

# Optimal Virtualized In-Network Processing with Applications to Aggregation and Multicast

M.Sc. Thesis Defense Talk

Matthias Rost

April 22nd, 2014

Technische Universität Berlin

## **Reviewer**

Prof. Anja Feldmann, Ph.D., Technische Universität Berlin

Prof. Dr. Andreas Bley, Universität Kassel

## **Supervisor**

Dr. Stefan Schmid, Technische Universität Berlin

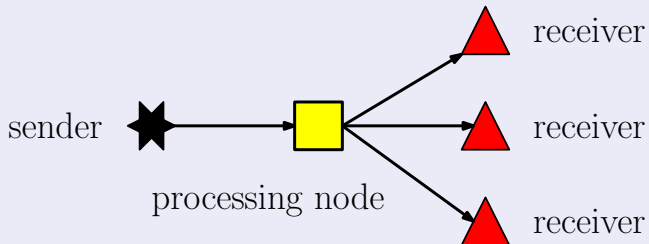
# Mindset

## Service Deployment $\neq$ VNet Embedding

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

## Communication Schemes: Multicast

processing = duplication + reroute



# Communication Schemes: Multicast

processing = duplication + reroute

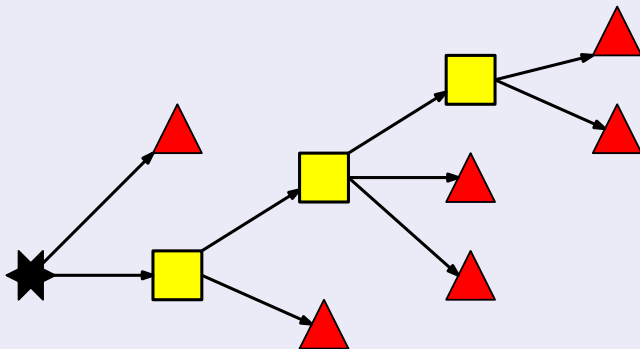
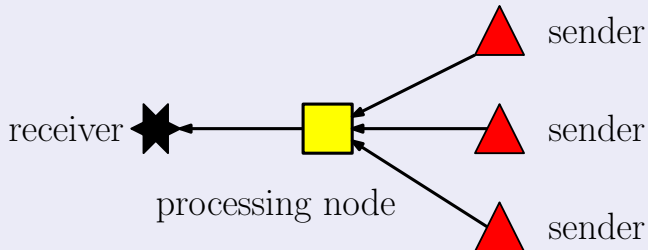


Figure: Hierarchy of processing nodes

## Communication Schemes: Aggregation

processing = merge + reroute



# Communication Schemes: Aggregation

processing = merge + reroute

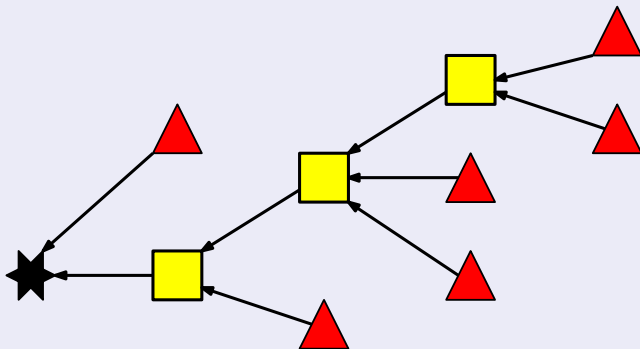


Figure: Hierarchy of processing nodes

# Problem Statement

## Enablers: Network Virtualization, e.g. SDN + NFV

- Routes can be selected arbitrarily
- Network functions can be placed on specific nodes

## 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

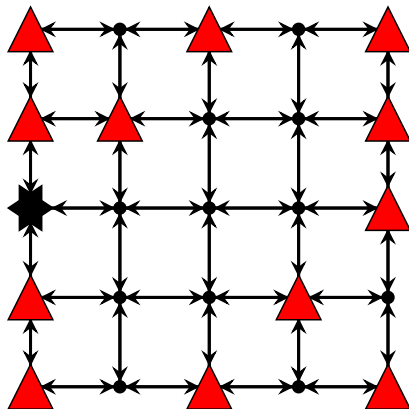
Aggregation scenario

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

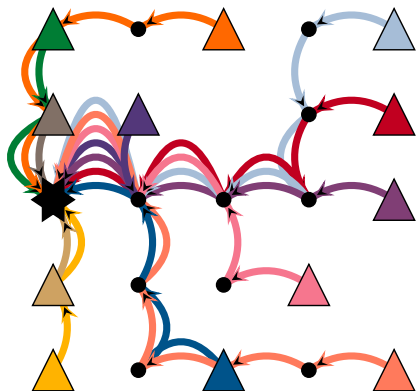


Figure: Unicast solution

# With in-network processing at all nodes

## Solution Method

- Steiner arborescence

## Solution uses

- 16 edges
- 9 processing nodes

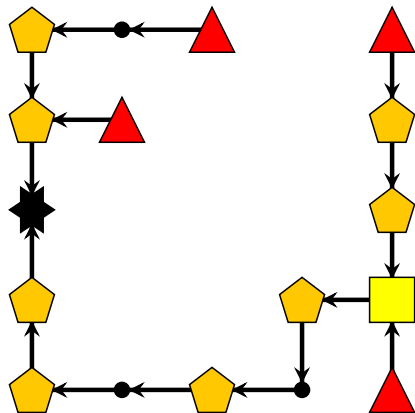
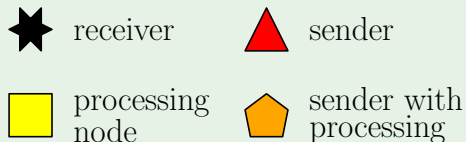
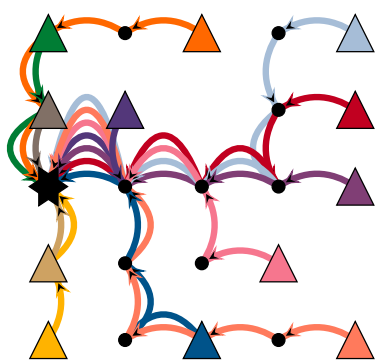
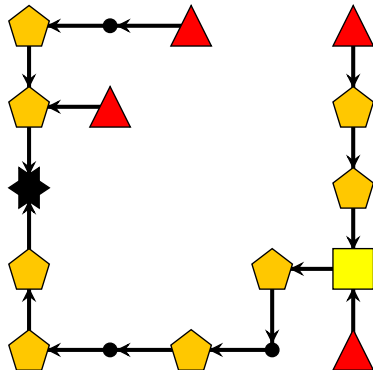


Figure: Aggregation tree

# How to Trade-off?



vs.



