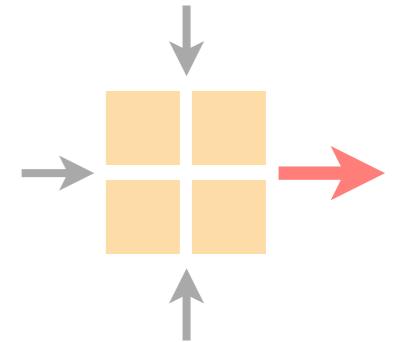


Advanced Topics in Communication Networks

Internet Routing and Forwarding



Laurent Vanbever
nsg.ee.ethz.ch

15 Sep 2020

Lecture starts at 14:15

In this lecture, you'll learn...

In this lecture, you'll learn
how to optimize the

- performance
- flexibility
- reliability

of large-scale network infrastructures

In this lecture, you'll learn
how to optimize the

- performance
- flexibility
- reliability

of large-scale network infrastructures

Why should *you* care, as a user?



Fortnite, Epic Games

Many gamers care about low latency!

Google

Fortnite Network Lag

X |  

About 917'000 results

Fortnite slow download speeds occur when you have low **Internet** connection speeds. ... That's because excellent **Internet** speed isn't all that's required to avoid **Fortnite lag** spikes. Low **latency** also comes into play. **Latency** refers to the time your computer takes to communicate with the **Fortnite** game server.

Dec 18, 2018

vortex.gg › cloud-gaming › fortnite-lag-how-to-fix ▾

[Ultimate Guide to Fix Fortnite Lag \(PS4, PC, Mac, and Xbox\)](#)



A photograph showing a person from behind, wearing a white t-shirt, a grey baseball cap, and black headphones. They are seated at a light-colored wooden desk, facing a computer monitor. The monitor displays a vibrant scene from a video game, possibly Fortnite, featuring a character in a village setting. On the desk, there is a black keyboard with blue backlit keys and a small blue figurine. The background is a plain, light-colored wall.

A new survey of 1,000 UK video gamers, which was commissioned by full fibre broadband ISP [Hyperoptic](#), has found that 44% of adult respondents named "*the internet lagging*" as the most infuriating aspect of online gaming, while 27% also complained about other players' slow online speeds.

The survey revealed that respondents are now spending an average of 10.5 hours online gaming a week (up from 8 hours pre-lockdown) and more than half said the

hobby has kept them entertained during this time (we assume the other half must have been either losing or playing a really boring game, since they aren't mentioned). Meanwhile 43% said they found playing online had helped them to stay connected with others.

How to reduce lag in PC games

By NICK GREENE

Last updated 26 Mar 2020

Actiontec

PRODUCTS ▾ SERVICES ▾ RESOURCES ▾ LEARN ▾ SUPPORT ▾

How to Reduce Latency or Lag in Gaming

7 Ways to Reduce Latency in Online Gaming

Sarah Pike
January 4, 2016

0
Comments

f | t | g

Source: Pixabay

Wh
you
fru
Int
my

Here's how to get lower ping for online gaming

Read this before spending hundreds on a new gaming router.

Ry Crist Aug. 28, 2019 5:00 a.m. PT

SHARE | 3

OUTFOX

About Us Forum Support Blog Download

LOG IN

START A FREE TRIAL

OUTSMART. OUTMANEUVER. OUTPERFORM.

