )
- 2  $\Delta_{\mathbf{e}_j} < 0$  possible if  $\text{alloc}(R_j) < c_s$

## $\Delta$ -Model Issue: LP Smearings!

### LP Relaxation Example

$\chi_{R_j}^+(\mathbf{e}_j) = 0.5$  for  $j \in \{1, 2\}$ :

$$-c_s + \text{alloc}(R_j) \leq \Delta_{\mathbf{e}_j} \leq \text{alloc}(R_j) + 0.5 \cdot c_s$$

### Implications

- 1  $\Delta_{\mathbf{e}_j} \leq 0$  is always feasible (when  $\chi_{R_j}^+(\mathbf{e}_j) = 0.5$  for  $j \in \{1, 2\}$ )
- 2  $\Delta_{\mathbf{e}_j} < 0$  possible if  $\text{alloc}(R_j) < c_s$

### This is really bad!

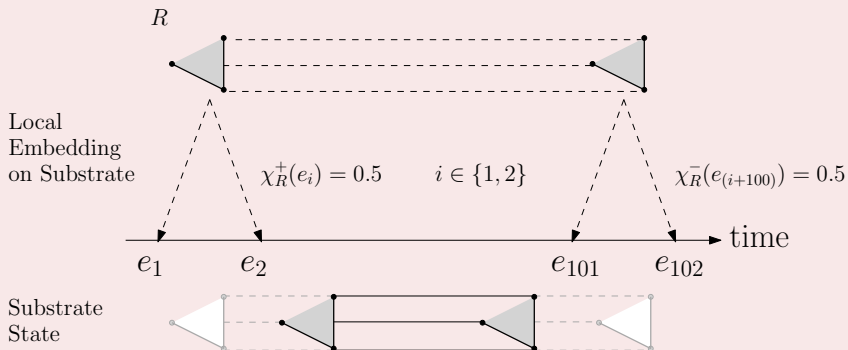
- 1 states do not 'materialize' well in LP relaxations:  
allocations will *never* be accounted for in the substrate's state
- 2 bounding is unable to reduce search space

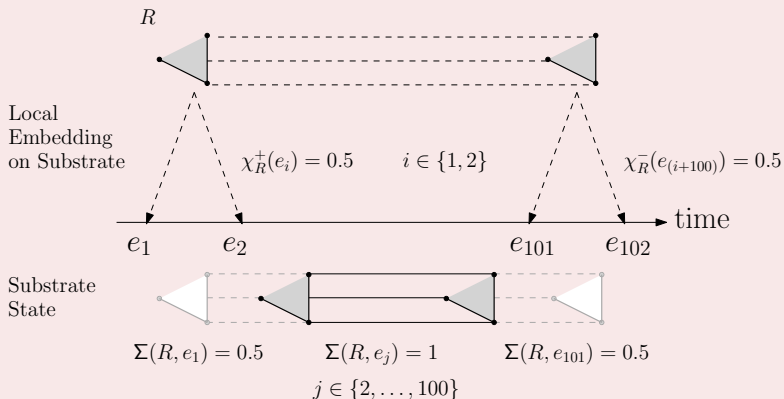
$\Sigma$ -Model

$\Sigma$ -Model: Intuition

## Requirement

Resource allocations must materialize in the substrate's state.

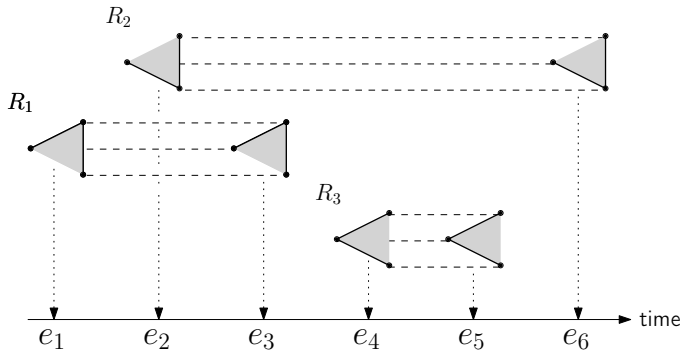


$\Sigma$ -Model: Intuition

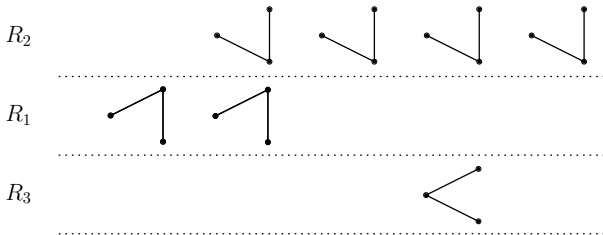
$$\forall R \in \mathcal{R}. \forall e_j \in \mathcal{E}.$$

$$\Sigma(R, e_j) = \sum_{j=1, \dots, i} \chi_R^+(e_j, R) - \sum_{j=1, \dots, i} \chi_R^-(e_j, R)$$

Local  
Embedding  
on Substrate



Allocations  
for each  
state and  
request



Substrate  
State





$c\Sigma$ -Model

# c $\Sigma$ -Model: Outline

## Computational Trade-Off

- The  $\Sigma$ -Model is provably stronger than the  $\Delta$ -Model.
- However, the  $\Sigma$ -Model uses (approximately)  $2 \cdot |\mathcal{R}|$  more variables!

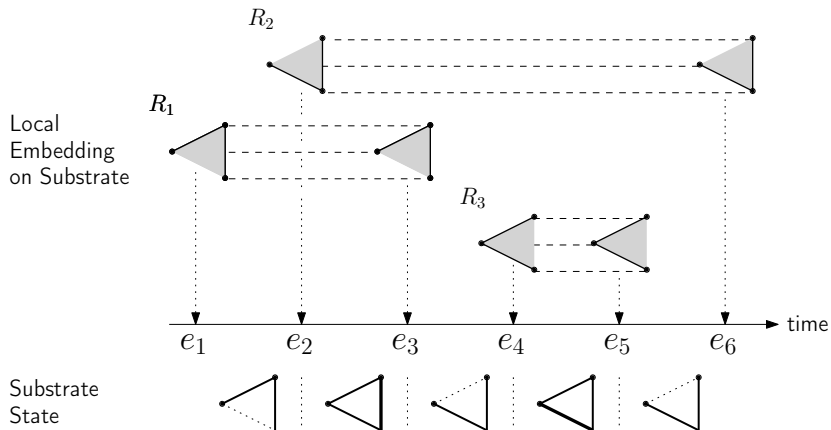
## $\Sigma$ -Model can be strengthened: c $\Sigma$ -Model

**Compactification** Consider only *partial* event order. Yields *state-space* and *symmetry reductions*.

**User cuts** Use temporal information to reduce *state-space* and strengthen formulation.

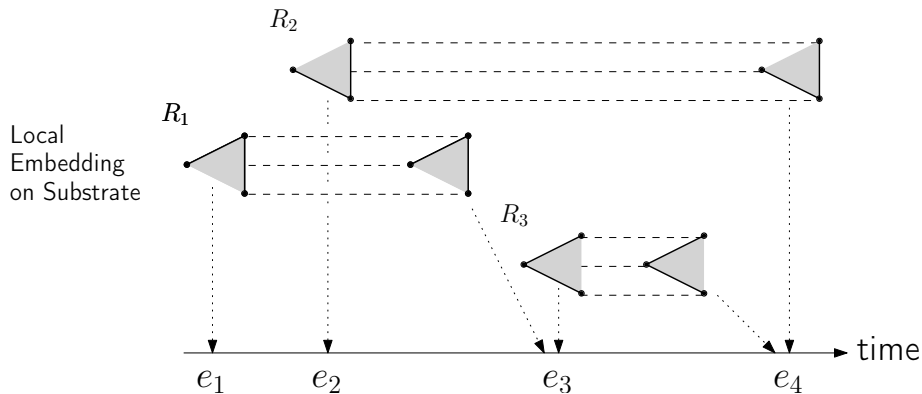
## c $\Sigma$ Optimization: State Compactification

# $c\Sigma$ -Model: State Compactification



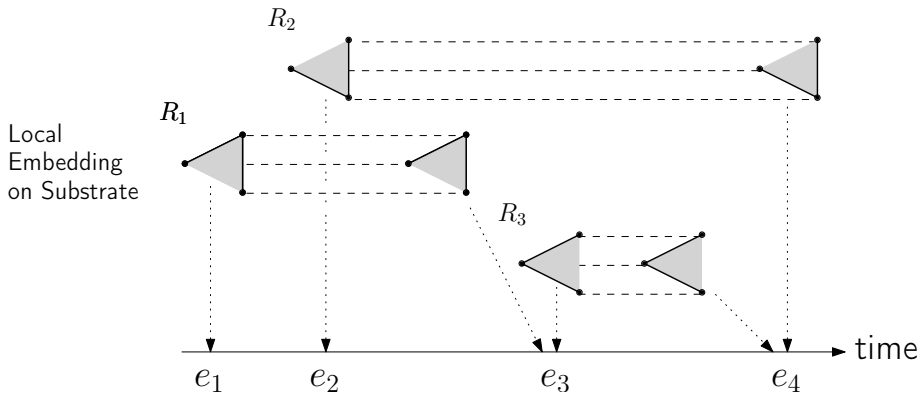
We only need to check feasibility after a request's start!

# c $\Sigma$ -Model: State Compactification



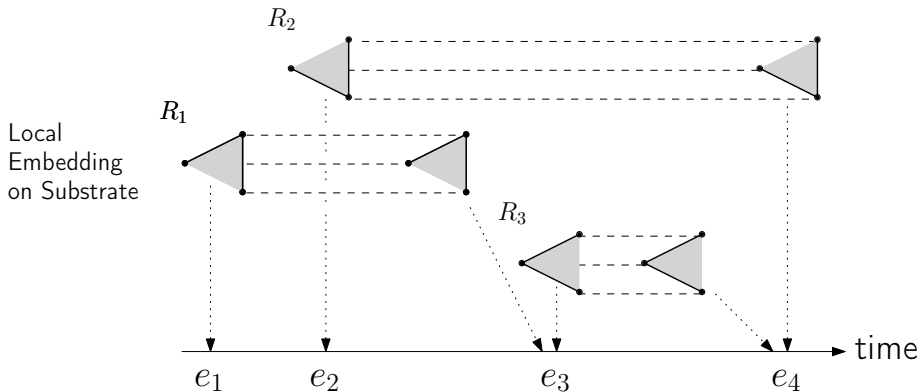
- consider only  $|\mathcal{R}| + 1$  event points
- injective mapping of request starts onto first  $|\mathcal{R}|$  event points
- mapping of request  $R$ 's end onto event  $e_j$ :  
 $R$  ends after  $e_{j-1}$  and before  $e_j$

# c $\Sigma$ -Model: State Compactification

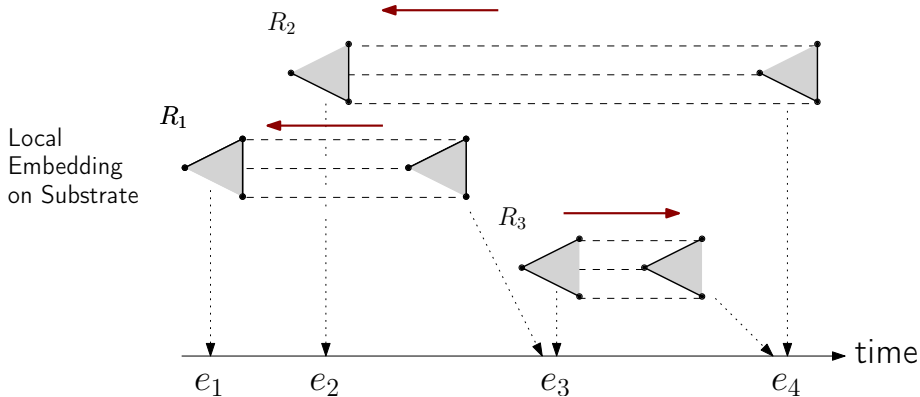


## State-space reduction

Number of states is halved  $\Rightarrow$  number of variables is halved.

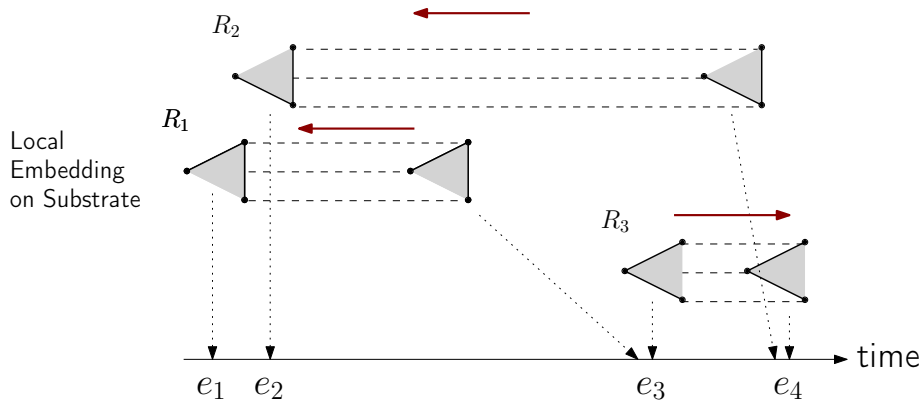
c $\Sigma$ -Model: State Compactification is Symmetry Reduction

# c $\Sigma$ -Model: State Compactification is Symmetry Reduction



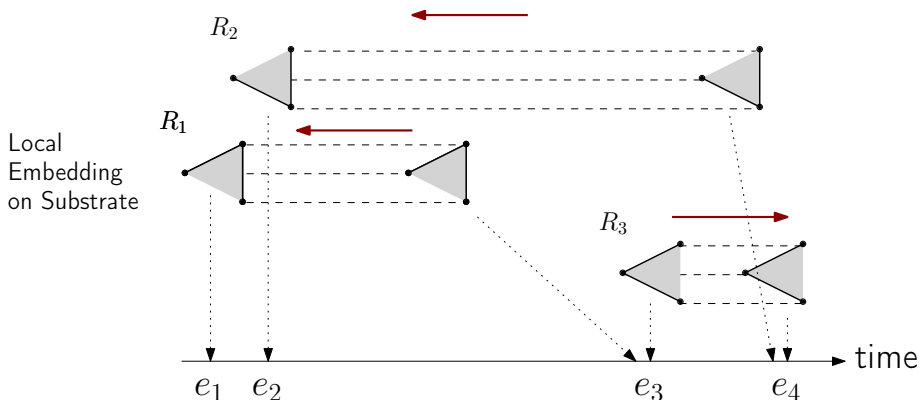


# $c\Sigma$ -Model: State Compactification is Symmetry Reduction



Same order as before!