# What we aim for

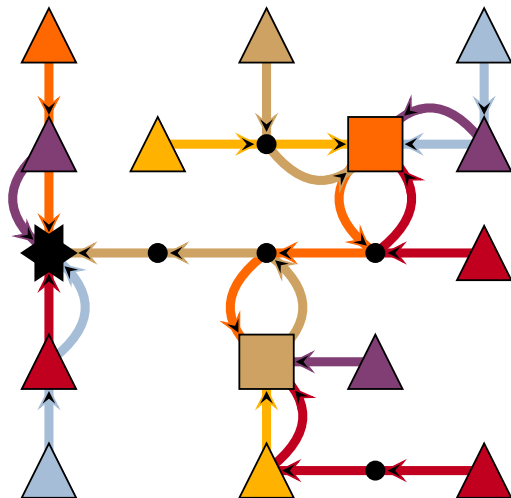
## Solution uses

- 26 edges
- 2 processing nodes

★ receiver

△ sender

□ processing node



# Solution Structure

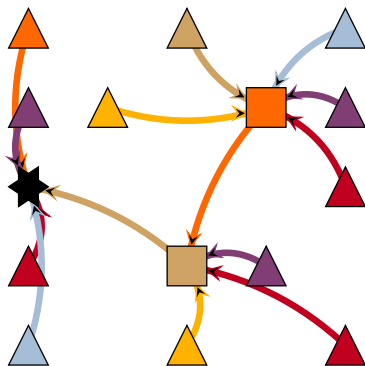


Figure: Virtual Arborescence

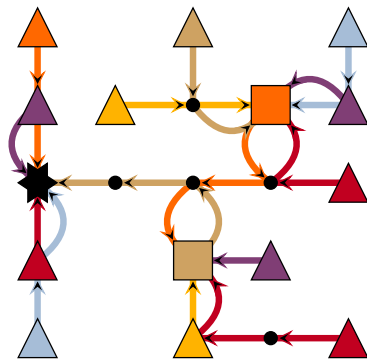


Figure: underlying routes

# New Model: Constrained Virtual Steiner Arborescence Problem

## Definition: CVSAP

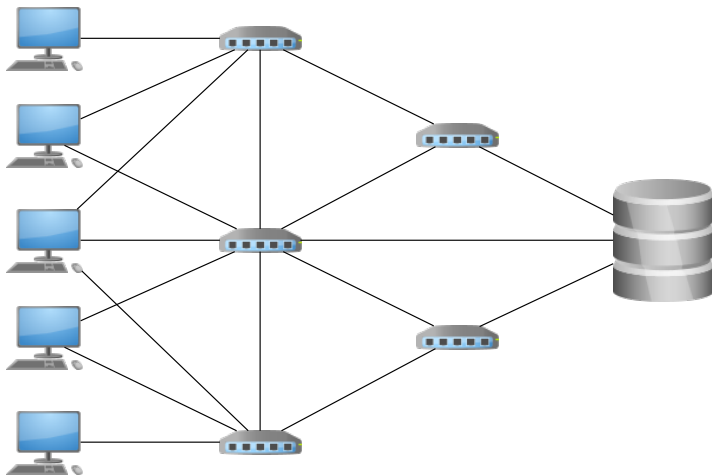
Find a Virtual Arborescence connecting senders to the single receiver, s.t.

- 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.

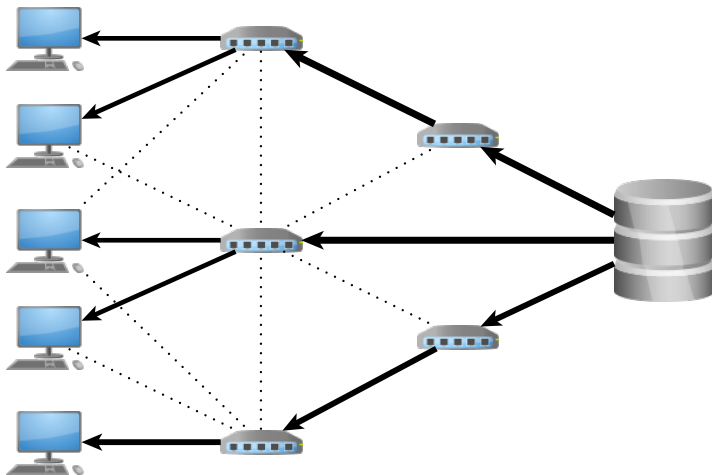
# Applications

# Service Replication



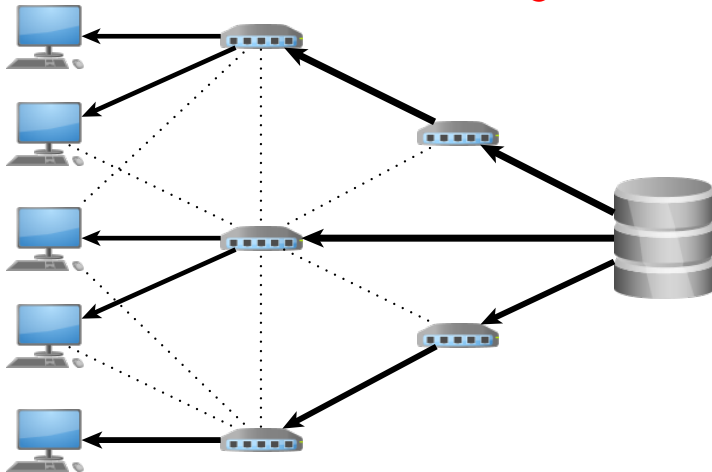


# Service Replication

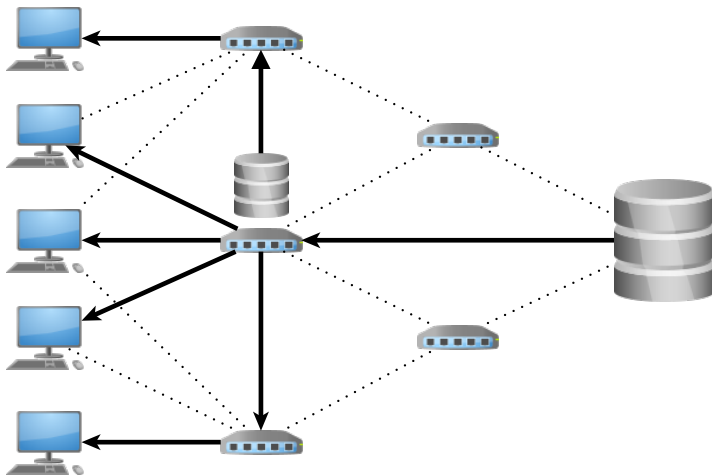


# Service Replication

What if backend links are congested?