THE ULTIMATE OPTIMIZED GAMING NETWORK

START A FREE TRIAL

GET HASTE NOW

HASTE

Features How Haste Works Supported Games Haste Pro Pricing

RECLAIM YOUR CONNECTION
DEFEAT PING AND STABILIZE YOUR CONNECTION

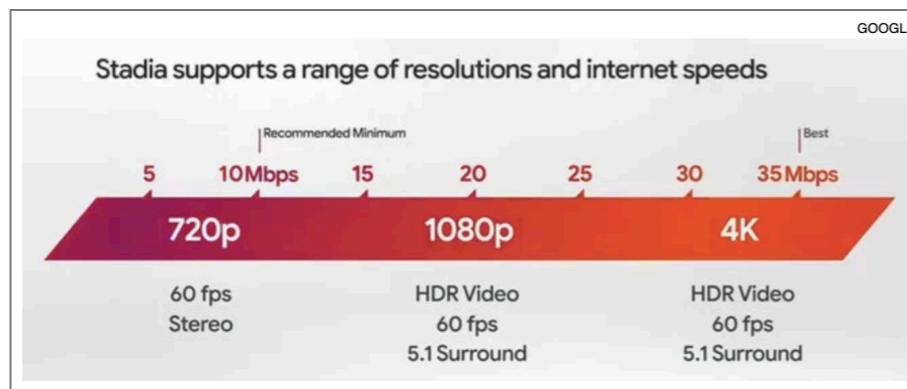
GET HASTE NOW

No Credit Card Required for Haste Free

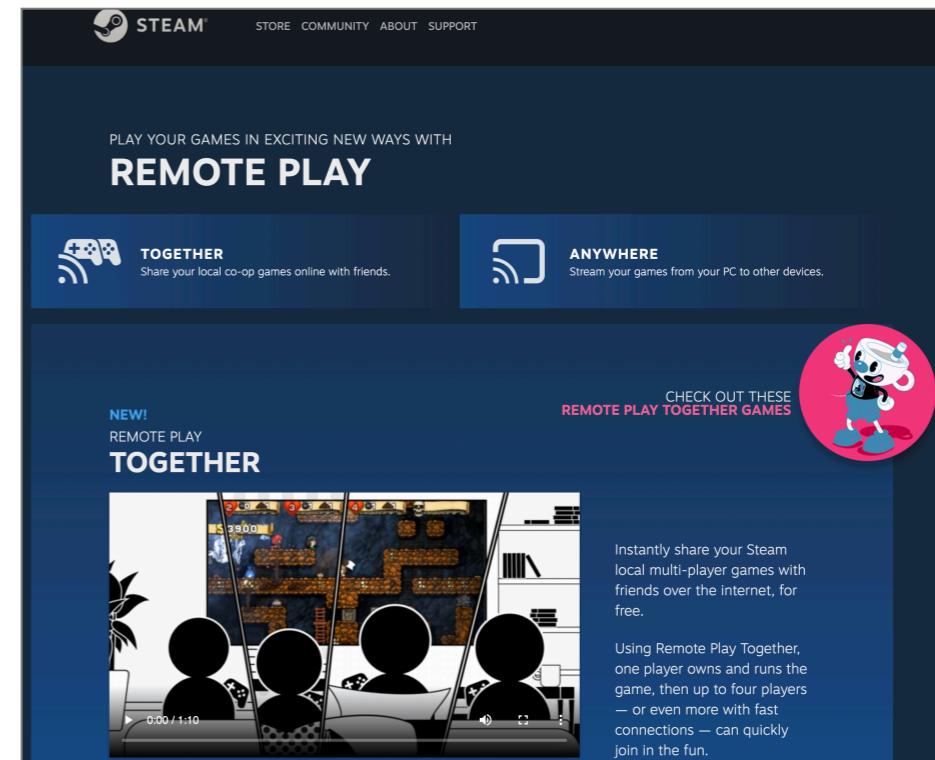
IMPROVES GAME CONNECTIONS

The problem of poor network performance will intensify as video games go streaming-based

