



Hashing on broken assumptions

Lorenzo Saino (@lorenzosaino)
Fastly

Problem:

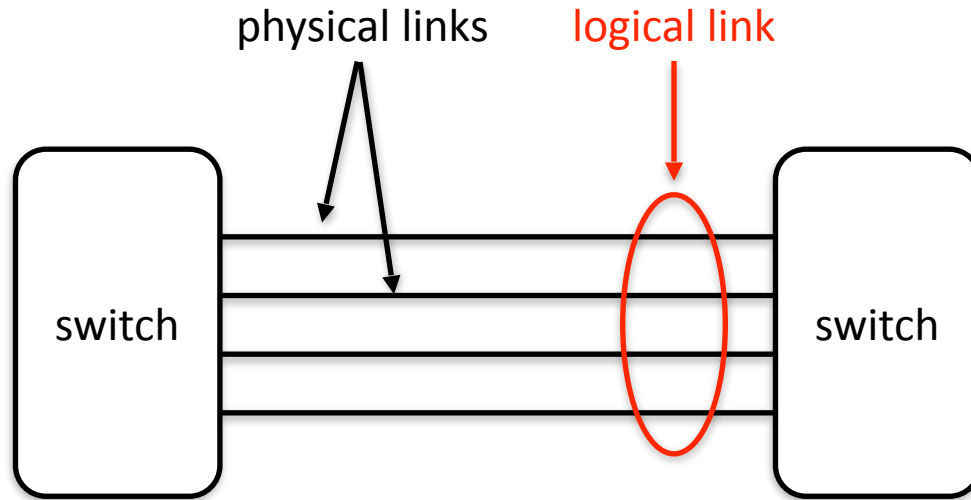
Spreading traffic across multiple links, paths, hosts

Solutions:

- Link Aggregation
- Equal Cost Multipath (ECMP)

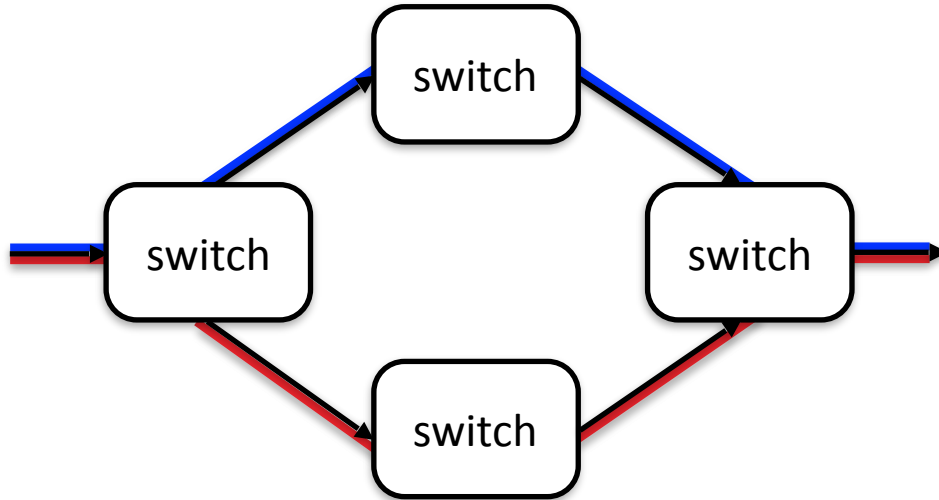
Link aggregation

Combine multiple physical links between network devices into one logical link

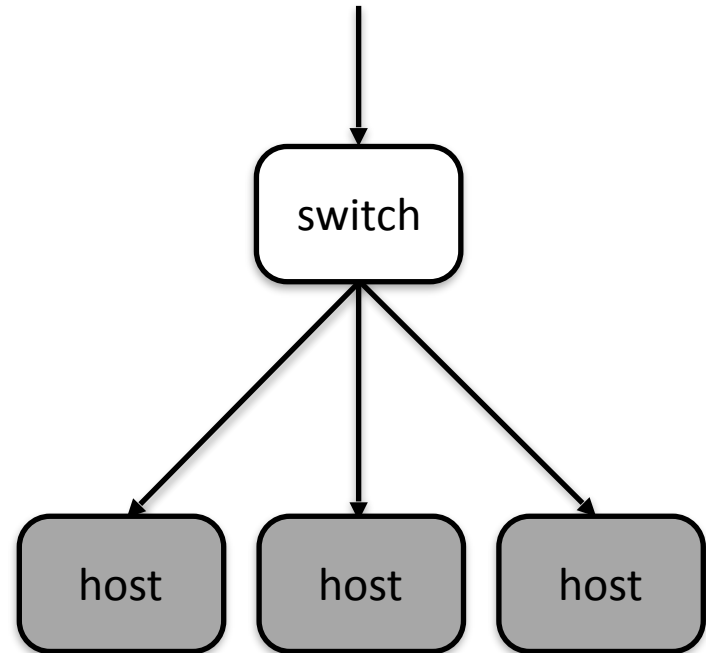


Equal Cost Multipath (ECMP)

Balance traffic across paths



Balance traffic across hosts



Requirements

Load balance

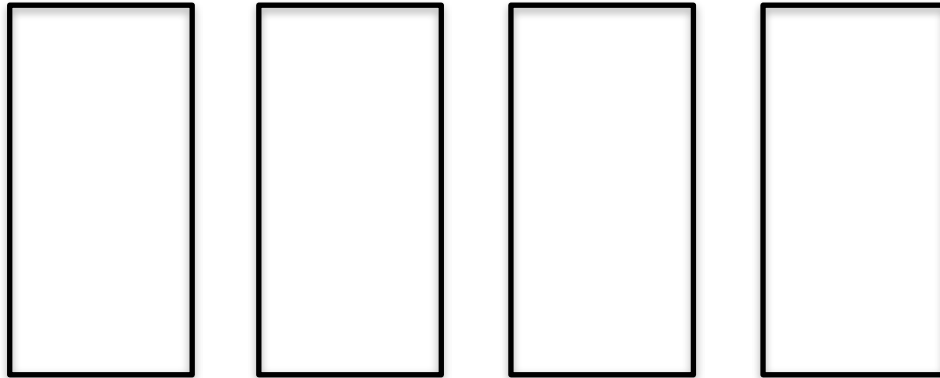
Traffic must be uniformly spread across next-hops