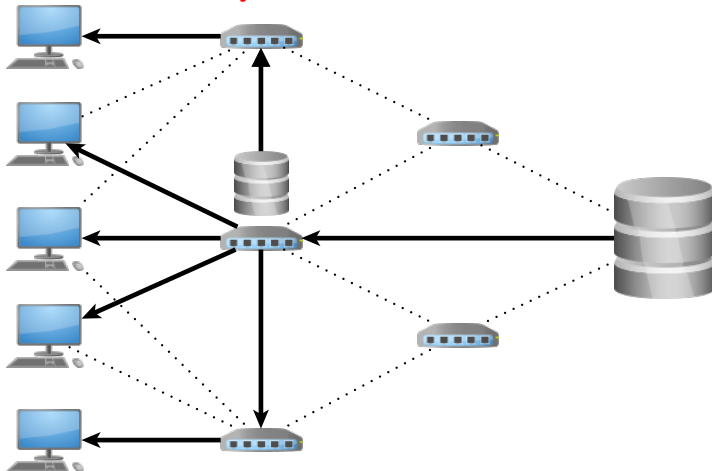


# Service Replication

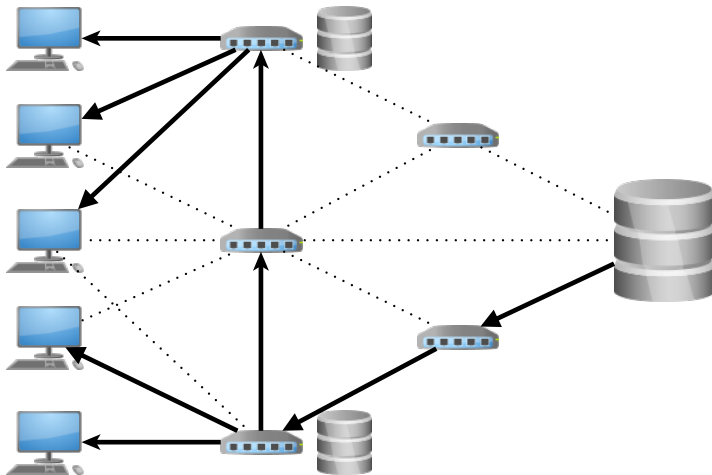


# Service Replication

What if only '3' users can be handled?



# Service Replication



# Applications

	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
	ISP	network analytics: GigaScope [3]	middleboxes / NFV + SDN
	data center	big data / map-reduce: Camdoop [2]	SDN

# 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

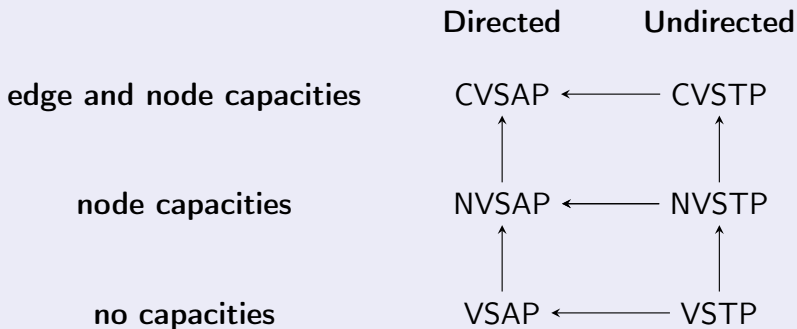
- FlowDecoRound
- MultipleShots
- GreedyDiving

### Combinatorial Heuristic

- GreedySelect

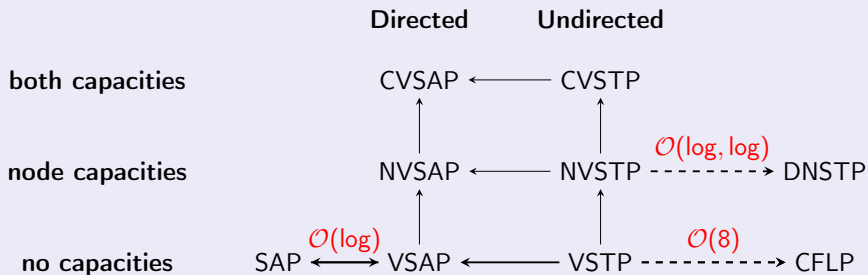
# Approximation Algorithms for Variants

# 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)

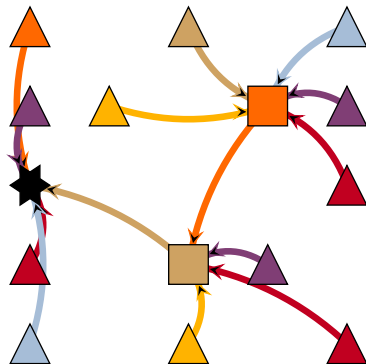
- successfully employed for solving related problems (STP, CFLP, ...)
- generates lower bounds on-the-fly



# Multi-Commodity Flow (MCF) Integer Program

## First approach: MCF IP

- explicitly represent virtual arborescence
- necessitates independent construction of paths for all processing nodes



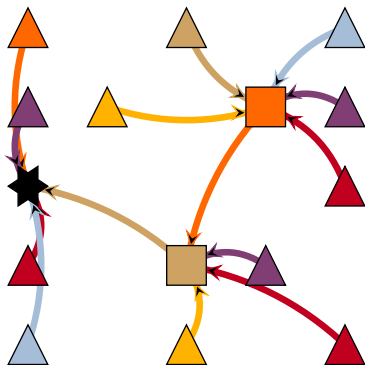
# Multi-Commodity Flow (MCF) Integer Program

## First approach: MCF IP

- explicitly represent virtual arborescence
- necessitates independent construction of paths for all processing nodes

## Does not scale well

- number of binary variables:  
 $\# \text{processing nodes} \cdot \# \text{edges}$



## Integer Program 1: A-CVSAP-MCF

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

# Single-Commodity Flow IP

## Single-commodity flow formulation

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

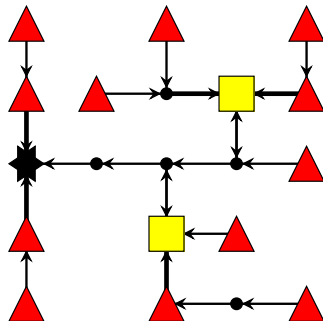


Figure: Single-commodity

# Multi- vs Single-Commodity

Example: 6000 edges and 200 Steiner sites

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

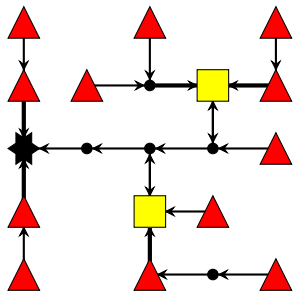


Figure: Single-commodity

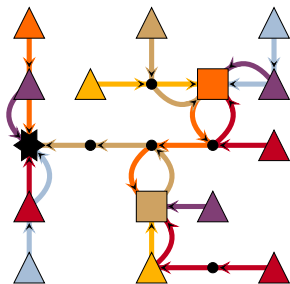


Figure: Multi-commodity

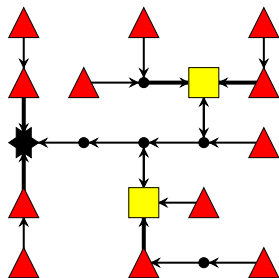
# VirtuCast Algorithm

# VirtuCast Algorithm

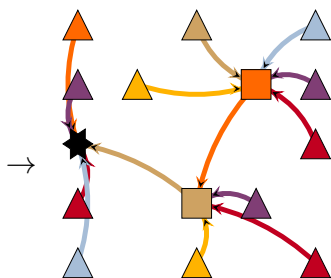
## Outline of VirtuCast

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

How to decompose?