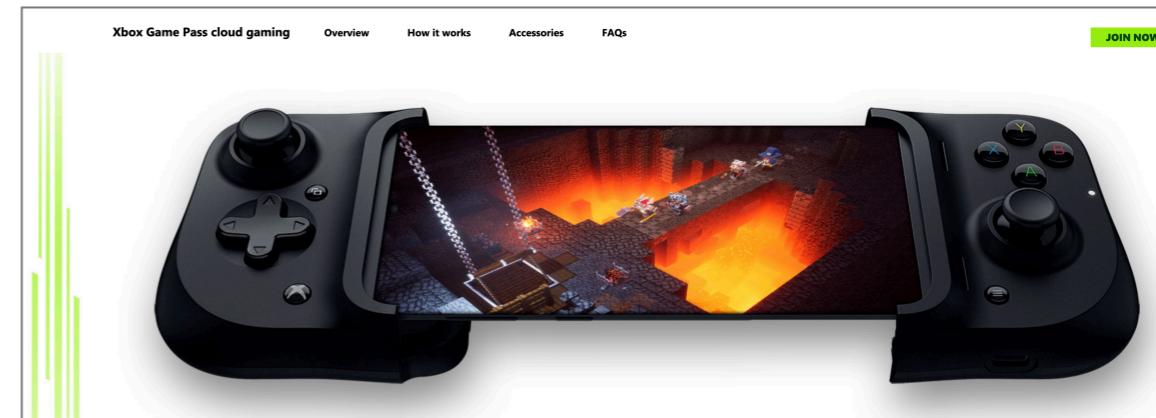
Google Stadia



Steam Remote Play



Xbox Game Pass Cloud Gaming



Audio/video-conferencing is also highly sensitive to network performance

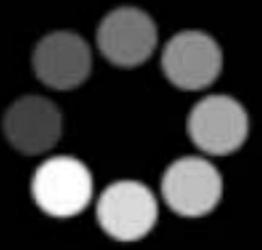


Google Hangouts



and *many* others...

And, needless to say,
so is video streaming...



Where does this come from?

By default, IP networks only offer
a (possibly lousy) best effort service

By default, IP networks only offer
a (possibly lousy) best effort service

best effort
service

IP routers do their best
to deliver packets to their destination
without making any promises

best effort
service

IP routers **do their best**
to deliver packets to their destination
without making any promises

For many applications and users
doing "its best" is just not good enough anymore

In this lecture, you'll learn
how to optimize the

- performance
- flexibility
- reliability

of large-scale network infrastructures

Why should *you* care, as a **network operator**?

IP networks carry always more critical services

IP networks carry always more critical services



Since the end of 2017,
All fixed-network services
(telephony, TV, and Internet)
run on IP technology

Same applies at ETH Zürich when you pick up the phone

The screenshot shows the ETH Zürich website interface. At the top left is the ETH Zürich logo. To its right, the text "Services & resources" is displayed. Below the logo, a horizontal navigation bar contains links: "News & events", "Organisation", "Employment & work", "Teaching", "Finance & controlling", and "IT Services". The "IT Services" link is underlined, indicating it is the current section. A breadcrumb navigation path is visible below the navigation bar: "Homepage > IT Services > IT Service Catalogue > Communication > Voice communication". On the left side, there is a "Subnavigation" button with a three-line icon. The main content area features a title "Voice communication (telephony)". Above the content, there is a "Description" section with a "Close" button. The text in this section reads: "The voice communication service (telephony) at ETH Zurich is exclusively operated by the ID ICT Networks division. It operates a comprehensive and reliable voice network. The service covers all aspects of voice communication." Below this, there are three expandable sections: "Customer Benefits" (with an "Open +" button), "Customer Groups / Cost / Order" (with an "Open +" button), and "Instructions / FAQ / How To" (with an "Open +" button).

ETH zürich Services & resources

News & events Organisation Employment & work Teaching Finance & controlling IT Services

Homepage > IT Services > IT Service Catalogue > Communication > Voice communication

Voice communication (telephony)

Description Close —

The voice communication service (telephony) at ETH Zurich is exclusively operated by the ID ICT Networks division. It operates a comprehensive and reliable voice network. The service covers all aspects of voice communication.

Customer Benefits Open +

Customer Groups / Cost / Order Open +

Instructions / FAQ / How To Open +

With great powers come great responsibilities

T WIRTSCHAFT

Unternehmen & Konjunktur Geld & Recht Karriere Börse

Swisscom-Panne: Bund kündigt Untersuchung an

Am Dienstagabend haben landesweit Notrufnummern, Internet und TV nicht funktioniert. Es ist die fünfte grosse Störung seit zwei Jahren.

Publiziert: 12.02.2020, 16:08



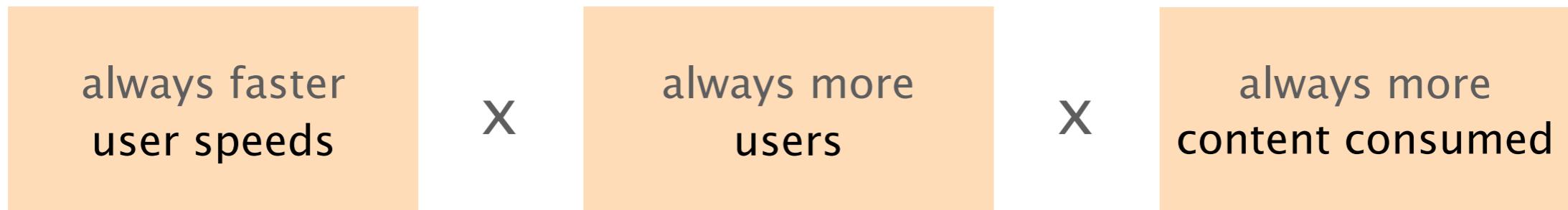
Kein Internet, keine Anrufe: Swisscom-Kunden konnten gestern wegen einer Störung nicht einmal den Notruf erreichen. (Bild: Franziska Rothenbühler)

