# Algorithmic Nuggets in Network Slicing and Resource Allocation

#### **Stefan Schmid**

Faculty of Computer Science University of Vienna

> IEEE ComSoc Autumn School on Network Slicing and Data Driven Communication

# Virtualization and Programmability: It's a great time to be a networking researcher!



Rhone and Arve Rivers, Switzerland

Credits: George Varghese.

# Teaching Goals: You will understand...

□ ... where flexibilities arise in current networks

□ ... how to exploit them algorithmically

In requirements and challenges of performance isolation

 $\Box$  ... where more research is urgently needed  $\bigcirc$ 

# Roadmap

#### A short recap:

- Networking preliminaries
- Opportunities and challenges of SDNs
- Network virtualization: today and tomorrow
- □ Algorithms I: Flexible routing
- Algorithms II: Flexible embedding and slicing
- Algorithms III: Flexible fast rerouting
- Algorithms IV: Toward self-adjusting networks
- Security aspects (if time permits)

# Roadmap

#### A short recap:

- Networking preliminaries
- Opportunities and challenges of SDNs
- Network virtualization: today and tomorrow
- Algorithms I: Flexible routing
- Algorithms II: Flexible embedding and slicing
- Algorithms III: Flexible fast rerouting
- □ Algorithms IV: Toward self-adjusting networks
  - Security aspects (if time permits)

#### How to design a datacenter network?



#### How to design a datacenter network?



## How to design a datacenter network?



#### Network Layer:

- To transport datagrams between hosts...
- ... routers execute a routing algorithm: compute shortest and policy-compliant paths...
- In based on hierarchical, location-dependent IP addresses!



IP address like **postal address**: changes when moving to new LAN.





But how do switches in a LAN know to which ports they need to forward a packet?





And given an IP address (e.g., of gateway router) how can a host

determine the required MAC destination address?

7



LANs are not about routing but broadcasting and learning. Layer-2 network = single broadcast domain.

Assume: host A wants to send a packet to host B via router R:



- Step 1: Host A creates IP datagram (IPsrc=A, IPdst=B)
- Step 2: A creates link-layer frame with MACdst=R (R's MAC address)



- Step 3: Frame sent from A to R
- Step 4: Frame received at R, passed up to IP



- Step 5: R forwards datagram with same IPsrc=A, IPdst=B
- Step 6: R creates link-layer frame with MACdst=B (B's MAC address), frame contains A-to-B IP datagram



- Step 5: R forwards datagram with same IPsrc=A, IPdst=B
- Step 6: R creates link-layer frame with MACdst=B (B's MAC address), frame contains A-to-B IP datagram





But what if host A does not know R's MAC address yet, but only the gateway R's IP address?



# *Given IP, find MAC!* The ARP Protocol



To learn R's MAC address: Host A broadcasts ARP query containing R's IP address ARP packet dstMAC = FF-FF-FF-FF-FF all nodes on LAN receive ARP query Router R receives ARP packet, replies to A with its (R's) MAC address Unicast to A's MAC address (why?)

# Layer 2: About broadcasting and learning!



# *The learned mappings are cached:* each IP node (host, router) on LAN maintains ARP table:

< IP address; MAC address; ۲۹>

Time-To-Live: avoid stale entries

ARP is "plug&play": nodes learn and create their ARP tables without intervention from administrator!

## Also the switches broadcast and learn



Another question: how does *a switch* know where to forward a frame? E.g., that A' is reachable via interface 4, B' reachable via interface 5?



switch with six interfaces (1,2,3,4,5,6)

# Also the switches broadcast and learn

TTL

3min



Another question: how does a switch know where to forward a frame? E.g., that A' is reachable via interface 4, B' reachable via interface 5?

# **MAC** learning:

Host's MAC

address

A'

B'

Each switch runs a MAC learning protocol and has a switch table:

Interface to

reach host

4

5



# **MAC Learning**

• Destination A' location unknown:

flood but learn A

 Now destination A location known:

send directly (and learn A')



Host's MAC address	Interface to reach host	TTL
A	1	10min
A'	4	10min

# **Intermediate Conclusion**

- Layer-2 networks are very flexible: locationindependent addresses, plug&play, self-learning, etc.: devices (and virtual machines!) can move (migrate)
- But: Layer-2 networks do not scale: despite caching, LAN-wide broadcasts needed once in a while (ARP, MAC learning, DHCP, etc.)!

# **Intermediate Conclusion**

- Layer-2 networks are very flexible: locationindependent addresses, plug&play, self-learning, etc.: devices (and virtual machines!) can move (migrate)
- But: Layer-2 networks do not scale: despite caching, LAN-wide broadcasts needed once in a while (ARP, MAC learning, DHCP, etc.)!
  - ? How large should a LAN be?

Flexibility vs Scalability tradeoff!

## So how to design a datacenter?



**Racks of Servers** 













Virtual Machine with IP address

No need to change IP!

#### A large LAN: High mobility...





Virtual Machine with IP address

No need to change IP!

A large LAN: High mobility... ... but high overhead due to learning and broadcasting.





Virtual Machine with IP address





Virtual Machine with IP address





Virtual Machine with IP address



A small LAN: A different mobility – overhead (scalability) tradeoff!

# What about isolation of tenants?


#### What about isolation of tenants?



Network virtualization: VLANs, VxLANs, tunneling, ... or SDN!

#### Isolation: Required on all levels!



## Roadmap

#### □ A short recap:

- Networking preliminaries
- Opportunities and challenges of SDNs
- On predictable performance
- Algorithms I: Flexible routing
- Algorithms II: Flexible embedding and slicing
- □ Algorithms III: Flexible fast rerouting
- □ Algorithms IV: Theory of demand-aware networks
- Security aspects (if time permits)

## Roadmap

#### A short recap:

- Networking preliminaries
- Opportunities and challenges of SDNs
- On predictable performance
- Algorithms I: Flexible routing
- Algorithms II: Flexible embedding and slicing
- Algorithms III: Flexible fast rerouting
- Algorithms IV: Theory of demand-aware networks
  - Security aspects (if time permits)

## SDN in a Nutshell

- Logically centralized control
- Innovation at speed of software development
- Direct control over paths
- Generalized match-action
- **Example benefits:** 
  - Flexible network virtualization
  - Traffic engineering











Kudos to: Per

- E.g., OSPF igodol
- Indirectly: changing weights igodol



Kudos to: Ped

- E.g., OSPF
- Indirectly: changing weights, equal splits ightarrow



Kudos to: Ped

Indirectly: changing weights, equal splits ightarrow



No link-weight assignment can attain ≤ 100% link utilization!