(a) IP solution



(b) Virtual Arborescence

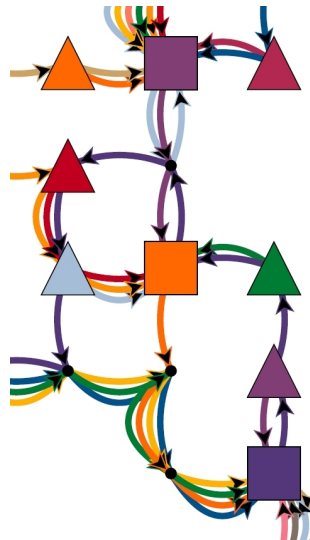
# Decomposing flow is non-trivial (5 pages proof)!

## 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**


---

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

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

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

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

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

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

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

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

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

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

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

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


---

## 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

Greedyly 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  $\bar{s}$ ,

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

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

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

## 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 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 Exact Algorithms

**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  $\Gamma_t \leftarrow \Gamma_t \setminus \{(P, f) \in \Gamma_t \mid \exists e \in P, u(e) = 0\}$ 
9   set  $\Gamma_t \leftarrow \Gamma_t \setminus \{(P, f) \in \Gamma_t \mid (\hat{V}_T + t, \hat{E}_T + (t, P_{|P|-1}) \text{ is not acyclic})\}$ 
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  $\hat{T} \triangleq (\hat{T} \setminus \hat{V}_T) \cup (\{s \in S \cap \hat{V}_T \mid \hat{\delta}_{\hat{E}_T}^+(s) = 0\})$ 
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: VCPrimConnect

## 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: VCPrimConnect

---

**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  $\bar{S} \triangleq U \cup S$ 
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  $P = \text{MinAllShortestPath}(G, c_E, \bar{S}, R)$ 
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  $\bar{T} \triangleq U \cup T$ 
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  $E_T \leftarrow E_T + (t, P_{|P^i|}^i)$  and  $\pi(t, P_{|P^i|}^i) \triangleq P^i$  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 Exact Algorithms

**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  $\hat{S}_0 \triangleq |S|$  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**


---

```

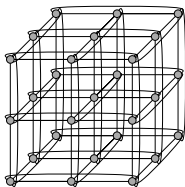
Input : Network  $G = (V_G, E_G, c_E, c_E, u_E)$ , Request
         $R_G = (r, S, T, u_r, c_S, u_S)$ ,
        LP relaxation solution  $(\hat{x}, \hat{r}) \in \mathcal{F}_{LP}$  to Exact Algorithms
Output: A Feasible Virtual Arborescence  $\mathcal{T}_G$  or null
1 set  $[S] \triangleq \{s \in S | \hat{x}_s \leq 0.01\}$  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  $\{(x_s = 1)\}$ 
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     $[\hat{E}] \triangleq \{e \in E_{\text{ext}} | |\hat{r}_e - [\hat{r}_e]| \leq 0.001\}$ 
23    addConstraintsLocally  $\{(f_e = [\hat{r}_e] | e \in [E]) \cup \{f_e = [\hat{r}_e] | e \in [\hat{E}]\}$ 
24    set  $\hat{E} \leftarrow \hat{E} \cup [E] \cup [\hat{E}]$ 
25    set  $\hat{E}_{\text{ext}} \triangleq E_{\text{ext}} \setminus \hat{E}$ 
26    if  $\hat{E}_{\text{ext}} \neq \emptyset$  then
27      choose  $\hat{e} \in \hat{E}_{\text{ext}}$  with  $|\hat{r}_{\hat{e}}| - \hat{r}_{\hat{e}}$  minimal
28      addConstraintsLocally  $\{(f_{\hat{e}} \geq [\hat{r}_{\hat{e}}])\}$ 
29      set  $\hat{E} \leftarrow \hat{E} + \hat{e}$ 
29  else
30    break
31 set  $\hat{r}_e \leftarrow [\hat{r}_e]$  for all  $e \in E_{\text{ext}} \setminus \hat{E}$ 
32 set  $\mathcal{T}_G^L \leftarrow \text{PartialDecompose}(G, R_G, (\hat{x}, \hat{r}))$ 
33 return  $\text{VCPPrimConnect}(G, R_G, \mathcal{T}_G^L)$ 

```

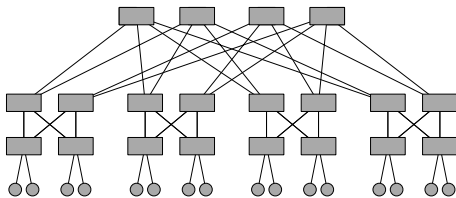
---

# Computational Evaluation

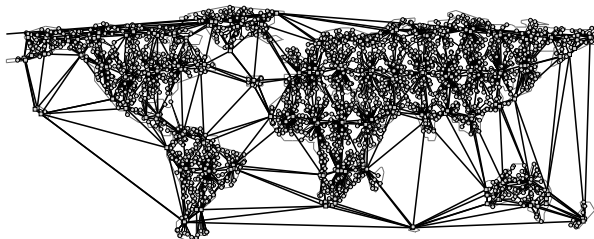
# Topologies



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>I Gen</b>	4000	16924	401	800