February 2020

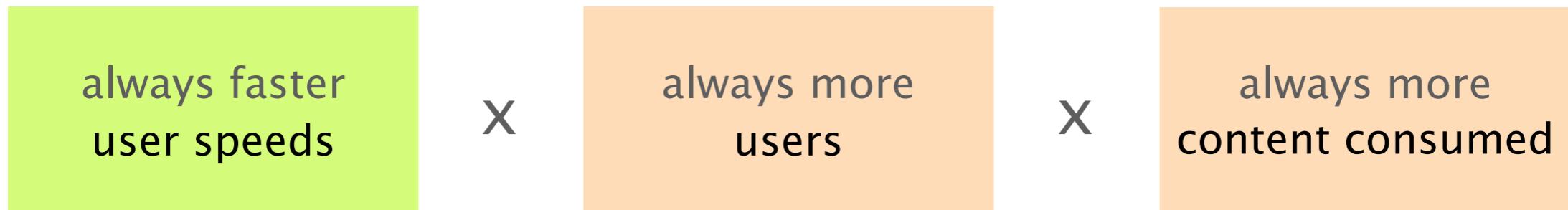
Emergency numbers
(117, 112, 144, 118)
were not reachable for >1.5h
due to a major outage
in Swisscom

IP networks carry always more traffic
critical services

IP networks carry always more traffic



IP networks carry always more traffic



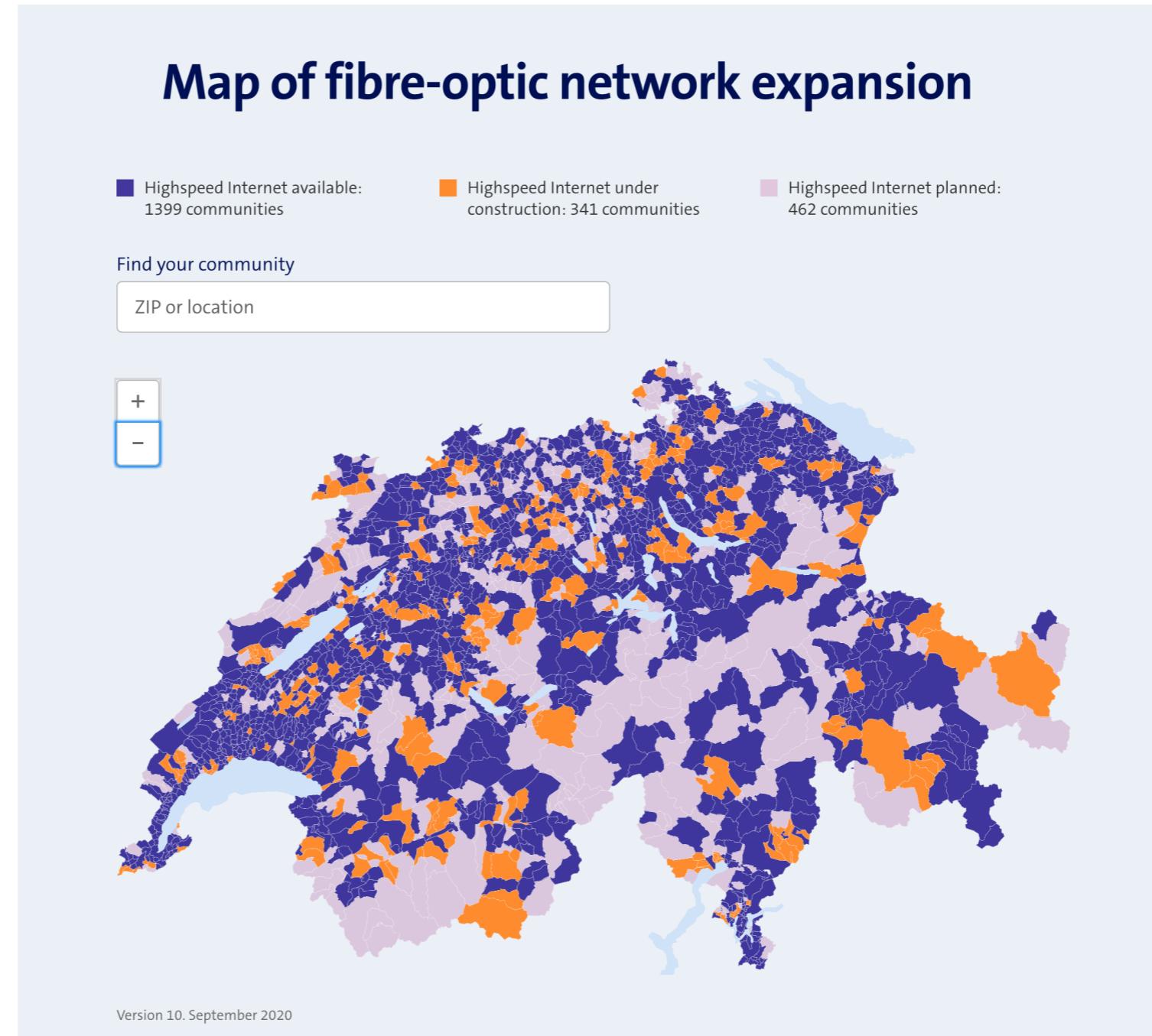