Kudos to: Pea

all link capacities 1

- E.g., OSPF igodol
- Indirectly: changing weights, equal splits igodol

# **Easy in SDN!**



Only 2 possible demands: 1. s1 
ightarrow d = 12. s2 
ightarrow d = 1

#### **SDN: Not necessarily simpler!**

Networking «Hello World»: MAC learning

#### Principle: for packet (*src,dst*) arriving at port *p*

- □ If *dst* unknown: broadcast packets to all ports
  - Otherwise forward directly to known port
- Also: if *src* unknown, switch learns: *src* is behind *p*

Networking «Hello World»: MAC learning

#### Principle: for packet (*src,dst*) arriving at port *p*

- □ If *dst* unknown: broadcast packets to all ports
  - Otherwise forward directly to known port
- Also: if *src* unknown, switch learns: *src* is behind *p*

#### Example

h1 sends to h2:

3 h3

Networking «Hello World»: MAC learning

#### Principle: for packet (*src,dst*) arriving at port *p*

- □ If *dst* unknown: broadcast packets to all ports
  - Otherwise forward directly to known port
- Also: if *src* unknown, switch learns: *src* is behind *p*

#### Example

h1 sends to h2: flood



Networking «Hello World»: MAC learning

#### Principle: for packet (*src,dst*) arriving at port *p*

- □ If *dst* unknown: broadcast packets to all ports
  - Otherwise forward directly to known port
- Also: if *src* unknown, switch learns: *src* is behind *p*

#### Example

h1 sends to h2: flood, learn (h1,p1)



Networking «Hello World»: MAC learning

#### Principle: for packet (*src,dst*) arriving at port *p*

- □ If *dst* unknown: broadcast packets to all ports
  - Otherwise forward directly to known port
- Also: if *src* unknown, switch learns: *src* is behind *p*

#### Example

- h1 sends to h2: flood, learn (h1,p1)
- □ h3 sends to h1: forward to p1



Networking «Hello World»: MAC learning

#### Principle: for packet (*src,dst*) arriving at port *p*

- If *dst* unknown: broadcast packets to all ports
  - □ Otherwise forward directly to known port
- Also: if *src* unknown, switch learns: *src* is behind *p*

dstmac=h1,fwd(1)
dstmac=h3,fwd(3)

#### Example

- h1 sends to h2: flood, learn (h1,p1)
- □ h3 sends to h1: forward to p1, learn (h3,p3)



Networking «Hello World»: MAC learning

#### Principle: for packet (*src,dst*) arriving at port *p*

- If *dst* unknown: broadcast packets to all ports
  - Otherwise forward directly to known port
- Also: if *src* unknown, switch learns: *src* is behind *p*

dstmac=h1,fwd(1)
dstmac=h3,fwd(3)

#### Example

- h1 sends to h2: flood, learn (h1,p1)
- □ h3 sends to h1: forward to p1, learn (h3,p3)
- h1 sends to h3: forward to p3



#### **From Traditional Networks to SDN**

How to implement this behavior in SDN?



Initial table: Send everything to controller



Pattern	Action
*	send to controller

Initial table: Send everything to controller



Pattern	Action
*	send to controller

□ When h1 sends to h2:



- When h1 sends to h2:
  - Controller learns that h1@p1, updates table, and floods

Principle: only send to ctrl if destination unknown



Pattern	Action
dstmac=h1	Forward(1)
*	send to controller

□ Now assume h2 sends to h1:

Principle: only send to ctrl if destination unknown



Pattern	Action
dstmac=h1	Forward(1)
*	send to controller

- Now assume h2 sends to h1:
  - Switch knows destination: message forwarded to h1
  - BUT: No controller interaction, does not learn about h2: no new rule for h2

Principle: only send to ctrl if destination unknown



Pattern	Action	
dstmac=h1	Forward(1)	
*	send to controller	h3 sends to h2

□ Now, when h3 sends to h2:

#### Controller **Example: SDN MAC Learning Done Wrong** h1 3 h3 Principle: only send to ctrl h2 if destination unknown OpenFlow switch Pattern Action Pattern Action dstmac=h3 Forward(3) dstmac=h1 Forward(1) h3 sends to h2 \* dstmac=h1 Forward(1) send to controller \* send to controller

- □ Now, when h3 sends to h2:
  - Dest unknown: goes to controller which learns about h3
  - And then floods

Principle: only send to ctrl if destination unknown



Pattern	Action
dstmac=h3	Forward(3)
dstmac=h1	Forward(1)
*	send to controller

□ Now, if h2 sends to h3 or h1:

Principle: only send to ctrl if destination unknown



Pattern	Action
dstmac=h3	Forward(3)
dstmac=h1	Forward(1)
*	send to controller

- □ Now, if h2 sends to h3 or h1:
  - Destinations known: controller does not learn about h2

Principle: only send to ctrl if destination unknown



Pattern	Action
dstmac=h3	Forward(3)
dstmac=h1	Forward(1)
*	send to controller

*Ouch!* Controller cannot learn about h2 anymore: whenever h2 is source, destination is known. All future requests to h2 will *all be flooded*: inefficient!