Table: Largest graph sizes

# Computational Setup

## Implementation

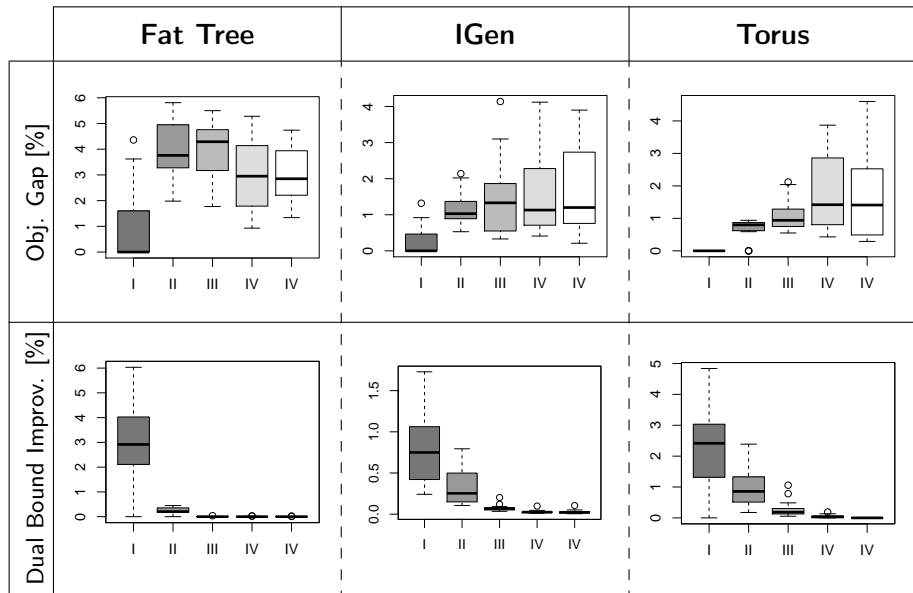
- *all* algorithms (except MCF-IP) are implemented in C/C++
- VirtuCast uses SCIP [1], many different parameters to consider
  - separation
  - branching
  - heuristics
  - ...
- MCF-IP is implemented using GMPL + CPLEX