# $c\Sigma$ -Model: State Compactification is Symmetry Reduction



## Theorem

Compactification is symmetry reduction.

Greedy Heuristic  $c \sum_A^G$

# Greedy Heuristic $c\Sigma_A^G$

## Outline

- 1 Order requests according to their earliest start time.
- 2 Iteratively try to embed requests as soon as possible using  $c\Sigma$ -Model
  - 1 If the request was embedded: fix start and end time.

# Greedy Heuristic $c\Sigma_A^G$

## Outline

- 1 Order requests according to their earliest start time.
- 2 Iteratively try to embed requests as soon as possible using  $c\Sigma$ -Model
  - 1 If the request was embedded: fix start and end time.

Theorem:  $c\Sigma_A^G$  is polynomial-time algorithm

There are maximally  $|\mathcal{R}|$  many possible orderings to consider.

### Important

All link allocations are re-computed in each iteration.

# Computational Evaluation

## Scenario: One day workload

- 20 requests (star-graphs) are to be embedded on  $4 \times 5$  grid
- Expected inter-arrival time of one hour [Poisson]
- Expected duration of 3.5 hours [Weibull: heavy-tailed]
- Node-mappings are fixed to concentrate on temporal aspects
- Link-mappings are not fixed
- Increasing temporal flexibility: 0, 30, 60,  $\dots$ , 300 minutes.

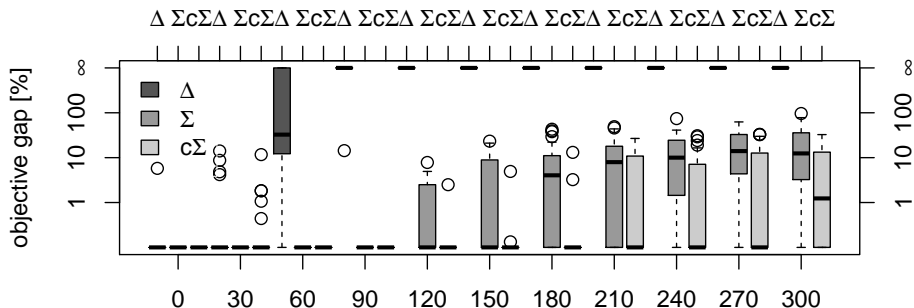
## Computational Setup

- 24 independently generated scenarios
- Limited runtime of one hour for MIPs [Gurobi]

## Task: Maximize revenue $\propto$ load $\cdot$ duration

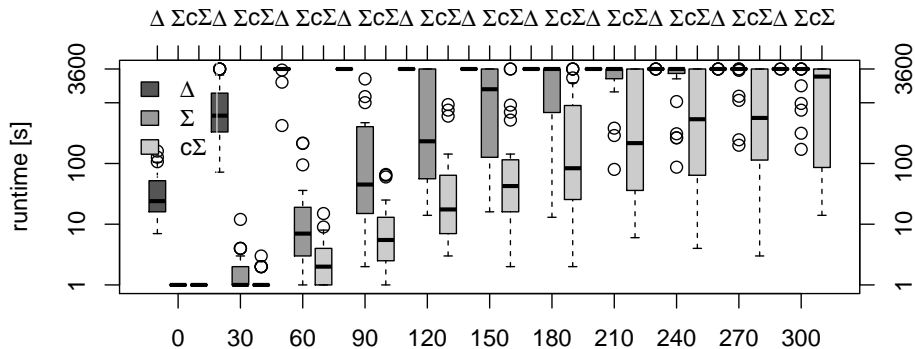
- 1 Decide which requests to embed (access control).
- 2 Find time-invariant embedding (routing of data).
- 3 Decide when to embed the requests.

## Objective Gap: MIP Formulations

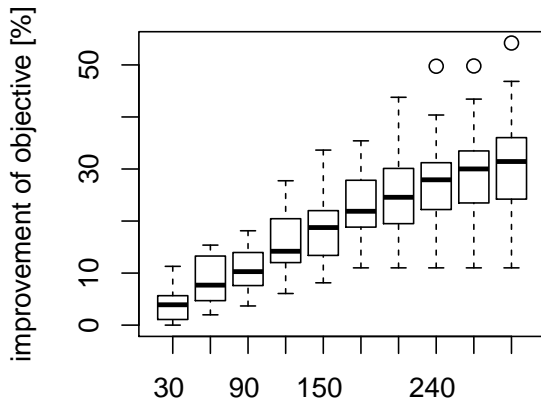


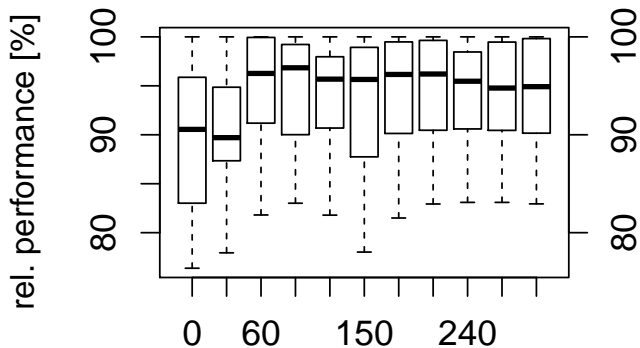


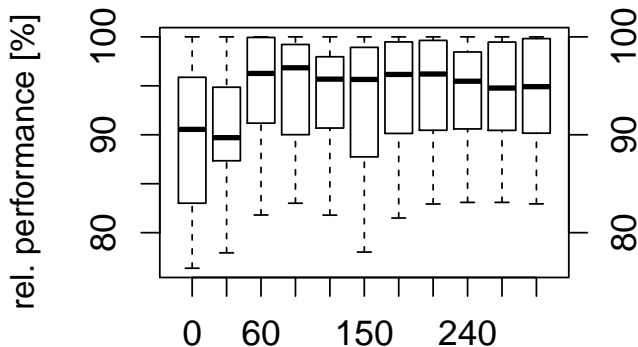
## Runtime: MIP Formulations



# Benefit of Flexibility



Performance of  $c\Sigma_A^G$ 

Performance of  $c\Sigma_A^G$ 

Fast: runtime of few seconds.

Conclusion: Modeling Opportunities / Discussion

# A possible application: Google B4

## Google's Software Defined WAN

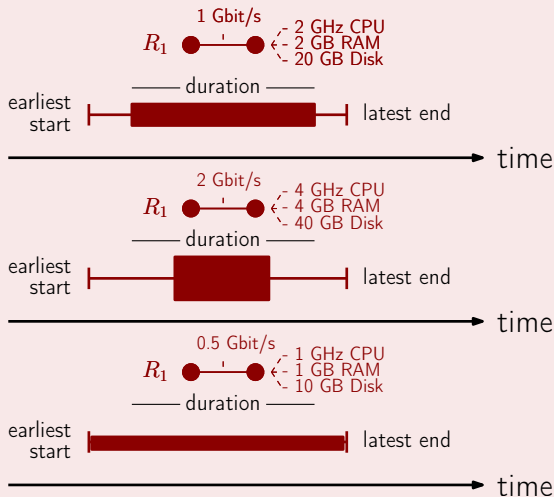
---



Goal: Utilize close to 100% of available bandwidth.

# A possible application: Google B4

## Interesting Extension: Delay-tolerant requests



# Important Notice

## We can schedule \*anything

- The way embeddings are computed can be changed arbitrarily.
- We can use the presented approach to e.g. schedule in-network processing service deployments!



# Takeaway message part I

- 1 The 'way' a MIP is formulated can have a significant impact on the solver's ability to *efficiently* solve it.
- 2 Our approach proposes a very general framework for offline scheduling of networking tasks.

# References I

- [1] A. Fischer, J. Botero, M. Beck, H. De Meer, and X. Hesselbach.  
Virtual network embedding: A survey.  
2013.
- [2] C. A. Floudas and X. Lin.  
Continuous-time versus discrete-time approaches for scheduling of chemical processes: a review.  
*Computers & Chemical Engineering*, 28(11), 2004.
- [3] T. A. Henzinger, A. V. Singh, V. Singh, T. Wies, and D. Zufferey.  
A marketplace for cloud resources.  
In *Proc. of the ACM EMSOFT '10*, 2010.
- [4] L. Mai, E. Kalyvianaki, and P. Costa.  
Exploiting time-malleability in cloud-based batch processing systems.  
In *Proceeding of the ACM SIGOPS LADIS Workshop '13*, 2013.
- [5] S. Jain et al.  
B4: experience with a globally-deployed software defined wan.  
In *Proc. of the ACM SIGCOMM '13*, 2013.

## References II

- [6] A. Singla, A. Singh, K. Ramachandran, L. Xu, and Y. Zhang.  
Proteus: A topology malleable data center network.  
In *Proc. of the ACM SIGCOMM HotNets Workshop '10*, 2010.

# The End

- 1 Abstract event point model
- 2  $\Delta$ -,  $\Sigma$ - and  $c\Sigma$ -Model
  - state-space reductions
  - symmetry reduction
- 3 Greedy heuristic  $c\Sigma_A^G$  based on  $c\Sigma$
- 4 Initial computational evaluation
  - $\Delta \ll \Sigma < c\Sigma$
  - $c\Sigma$ : near optimal solutions within one hour
  - $c\Sigma_A^G$  only approx. 5-10% off optimum

# Future Work / Discussion

## Modeling

- Consider flexible duration of requests.
- Consider delay-tolerant VNets.
- Consider more complex scenarios, e.g. migrations.

## Algorithmic

- Incorporate other heuristical embedding approaches.
- Develop local-search algorithms for the TVNEP.

## Related Work

### Chemical plants [2]

Utilize similar event abstraction, but no resource sharing.

### Business Perspective [3]

Marketplace based on temporal flexibilities.

### MapReduce [4]

Consider temporally predictable jobs (MapReduce-like) and allow for temporally interleaved resource sharing.

### VNet Survey [1]

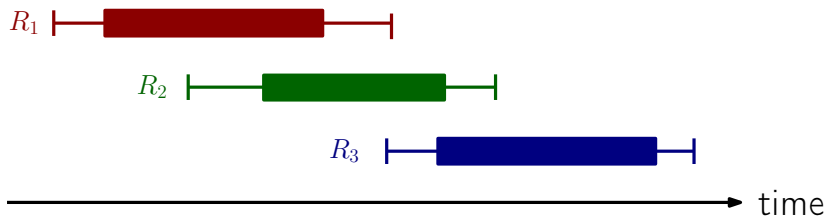
There is no comparable work on TVNEP.

### Google B4 [5]

Software-defined network (wide-area) connecting data centers. Only some dozen locations.

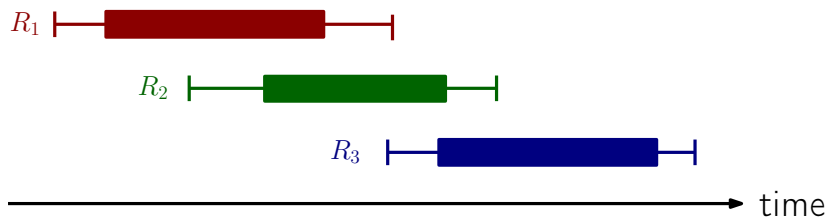
# Optimization: Temporal Dependency Graph

# Temporal Dependency Graph





# Temporal Dependency Graph



Latest possible point in time for  $R_1$  to start is less than the earliest point in time at which  $R_2$  can start.

$\Rightarrow$  We know that  $R_1$  must start before  $R_2$ .

# Temporal Dependency Graph

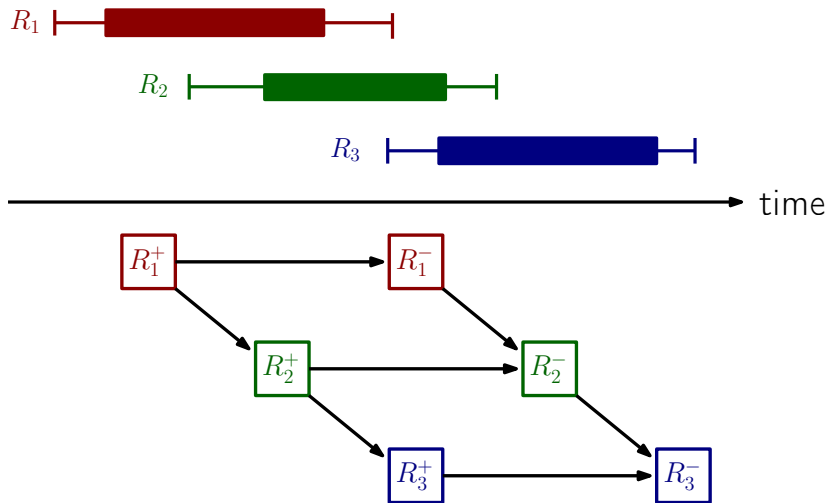


Figure: Temporal Dependency Graph

# Temporal Dependency Graph (Formal)

## Definition

- $G_{dep}(\mathcal{R}) = (V_{dep}, E_{dep})$
- $V_{dep} = \mathcal{R} \times \{start, end\}$
- $E_{dep} = \{(v, w) \in V_{dep}^2 \mid latest(v) < earliest(w)\}$

$$earliest((R, t) \in V_{dep}) = \begin{cases} t_R^s & , \text{ if } t = start \\ t_R^s + \mathbf{d}_R & , \text{ if } t = end \end{cases}$$

$$latest((R, t) \in V_{dep}) = \begin{cases} t_R^e - \mathbf{d}_R & , \text{ if } t = start \\ t_R^e & , \text{ if } t = end \end{cases}$$

# Weighted Temporal Dependency Graph

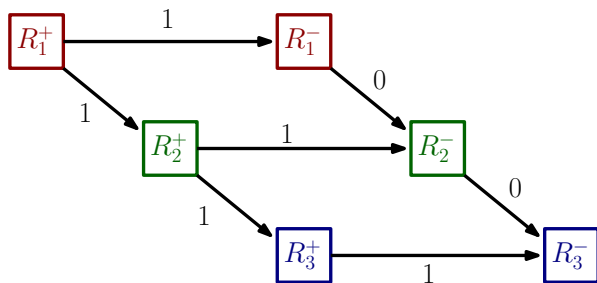


Figure: Temporal Dependency Graph with weights

By computing maximal distances (in polynomial time) we obtain:

- Start of  $R_1$ :  $e_1$
- Start of  $R_2$ :  $e_2$
- Start of  $R_3$ :  $e_3$
- End of  $R_1$ :  $e_2, e_3, e_4$
- End of  $R_2$ :  $e_3, e_4$
- End of  $R_3$ :  $e_4$

## Overview $c\Sigma$ -Model

## Access Control & Resource Mapping

### Variables

Access Control  $x_{\mathcal{R}} : \mathcal{R} \rightarrow \mathbb{B}$

Node Mapping  $\forall R \in \mathcal{R}. x_V : \mathbf{V}_R \times \mathbf{V}_S \rightarrow \mathbb{B}$

Link Mapping  $\forall R \in \mathcal{R}. x_E : \mathbf{E}_R \times \mathbf{E}_S \rightarrow [0, 1]$

Node mapping:  $\forall R \in \mathcal{R}. \forall N_V \in \mathbf{V}_R.$

$$x_{\mathcal{R}}(R) = \sum_{N_S \in \mathbf{V}_S} x_V(N_V, N_S)$$

Link mapping:  $\forall R \in \mathcal{R}. \forall L_V = (N_V^+, N_V^-) \in \mathbf{E}_R. \forall N_S \in \mathbf{V}_S$

$$\sum_{L_S \in \delta^+(N_S)} x_E(L_V, L_S) - \sum_{L_S \in \delta^-(N_S)} x_E(L_V, L_S) = x_V(N_V^-, N_S) - x_V(N_V^+, N_S)$$

Macro  $alloc_V(R, N_S): \forall R \in \mathcal{R}. \forall N_S \in \mathbf{V}_S$

$$alloc_V(R, N_S) = \sum_{N_V \in \mathbf{V}_R} c_{\mathcal{R}}(N_V) \cdot x_V(N_V, N_S)$$

Macro  $alloc_E(R, L_S): \forall R \in \mathcal{R}. \forall L_S \in \mathbf{E}_S$

$$alloc_E(R, L_S) = \sum_{L_V \in \mathbf{E}_R} c_{\mathcal{R}}(L_V) \cdot x_E(L_V, L_S)$$

## Access Control & Resource Mapping

### Mapping onto Event Points

#### Variables

- $\forall R \in \mathcal{R}. \chi_R^+ : \mathcal{E} \rightarrow \mathbb{B}$
- $\forall R \in \mathcal{R}. \chi_R^- : \mathcal{E} \rightarrow \mathbb{B}$

Mapping each start / end:  $\forall R \in \mathcal{R}.$

$$\sum_{\mathbf{e}_i \in \{\mathbf{e}_1, \dots, \mathbf{e}_{|\mathcal{R}|}\}} \chi_R^+(\mathbf{e}_i) = 1 \qquad \sum_{\mathbf{e}_i \in \{\mathbf{e}_2, \dots, \mathbf{e}_{|\mathcal{R}|+1}\}} \chi_R^-(\mathbf{e}_i) = 1$$

Mapping start injectively:  $\forall \mathbf{e}_i \in \{\mathbf{e}_1, \dots, \mathbf{e}_{|\mathcal{R}|}\}.$

$$\sum_{R \in \mathcal{R}} (\chi_R^+(\mathbf{e}_i)) = 1$$

# Access Control & Resource Mapping

## Mapping onto Event Points

### Guaranteeing State Feasibility

#### Variables

$$alloc_V : \mathcal{R} \times \mathcal{S} \times \mathbf{V}_S \rightarrow \mathbb{R}_{\geq 0} \quad alloc_E : \mathcal{R} \times \mathcal{S} \times \mathbf{E}_S \rightarrow \mathbb{R}_{\geq 0}$$

Computing allocations at states:  $\forall R \in \mathcal{R}. \forall s_i \in \mathcal{S}. \forall N_s \in \mathbf{V}_S / \forall L_s \in \mathbf{E}_S$ .

- $alloc_V(R, s_i, N_s) \geq alloc_V(R, N_s) - c_S(N_s) \cdot (1 - \Sigma(R, e_i))$
- $alloc_E(R, s_i, L_s) \geq alloc_E(R, L_s) - c_S(L_s) \cdot (1 - \Sigma(R, e_i))$

Ensuring feasibility:  $\forall s_i \in \mathcal{S}. \forall N_s \in \mathbf{V}_S / L_s \in \mathbf{E}_S$ .

- $c_S(N_s) \geq \sum_{R \in \mathcal{R}} alloc_V(R, s_i, N_s)$
- $c_S(L_s) \geq \sum_{R \in \mathcal{R}} alloc_E(R, s_i, L_s)$



# Access Control & Resource Mapping

## Mapping onto Event Points

## Guaranteeing State Feasibility

## Guaranteeing Temporal Feasibility

### Variables

$$\forall R \in \mathcal{R}. t_R^+, t_R^- \in \mathbb{R}_{\geq 0} \quad \forall e_i \in \mathcal{E}. t_{e_i} \in \mathbb{R}_{\geq 0}$$

$$\forall R \in \mathcal{R}.$$

$$d_R = t_R^- - t_R^+$$

Setting start times:  $\forall R \in \mathcal{R}. \forall e_i \in \{\mathbf{e}_1, \dots, \mathbf{e}_{|\mathcal{R}|}\}.$

$$t_R^+ \leq t_{e_i} + (1 - \sum_{j=1, \dots, i} \chi_R^+(e_j, R)) \cdot T \quad t_R^+ \geq t_{e_i} - (1 - \sum_{j=i, \dots, |\mathcal{E}|} \chi_R^+(e_j, R)) \cdot T$$

Setting end times:  $\forall R \in \mathcal{R}. \forall e_i \in \{\mathbf{e}_2, \dots, \mathbf{e}_{|\mathcal{R}|+1}\}.$

$$t_R^- \leq t_{e_i} + (1 - \sum_{j=2, \dots, i} \chi_R^-(e_j, R)) \cdot T \quad t_R^- \geq t_{e_{i-1}} - (1 - \sum_{j=i, \dots, |\mathcal{E}|} \chi_R^-(e_j, R)) \cdot T$$

Access Control & Resource Mapping

Mapping onto Event Points

Guaranteeing State Feasibility

Guaranteeing Temporal Feasibility

Temporal Dependency Graph User Cuts

$\forall v \in V_{dep}$ .

$$\sum_{i=|dist_{\max}^+(v)|+1}^{|\mathcal{R}|+1-|dist_{\max}^-(v)|} \chi_{Event}(\mathbf{e}_i, v) = 1$$

$\forall v \in V_{dep} \cdot \forall w \in dist_{\max}^-(v) \cdot \forall \mathbf{e}_i \in \mathcal{E}, dist_{\max}(v, w) + 1 \leq i \leq |\mathcal{R}|$ .

$$\sum_{j=1}^i \chi_{Event}(\mathbf{e}_j, w) \leq \sum_{\mathbf{e}_j \in \mathcal{E}} \chi_{Event}(\mathbf{e}_j, v)$$

with  $j \leq i - dist_{\max}^-(v, w)$

Access Control & Resource Mapping

Mapping onto Event Points

Guaranteeing State Feasibility

Guaranteeing Temporal Feasibility

Temporal Dependency Graph User Cuts

### Some further optimizations

- Big-M constants are chosen as *tight* as possible
- virtual links can be aggregated if their virtual source or their virtual destination is the same

# Applications

## Data center

- e.g. MapReduce cycles through different phases, traffic only during 30-60% of execution [6]
- price incentives for customers and providers to allow for / harness temporal flexibility [4]

## Wide area networks

- Google uses SDN in the WAN to connect data centers [5]
- scheduling of bandwidth-intensive synchronizations
  - is necessary to achieve good utilization and resource isolation
  - is enabled by central SDN control

## Part II: Virtualized In-Network Processing

# Mindset

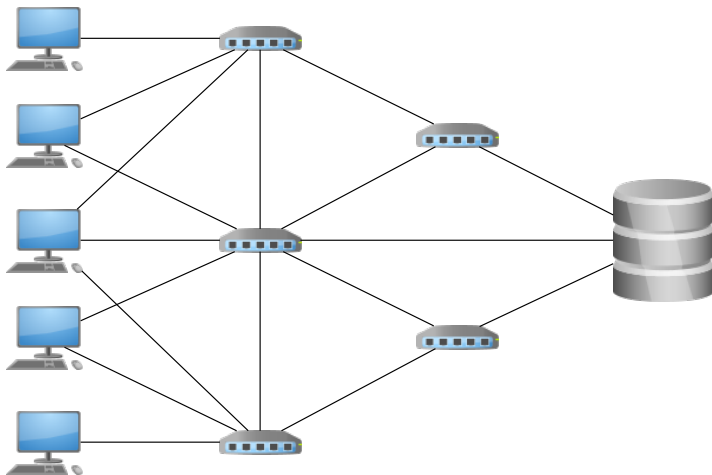
## Service Deployment $\neq$ VNet Embedding

- Customer requests communication *service* between locations.
- Service provider *finds* an appropriate topology.

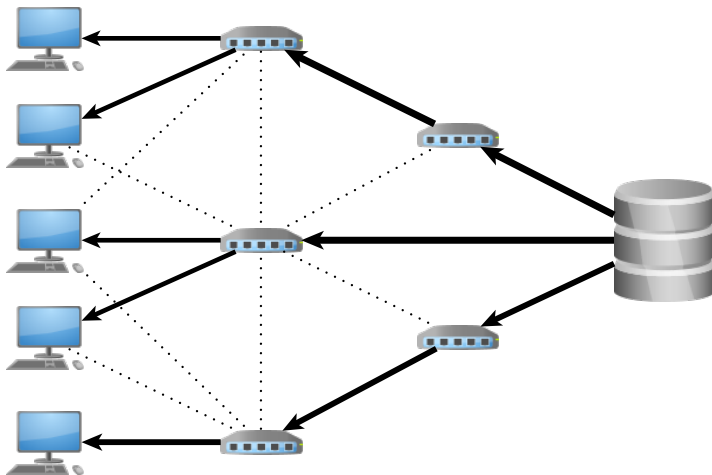
## Publications

- Joint work with Stefan Schmid: OPODIS 2013, arXiv 2013
- My thesis January 2014

# Service Replication



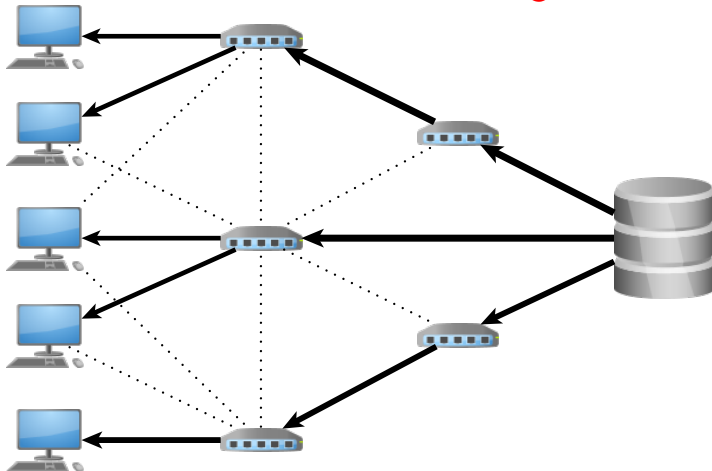
# Service Replication



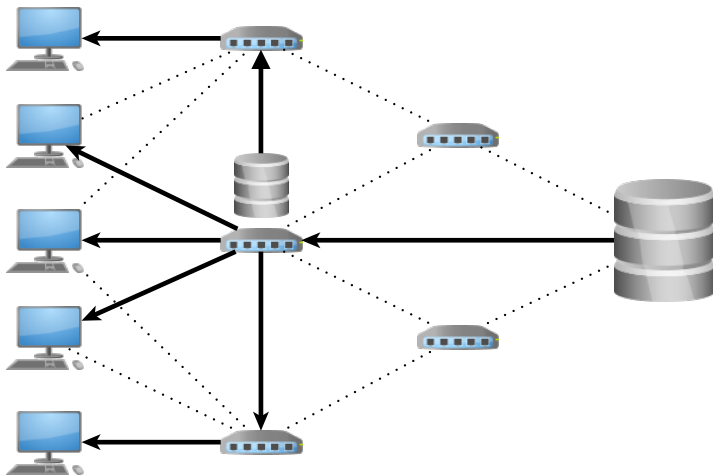


# Service Replication

What if backend links are congested?

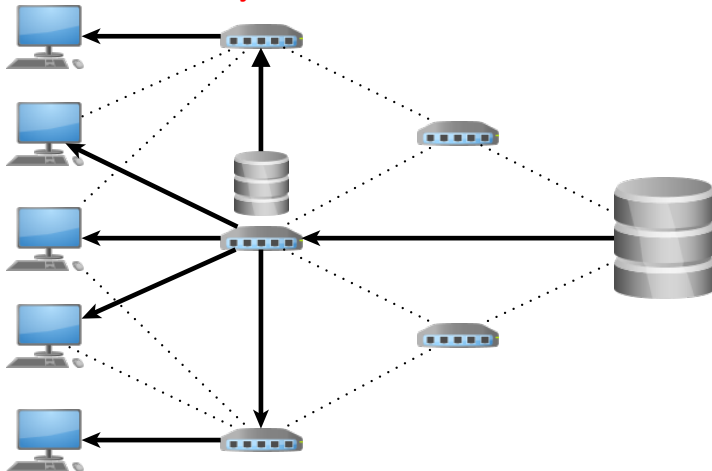


# Service Replication

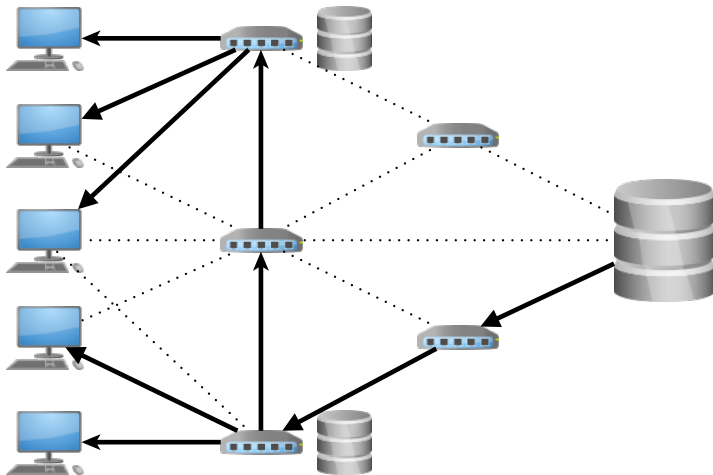


# Service Replication

What if only '3' users can be handled?