Stateless-but-sticky path pinning

All packets of a flow must take the same path

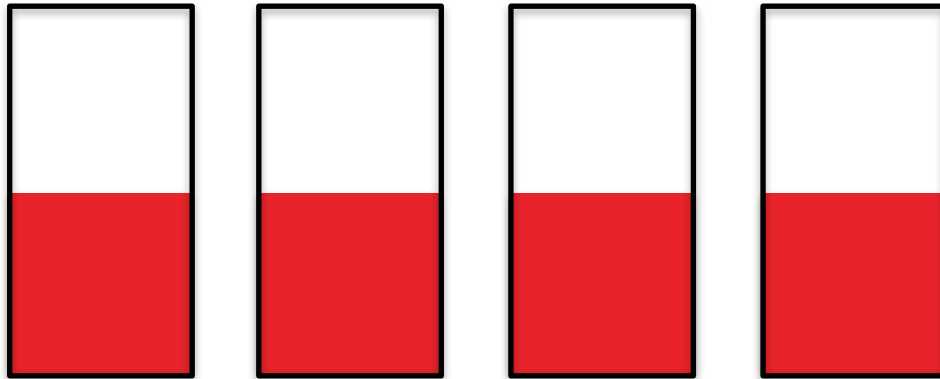
Load imbalance

Load imbalance reduces system capacity



Load imbalance

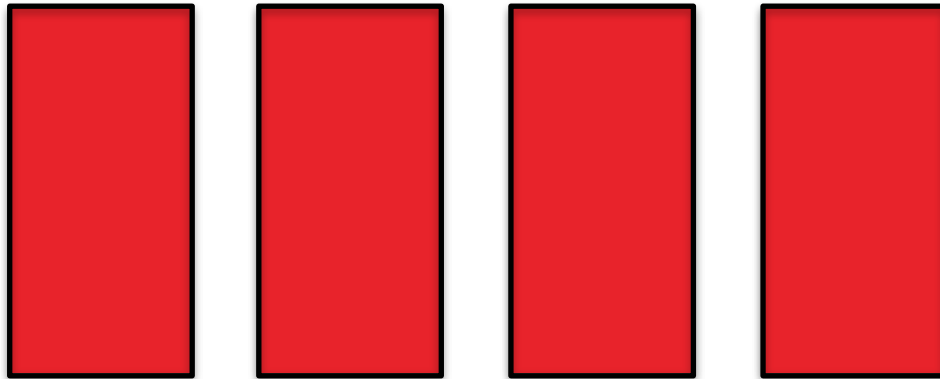
Load imbalance reduces system capacity



Perfect load balance

Load imbalance

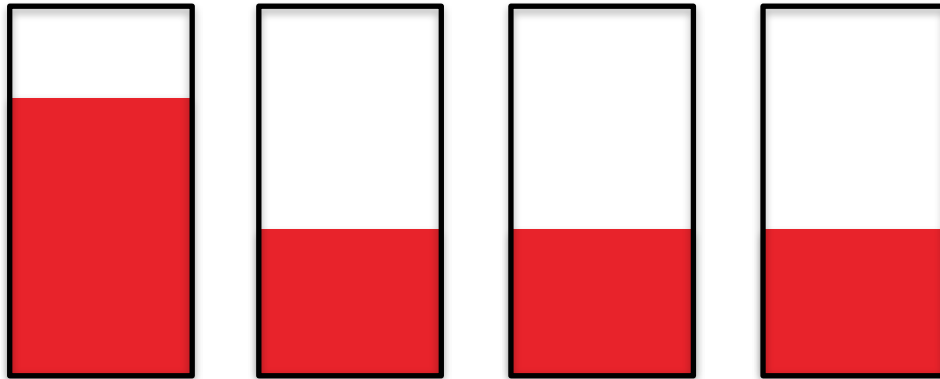
Load imbalance reduces system capacity



All resources fully utilized

Load imbalance

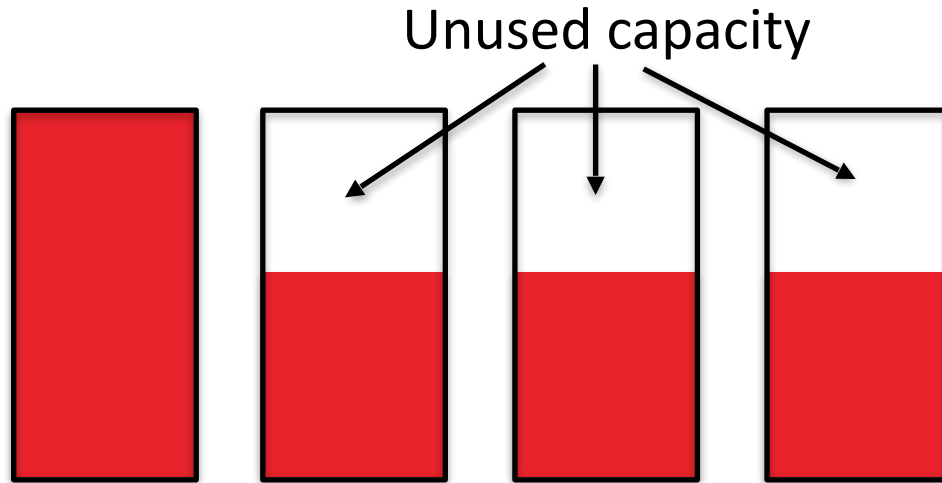