IP networks carry always more traffic

Expected growth in throughput by 2023

broadband	>x2	110 Mbps
cellular	>x3	44
Wi-Fi	>x3	92

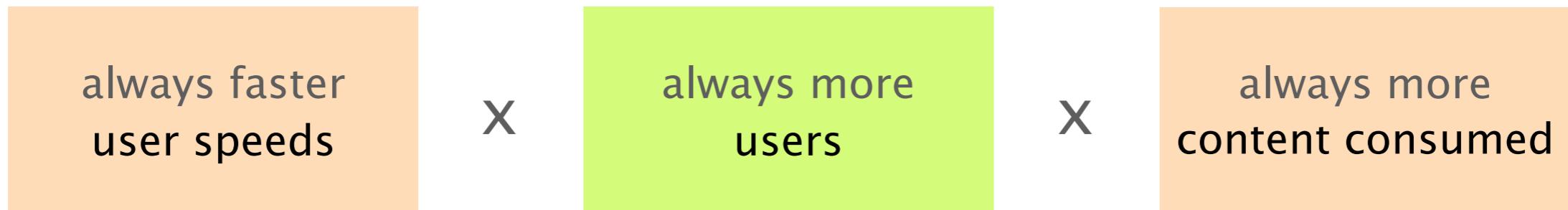
Source: Cisco Annual Internet Report (2018–2023)

Up to 60% of Swiss homes and business
will have access to Fiber connectivity by 2025