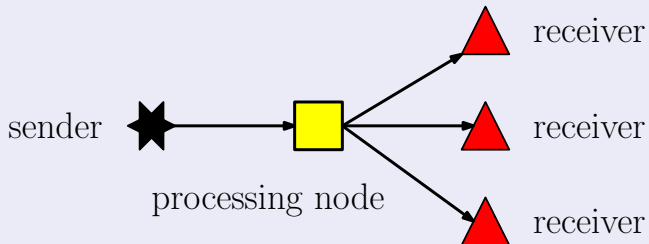
# Service Replication



# Generalized Communication Services

# Communication Services: Multicast

processing = duplication + reroute



# Communication Services: Multicast

processing = duplication + reroute

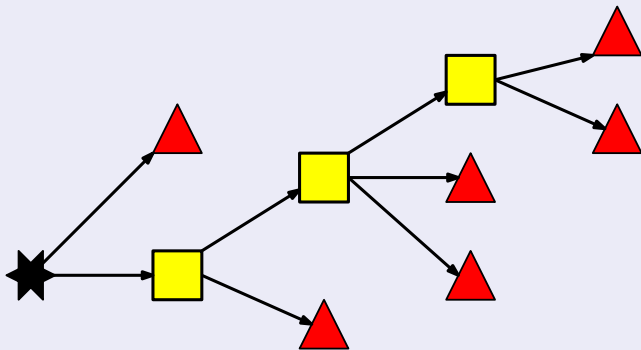
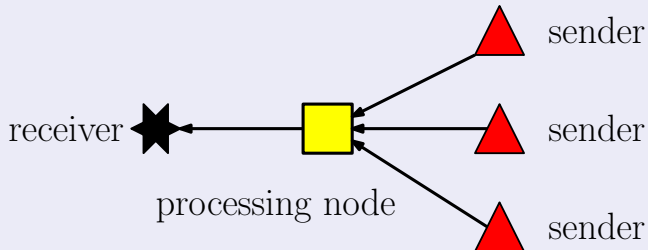


Figure: Hierarchy of processing nodes

# Communication Services: Aggregation

processing = merge + reroute





# Communication Services: Aggregation

processing = merge + reroute

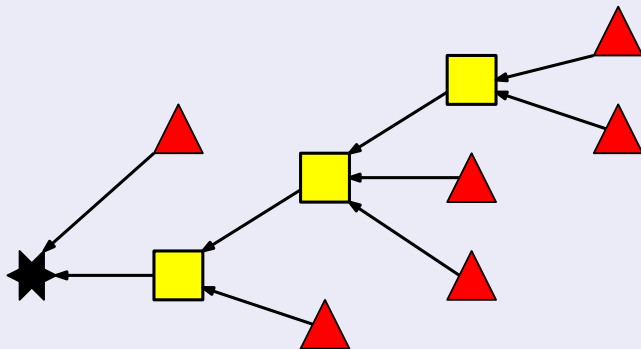


Figure: Hierarchy of processing nodes

# Problem Statement

## Questions

- How to compute *virtual* aggregation / multicasting trees?
- Where to place in-network processing functionality?
- How to trade-off between traffic and processing?

# Introductory Example

# Introductory Example

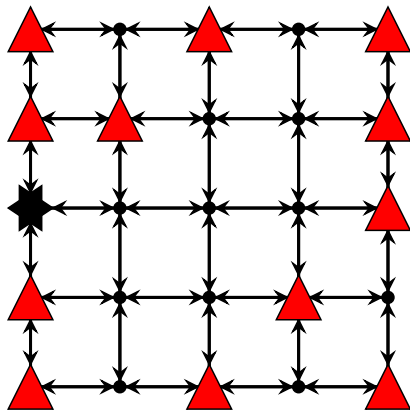
Aggregation scenario

grid graph: 14 senders, one receiver

Virtualized links

data can be routed arbitrarily

★ receiver    ▲ sender



# Without in-network processing: Unicast

## Solution Method

- minimal cost flow

## Solution uses

- 41 edges
- 0 processing nodes



receiver



sender

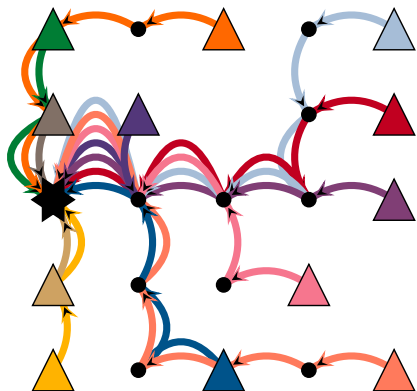


Figure: Unicast solution

# With in-network processing at all nodes

## Solution Method

- Steiner arborescence

## Solution uses

- 16 edges
- 9 processing nodes

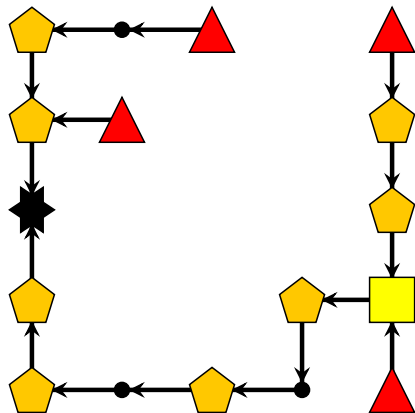
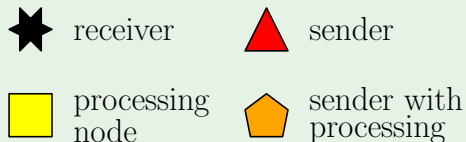
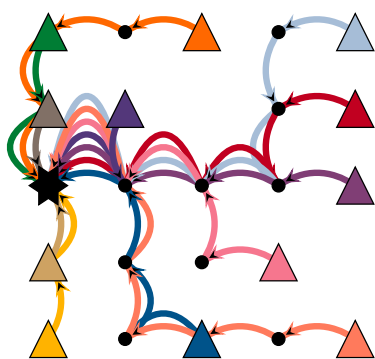
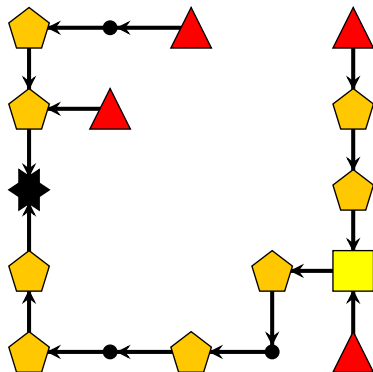


Figure: Aggregation tree

# How to Trade-off?



vs.



# What we aim for

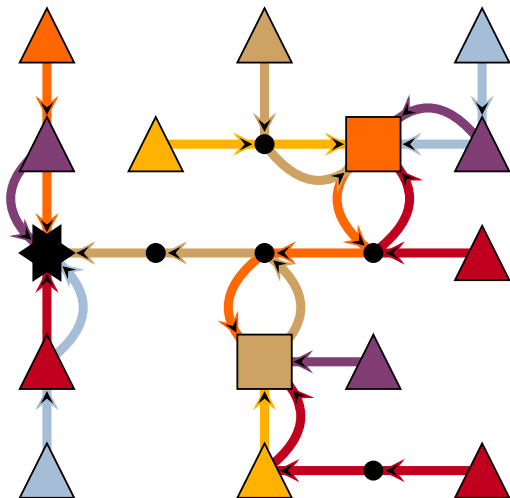
## Solution uses

- 26 edges
- 2 processing nodes

★ receiver

△ sender

□ processing node





# Solution Structure

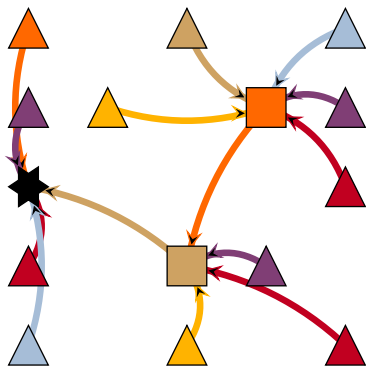


Figure: Virtual Arborescence

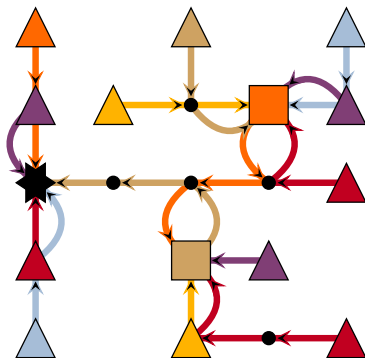


Figure: underlying routes

# New Model: Constrained Virtual Steiner Arborescence Problem

## Definition: CVSAP

Find a Virtual Arborescence connecting senders to the single receiver, s.t.

- 1 bandwidth of substrate is not exceeded,
- 2 inner nodes are capable of processing flow,
- 3 the processing nodes' capacities are not exceeded,

minimizing the joint cost for bandwidth allocations and function placement.

# Applications Overview

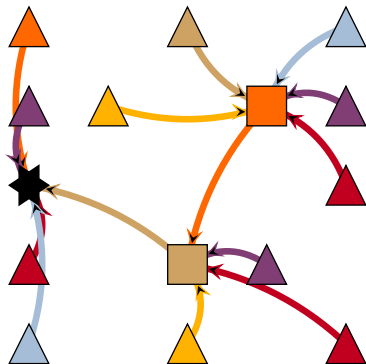
|             | Network        | Application                                  | Technology, e.g.        |
|-------------|----------------|--|-------------------------|
| multicast   | ISP            | service replication / cache placement [8, 9] | middleboxes / NFV + SDN |
|             | backbone       | optical multicast [5]                        | ROADM + SDH             |
|             | all            | application-level multicast [12]             | different               |
| aggregation | sensor network | value & message aggregation [4, 6]           | source routing          |
|             | data center    | big data / map-reduce: Camdoop [2]           | SDN                     |
|             | ISP            | network analytics: Gigascope [3]             | middleboxes / NFV + SDN |

# MIP Formulations

# Multi-Commodity Flow (MCF) Integer Program

## First approach: MCF IP

- explicitly represent virtual arborescence by the explicit construction of paths for all processing nodes
- enforces loop-freeness of connections



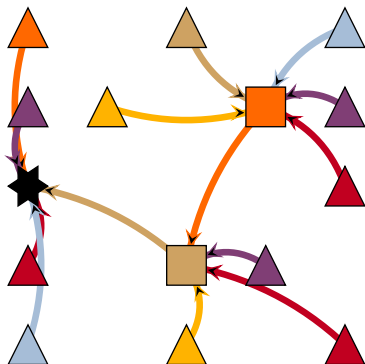
# Multi-Commodity Flow (MCF) Integer Program

## First approach: MCF IP

- explicitly represent virtual arborescence by the explicit construction of paths for all processing nodes
- enforces loop-freeness of connections

## Does not scale well

- number of binary variables:  
 $\# \text{processing nodes} \cdot \# \text{edges}$
- 'not so good relaxations' due to weak loop-freeness formulation



## Integer Program 1: A-CVSAP-MCF

$$\begin{aligned}
 \text{minimize} \quad & C_{\text{MCF}} = \sum_{e \in E_G} c_e (f_e + \sum_{s \in S} f_{s,e}) && \text{(MCF-OBJ)} \\
 & + \sum_{s \in S} c_s \cdot x_s \\
 \text{subject to} \quad & f^T(\delta_{\text{EMCF}}^+(v)) = f^T(\delta_{\text{EMCF}}^-(v)) + |\{v\} \cap T| && \forall v \in V_G \quad \text{(MCF-1)} \\
 & f^s(\delta_{\text{EMCF}}^+(v)) = f^s(\delta_{\text{EMCF}}^-(v)) + \delta_{s,v} \cdot x_s && \forall s \in S, v \in V_G \quad \text{(MCF-2)} \\
 & f_e^T + \sum_{s \in S} f_e^s \leq \begin{cases} \mathbf{u}_s x_s, & e = (s, o^-), s \in S \\ \mathbf{u}_r, & e = (r, o^-) \\ \mathbf{u}_e, & e \in E_G \end{cases} && \forall e \in E_{\text{MCF}} \quad \text{(MCF-3)} \\
 & -|S|(1 - f_{\bar{s}, o^-}^s) \leq p_s - p_{\bar{s}} - 1 && \forall s, \bar{s} \in S \quad \text{(MCF-4)} \\
 & f_{(\bar{s}, o^-)}^s \leq x_{\bar{s}} && \forall s \in S, \bar{s} \in S - s \quad \text{(MCF-5*)} \\
 & f_{s, o^-}^s = 0 && \forall s \in S \quad \text{(MCF-6*)} \\
 & f_{\bar{s}, o^-}^s + f_{s, o^-}^{\bar{s}} \leq 1 && \forall s, \bar{s} \in S \quad \text{(MCF-7*)} \\
 & x_s \in \{0, 1\} && \forall s \in S \quad \text{(MCF-8)} \\
 & f_e^T \in \mathbb{Z}_{\geq 0} && \forall e \in E_{\text{MCF}} \quad \text{(MCF-9)} \\
 & f_e^s \in \{0, 1\} && \forall s \in S, e \in E_{\text{MCF}} \quad \text{(MCF-10)} \\
 & p \in [0, |S| - 1] && \forall s \in S \quad \text{(MCF-11)}
 \end{aligned}$$

# Single-Commodity Flow IP

## Single-commodity flow formulation

- computes *aggregated* flow on edges independently of the origin
- does not represent virtual arborescence

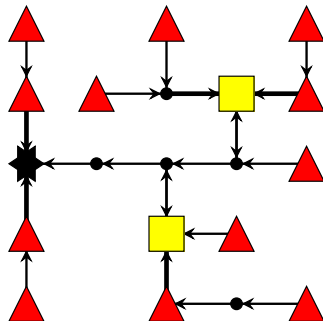


Figure: Single-commodity



# Multi- vs Single-Commodity

Example: 6000 edges and 200 Steiner sites

- Single-commodity: 6000 integer variables
- Multi-commodity: 1,200,000 binary variables

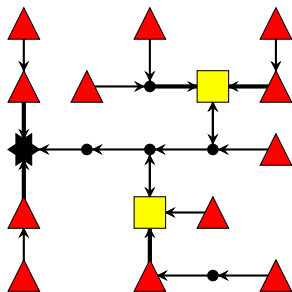


Figure: Single-commodity

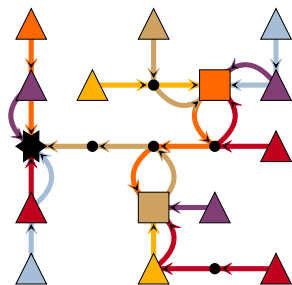


Figure: Multi-commodity

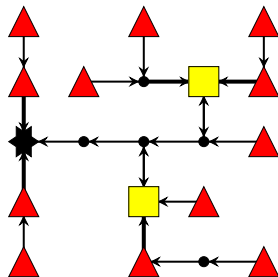
# VirtuCast Algorithm

# VirtuCast Algorithm

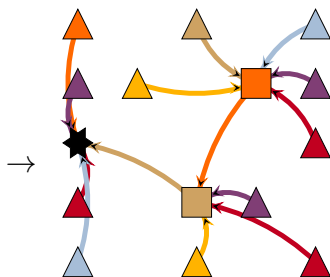
## Outline of VirtuCast

- 1 Solve single-commodity flow IP formulation.
- 2 Decompose IP solution into Virtual Arborescence.