## Roadmap

#### □ A short recap:

- Networking preliminaries
- Opportunities and challenges of SDNs
- On predictable performance
- Algorithms I: Flexible routing
- Algorithms II: Flexible embedding and slicing
- □ Algorithms III: Flexible fast rerouting
- □ Algorithms IV: Theory of demand-aware networks
- Security aspects (if time permits)

## Roadmap

- □ A short recap:
  - Networking preliminaries
  - Opportunities and challenges of SDNs
  - On predictable performance
- □ Algorithms I: Flexible routing
- Algorithms II: Flexible embedding and slicing
- Algorithms III: Flexible fast rerouting
- □ Algorithms IV: Theory of demand-aware networks
- Security aspects (if time permits)

#### **On Predictable Performance...**



An Experiment: 2 vSDNs with bw guarantee!


An Experiment: 2 vSDNs with bw guarantee!



An Experiment: 2 vSDNs with bw guarantee!



An Experiment: 2 vSDNs with bw guarantee!



#### **The Many Faces of Performance Interference**



#### Performance also depends on hypervisor type...

(multithreaded or not, which version of Nagle's algorithm, etc.)

#### ... number of tenants...



#### **Further Reading**

On The Impact of the Network Hypervisor on Virtual Network Performance Andreas Blenk, Arsany Basta, Wolfgang Kellerer, and Stefan Schmid. IFIP Networking, Warsaw, Poland, May 2019.

### Roadmap

#### □ A short recap:

- Networking preliminaries
- Opportunities and challenges of SDNs
- On predictable performance
- Algorithms I: Flexible routing
- Algorithms II: Flexible embedding and slicing
- □ Algorithms III: Flexible fast rerouting
- □ Algorithms IV: Theory of demand-aware networks
- Security aspects (if time permits)

### Roadmap

#### □ A short recap:

- Networking preliminaries
- Opportunities and challenges of SDNs
- On predictable performance
- Algorithms I: Flexible routing
- Algorithms II: Flexible embedding and slicing
- Algorithms III: Flexible fast rerouting
- □ Algorithms IV: Theory of demand-aware networks
- Security aspects (if time permits)

## A Nugget: Flexible Traffic Steering

A novelty:

- □ Traditionally: routes form simple paths (e.g., shortest paths)
- Now: routing through middleboxes may require more general paths, with loops: a walk



Computing shortest routes through waypoints is non-trivial!
Assume unit capacity and

demand for simplicity!



Computing shortest routes through waypoints is non-trivial!
Assume unit capacity and

demand for simplicity!



Greedy fails: choose shortest path from s to w...

Computing shortest routes through waypoints is non-trivial!
Assume unit capacity and

demand for simplicity!



Greedy fails: ... now need long path from w to t

Computing shortest routes through waypoints is non-trivial!
Assume unit capacity and

demand for simplicity!

Greedy fails: ... now need long path from w to t

Computing shortest routes through waypoints is non-trivial!
Assume unit capacity and

demand for simplicity!

A better solution: jointly optimize the two segments!

W

## **Relationship to Shortest Disjoint Paths**

If capacities are 1, segments need to be edgedisjoint: A **disjoint paths problem** 



- A well-known combinatorial problem!
- NP-hard on directed networks
- Feasibility in P on undirected networks for small (constant) number of flows
- Polytime randomized algorithm for 2 disjoint paths (recent result!)

#### NP-hard on *Directed* Networks: Reduction from Disjoint Paths Problem

**Reduction:** From joint shortest paths  $(s_1,t_1),(s_2,t_2)$ to shortest walk (s,w,t) problem



#### NP-hard on *Directed* Networks: Reduction from Disjoint Paths Problem

**Reduction:** From joint shortest paths  $(s_1,t_1),(s_2,t_2)$ to shortest walk (s,w,t) problem



#### NP-hard on *Directed* Networks: Reduction from Disjoint Paths Problem

**Reduction:** From joint shortest paths  $(s_1,t_1),(s_2,t_2)$ to shortest walk (s,w,t) problem





#### What about waypoint routes on undirected networks?

## What about waypoint routes on undirected networks? (2)

Idea: Reduce it to disjoint paths problem!

For a single waypoint, can even compute *shortest route (walk)*!
 Recall: there is a randomized polytime algorithm for 2 disjoint paths

Step 1: replace weights with parallel links



Step 2: compute 2 disjoint paths (A,W) and (W,B)



## What about waypoint routes on undirected networks? (2)

Idea: Reduce it to disjoint paths problem!

For a single waypoint, can even compute *shortest route (walk)*!
 Recall: there is a randomized polytime algorithm for 2 disjoint paths

Good news: For a single waypoint, *shortest* paths can be computed even *faster*!

**Step 2:** compute 2 disjoint paths (A,W) and (W,B)



Suurballe's algorithm: finds two (edge-)disjoint shortest paths between same endpoints:



Suurballe's algorithm: finds two (edge-)disjoint shortest paths between same endpoints:



Step 1: replace capacities with parallel edges: paths will become edge-disjoint

$$w \xrightarrow{2} s \xrightarrow{2} t \implies w \equiv s \equiv t$$

□ Step 2: Reduction to Suurballe's algorithm:



□ Step 2: Reduction to Suurballe's algorithm:



**Step 2:** Reduction to Suurballe's algorithm:



**Step 2:** Reduction to **Suurballe's algorithm**:



#### Wait A Moment...!?

Can we not use Suurballe as well to solve 2 disjoint paths?



#### Wait A Moment...!?

**No!** Solves a much easier problem: 2 routes from  $\{s_1, s_2\}$  to  $\{t_1, t_2\}$ .



#### **Reduction** Waypoint Routing ⇒ Suurballe

#### **Reduction** 2 Disjoint Paths ⇒ Suurballe



#### **Remarks: Under the rug...**

Remark 1: Suurballe is actually for directed substrate graphs, so need gadget to transform problem in right form:



#### Remark 2: Suurballe: for vertex disjoint

□ Suurballe & Tarjan: edge disjoint

#### **Remark: There are complex requests...**

#### **IETF Draft**:



- ❑ Service chain for mobile operators
- Load-balancers are used to route (parts of) the traffic through cache

## Can be seen as a virtual network...

Credits: https://tools.ietf.org/html/draft-ietf-sfc-use-case-mobility-06

#### **Further Reading**

Walking Through Waypoints

Saeed Akhoondian Amiri, Klaus-Tycho Foerster, and Stefan Schmid.