## Objective

Solve instances within reasonable time: 1 hour runtime limit

# VirtuCast + LP-based Heuristics

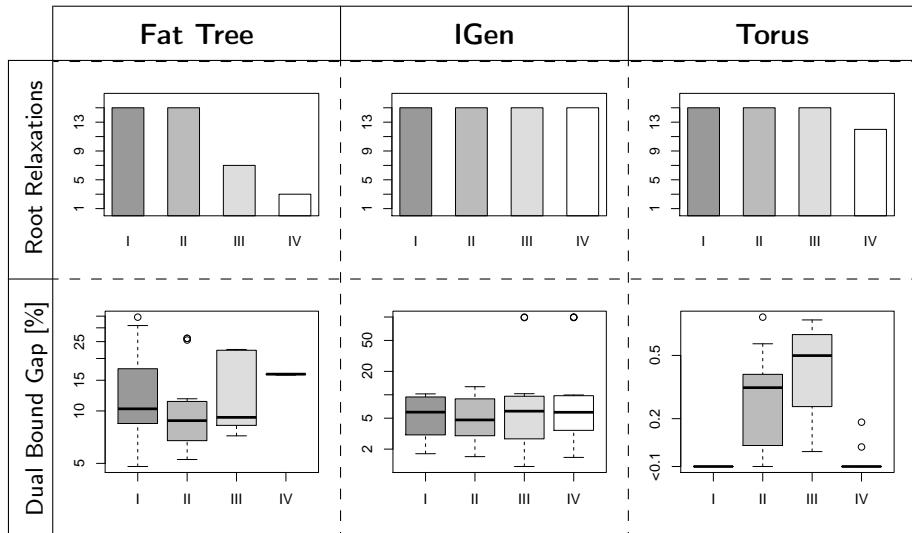
# VirtuCast + LP-based Heuristics





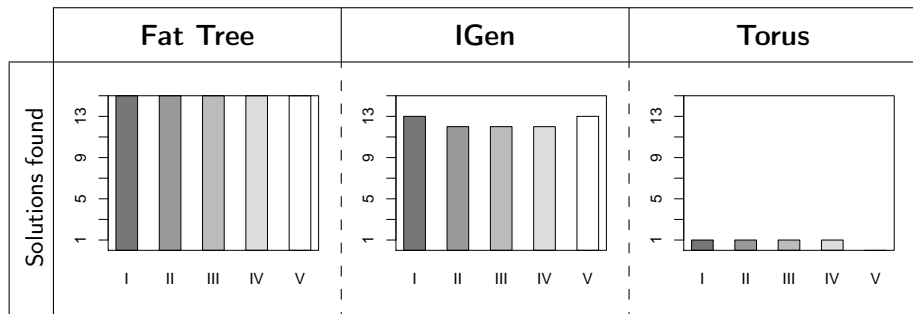
MCF-IP

## MCF-IP: Performance

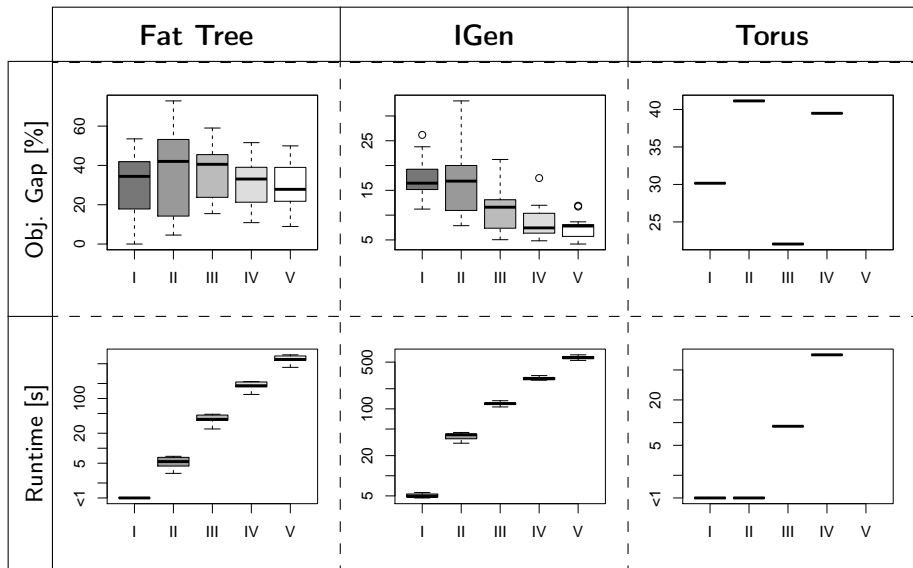


GreedySelect

# GreedySelect: Efficacy

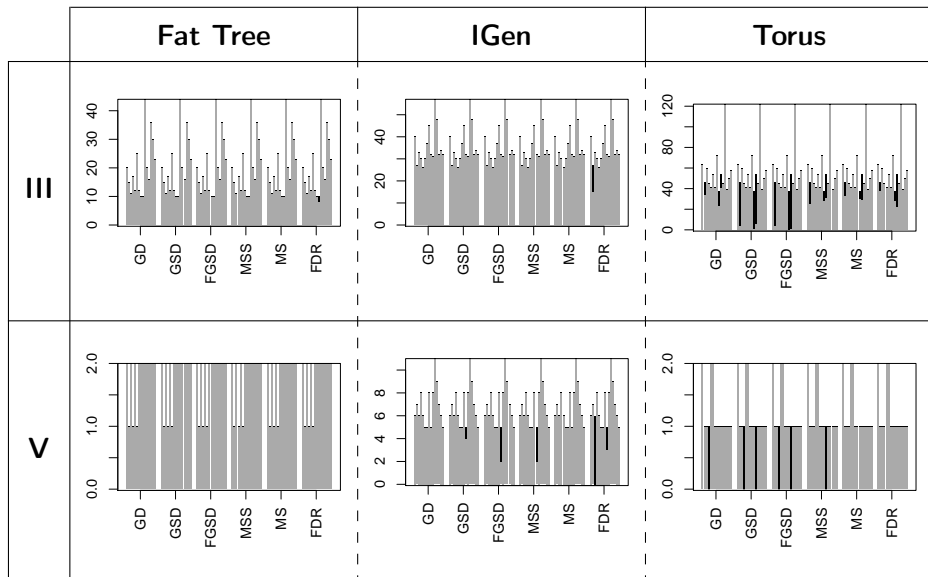


## GreedySelect: Performance

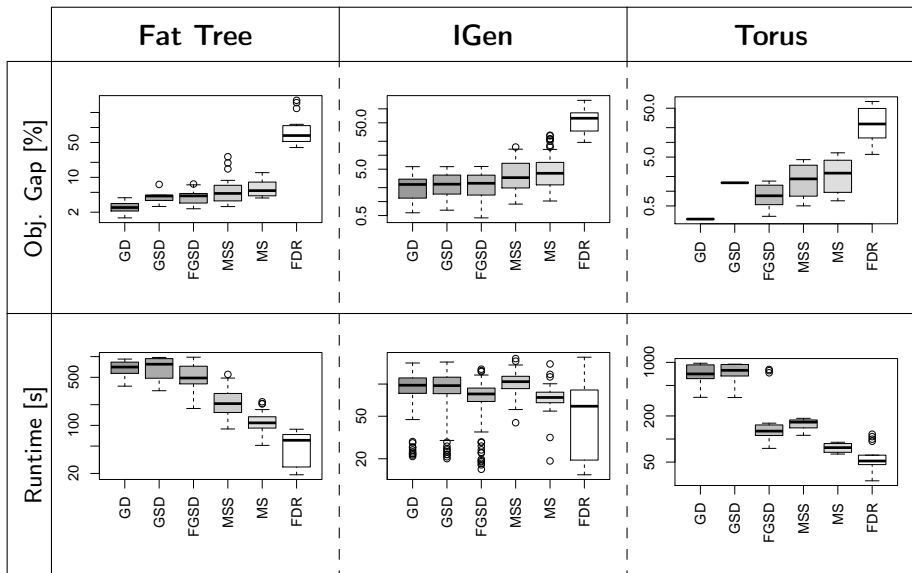


## LP-based Heuristics

## LP-based Heuristics: Efficacy



## LP-based Heuristics: Performance on graph size V





Conclusion

## 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.

# Future Work

## Model Extensions

- prize-collecting variants
- concurrent multicast / aggregation sessions

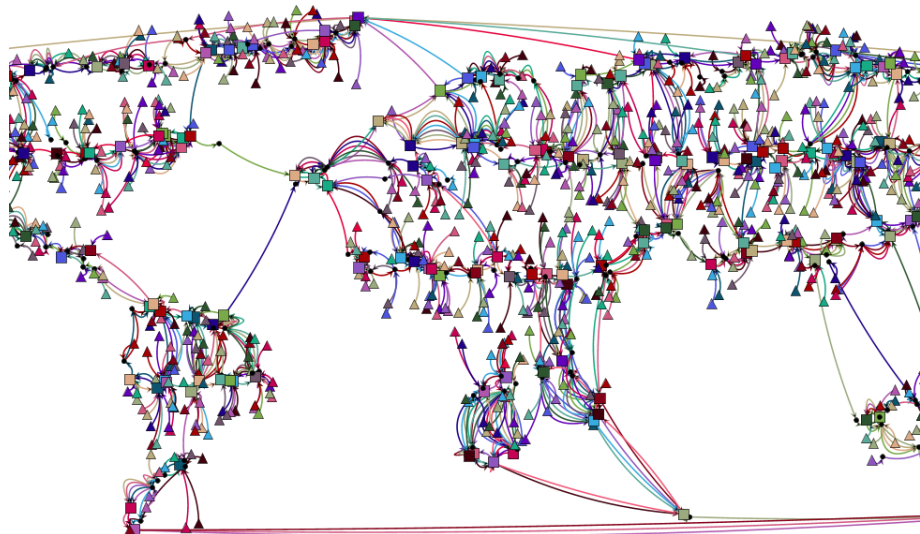
## Application Modeling

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

## IP formulation

- try to derive further cuts / facets

# 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.

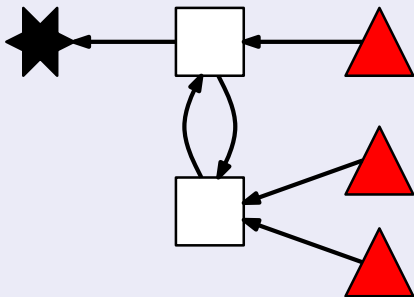
# References III

- [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.



# Example

## Scenario



sender



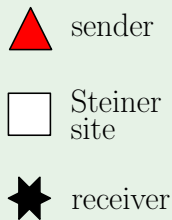
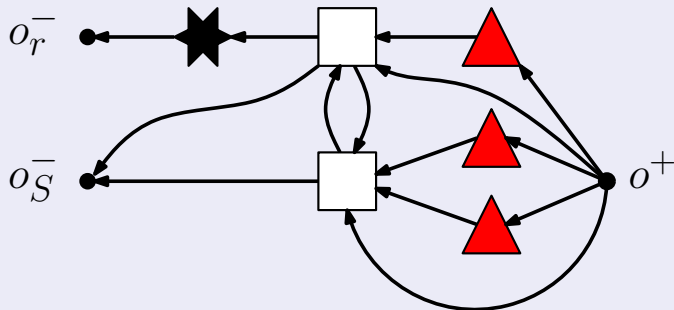
Steiner  
site



receiver

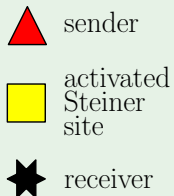
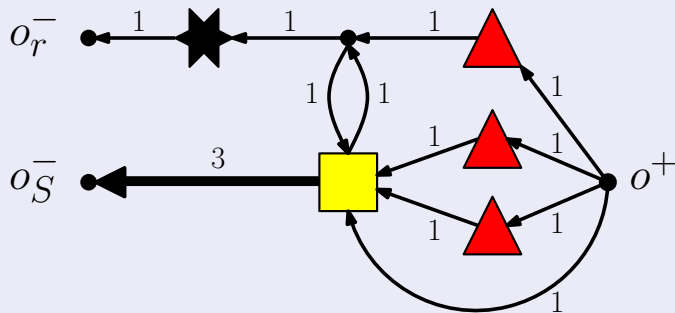
# Example