How to decompose?



(a) IP solution



(b) Virtual Arborescence

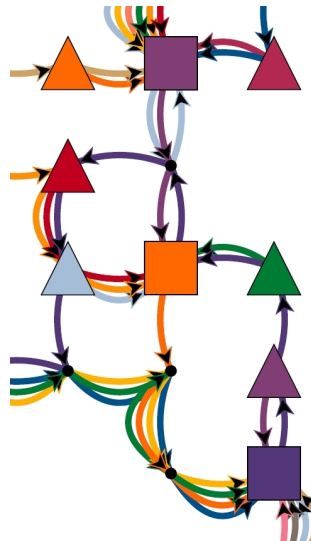
# Decomposing flow is non-trivial (5 pages proof)!

## Flow solution ...

- contains cycles and
- represents *arbitrary* hierarchies.

## Main Results

- decomposition is *always* feasible
- constructive proof:  
polynomial time algorithm



## Integer Program 2: IP-A-CVSAP

$$\text{minimize} \quad C_{IP}(x, f) = \sum_{e \in E_G} c_e f_e + \sum_{s \in S} c_s x_s \quad (\text{IP-OBJ})$$

$$\text{subject to} \quad f(\delta_{E_{\text{ext}}}^+(v)) = f(\delta_{E_{\text{ext}}}^-(v)) \quad \forall v \in V_G \quad (\text{IP-1})$$

$$f(\delta_{E_{\text{ext}}}^+(W)) \geq x_s \quad \forall W \subseteq V_G, s \in W \cap S \neq \emptyset \quad (\text{IP-2})$$

$$f(\delta_{E_{\text{ext}}}^+(W)) \geq 1 \quad \forall W \subseteq V_G, T \cap W \neq \emptyset \quad (\text{IP-3}^*)$$

$$f_e \geq x_s \quad \forall e = (s, o_s^-) \in E_{\text{ext}}^{S^-} \quad (\text{IP-4}^*)$$

$$f_e \leq u_s x_s \quad \forall e = (s, o_s^-) \in E_{\text{ext}}^{S^-} \quad (\text{IP-5})$$

$$f_{(r, o_r^-)} \leq u_r \quad (\text{IP-6})$$

$$f_e \leq u_e \quad \forall e \in E_G \quad (\text{IP-7})$$

$$f_e = 1 \quad \forall e \in E_{\text{ext}}^{T^+} \quad (\text{IP-8})$$

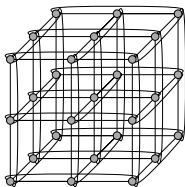
$$f_e = x_s \quad \forall e = (o^+, s) \in E_{\text{ext}}^{S^+} \quad (\text{IP-9})$$

$$x_s \in \{0, 1\} \quad \forall s \in S \quad (\text{IP-10})$$

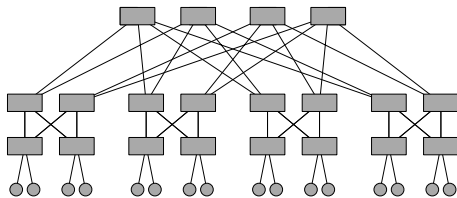
$$f_e \in \mathbb{Z}_{\geq 0} \quad \forall e \in E_{\text{ext}} \quad (\text{IP-11})$$

# Computational Evaluation

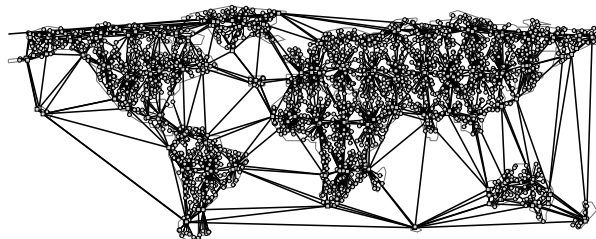
# Topologies



3D torus



Fat tree



An ISP topology generated by IGen with 2400 nodes.

# Instances

## Generation Parameters

- five graph sizes I-V
- 15 instances per graph size: different Steiner costs, different edge capacities

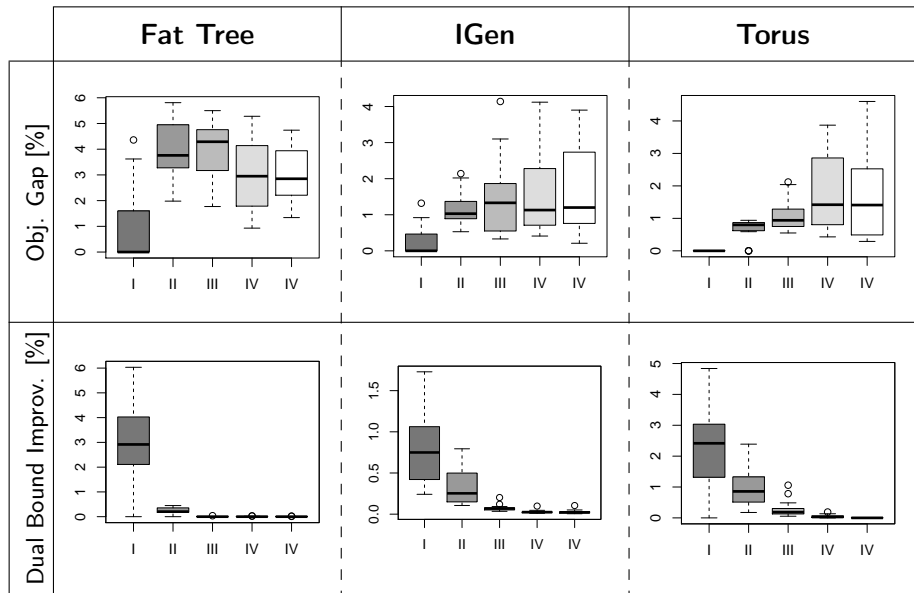
|                 | Nodes | Edges | Processing Locations | Senders |
|-----------------|-------|-------|----------------------|---------|
| <b>Fat tree</b> | 1584  | 14680 | 720                  | 864     |
| <b>3D torus</b> | 1728  | 10368 | 432                  | 864     |
| <b>IGen</b>     | 4000  | 16924 | 401                  | 800     |

Table: Largest graph sizes



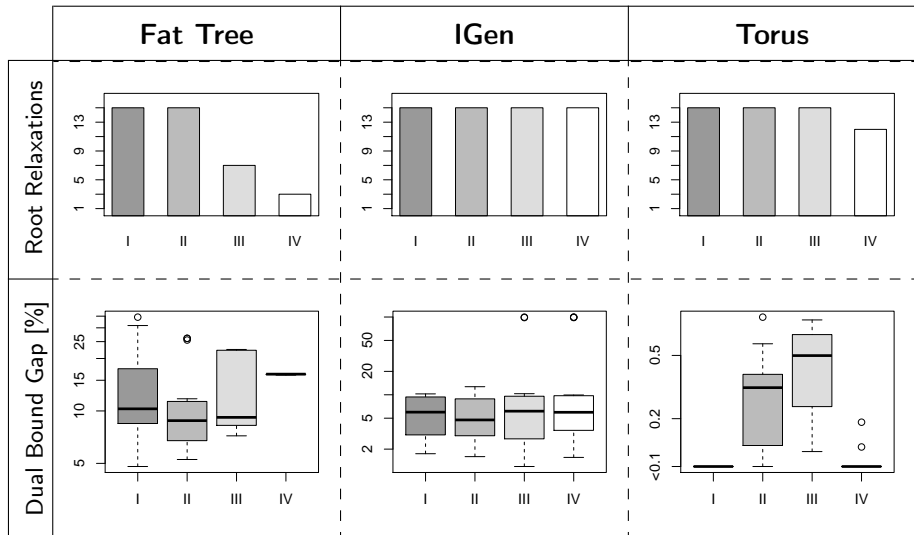
# VirtuCast + LP-based Heuristics

# VirtuCast + LP-based Heuristics



MCF-IP

## MCF-IP: Performance



## Discussion: Applicability to Aggregation Applications

# Applicability of Aggregation Model

Sensor Networks

- commutative
- associative

MapReduce

good aggregation  
in general possible

GigaScope

aggregation factor  
of n:1 unrealistic

← fully applicable

→ not really applicable

# Applicability of Aggregation Model: Discussion

## GigaScope Evaluation Trees

- In GigaScope queries are broken down into evaluation trees
- CVSAP could be applied to subtrees

## Simple Restriction on e.g. OpenFlow statistics

- Byte and packet counter can easily be aggregated

# Future Work

## Model Extensions

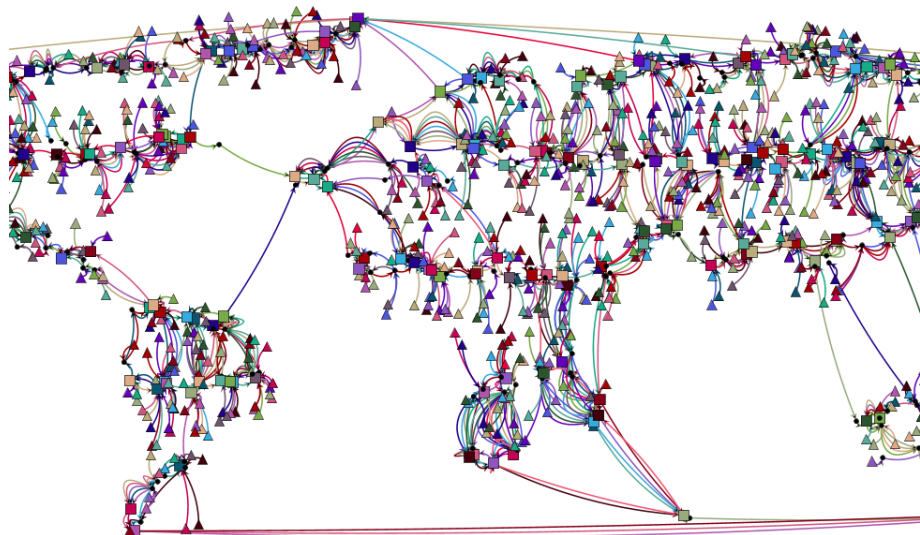
- prize-collecting variants
- concurrent multicast / aggregation sessions

## Application Modeling

- Stratosphere II: Big Data
- UNIFY EU Project: flow analytics



# Thanks



# References I

- [1] T. Achterberg.  
SCIP: solving constraint integer programs.  
*Mathematical Programming Computation*, 1(1):1–41, 2009.
- [2] P. Costa, A. Donnelly, A. Rowstron, and G. O. Shea.  
Camdoop: Exploiting In-network Aggregation for Big Data Applications.  
In *Proc. USENIX Symposium on Networked Systems Design and Implementation (NSDI)*, 2012.
- [3] C. Cranor, T. Johnson, O. Spataschek, and V. Shkapenyuk.  
Gigascop: A Stream Database for Network Applications.  
In *Proc. ACM SIGMOD International Conference on Management of Data*, pages 647–651, 2003.
- [4] M. Ding, X. Cheng, and G. Xue.  
Aggregation tree construction in sensor networks.  
In *Vehicular Technology Conference, 2003. VTC 2003-Fall. 2003 IEEE 58th*, volume 4, pages 2168–2172. IEEE, 2003.
- [5] C. Hermsmeyer, E. Hernandez-Valencia, D. Stoll, and O. Tamm.  
Ethernet aggregation and core network models for efficient and reliable iptv services.  
*Bell Labs Technical Journal*, 12(1):57–76, 2007.

# References II

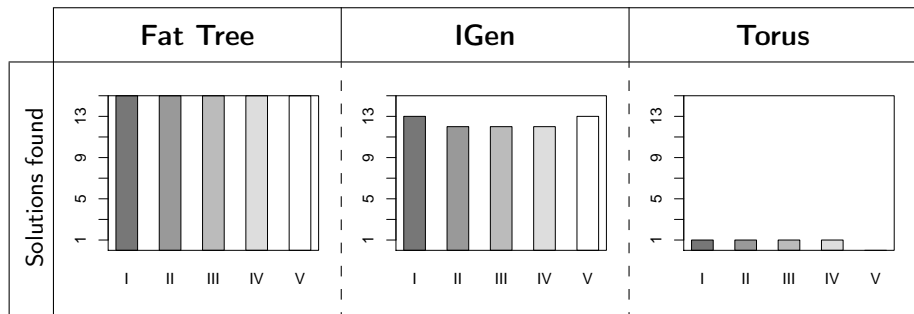
- [6] B. Krishnamachari, D. Estrin, and S. Wicker.  
Modelling data-centric routing in wireless sensor networks.  
In *IEEE infocom*, volume 2, pages 39–44, 2002.
- [7] M. Molnár.  
Hierarchies to Solve Constrained Connected Spanning Problems.  
Technical Report Irimm-00619806, University Montpellier 2, LIRMM, 2011.
- [8] S. Narayana, W. Jiang, J. Rexford, and M. Chiang.  
Joint Server Selection and Routing for Geo-Replicated Services.  
In *Proc. Workshop on Distributed Cloud Computing (DCC)*, 2013.
- [9] C. Oliveira and P. Pardalos.  
Streaming cache placement.  
In *Mathematical Aspects of Network Routing Optimization*, Springer Optimization and Its Applications, pages 117–133. Springer New York, 2011.
- [10] M. Rost and S. Schmid.  
The Constrained Virtual Steiner Arborescence Problem: Formal Definition,  
Single-Commodity Integer Programming Formulation and Computational Evaluation.  
Technical report, arXiv, 2013.