13th Latin American Theoretical Informatics Symposium (LATIN), Buenos Aires, Argentina, April 2018.

Charting the Algorithmic Complexity of Waypoint Routing Saeed Akhoondian Amiri, Klaus-Tycho Foerster, Riko Jacob, and Stefan Schmid.

ACM SIGCOMM Computer Communication Review (CCR), 2018.

### Roadmap

#### □ A short recap:

- Networking preliminaries
- Opportunities and challenges of SDNs
- On predictable performance
- Algorithms I: Flexible routing
- Algorithms II: Flexible embedding and slicing
- □ Algorithms III: Flexible fast rerouting
- □ Algorithms IV: Theory of demand-aware networks
- Security aspects (if time permits)

## Roadmap

#### □ A short recap:

- Networking preliminaries
- Opportunities and challenges of SDNs
- On predictable performance
- □ Algorithms I: Flexible routing

#### Algorithms II: Flexible embedding and slicing

- □ Algorithms III: Flexible fast rerouting
- □ Algorithms IV: Theory of demand-aware networks
- Security aspects (if time permits)
- A fundamental resource allocation problem
- **2 dimensions** of flexibility:
  - Mapping of virtual nodes (to physical nodes)
  - Mapping of virtual links (to paths)



- A fundamental resource allocation problem
- 2 dimensions of flexibility:



- A fundamental resource allocation problem
- 2 dimensions of flexibility:

















## **Therefore: Mapping Virtual Links is Challenging**

Bad news: The Virtual Network Embedding Problem is hard even if endpoints are already mapped and given.

But maybe at least mapping nodes is simple?

Let's start simple again: assume paths are trivial, e.g., the physical network (host graph) is a line



Let's start simple again: assume paths are trivial, e.g., the physical network (host graph) is a line



Let's start simple again: assume paths are trivial, e.g., the physical network (host graph) is a line





## Therefore: VNEP is Hard "in Both Dimensions"!

- □ We have seen examples that:
  - mapping virtual links is hard (even if nodes are given)
  - mapping virtual nodes is hard (even if links are trivial)
- Remark: the VNEP can also be seen as a generalization of the Subgraph Isomorphism Problem (SIP)

Known? Why is SIP NP-hard?

## Therefore: VNEP is Hard "in Both Dimensions"!

#### We have seen examples that:

- mapping virtual links is hard (even if nodes are given)
- mapping virtual nodes is hard (even if links are trivial)
- Remark: the VNEP can also be seen as a generalization of the Subgraph Isomorphism Problem (SIP)
  - The SIP problem: Given two graphs G,H, determine whether G contains a subgraph that is isomorphic to H?
  - □ NP-hard: "does G contain an n-node cycle?" is a Hamilton cycle problem (each node visited exactly once), a solution to "does G contain a k-clique?" solves maximum clique problem, etc.



## Therefore: VNEP is Hard "in Both Dimensions"!

#### □ We have seen examples that:

- mapping virtual links is hard (even if nodes are given)
- mapping virtual nodes is hard (even if links are trivial)
- Remark: the VNEP can also be seen as a generalization of the Subgraph Isomorphism Problem (SIP)
  - The SIP problem: Given two graphs G,H, determine whether G contains a subgraph that is isomorphic to H?

So if SIP is hard, why is VNEP hard? **a problem** (each node visited and a problem, etc.)





### **NP-Hardness: From SIP to VNEP**

- ❑ Observe: VNEP is a generalization of SIP
- **For example:**



## **Remark: Graph Minors**

Note: It is sometimes possible to embed a guest graph G on a host graph H, even though G is *not* a minor of H:



Assume planar host graph H: K<sub>5</sub> and K<sub>3,3</sub> minor-free...



... but it is possible to embed non-planar guest graph G=K<sub>5</sub>!

?

#### □ Recall: Mixed Integer Program (MIP)

- Linear objective function (e.g., minimize embedding footprint)
- Linear constraints (e.g., do not violate capacity constraints)

One that provides

good relaxations!

- Recall: Mixed Integer Program (MIF)
  - Linear objective function (e.g., minimize
  - Linear constraints (e.g., do not violate capacity constraints

#### □ Recall: Mixed Integer Program (MIP)

- Linear objective function (e.g., minimize embedding footprint)
- Linear constraints (e.g., do not violate capacity constraints)



#### □ Recall: Mixed Integer Program (MIP)

- Linear objective function (e.g., minimize embedding footprint)
- Linear constraints (e.g., do not violate capacity constraints)



#### □ Recall: Mixed Integer Program (MIP)

- Linear objective function (e.g., minimize embedding footprint)
- Linear constraints (e.g., do not violate capacity constraints)



#### □ Recall: Mixed Integer Program (MIP)

- Linear objective function (e.g., minimize embedding footprint)
- Linear constraints (e.g., do not violate capacity constraints)



#### □ Recall: Mixed Integer Program (MIP)

- Linear objective function (e.g., minimize embedding footprint)
- Linear constraints (e.g., do not violate capacity constraints)

#### □ Solved, e.g., with branch-and-bound search tree

Usual trick: Relax! Solve LP (fast!), and if **relaxed solution** (more general!) **not better** then best solution so far: skip it!

Bottomline: If MIP provides «good relaxations», large parts of the search space can be pruned.

A typical MIP formulation:

- Introduce binary variables map(v,s) to map virtual nodes v to substrate node s
- Introduce flow variables for paths (say splittable for now)
- Ensure flow conservation: all flow entering a node must leave the node, unless it is the source or the destination





 $\sum_{u \to v} f_{uv} = \sum_{v \to w} f_{vw}$ 



What does this formula do and why is it correct?









## What will happen in this example?



### What will happen in this example?


#### What will happen in this example?



Minimal flow = 0: fulfills flow conservation! Relaxation useless: does not provide any lower bound or indication of good mapping!