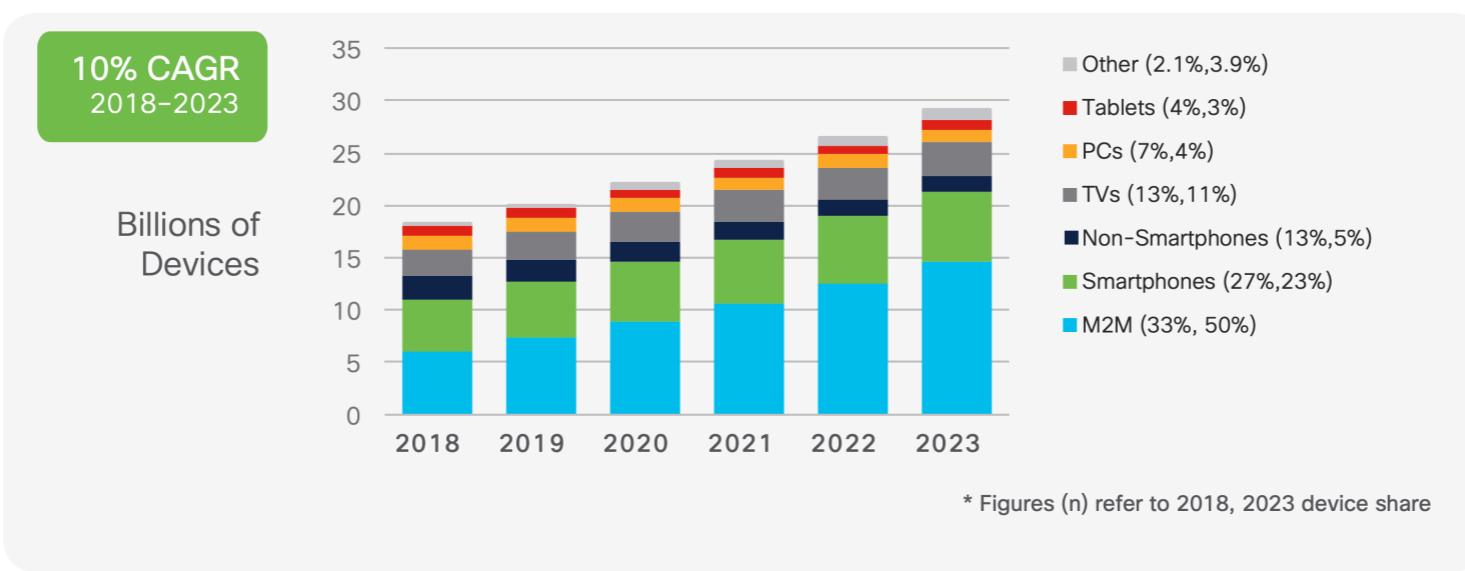
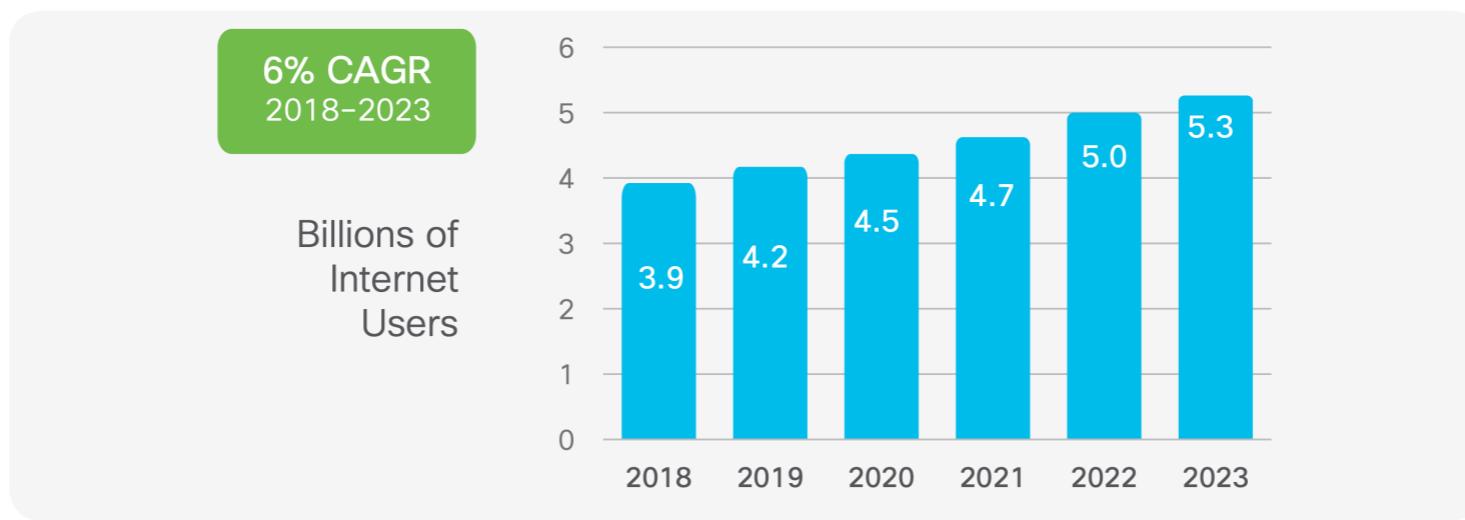