Load imbalance reduces system capacity



Load imbalance

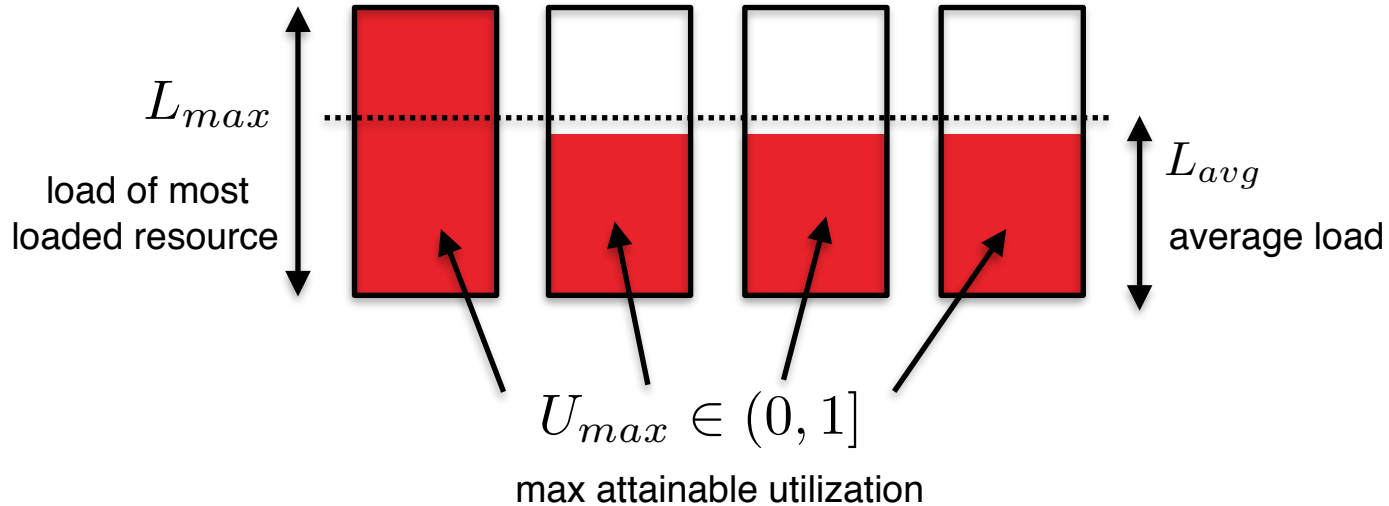
Load imbalance

Load imbalance reduces system capacity



Cannot take any additional load

Quantifying impact of load imbalance



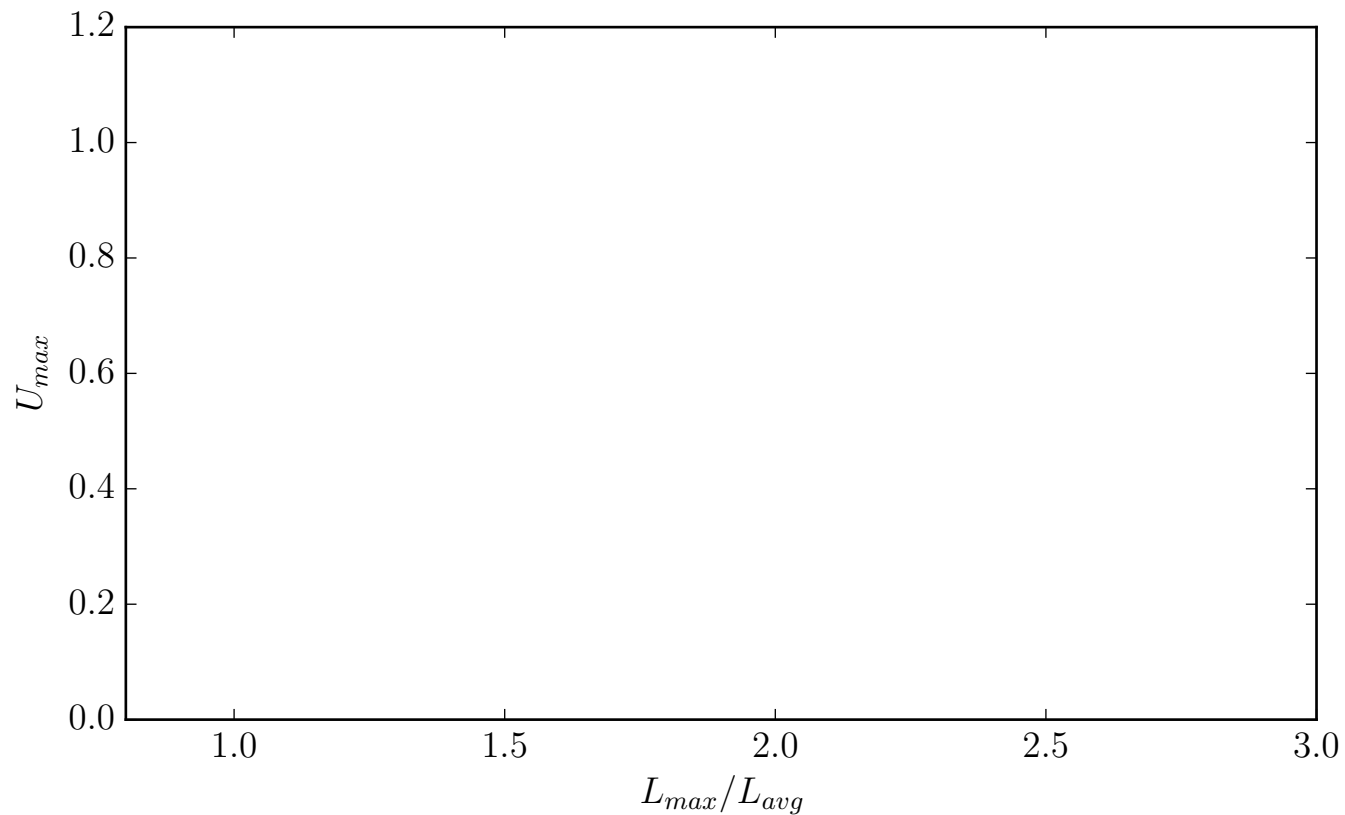
Load imbalance:

$$\frac{L_{max}}{L_{avg}} = [1, +\infty)$$

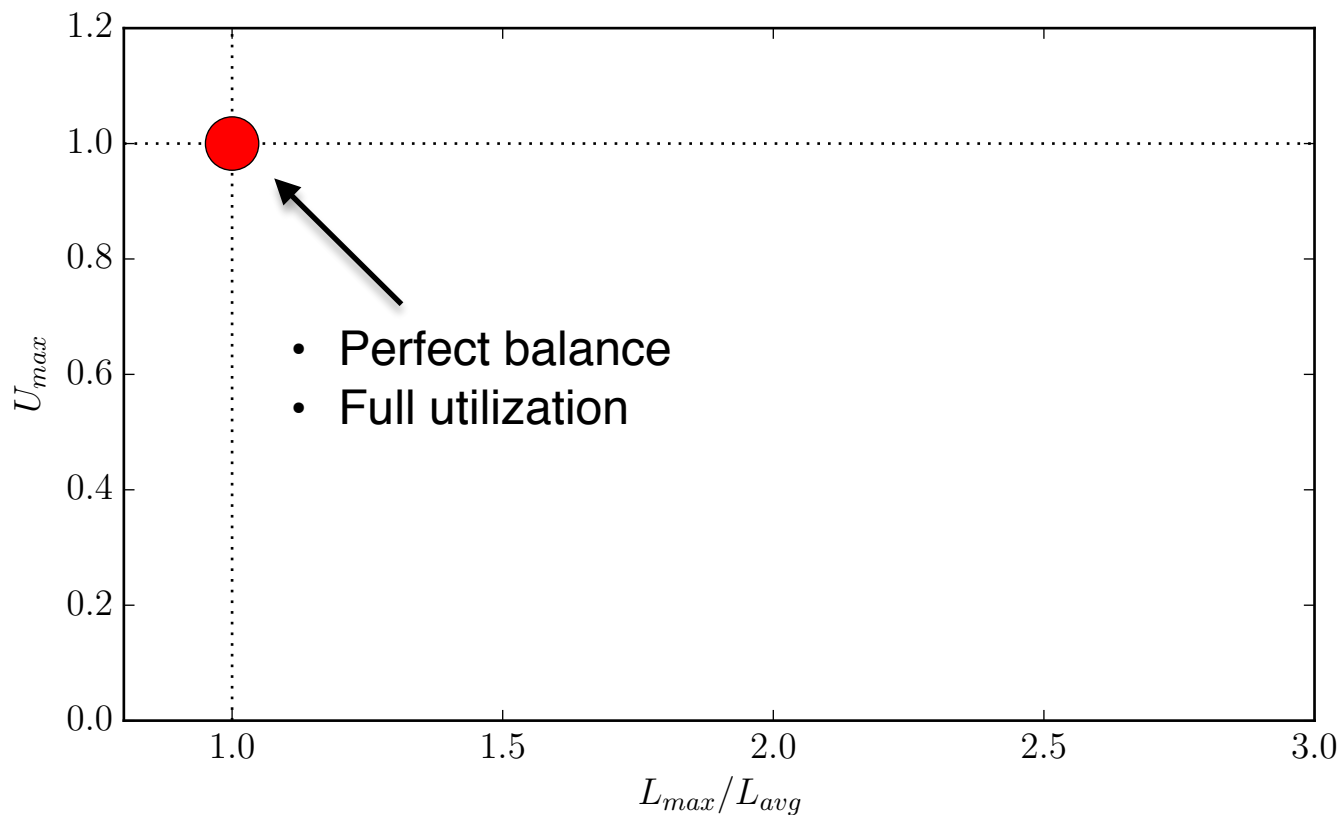
Max attainable utilization:

$$U_{max} = \left(\frac{L_{max}}{L_{avg}} \right)^{-1} = \frac{L_{avg}}{L_{max}}$$