### Remark: What about using randomized rounding?

An approximation approach which:

- ... computes a solution to the linear relaxation of the IP,
- ... decomposes this solution into convex combinations of elementary solutions, and
- Image: Image: Image: constraint of the selementary solutions based on their weight.

### What about using randomized rounding?

- Problem 1: relaxed solutions may not be very meaningful
  - see example for splittable flows before
- Problem 2: also for unsplittable flows, if using a standard Multi-Commodity Flow (MCF) formulation of VNEP, the integrality gap can be huge
  - Tree-like VNets are still ok
  - VNets with cycles: randomized rounding not applicable, since problem not decomposable

The linear solutions can be decomposed into convex combinations of valid mappings.

### **Non-Decomposability**

Relaxations of classic MCF formulation cannot be decomposed into convex combinations of valid mappings (so we need different formulations!)



### **Non-Decomposability**









Wait a minute! These problems need to be solved! And they often can, even with guarantees.

That's all Folks

### **Theory vs Practice:** In Practice There is Hope!

Guest graphs may not be general graphs, but e.g., virtual clusters: very simple and symmetric, used in context of batch processing

k VMs/compute-units/tasks/... 







Consider guest graph:





1. Place logical switch (try all options)



- 1. Place logical switch (try all options)
- 2. Extend network with artificial source s and sink t



- 1. Place logical switch (try all options)
- 2. Extend network with artificial source s and sink t
- 3. Add capacities (recall that b=1, so each virtual node
- needs one unit of capacity)



Then: Compute min-cost max flow of size n from s to t (e.g., successive shortest paths): due to capacity constraints at most size n.



Then: Compute min-cost max flow of size n from s to t (e.g., successive shortest paths): due to capacity constraints at most size n.

... and assign virtual nodes (and edges) accordingly!





In fact: this physical network is even *a tree*! For trees with servers at leaves, even simpler algorithms are possible. Ideas?









..., n} virtual nodes in left subtree...

..., n} virtual nodes in right subtree...





To compute cost of embedding x nodes in T, place y nodes on the left, x-y on the right subtree, and compute cost due to links across root.

# **Remark on Virtual Cluster Abstraction**

**Two interpretations:** 

- □ Logical switch at unique location
- □ Logical switch can be distributed («Hose model»)



# **Remark on Virtual Cluster Abstraction**

**Two interpretations:** 

- Logical switch at unique location
- □ Logical switch can be distributed («Hose model»)



### **Remark on Virtual Cluster Abstraction**

**Two interpretations:** 

- Logical switch at unique location
- Logical switch can be distributed («Hose model»)



But **embedding costs** can be different if we do not insist on placing the logical switch explicitly! **Not on trees** though, and **not in uncapacitated networks**: without routing restrictions, optimal routing paths form a tree (**SymG=SymT** a.k.a. **VPN Conjecture**).





Impossible: need at least 5 units of flow from/to node where *star center is mapped*. However, capacity of incident links is only 4.
























### All virtual links mapped to routes!



Can demand really be satisfied given link capacity of 2?!





## What about requests "with flexibilities"?

## **Requests can be more complex**



We have seen: Already non-trivial!

But what if requests allow for alternatives and different decompositions?



## **Requests can be more complex**



We have seen: Already non-trivial!

But what if requests allow for alternatives and different decompositions?



Known as PR (Processing and Routing) Graph: allows to model different choices and implementations!

























Su

This problem can be solved using mincost unsplittable multi-commodity flow (approximation) algorithms (e.g., randomized rounding).



This problem can be solved using mincost unsplittable multi-commodity flow (approximation) algorithms (e.g., randomized rounding).

Su

But note: cannot keep track of dependencies across stages (e.g., allocation on links or nodes): may yield oversubscription.

## What about requests which arrive over time?

**Primal and Dual** 



Fig. 1: (I) The primal covering LP. (II) The dual packing LP.

## $\overline{\nabla}$

#### Algorithm

Algorithm 1 The General Integral (all-or-nothing) Packing Online Algorithm (GIPO).

Upon the *j*th round:

- 1.  $f_{j,\ell} \leftarrow \operatorname{argmin}\{\gamma(j,\ell) : f_{j,\ell} \in \Delta_j\}$  (oracle procedure)
- 2. If  $\gamma(j, \ell) < b_j$  then, (accept)
  - (a)  $y_{j,\ell} \leftarrow 1$ .
  - (b) For each row e : If  $A_{e,(j,\ell)} \neq 0$  do

$$x_{e} \leftarrow x_{e} \cdot 2^{A_{e,(j,\ell)}/c_{e}} + \frac{1}{w(j,\ell)} \cdot (2^{A_{e,(j,\ell)}/c_{e}} - 1)$$

(c)  $z_j \leftarrow b_j - \gamma(j, \ell)$ . 3. Else, (reject)

- b. Else, (leject)
  - (a)  $z_j \leftarrow 0$ .

**Primal and Dual** 



**Primal and Dual** 

(a)  $z_i \leftarrow 0$ .

 $\min Z_j^T \cdot \mathbf{1} + X^T \cdot C \ s.t.$  $\max B_j^T \cdot Y_j \ s.t.$  $A_j \cdot Y_j \le C$  $D_j \cdot Y_j \le \mathbf{1}$  $Y_j \ge \mathbf{0}$  $Z_j^T \cdot D_j + X^T \cdot A_j \ge B_j^T$  $X, Z_i \geq \mathbf{0}$ (II)s.t. constraints (I)Fig. 1: (I) The primal covering LP. (II) The dual packing LP. Algorithm Algorithm 1 The General Integral (all-or-nothing) Packing Online Algorithm (GIPO). Upon the *j*th round: 1.  $f_{j,\ell} \leftarrow \operatorname{argmin}\{\gamma(j,\ell) : f_{j,\ell} \in \Delta_j\}$  (oracle procedure) 2. If  $\gamma(j, \ell) < b_j$  then, (accept) (a)  $y_{i,\ell} \leftarrow 1$ . (b) For each row e : If  $A_{e,(j,\ell)} \neq 0$  do  $x_e \leftarrow x_e \cdot 2^{A_{e,(j,\ell)}/c_e} + \frac{1}{w(j,\ell)} \cdot (2^{A_{e,(j,\ell)}/c_e} - 1).$ (c)  $z_i \leftarrow b_i - \gamma(j, \ell)$ . 3. Else, (reject)

Buchbinder&Naor

**Primal and Dual** 



Fig. 1: (I) The primal covering LP. (II) The dual packing LP.