## Extended Graph

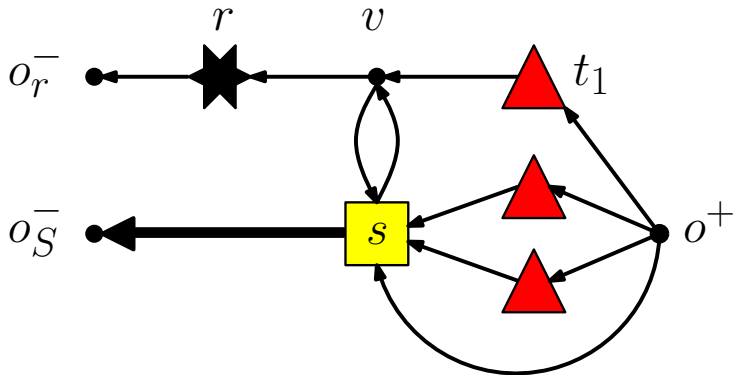


## Example

## Solution

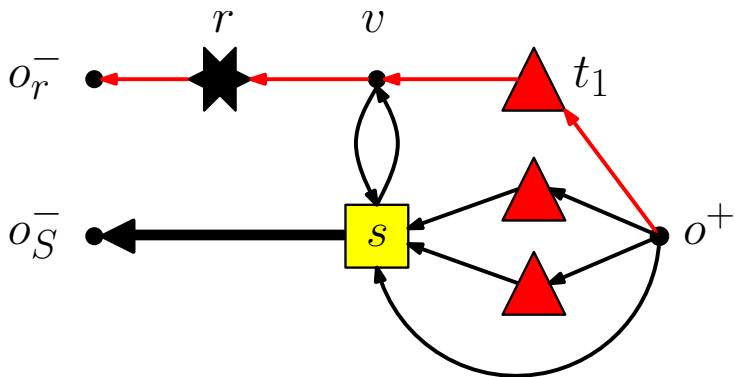


## Decomposition Example I



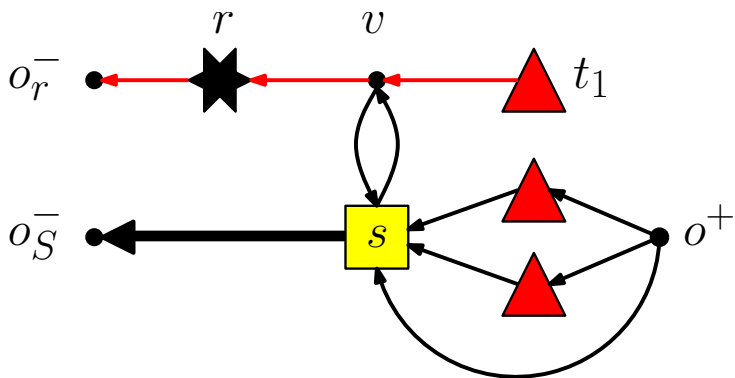
## Decomposition Example I

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



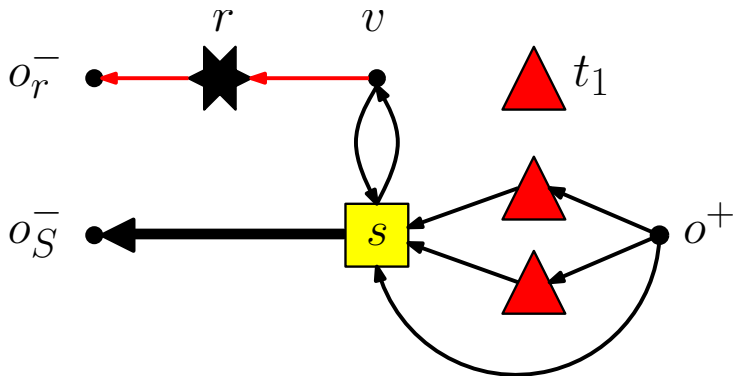
## Decomposition Example I

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



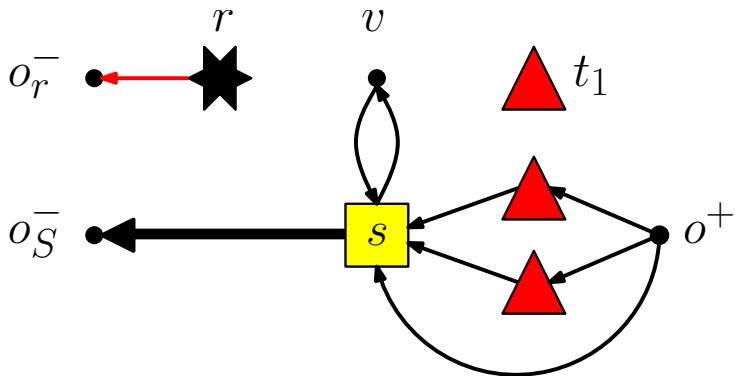
## Decomposition Example I

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



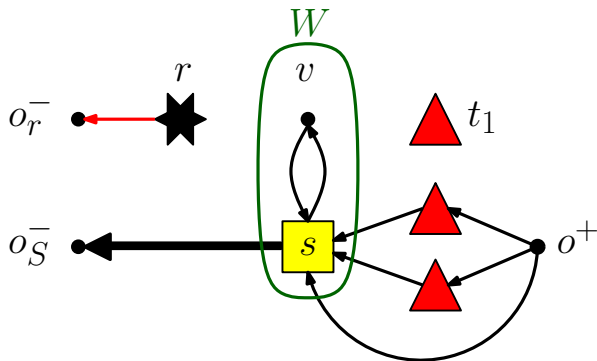
## Decomposition Example I

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





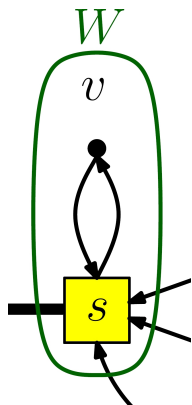
# Redirecting Flow



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



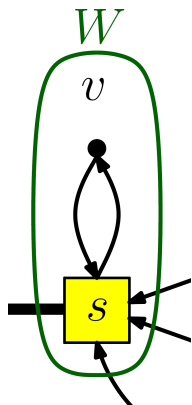
Redirection towards  $o_s^-$  is possible!