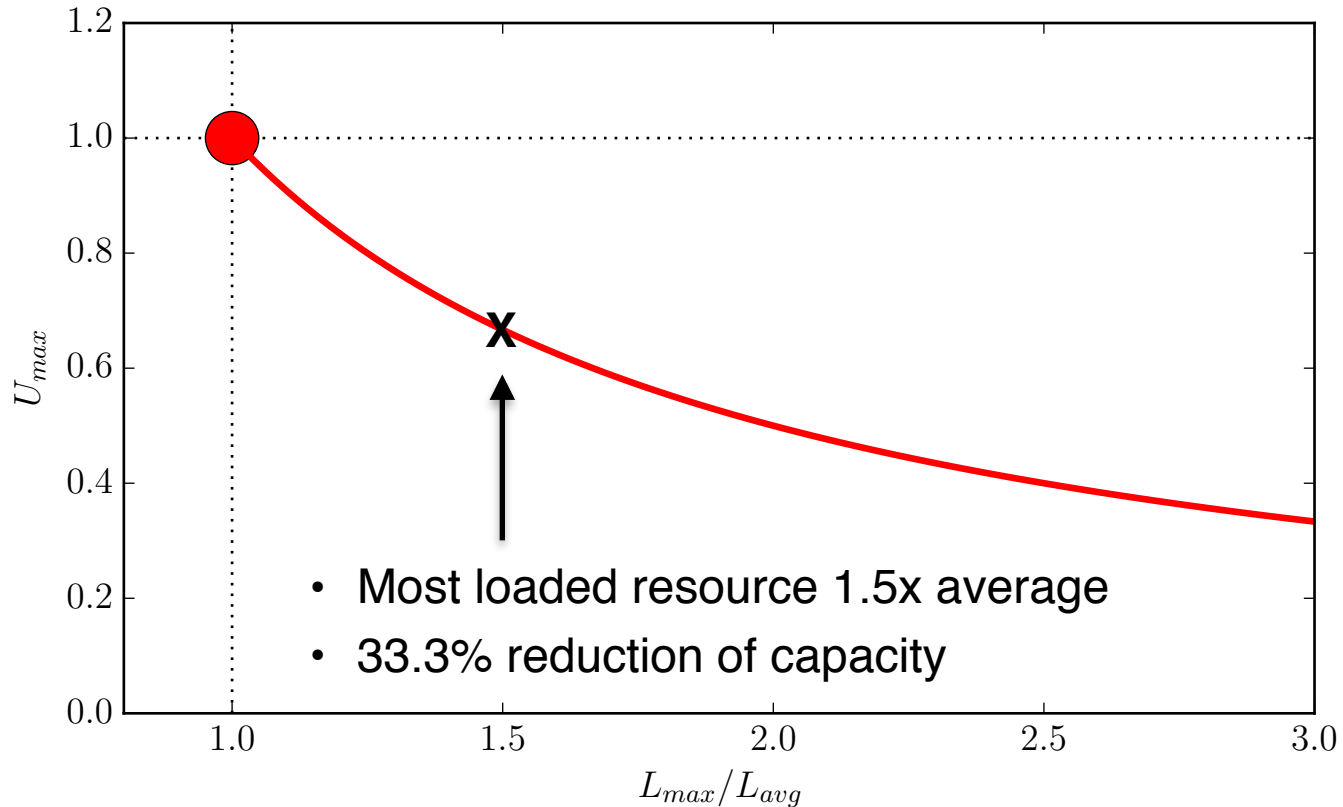
Quantifying impact of load imbalance



Quantifying impact of load imbalance



Quantifying impact of load imbalance



What happens without path pinning?

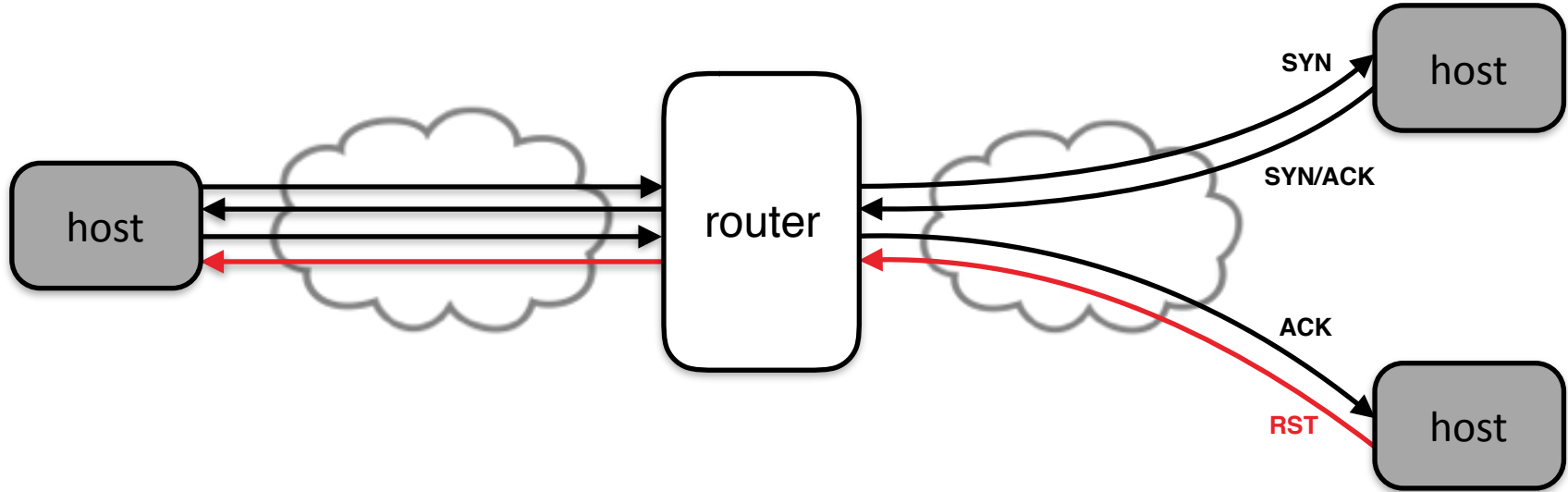
Same endpoints, different paths:

- Out-of-order packets
- Frequent drops of TCP congestion window (CWND)
- Poor throughput performance

Different endpoints:

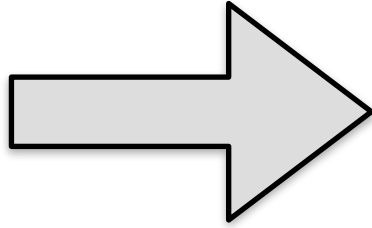
- TCP resets

TCP resets



Requirements:

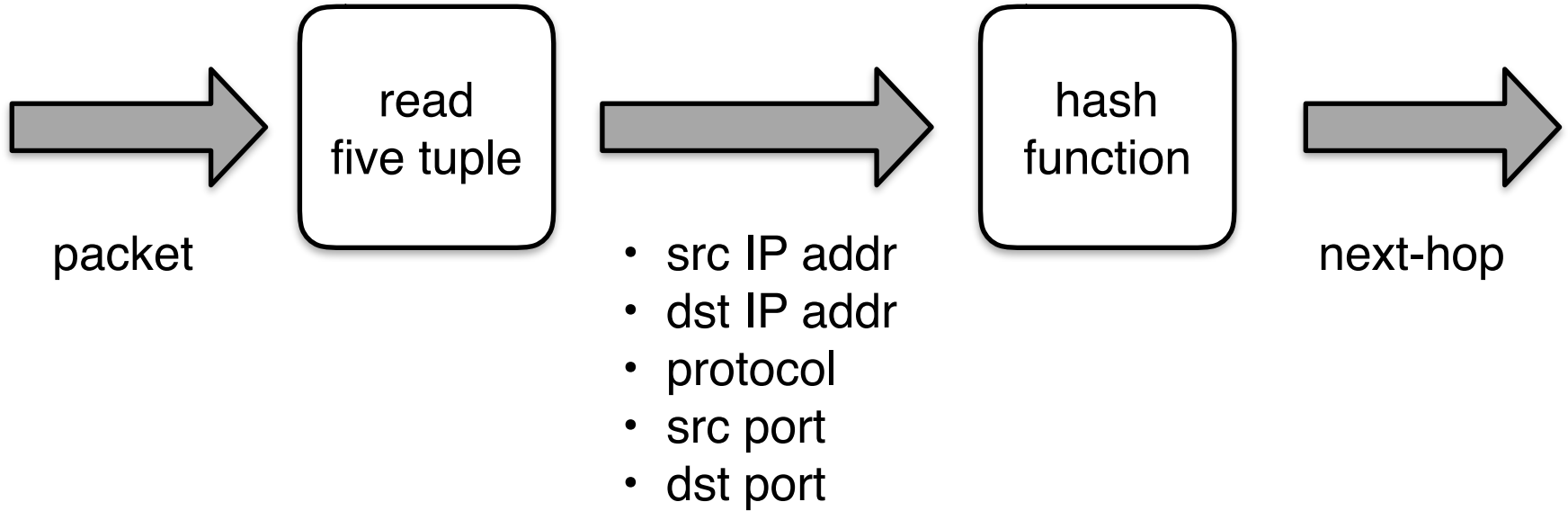
- Load balance
- Path pinning



Solution:

Flow-level hashing

Flow-level hashing



Assumptions

Load balance

Hashing uniformly spread traffic across next-hops

Path pinning

Hashing pins packets of a flow to the same path

Do these assumptions hold?