#### Algorithm

Algorithm 1 The General Integral (all-or-nothing) Packing Online Algorithm (GIPO).

Upon the jth round:

- 1.  $f_{j,\ell} \leftarrow \operatorname{argmin}\{\gamma(j,\ell) : f_{j,\ell} \in \Delta_j\}$  (oracle procedure)
- 2. If  $\gamma(j, \ell) < b_j$  then, (accept)
  - (a)  $y_{j,\ell} \leftarrow 1$ .
  - (b) For each row e : If  $A_{e,(j,\ell)} \neq 0$  do

$$x_e \leftarrow x_e \cdot 2^{A_{e,(j,\ell)}/c_e} + \frac{1}{w(j,\ell)} \cdot (2^{A_{e,(j,\ell)}/c_e} - 1).$$

(c)  $z_j \leftarrow b_j - \gamma(j, \ell)$ . 3. Else, (reject)

- - (a)  $z_j \leftarrow 0$ .

 $\begin{array}{|c|c|c|c|c|c|}\hline \min Z_j^T \cdot \mathbf{1} + X^T \cdot C & s.t. \\ Z_j^T \cdot D_j + X^T \cdot A_j \ge B_j^T \\ & X, Z_j \ge \mathbf{0} \\ & (I) \end{array} \qquad \begin{array}{|c|c|c|c|c|} \max B_j^T \cdot Y_j & s.t. \\ A_j \cdot Y_j \le C \\ D_j \cdot Y_j \le \mathbf{1} \\ & Y_j \ge \mathbf{0} \\ & (II) \end{array}$ 

Fig. 1: (I) The primal covering LP. (II) The dual packing LP.

# $\bigtriangledown$

#### Algorithm

**Primal and Dual** 

Algorithm 1 The General Integral (all-or-nothing) Packing Online Algorithm (GIPO).

Upon the jth round:

1. 
$$f_{j,\ell} \leftarrow \operatorname{argmin}\{\gamma(j,\ell) : f_{j,\ell} \in \Delta_j\}$$
 (oracle procedure)  
2. If  $\gamma(j,\ell) < b_j$  then, (accept)  
(a)  $y_{j,\ell} \leftarrow 1$ .  
(b) For each row  $e$  : If  $A_{e,(j,\ell)} \neq 0$  do  
 $x_e \leftarrow x_e \cdot 2^{A_{e,(j,\ell)}/c_e} + \frac{1}{w(j,\ell)} \cdot (2^{A_{e,(j,\ell)}/c_e} - 1)$ .  
(c)  $z_j \leftarrow b_j - \gamma(j,\ell)$ .  
3. Else, (reject)  
(a)  $z_j \leftarrow 0$ .

**Primal and Dual** 



Fig. 1: (I) The primal covering LP. (II) The dual packing LP.

## $\nabla$

#### Algorithm

Algorithm 1 The General Integral (all-or-nothing) Packing Online Algorithm (GIPO).

Upon the jth round:

(a)  $z_i \leftarrow 0$ .

1. 
$$f_{j,\ell} \leftarrow \operatorname{argmin}\{\gamma(j,\ell) : f_{j,\ell} \in \Lambda_j\}$$
 (oracle procedure)  
2. If  $\gamma(j,\ell) < b_j$  then, (accept)  
(a)  $g_{j,\ell} \leftarrow 1$ .  
(b) For each row  $e$  : If  $A_{e,(j,\ell)} \neq 0$  do  
 $x_e \leftarrow x_e \cdot 2^{A_{e,(j,\ell)}/c_e} + \frac{1}{w(j,\ell)} \cdot (2^{A_{e,(j,\ell)}/c_e} - 1)$ .  
(c)  $z_j \leftarrow b_j - \gamma(j,\ell)$ .  
3. Else, (reject)

Primal and Dual



Fig. 1: (I) The primal covering LP. (II) The dual packing LP.

## $\overline{\nabla}$

#### Algorithm

Algorithm 1 The General Integral (all-or-nothing) Packing Online Algorithm (GIPO).

Upon the jth round:

1. 
$$f_{j,\ell} \leftarrow \operatorname{argmin}\{\gamma(j,\ell) : f_{j,\ell} \in \Delta_j\}$$
 (oracle procedure)

2. If 
$$\gamma(j, \ell) < b_j$$
 then, (accept)

(a) 
$$y_{j,\ell} \leftarrow 1$$
.

(b) For each row 
$$e : \text{If } A_{e,(j,\ell)} \neq 0 \text{ do}$$

$$x_e \leftarrow x_e \cdot 2^{A_{e,(j,\ell)}/c_e} + \frac{1}{w(j,\ell)} \cdot (2^{A_{e,(j,\ell)}/c_e} - 1).$$

(c) 
$$z_j \leftarrow b_j - \gamma(j, \ell)$$
.

3. Else, (reject)

(a) 
$$z_j \leftarrow 0$$
.

If cheap: accept and - update primal variables (always feasible solution)
### **The Buchbinder-Naor Approach**

**Primal and Dual** 



Fig. 1: (I) The primal covering LP. (II) The dual packing LP.

# $\overline{\nabla}$

#### Algorithm

Algorithm 1 The General Integral (all-or-nothing) Packing Online Algorithm (GIPO).

Upon the jth round:

1. 
$$f_{j,\ell} \leftarrow \operatorname{argmin}\{\gamma(j,\ell) : f_{j,\ell} \in \Delta_j\}$$
 (oracle procedure)

2. If  $\gamma(j, \ell) < b_j$  then, (accept)

(a)  $y_{j,\ell} \leftarrow 1$ .