# References III

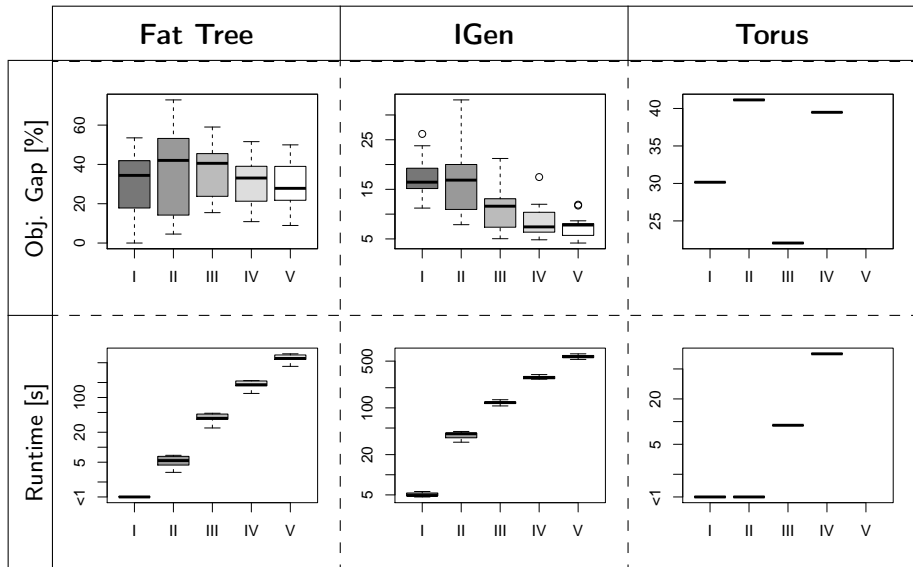
- [11] M. Rost and S. Schmid.  
Virtucast: Multicast and aggregation with in-network processing.  
In R. Baldoni, N. Nisse, and M. Steen, editors, *Principles of Distributed Systems*, volume 8304 of *Lecture Notes in Computer Science*, pages 221–235. Springer International Publishing, 2013.
- [12] S. Shi.  
A proposal for a scalable internet multicast architecture.  
In *Washington Universtiy*, 2001.

GreedySelect

# GreedySelect: Efficacy



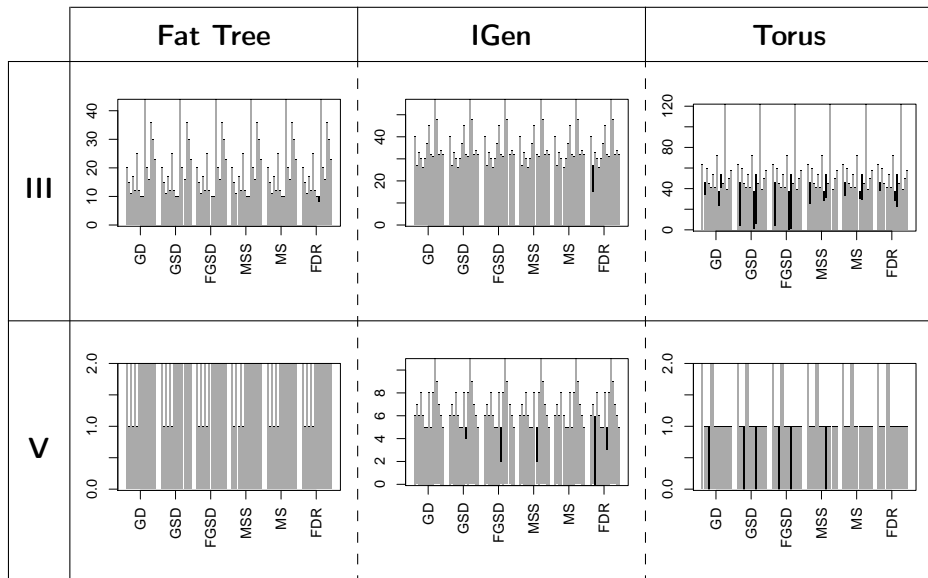
# GreedySelect: Performance



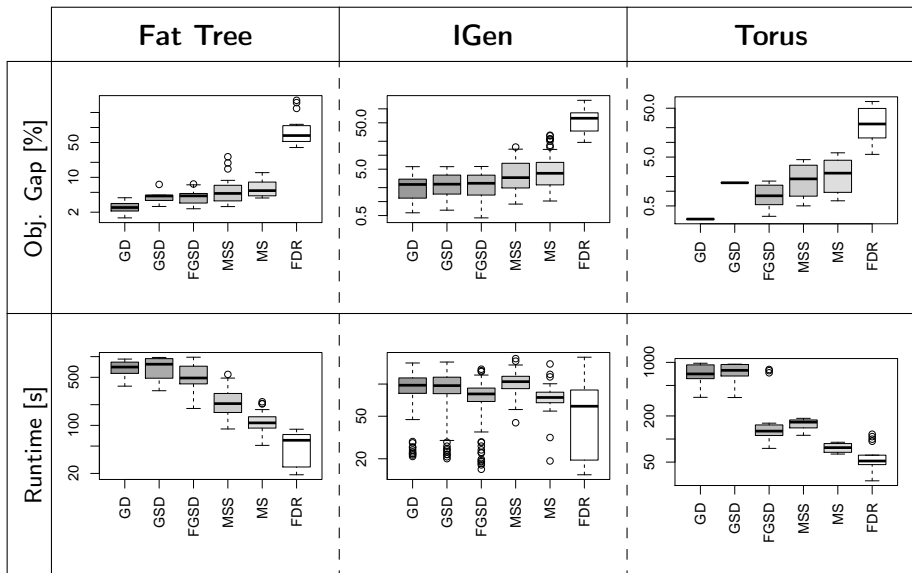
## LP-based Heuristics



# LP-based Heuristics: Efficacy



# LP-based Heuristics: Performance on graph size V



## Publications

Matthias Rost, Stefan Schmid: OPODIS 2013 & arXiv [11, 10]

Applications → Concise definition of CVSAP

## Algorithmic Study

## Inapproximability

## Approximations

- NVSTP
- VSTP
- VSAP

## Exact Algorithms

- multi-commodity flow
- single-commodity flow  
→ VirtuCast

## Heuristics

- FlowDecoRound
- MultipleShots
- GreedyDiving
- GreedySelect

Extensive explorative Computational Evaluation

## Related Work

### Molnar: Constrained Spanning Tree Problems [7]

- Shows that optimal solution is a 'spanning hierarchy' and not a DAG.

### Oliveira et. al: Flow Streaming Cache Placement Problem [9]

- Consider a weaker variant of multicasting CVSAP without bandwidth
- Give weak approximation algorithm

### Shi: Scalability in Overlay Multicasting [12]

- Provided heuristic and showed improvement in scalability.

# Solution Approaches

# Solution Approaches

Wishful thinking: there exists a

- polynomial time algorithm
- solving CVSAP to optimality
- considering all constraints.

# Solution Approaches

Wishful thinking: there exists a

- polynomial time algorithm
- solving CVSAP to optimality
- considering all constraints.

Theorem: Inapproximability of CVSAP

Finding a feasible solution is NP-complete!

# Solution Approaches

Wishful thinking: there exists a

- polynomial time algorithm
- solving CVSAP to optimality
- considering all constraints.

Theorem: Inapproximability of CVSAP

Finding a feasible solution is NP-complete!

## Approximations

- polynomial
- quality guarantee
- weaker models

## Exact Algorithms

- non-polynomial
- optimality
- full model

## Heuristics

- polynomial
- no solution guarantee
- full model



# Comprehensive algorithmic study

## Algorithms

### Approximations

- NVSTP
- VSTP
- VSAP

### Exact Algorithms

- multi-commodity flow
- single-commodity flow  
→ VirtuCast

### LP-based Heuristics

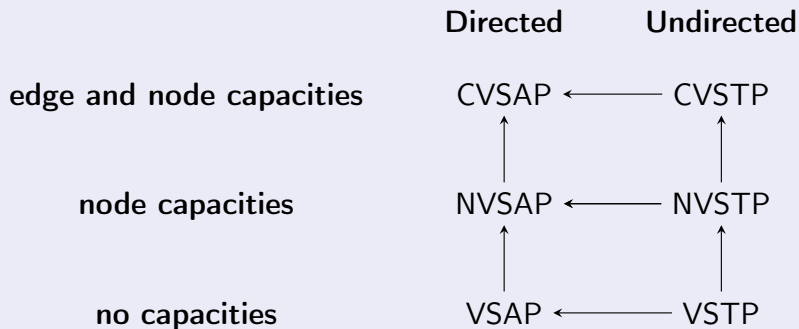
- FlowDecoRound
- MultipleShots
- GreedyDiving

### Combinatorial Heuristic

- GreedySelect

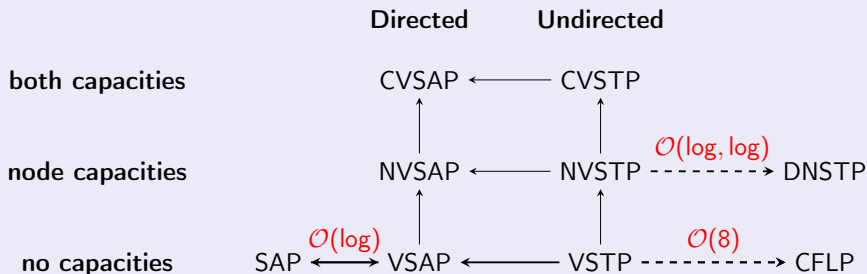
# Approximation Algorithms for Variants

# Variants



# Approximation via related problems

## Results



## Bottom Line

- Better understanding of how to incorporate *virtualized links*.
- Obtained lower bounds via reductions.

# Exact Algorithms for CVSAP

# Overview

## Why exact algorithms matter

- allow trading-off runtime with solution quality
- baseline for heuristics

## Choice: Integer Programming (IP)

- successfully employed for solving related problems (STP, CFLP, ...)
- generates lower bounds on-the-fly

## Combinatorial Heuristic: GreedySelect

# Combinatorial Heuristics

## On Chickens and Eggs

- How and when to place processing nodes?
- How and when to reserve bandwidth for routes?
- How to react to infeasibilities?

## Our Approach

- Place processing functionality and reserve bandwidth jointly.
- Try to avoid infeasibilities by proactive routing decisions.



# GreedySelect Heuristic

Greedly either ...

- connect a single node to the connected component of the receiver or
- connect multiple nodes to an inactive processing node

minimizing the averaged discounted cost per connected node.

Selecting processing node + terminals + paths :  $\mathcal{O}(|V| \cdot |E| + |V|^2 \log |V|)$

compute  $\mathcal{P}_{\bar{s}} \triangleq (\bar{s} \in \bar{S}, T' \subseteq \bar{T}, \mathcal{P}_{T'} = \{P_{t,\bar{s}} | t \in T'\})$ ,

such that  $P_{t,\bar{s}}$  connects  $t$  to  $\bar{s}$ ,

$$u^{\bar{s}}(e) - |\mathcal{P}_{T'}[e]| \geq 0 \text{ for all } e \in E_G,$$

$$2 \leq |T'| \leq u_{r,S}(\bar{s})$$

minimizing  $c_{\bar{s},T'} \triangleq \left( \sum_{t \in T'} (c_E(P_{t,\bar{s}}) - c_E(P_{t,R})) + c_E(P_{\bar{s},R}) + c_S(\bar{s}) \right) / |T'|$

## LP-based Heuristics

# Overview

## Linear Relaxations

- The linear relaxation of an IP is obtained by relaxing the integrality constraints of the variables, thereby obtaining a Linear Program (LP).
- Solutions to linear relaxations are readily available when using branch-and-bound to solve an IP.
- May provide useful information to guide the construction of a solution.

## Usage

- LP-based heuristics are employed within the VirtuCast *solver* to improve the bounding process.
- Yield polynomial time heuristics when used together with the root relaxation.