Assumptions

Load balance

Hashing uniformly spread traffic across next-hops

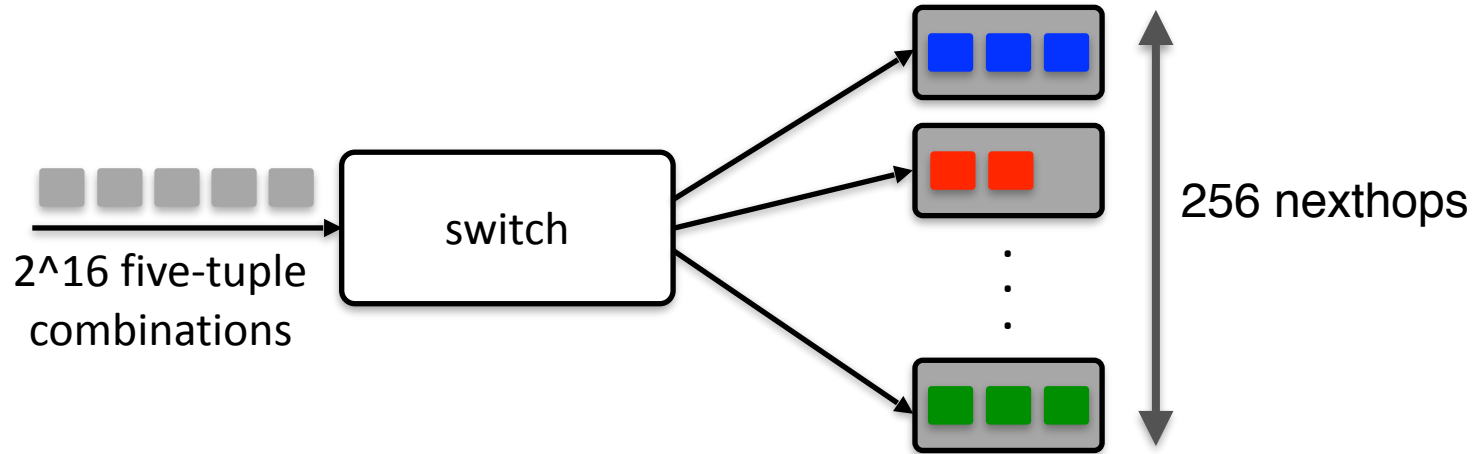
Path pinning

Hashing pins packets of a flow to the same path

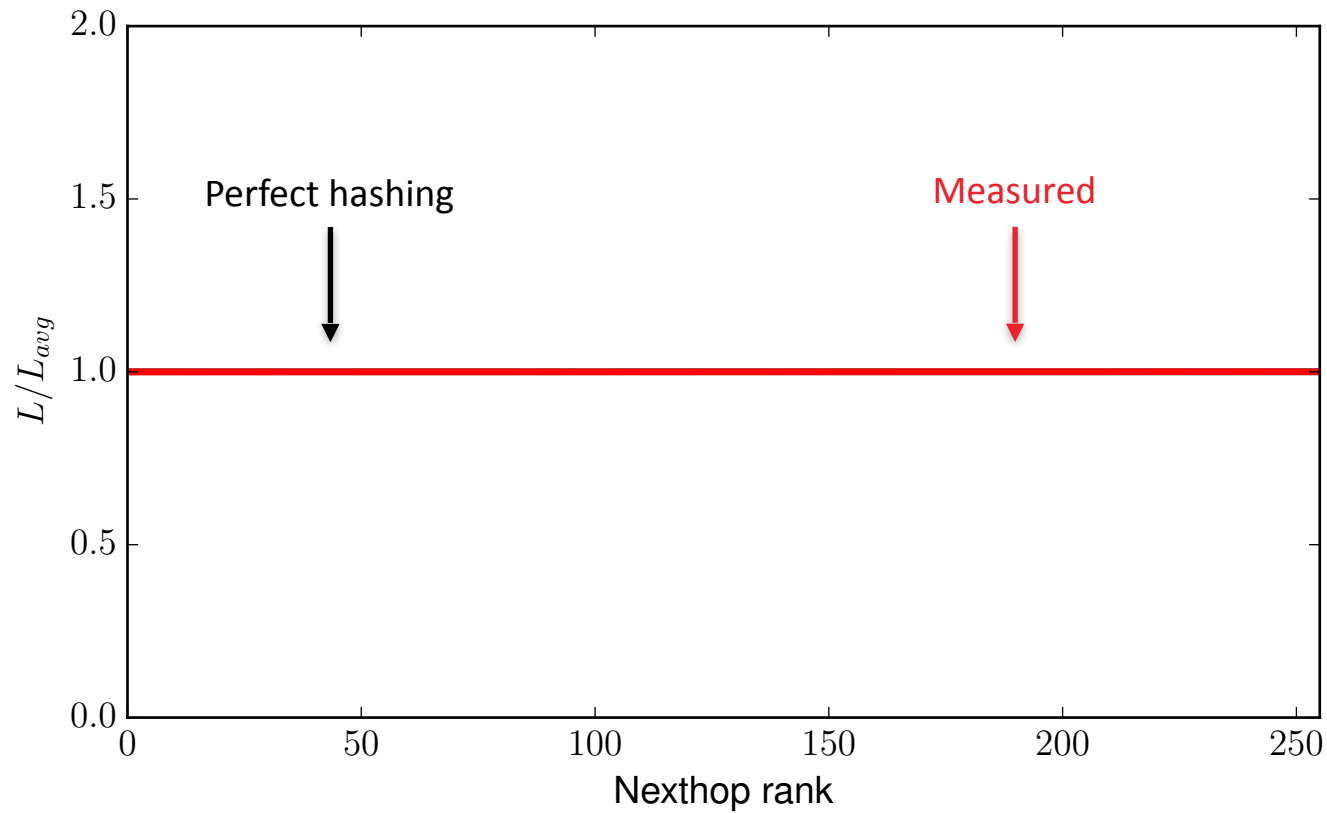
Hashing quality

Two switch models:

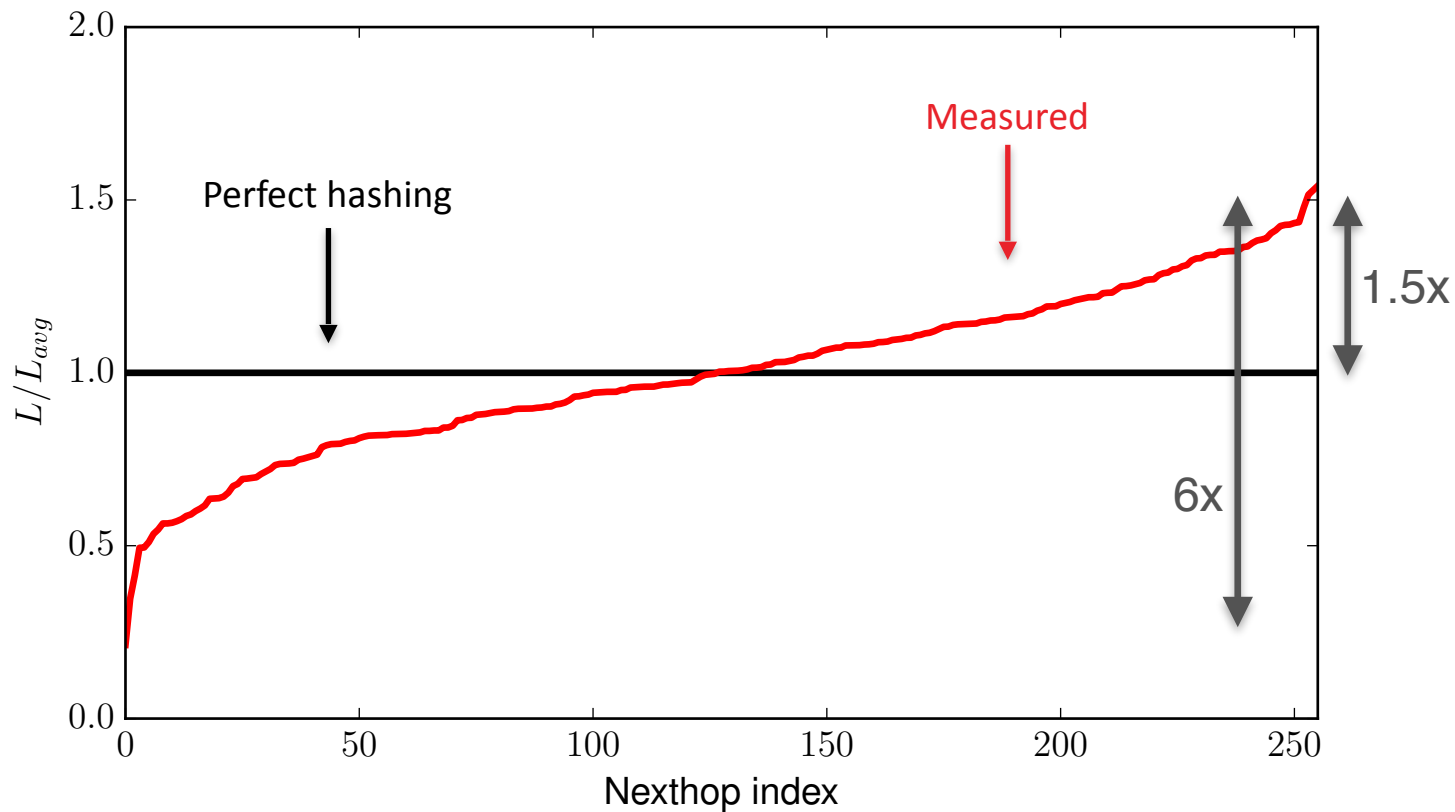
- Switch A
- Switch B



Switch A



Switch B



Switch B

Only a subset of next-hops are actually supported

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
	18		20		22		24		26		28		30		32
	34		36		38		40				44				48
			52				56				60				64
							72								80
							88								96
							104								112
							120	X	X	X	X	X	X		128
															...