(b) For each row 
$$e$$
 : If  $A_{e,(j,\ell)} \neq 0$  do

$$x_{e} \leftarrow x_{e} \cdot 2^{A_{e,(j,\ell)}/c_{e}} + \frac{1}{w(j,\ell)} \cdot (2^{A_{e,(j,\ell)}/c_{e}} - 1)$$



## **The Buchbinder-Naor Approach**

How to do this

optimally?!

Fig. 1: (I) The primal covering LP. (II) The dual packing LP.

# $\bigtriangledown$

### Algorithm

**Primal and Dual** 

Algorithm 1 The General Integral (all-or-nothing) Packing Online Algorithm (GIPO).

Upon the jth round:

1. 
$$f_{j,\ell} \leftarrow \operatorname{argmin}\{\gamma(j,\ell) : f_{j,\ell} \in \Delta_j\}$$
 (oracle procedure)

2. If 
$$\gamma(j, \ell) < b_j$$
 then, (accept)

(a) 
$$y_{j,\ell} \leftarrow 1$$
.

(b) For each row 
$$e$$
 : If  $A_{e,(j,\ell)} \neq 0$  do

$$x_{e} \leftarrow x_{e} \cdot 2^{A_{e,(j,\ell)}/c_{e}} + \frac{1}{w(j,\ell)} \cdot (2^{A_{e,(j,\ell)}/c_{e}})$$

(c)  $z_j \leftarrow b_j - \gamma(j, \ell)$ . 3. Else, (reject)

(a)  $z_j \leftarrow 0$ .

## **Further Reading**

Virtual Network Embedding Approximations: Leveraging Randomized Rounding Matthias Rost and Stefan Schmid. IEEE/ACM Transactions on Networking (**TON**), 2019. Parametrized Complexity of Virtual Network Embeddings: Dynamic & Linear Programming Approximations Matthias Rost, Elias Döhne, and Stefan Schmid. ACM SIGCOMM Computer Communication Review (**CCR**), January 2019. Charting the Complexity Landscape of Virtual Network Embeddings Matthias Rost and Stefan Schmid. **IFIP Networking**, Zurich, Switzerland, May 2018. Competitive and Deterministic Embeddings of Virtual Networks Guy Even, Moti Medina, Gregor Schaffrath, and Stefan Schmid. Journal Theoretical Computer Science (**TCS**), Elsevier, 2013.

- Networking preliminaries
- Opportunities and challenges of SDNs
- On predictable performance
- Algorithms I: Flexible routing
- Algorithms II: Flexible embedding and slicing
- □ Algorithms III: Flexible fast rerouting
- □ Algorithms IV: Theory of demand-aware networks
- Security aspects (if time permits)

- Networking preliminaries
- Opportunities and challenges of SDNs
- On predictable performance
- Algorithms I: Flexible routing
- Algorithms II: Flexible embedding and slicing
- Algorithms III: Flexible fast rerouting
- □ Algorithms IV: Theory of demand-aware networks
- Security aspects (if time permits)





Ctrl

Control

Programs ()

Control

Programs

The Crux: How to define conditional rules which have local failure knowledge only?

**Failover via** 

control plane

too slow.

0

OpenFlow etc. allow to preconfigure conditional failover rules: 1st line of defense!

Open problem: How many link failures can be tolerated in kconnected network without going through controller?



OpenFlow etc. allow to preconfigure conditional failover rules: 1st line of defense!

# Solution: Use Arborescences (Chiesa et al.)



### **Basic principle:**

- Route along fixed arborescence ("directed spanning tree") towards the destination d
- If packet hits a failed edge at vertex v, reroute along a different arborescence

The Crux: which arborescence to choose next? Influences resiliency!



## **Simple Example: Hamilton Cycle**

Chiesa et al.: if *k*-connected graph has *k* arc disjoint Hamilton Cycles, *k*-1 resilient routing can be constructed!



### **Example: 3-Resilient Routing Function for 2-dim Torus**



### **Example: 3-Resilient Routing Function for 2-dim Torus**









# **Example: 3-Resilient Routing Function for 2-dim Torus**

**Failover:** In order to reach destination d: go along 1<sup>st</sup> directed HC, if hit failure, reverse direction, if again failure switch to 2<sup>nd</sup> HC, if again failure reverse direction: no more failures possible!



**Arc-Disjoint Arborescences** 



- Networking preliminaries
- Opportunities and challenges of SDNs
- On predictable performance
- Algorithms I: Flexible routing
- Algorithms II: Flexible embedding and slicing
- □ Algorithms III: Flexible fast rerouting
- □ Algorithms IV: Theory of demand-aware networks
- Security aspects (if time permits)

- Networking preliminaries
- Opportunities and challenges of SDNs
- On predictable performance
- Algorithms I: Flexible routing
- Algorithms II: Flexible embedding and slicing
- □ Algorithms III: Flexible fast rerouting
- □ Algorithms IV: Theory of demand-aware networks
  - Security aspects (if time permits)

## **Extra slides!**

- Networking preliminaries
- Opportunities and challenges of SDNs
- On predictable performance
- Algorithms I: Flexible routing
- Algorithms II: Flexible embedding and slicing
- Algorithms III: Flexible fast rerouting
- □ Algorithms IV: Toward self-adjusting networks
  - Security aspects (if time permits)

- Networking preliminaries
- Opportunities and challenges of SDNs
- On predictable performance
- Algorithms I: Flexible routing
- Algorithms II: Flexible embedding and slicing
- Algorithms III: Flexible fast rerouting
- □ Algorithms IV: Toward self-adjusting networks
  - **Security aspects (if time permits)**

## **Programmable and Virtualized Networks**



## **Virtual Switches**



Virtual switches reside in the **server's virtualization layer** (e.g., Xen's Dom0). Goal: provide connectivity and isolation.