# FlowDecoRound Heuristic

- computes a *flow* decomposition and connects nodes randomly according to the decomposition
- processing nodes are activated if another node connects to it, must be connected itself
- failsafe: shortest paths

---

**Algorithm 1:** FlowDecoRound
 

---

**Input** : Network  $G = (V_G, E_G, c_E, u_E)$ , Request  $R_G = (r, S, T, u_r, c_S, u_S)$ , LP relaxation solution  $(\hat{r}, \hat{f}) \in \mathcal{F}_{LP}$  to MIP Formulations  
**Output**: A Feasible Virtual Arborescence  $\hat{T}_G$  or null

```

1 set  $\hat{S} \triangleq \emptyset$  and  $\hat{T} \triangleq \emptyset$  and  $U = T$ 
2 set  $\hat{V}_T \triangleq \{r\}$ ,  $\hat{E}_T \triangleq \emptyset$  and  $\hat{\pi} : \hat{E}_T \rightarrow \mathcal{P}_G$ 
3 set  $u(e) \triangleq \begin{cases} u_E(e) & , \text{ if } e \in E_G \\ u_r(r) & , \text{ if } e = (r, o_r^-) \\ u_S(s) & , \text{ if } e = (s, o_S^-) \in E_{\text{ext}}^S \\ 1 & , \text{ else} \end{cases}$  for all  $e \in E_{\text{ext}}$ 
4 while  $U \neq \emptyset$  do
5   choose  $t \in U$  uniformly at random and set  $U \leftarrow U - t$ 
6   set  $\Gamma_t \triangleq \text{MinCostFlow}(G_{\text{ext}}, \hat{r}, \hat{f}(o^+, t), t, \{o_S^-, o_r^-\})$ 
7   set  $\hat{r} \leftarrow \hat{r} - \sum_{(P,f) \in \Gamma_t, e \in P} f$ 
8   set  $\Gamma_t \leftarrow \Gamma_t \setminus \{(P, f) \in \Gamma_t | \exists e \in P, u(e) = 0\}$ 
9   set  $\Gamma_t \leftarrow \Gamma_t \setminus \{(P, f) \in \Gamma_t | (\hat{V}_T + t, \hat{E}_T + (t, P_{|P|-1}) \text{ is not acyclic})\}$ 
10  if  $\Gamma_t \neq \emptyset$  then
11    choose  $(P, f) \in \Gamma_t$  with probability  $f / (\sum_{(P_i, f_i) \in \Gamma_t} f_i)$ 
12    if  $P_{|P|-1} \notin \hat{V}_T$  then
13      set  $U \leftarrow U + P_{|P|-1}$  and  $\hat{V}_T \leftarrow \hat{V}_T + P_{|P|-1}$ 
14      set  $\hat{V}_T \leftarrow \hat{V}_T + t$  and  $\hat{E}_T \leftarrow \hat{E}_T + (t, P_{|P|-1})$ 
15      and  $\hat{\pi}(t, P_{|P|-1}) \triangleq P$ 
16      set  $u(e) \leftarrow u(e) - 1$  for all  $e \in P$ 
17  set  $u(e) \leftarrow 0$  for all  $e = (s, o_S^-) \in E_{\text{ext}}^S$  with  $s \in S \wedge s \notin \hat{V}_T$ 
18  set  $\hat{T} \triangleq (\hat{T} \setminus \hat{V}_T) \cup (\{s \in S \cap \hat{V}_T | \hat{\delta}_{\hat{E}_T}^+(s) = 0\})$ 
19  for  $t \in \hat{T}$  do
20    choose  $P \leftarrow \text{ShortestPath}(G_{\text{ext}}^u, c_E, t, \{o_S^-, o_r^-\})$ 
21    such that  $(\hat{V}_T + t, \hat{E}_T + (t, P_{|P|-1}))$  is acyclic
22    if  $P = \emptyset$  then
23      return null
24    set  $\hat{V}_T \leftarrow \hat{V}_T + t$  and  $\hat{E}_T \leftarrow \hat{E}_T + (t, P_{|P|-1})$  and  $\hat{\pi}(t, P_{|P|-1}) \triangleq P$ 
25    set  $u(e) \leftarrow u(e) - 1$  for all  $e \in P$ 
26  for  $e \in \hat{E}_T$  do
27    set  $P \triangleq \hat{\pi}(e)$ 
28    set  $\hat{\pi}(e) \leftarrow (P_1, \dots, P_{|P|-1})$ 
29  set  $\hat{T}_G \triangleq \text{Virtual Arborescence}(\hat{V}_T, \hat{E}_T, r, \hat{\pi})$ 
30  return PruneSteinerNodes( $\hat{T}_G$ )
  
```

---

# Intermezzo: VCPPrimConnect

## Important Observation

If all placed processing nodes are already connected, all senders can be assigned *optimally* using a minimum cost flow.

## Outline

- 1 connect all opened processing nodes in tree via a adaption of Prim's MST algorithm
- 2 assign all sending nodes using min-cost flow

---

### Algorithm 2: VCPPrimConnect

---

**Input** : Network  $G = (V_G, E_G, c_E, c_E, u_E)$ , Request

$R_G = (r, S, T, u_r, c_S, u_S)$ ,

Partial Virtual Arborecence  $\mathcal{T}_G^P = (V_T^P, E_T^P, r, \pi^P)$

**Output**: Feasible Virtual Arborecence  $\mathcal{T}_G = (V_T, E_T, r, \pi)$  or null

```

1 set  $U \triangleq \{v | v \in V_T^P \setminus \{r\}, \delta_{E_T^P}^+(v) = 0\}$ 
2 set  $\bar{S} \triangleq U \cap S$ 
3 set  $V_T \triangleq V_T^P$ ,  $E_T \triangleq E_T^P$  and  $\pi(u, v) = \pi^P(u, v)$  for all  $(u, v) \in E_T$ 
4 set  $u(e) \triangleq u_E(e) - |\pi(E_T)[e]|$  for all  $e \in E_G$ 
5 while  $\bar{S} \neq \emptyset$  do
6   compute  $R \leftarrow \{r' | r' \in \{r\} \cup (V_T \cap S), r' \text{ reaches } r \text{ in } \mathcal{T}_G, \delta_{E_T}^-(r') <$ 
       $u_{r,S}(r')\}$ 
7   compute  $P = \text{MinAllShortestPath}(G^u, c_E, \bar{S}, R)$ 
8   if  $P = \text{null}$  then
9     return null
10  end
11  set  $\bar{S} \leftarrow \bar{S} - P_1$ 
12  set  $E_T \leftarrow E_T + (P_1, P_{|P_1|})$  and  $\pi(P_1, P_{|P_1|}) \triangleq P$ 
13  set  $u(e) \leftarrow u(e) - 1$  for all  $e \in P$ 
14 end
15 set  $\bar{T} \triangleq U \cap T$ 
16 set  $u_V(r') \triangleq u_{r,S}(r') - \delta_{E_T}^-(r')$  for all  $r' \in \{r\} \cup (V_T \cap S)$ 
17 compute  $\Gamma = \{P^i\} \leftarrow \text{MinCostAssignment}(G, c_E, u, u_V, \bar{T}, \{r\} \cup V_T \cap S)$ 
18 if  $\Gamma = \emptyset$  then
19   return null
20 end
21 set  $E_T \leftarrow E_T + (t, P_{|P^i|}^t)$  and  $\pi(t, P_{|P^i|}^t) \triangleq P^i$  for all  $P^i \in \Gamma$ 
22 return  $\mathcal{T}_G \triangleq (V_T, E_T, r, \pi)$ 

```

---

# MultipleShots

- treats node variables as probabilities and iteratively places processing functionality accordingly
- tries to generate a feasible solution in each round via VCPPrimConnect

---

**Algorithm 3: MultipleShots**


---

**Input** : Network  $G = (V_G, E_G, c_E, u_E)$ , Request

 $R_G = (r, S, T, u_r, c_S, u_S)$ .

 LP relaxation solution  $(\hat{x}, \hat{r}) \in \mathcal{F}_{LP}$  to MIP Formulations

**Output**: A Feasible Virtual Arborescence  $\hat{T}_G$  or null

```

1 set  $[S] \triangleq \{s \in S \mid \hat{x}_s \leq 0.01\}$  and  $[S'] \triangleq \{s \in S \mid \hat{x}_s \geq 0.99\}$ 
2 addConstraintsLocally( $\{x_s = 0 \mid s \in [S]\} \cup \{x_s = 1 \mid s \in [S']\}$ )
3 set  $\hat{S}_0 \triangleq [S] \cup$  and  $\hat{S}_1 \triangleq [S']$ 
4 disableGlobalPrimalBound()
5 repeat
6    $(\hat{x}, \hat{r}) \leftarrow \text{solveSeparateSolve}()$ 
7   if infeasibleLP() return null
8   set  $[S] \triangleq \{s \in S \mid \hat{x}_s \leq 0.01\}$  and  $[S'] \triangleq \{s \in S \mid \hat{x}_s \geq 0.99\}$ 
9   addConstraintsLocally( $\{x_s = 0 \mid s \in [S]\} \cup \{x_s = 1 \mid s \in [S']\}$ )
10  set  $\hat{S}_0 \leftarrow \hat{S}_0 \cup [S]$  and  $\hat{S}_1 \leftarrow \hat{S}_1 \cup [S']$ 
11  set  $\hat{S} \triangleq S \setminus (\hat{S}_0 \cup \hat{S}_1)$ 
12  if  $\hat{S} \neq \emptyset$  then
13    repeat
14      set  $S_1 \triangleq \hat{S}$ 
15      remove  $s$  from  $S_1$  with probability  $1 - \hat{x}_s$  for all  $s \in S_1$ 
16      if  $S_1 = \emptyset$  and  $|S \setminus (\hat{S}_0 \cup \hat{S}_1)| < 10$  then
17        set  $S_1 \leftarrow S \setminus (\hat{S}_0 \cup \hat{S}_1)$ 
18    until  $S_1 \neq \emptyset$ 
19    addConstraintsLocally( $\{x_s = 1 \mid s \in S_1\}$ )
20    set  $\hat{S}_1 \leftarrow \hat{S}_1 \cup S_1$ 
21   $\hat{T}_G^p \triangleq (\hat{V}_T^p, \hat{E}_T^p, r, \emptyset)$  where  $\hat{V}_T^p \triangleq \{r\} \cup T \cup \hat{S}_1$  and  $\hat{E}_T \triangleq \emptyset$ 
22  set  $\hat{T}_G \triangleq \text{VCPPrimConnect}(G, R_G, \hat{T}_G^p)$ 
23  if  $\hat{T}_G \neq \text{null}$  then
24    return PruneSteinerNodes( $\hat{T}_G$ )
25 until  $\hat{S}_0 \cup \hat{S}_1 = S$ 
26 return null

```

---

# GreedyDiving

- aims at generating a feasible *IP* solution
- iteratively bounds at least a single variable from below, first fixing node variables
- complex failsafe:  
PartialDecompose + VCPPrimConnect

---

**Algorithm 4: GreedyDiving**


---

```

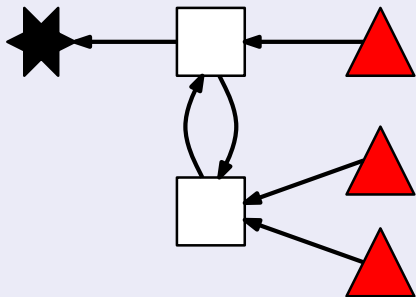
Input : Network  $G = (V_G, E_G, c_E, u_E)$ , Request
         $R_G = (r, S, T, u_r, c_S, u_S)$ ,
        LP relaxation solution  $(\hat{x}, \hat{r}) \in \mathcal{F}_{LP}$  to MIP Formulations
Output: A Feasible Virtual Arborescence  $\mathcal{T}_G$  or null
1 set  $|S| \triangleq \{s \in S | \hat{x}_s \leq 0.01\}$  and  $|S| \triangleq \{s \in S | \hat{x}_s \geq 0.99\}$ 
2 addConstraintsLocally  $\{(x_s = 0 | s \in |S|) \cup \{x_s = 1 | s \in |S|\}\}$ 
3 set  $\hat{S} \triangleq |S| \cup |S|$  and  $\hat{E} \triangleq \emptyset$ 
4 do
5    $(\hat{x}', \hat{r}') \leftarrow \text{solveSeparateSolve}()$ 
6   if infeasibleLP() and  $S = \hat{S}$  then
7     break
8   else if infeasibleLP() or objectiveLimit() then
9     return null
10  set  $(\hat{x}, \hat{r}) \leftarrow (\hat{x}', \hat{r}')$ 
11  if  $\hat{S} \neq S$  then
12    set  $|S| \triangleq \{s \in S | \hat{x}_s \leq 0.01\}$  and  $|S| \triangleq \{s \in S | \hat{x}_s \geq 0.99\}$ 
13    addConstraintsLocally  $\{(x_s = 0 | s \in |S|) \cup \{x_s = 1 | s \in |S|\}\}$ 
14    set  $S \leftarrow S \cup |S| \cup |S|$ 
15    set  $\hat{S} \triangleq S \setminus \hat{S}$ 
16    if  $\hat{S} \neq \emptyset$  then
17      choose  $\hat{s} \in \hat{S}$  with  $c_S(\hat{s})/\hat{x}_s$  minimal
18      addConstraintsLocally  $\{(x_{\hat{s}} = 1)\}$ 
19      set  $S \leftarrow S + \hat{s}$ 
20  else if  $\hat{E} \neq E_{\text{ext}}$  then
21    set  $|E| \triangleq \{e \in E_{\text{ext}} | |\hat{r}_e - \hat{r}_e| \leq 0.001\}$ ,
22     $|E| \triangleq \{e \in E_{\text{ext}} | |\hat{r}_e - \hat{r}_e| \leq 0.001\}$ 
23    addConstraintsLocally  $\{(f_e = \lceil \hat{r}_e \rceil | e \in |E|) \cup \{f_e = \lfloor \hat{r}_e \rfloor | e \in |E|\}\}$ 
24    set  $\hat{E} \leftarrow \hat{E} \cup |E| \cup |E|$ 
25    set  $\hat{E} \triangleq E_{\text{ext}} \setminus \hat{E}$ 
26    if  $\hat{E} \neq \emptyset$  then
27      choose  $\hat{e} \in \hat{E}$  with  $|\hat{r}_{\hat{e}}| - \hat{r}_{\hat{e}}$  minimal
28      addConstraintsLocally  $\{(f_{\hat{e}} \geq \lceil \hat{r}_{\hat{e}} \rceil)\}$ 
29      set  $\hat{E} \leftarrow \hat{E} + \hat{e}$ 
29  else
30    break
31 set  $\hat{r}_e \leftarrow \lceil \hat{r}_e \rceil$  for all  $e \in E_{\text{ext}} \setminus \hat{E}$ 
32 set  $\mathcal{T}_G^D \leftarrow \text{PartialDecompose}(G, R_G, (\hat{x}, \hat{r}'))$ 
33 return  $\text{VCPPrimConnect}(G, R_G, \mathcal{T}_G^D)$ 

```

---

# Example

## Scenario



sender



Steiner  
site

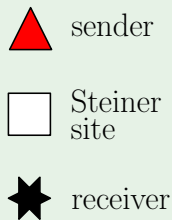
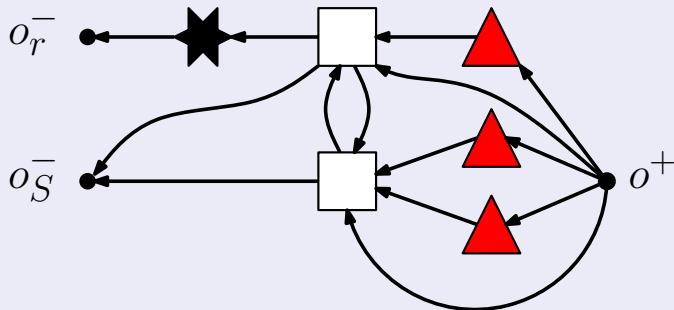


receiver



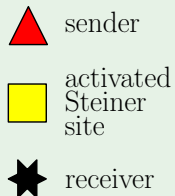
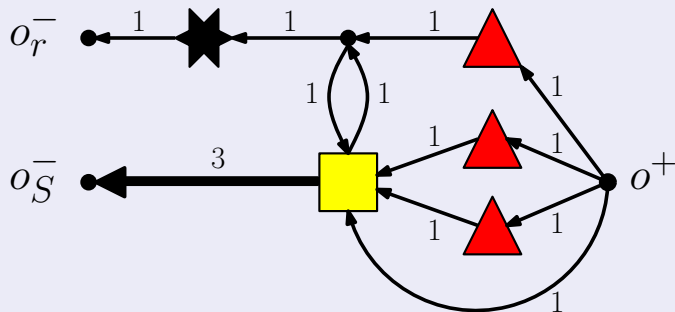
# Example