6 next-hops don't get any traffic

Assumptions

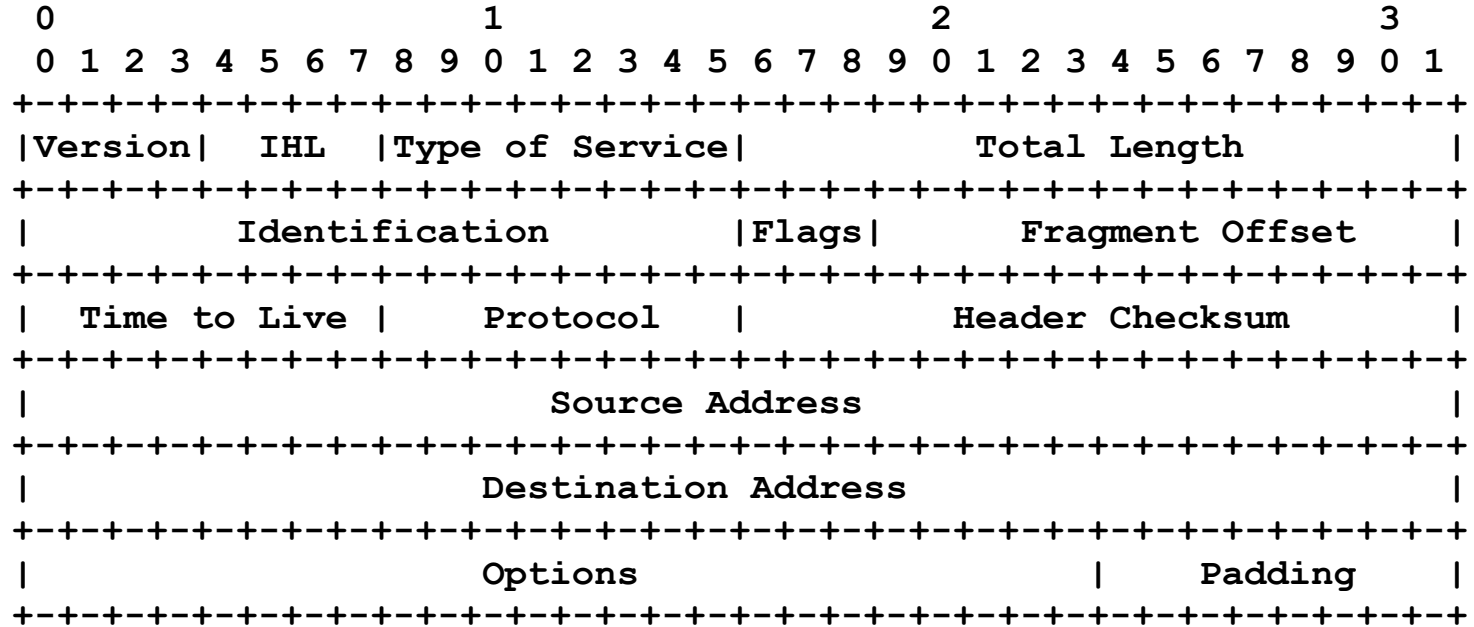
Load balance

Hashing uniformly spread traffic across next-hops

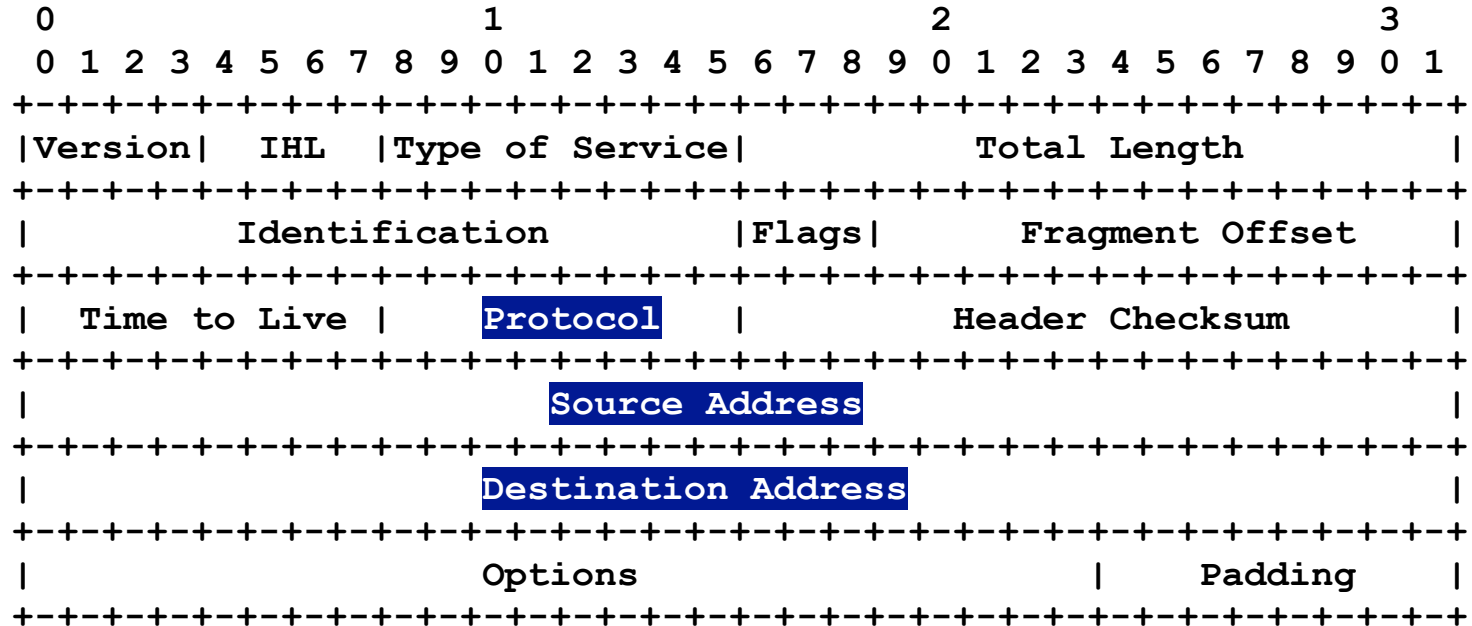
Path pinning

Hashing pins packets of a flow to the same path

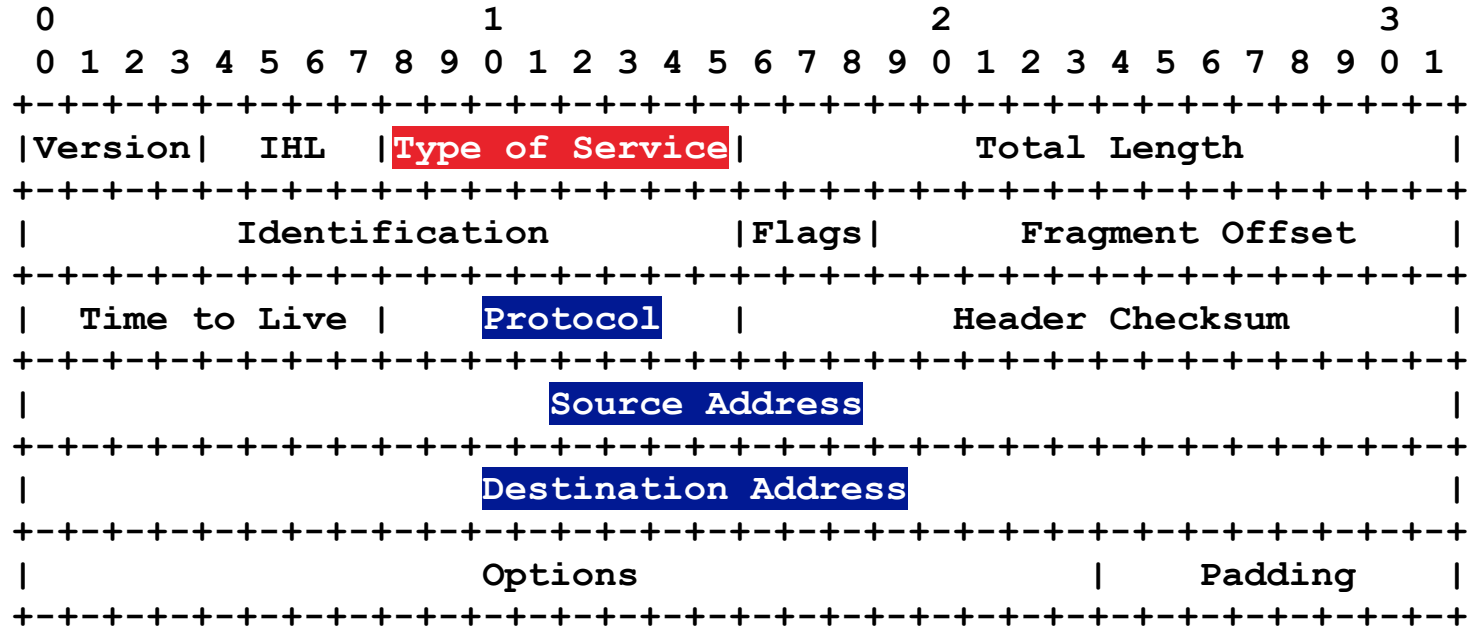
Hashing on IPv4 TOS field



Hashing on IPv4 TOS field



Hashing on IPv4 TOS field



Hashing on IPv4 TOS field

RFC 1812 - Requirements for IP Version 4 Routers

explicitly permits to involve the second-to-last bit of the TOS/DS octet in routing decisions