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



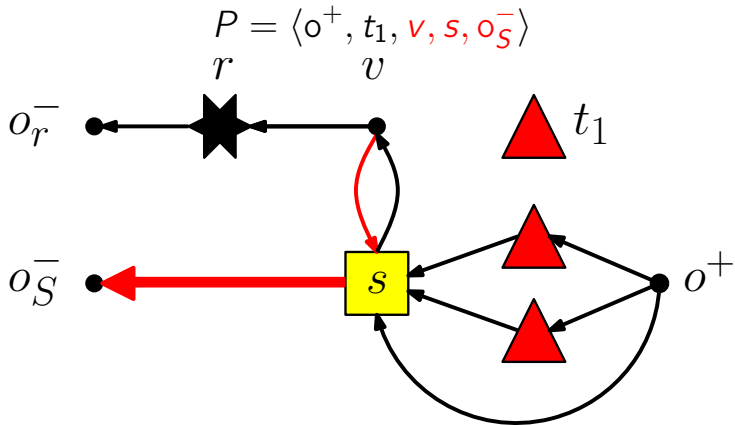
Redirection towards  $o_s^-$  is possible!

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

## Reasoning

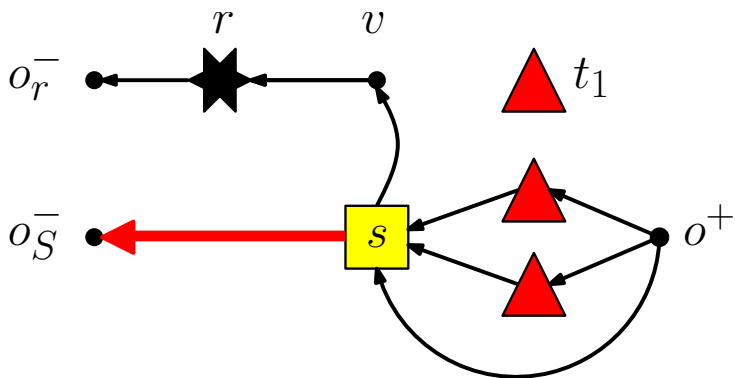
- ① Flow preservation holds within  $W$ .
- ②  $s$  could reach  $o_r^-$  via  $v$  before the reduction of flow.
- ③  $v$  receives at least one unit of flow.
- ④ Flow leaving  $v$  must eventually terminate at  $o_s^-$ .

## Decomposition Example II

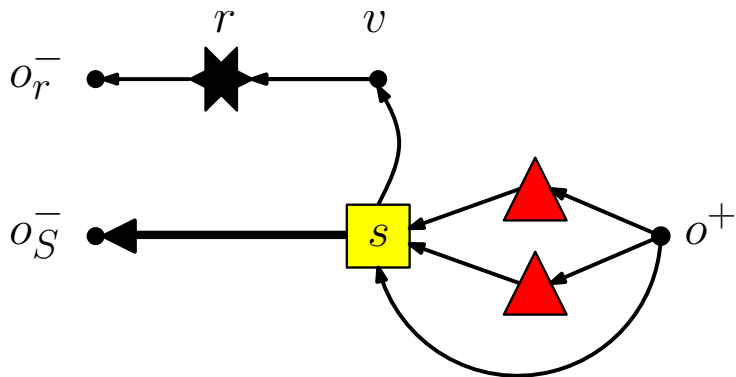


## Decomposition Example II

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



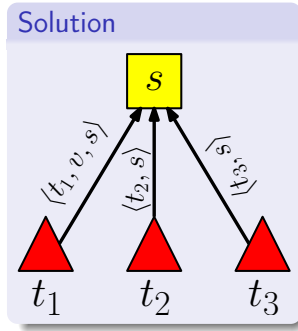
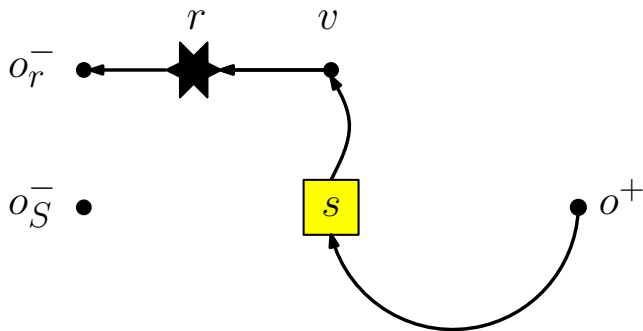
## Decomposition Example II



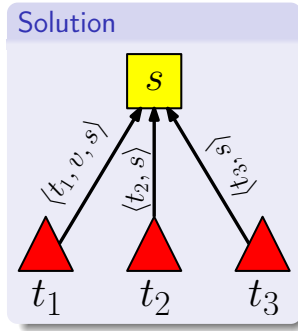
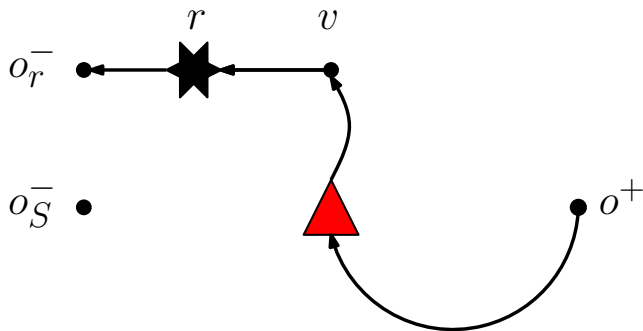
Solution



## Decomposition Example II

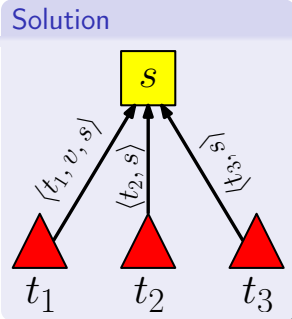
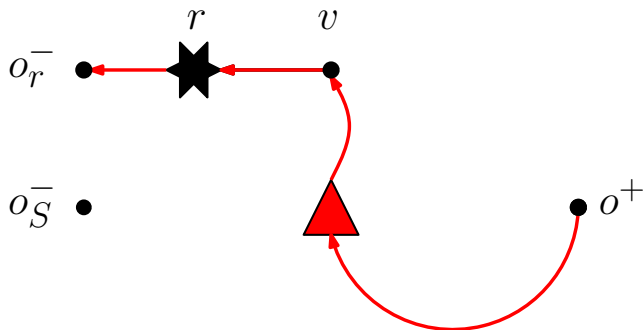


## Decomposition Example II



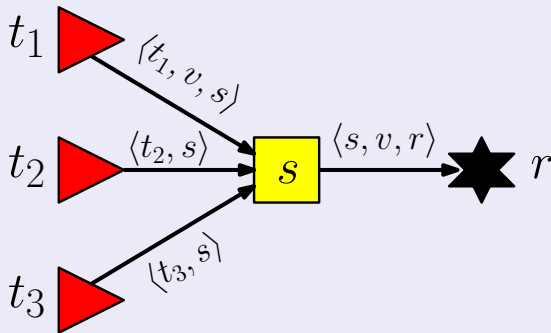


## Decomposition Example II



# 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.