## Extended Graph



## Example

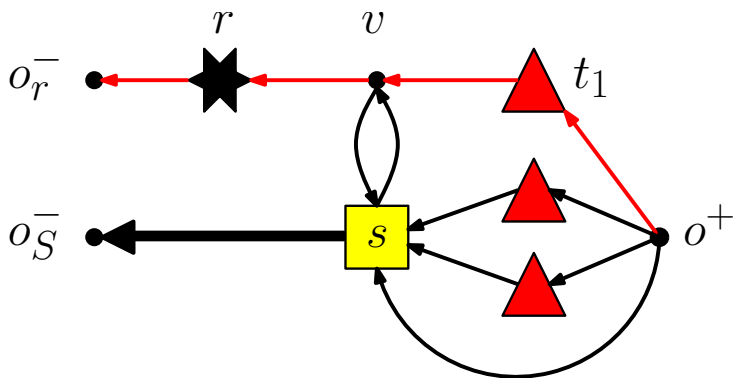
## Solution





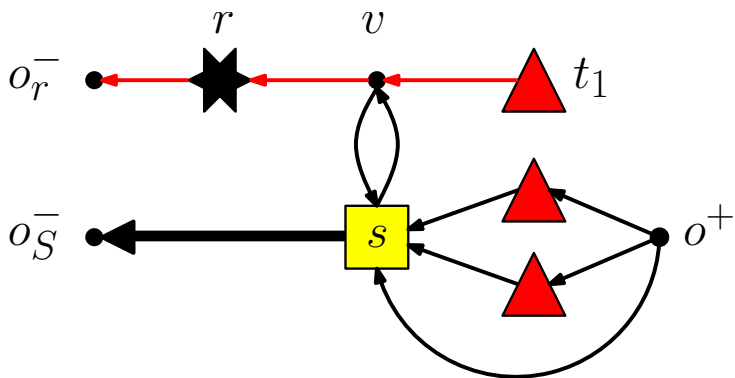
## Decomposition Example I

$$P = \langle o^+, t_1, v, r, o_r^- \rangle$$



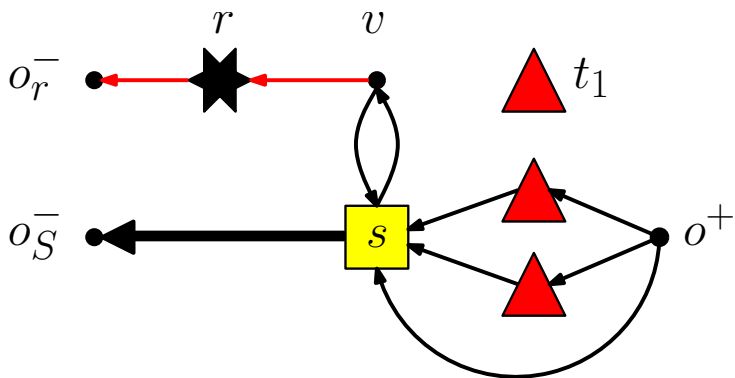
## Decomposition Example I

$$P = \langle o^+, t_1, v, r, o_r^- \rangle$$



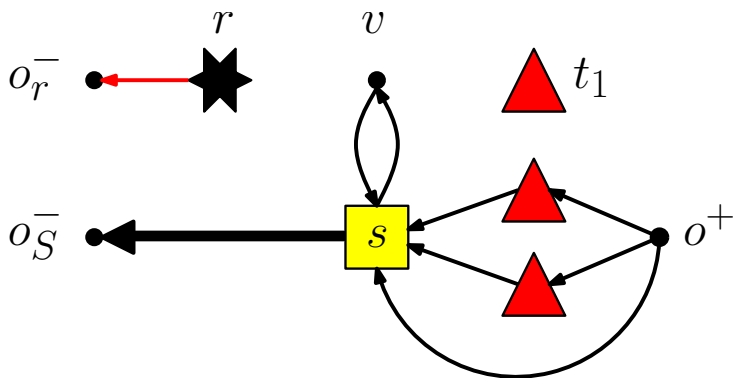
## Decomposition Example I

$$P = \langle o^+, t_1, v, r, o_r^- \rangle$$

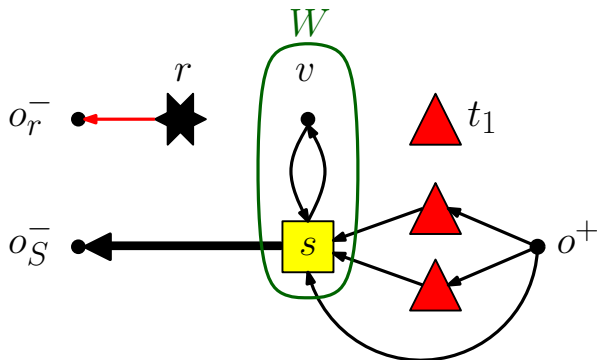


## Decomposition Example I

$$P = \langle o^+, t_1, v, r, o_r^- \rangle$$



# Redirecting Flow

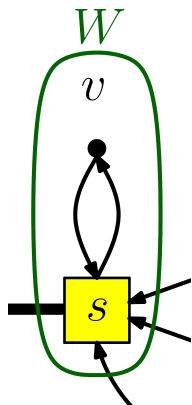


Violation of Connectivity Inequality

$$f(\delta_{E_{\text{ext}}}^+(W)) \geq x_s \quad \forall W \subseteq V_G, s \in W \cap S \neq \emptyset$$



# Redirecting Flow



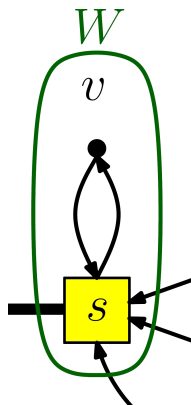
Redirection towards  $o_s^-$  is possible!

There exists a path from  $v$  towards  $o_s^-$  in  $W$ .

## Reasoning

- ① Flow preservation holds within  $W$ .
- ②  $s$  could reach  $o_r^-$  via  $v$  before the reduction of flow.
- ③  $v$  receives at least one unit of flow.
- ④ Flow leaving  $v$  must eventually terminate at  $o_s^-$ .

# Redirecting Flow



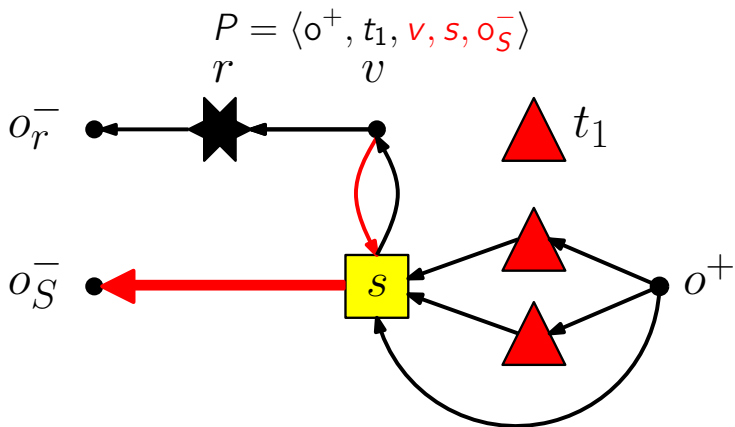
Redirection towards  $o_s^-$  is possible!

There exists a path from  $v$  towards  $o_s^-$  in  $W$ .

## Reasoning

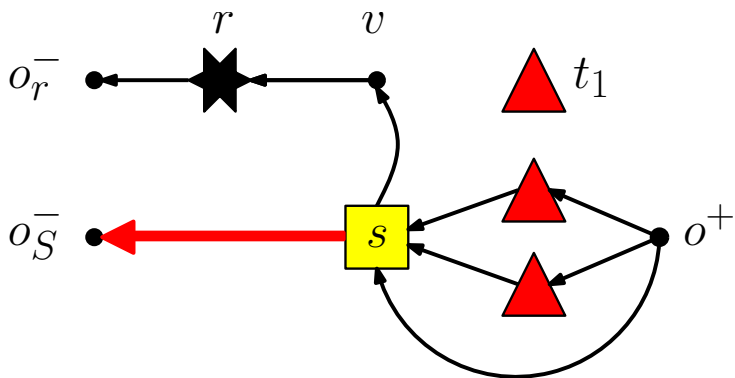
- ① Flow preservation holds within  $W$ .
- ②  $s$  could reach  $o_r^-$  via  $v$  before the reduction of flow.
- ③  $v$  receives at least one unit of flow.
- ④ Flow leaving  $v$  must eventually terminate at  $o_s^-$ .

## Decomposition Example II

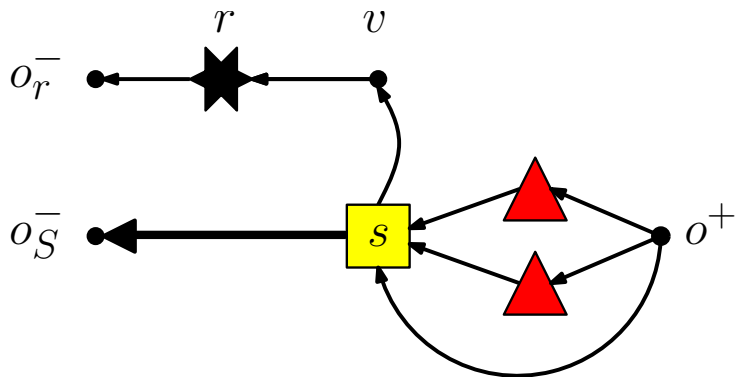


## Decomposition Example II

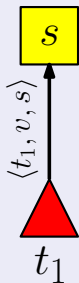
$$P = \langle o^+, t_1, v, s, o_S^- \rangle$$



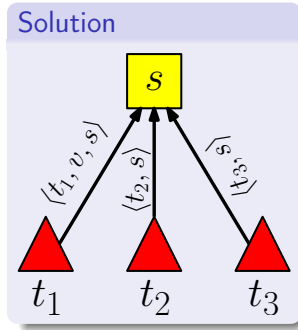
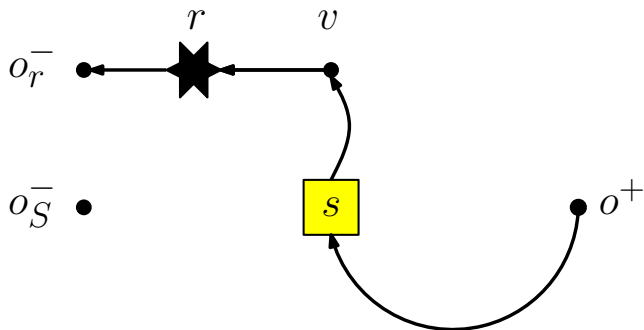
## Decomposition Example II



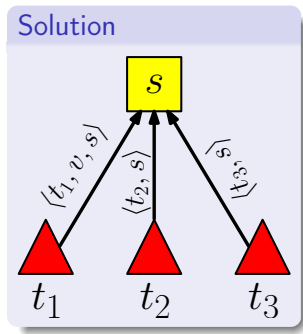
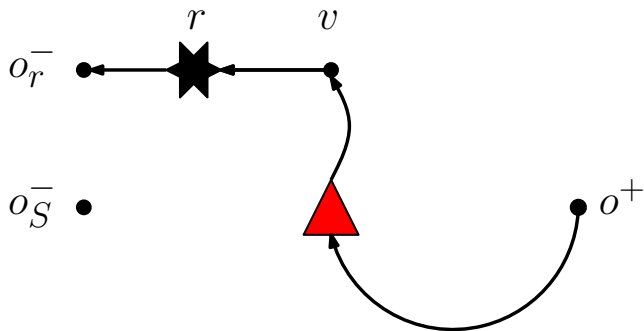
Solution



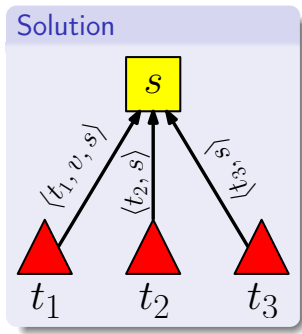
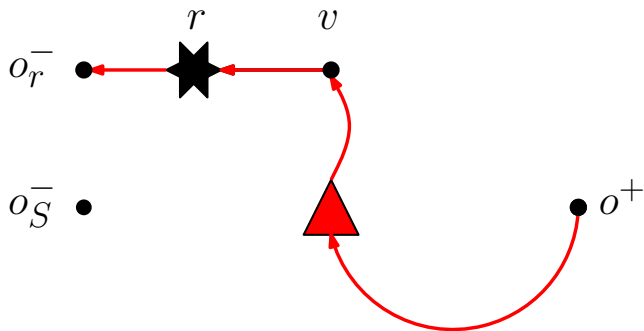
## Decomposition Example II



## Decomposition Example II



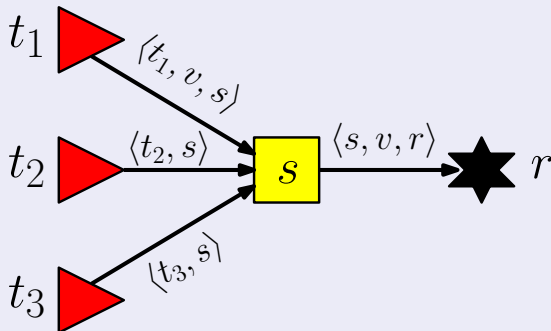
## Decomposition Example II





# Decomposition Example II

## Final Solution



## Related Work

### Molnar: Constrained Spanning Tree Problems [7]

- Shows that optimal solution is a 'spanning hierarchy' and not a DAG.

### Oliveira et. al: Flow Streaming Cache Placement Problem [9]

- Consider a weaker variant of multicasting CVSAP without bandwidth
- Give weak approximation algorithm

### Shi: Scalability in Overlay Multicasting [12]

- Provided heuristic and showed improvement in scalability.