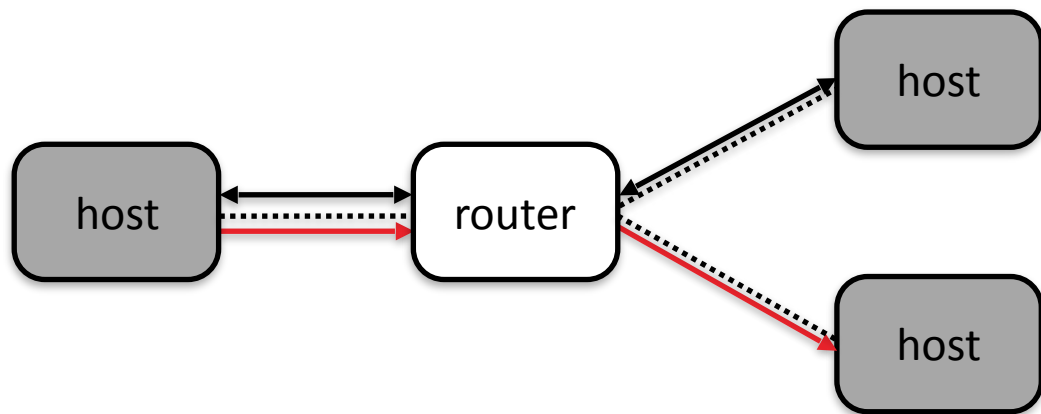
RFC 2474 - Definition of the Differentiated Services Field

deprecates the IPv4 Type of Service field redefines it as the Differentiated Services field

RFC 3168 - The Addition of Explicit Congestion Notification (ECN) to IP

reserves the last two bits of the DS octet for ECN

Hashing on IPv4 TOS field



TCP handshake:

- Hosts negotiate ECN support
- ECN-capable bits unset

Flow data:

- ECN-capable bits set

Scenario

- Hosts are ECN capable
- Router uses IPv4 TOS for hash computation (RFC 1812)

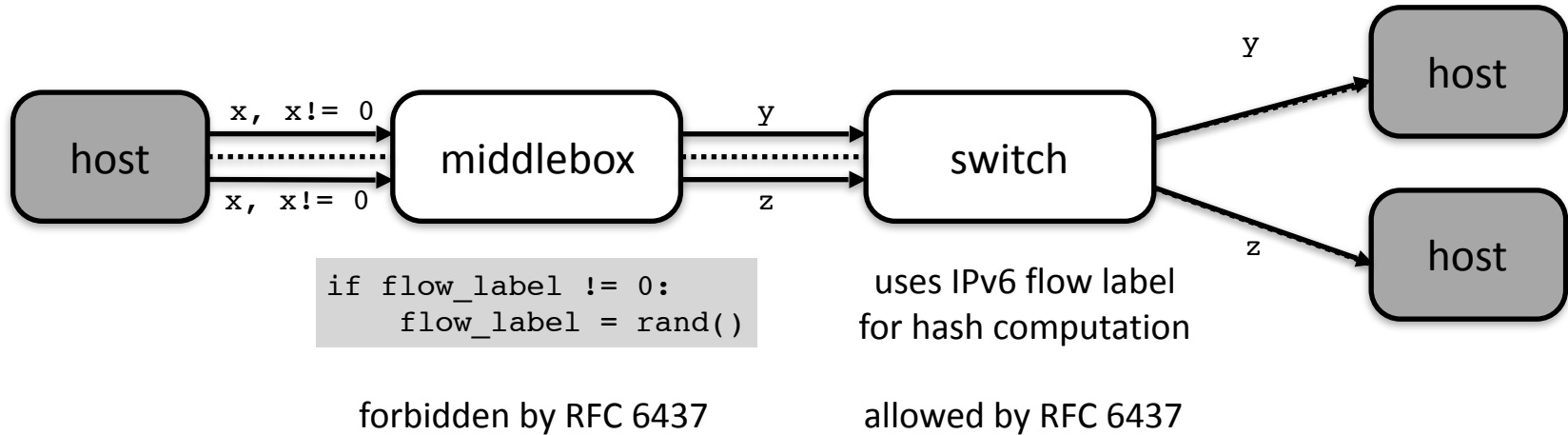


TCP handshake

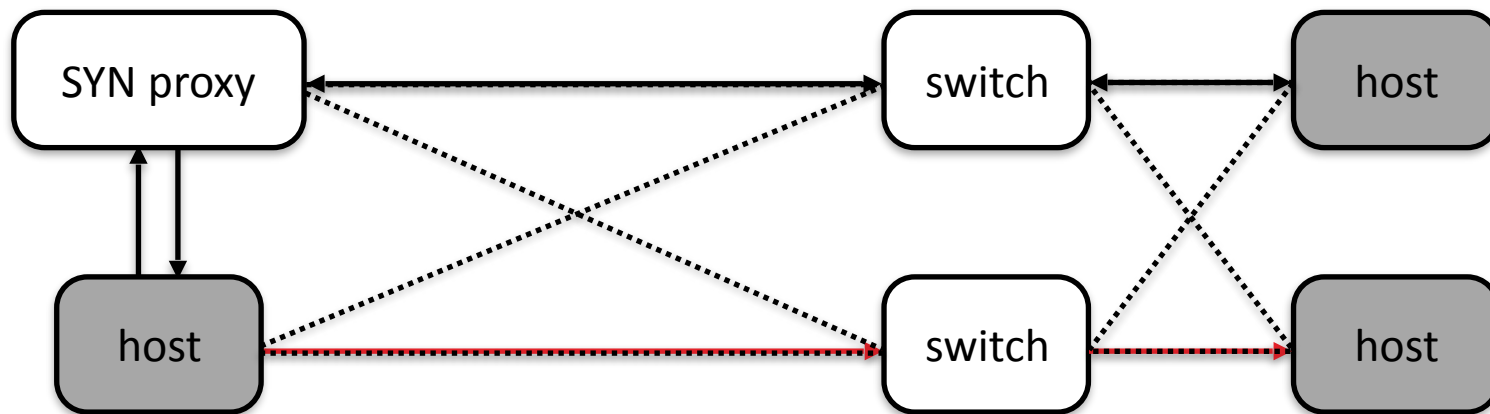


flow data

IPv6 flow label rewrite



SYN proxies



Switches:

- use ingress interface for hash computation, or
- use different hash function seeds



TCP handshake



flow data

Conclusions

Load balancing

There are devices that do not hash traffic uniformly

Path pinning

Hashing on fields other than five tuples breaks ECMP

- Ingress port
- IPv4 TOS
- IPv6 flow label

Recommendations

Operators:

- Ensure that your network devices hash flows uniformly or that could cost you money
- Disable additional inputs if you do not need extra entropy

Vendors:

- Disable hashing inputs other than five-tuple by default
- Make hash input fields configurable
- Make hash seed configurable

FIN