# Increasing Complexity: # Parsed Protocols

Number of parsed high-level protocols constantly increases:



# Increasing Complexity: Introduction of middlebox functionality



Increasing workloads and advancements in network virtualization drive virtual switches to implement middlebox functions such as load-balancing, DPI, firewalls, etc.







RARP

**IGMP** 

Unified packet parsing allows parse more and more protocols efficiently: in a single pass!

# **Complexity: The Enemy of Security!**

- Data plane security not well-explored (in general, not only virtualized): most security research on control plane
  - Two conjectures:
    - 1. Virtual switches increase the attack surface.

2. Impact of attack larger than with traditional data planes.



# The Attack Surface: Closer...

Attack surface becomes closer:

- Packet parser typically integrated into the code base of virtual switch
- First component of the virtual switch to process network packets it receives from the network interface
  - May process attacker-controlled packets!



## The Attack Surface: ... More Complex ...

Ethernet **PBB** LLC **IPv6 EXT HDR** VLAN **TUNNEL-ID MPLS** IPv6 ND IPv4 **IPv6 EXT HDR** ICMPv4 **IPv6HOPOPTS** TCP **IPv6ROUTING** UDP **IPv6Fragment** ARP **IPv6DESTOPT** SCTP **IPv6ESP** IPv6 IPv6 AH ICMPv6 RARP **IGMP IPv6 ND** GRE LISP **VXLAN** 



## ... Elevated Priviledges and Collocation ...

Collocated (at least partially) with hypervisor's Dom0 kernel space, guest VMs, image management, block storage, identity management, ...




## ... Elevated Priviledges and Collocation ...

Collocated (at least partially) with hypervisor's Dom0 kernel space, guest VMs, image management, block storage, identity management, ...





I ... the controller itself.

## ... Centralization ...



I ... the controller itself.



1. Rent a VM in the cloud (cheap)



2. Send **malformed MPLS packet** to virtual switch (**unified parser** parses label stack packet **beyond the threshold**)



3. Stack buffer overflow in (unified) MPLS parsing code: enables remote code execution



4. Send malformed packet to server (virtual switch) where controller is located (use existing communication channel)



#### 5. Spread

## **A New Threat Model**

### I Limited skills required

- Use standard fuzzer to find crashes
- Construct malformed packet
- Build ROP chain
- Limited resources
  - I rent a VM in the cloud



No physical access needed

No need to be a state-level attacker to compromise the dataplane (and beyond)!

Similar problems in NFV: need even more complex parsing/processing. And are often built on top of OvS.

# Conclusion

#### Programmability and virtualization: opportunities but also challenges

- **E.g.**,: faster innovation, flexibilities in resource allocation, etc.
- But, e.g.: performance isolation needs to be ensured across all involved resources, resulting resource allocation problems hard
- Algorithmic techniques for flexible resource allocation: from waypoint routing over virtual network embedding to online admission control
- Vision of demand-aware and self-adjusting networks: depends on the structure of demand!
- ❑ Security: even more opportunities and challenges ☺

## References

ble	On The Impact of the Network Hypervisor on Virtual Network Performance
	Andreas Blenk, Arsany Basta, Wolfgang Kellerer, and Stefan Schmid.
tal	IFIP Networking, Warsaw, Poland, May 2019.
for	Wowneint Pouting in Creatial Notworks
Pre	Waypoint Routing in Special Networks
	Saeed Aknoondian Amiri, Klaus-Tycho Foerster, Riko Jacob, Manmoud Parnam, and Stefan Schmid.
	IFIP Networking, Zurich, Switzerland, May 2018.
	Walking Through Waypoints
	Saeed Akhoondian Amiri, Klaus-Tycho Foerster, and Stefan Schmid.
ng	13th Latin American Theoretical Informatics Symposium (LATIN), Buenos Aires, Argentina, April 2018.
utii	
k Way Roi	Charting the Algorithmic Complexity of Waypoint Routing
	Saeed Akhoondian Amiri, Klaus-Tycho Foerster, Riko Jacob, and Stefan Schmid.
	ACM SIGCOMM Computer Communication Review ( <b>CCR</b> ), 2018.
	Virtual Network Embedding Approximations: Leveraging Randomized Rounding
	Matthias Rost and Stefan Schmid.
	IEEE/ACM Transactions on Networking ( <b>TON</b> ), 2019.
	······
Virtual Networ Embedding	Parametrized Complexity of Virtual Network Embeddings: Dynamic & Linear Programming
	Approximations
	Matthias Rost, Elias Döhne, and Stefan Schmid.
	ACM SIGCOMM Computer Communication Review (CCR), January 2019.
	Charting the Complexity Landscape of Virtual Network Embeddings (Best Paper Award)
	Matthias Rost and Stefan Schmid.
	IFIP Networking, Zurich, Switzerland, May 2018.
e	Competitive and Deterministic Embeddings of Virtual Networks
ast	Guy Even, Moti Medina, Gregor Schaffrath, and Stefan Schmid.
ai n	Journal Theoretical Computer Science ( <b>TCS</b> ), Elsevier, 2013.
	Improved Fast Rerouting Using Postprocessing
	Klaus-Tycho Foerster, Andrzej Kamisinski, Yvonne-Anne Pignolet, Stefan Schmid, and Gilles Tredan.
μ υ	38th International Symposium on Reliable Distributed Systems (SRDS), Lyon, France, October 2019.
an c	
a Pl	MTS: Bringing Multi-Tenancy to Virtual Switches
ecu	Kashyap Thimmaraju, Saad Hermak, Gabor Retvari, and Stefan Schmid.
ΔŇ	USENIX Annual Technical Conference (ATC), Renton, Washington, USA, July 2019.
	Taking Control of SDN-based Cloud Systems via the Data Plane
	Kashyap Thimmaraju, Bhargava Shastry, Tobias Fiebig, Felicitas Hetzelt, Jean-Pierre Seifert, Ania
	Feldmann, and Stefan Schmid.
	ACM Symposium on SDN Research (SOSR), Los Angeles, California, USA, March 2018.