Source: Swisscom

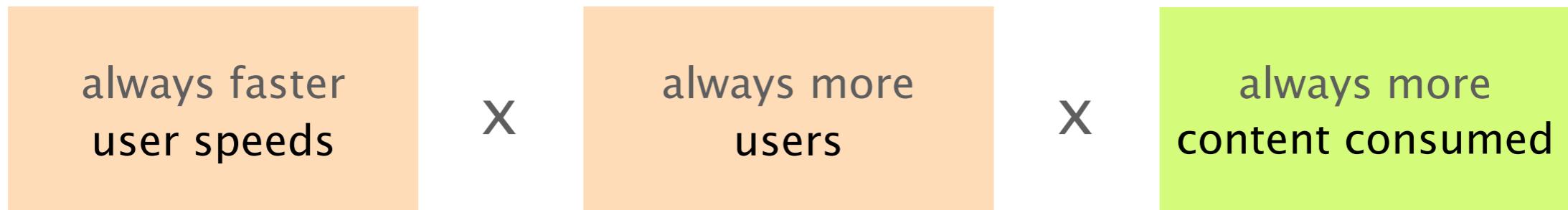
IP networks carry always more traffic



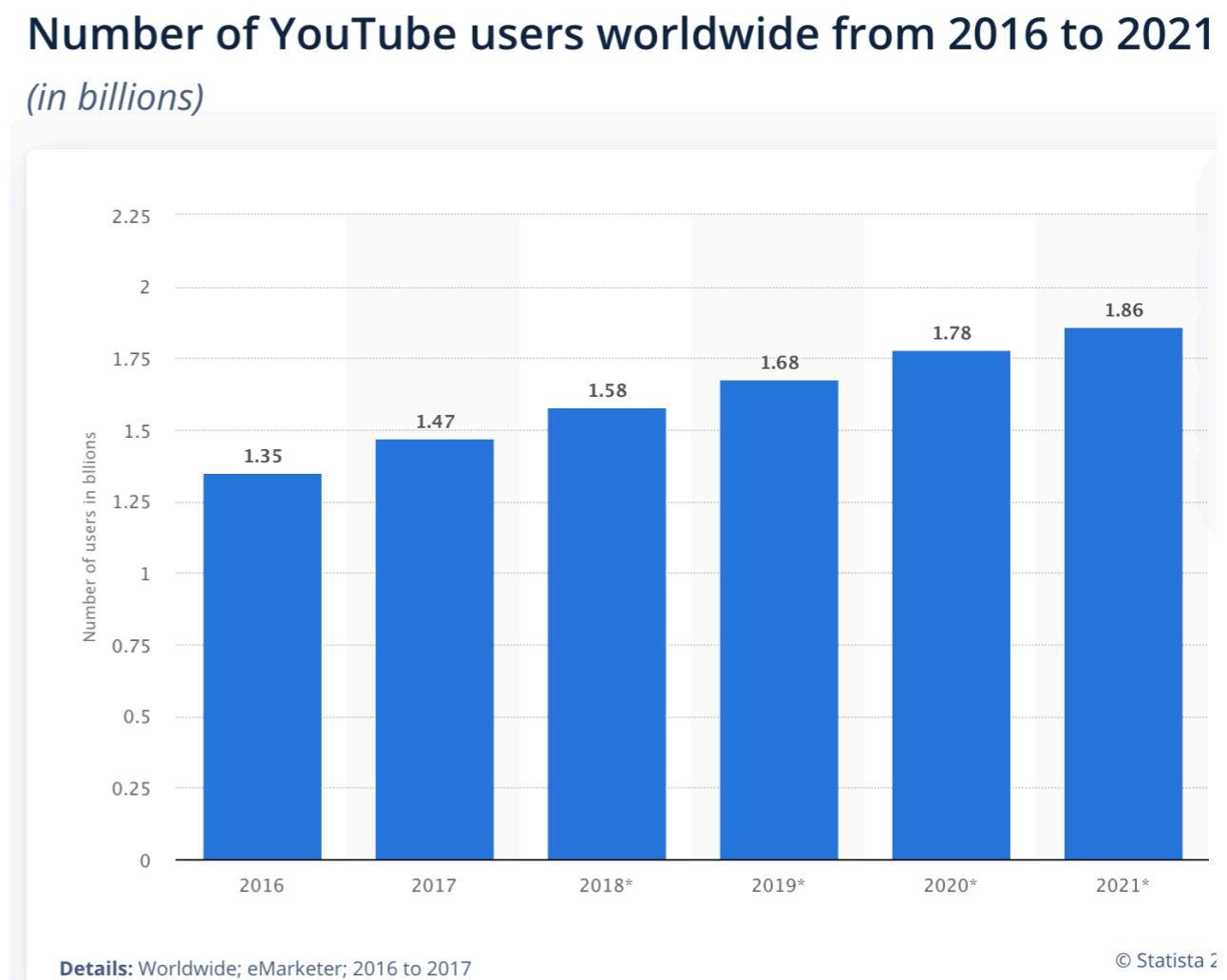
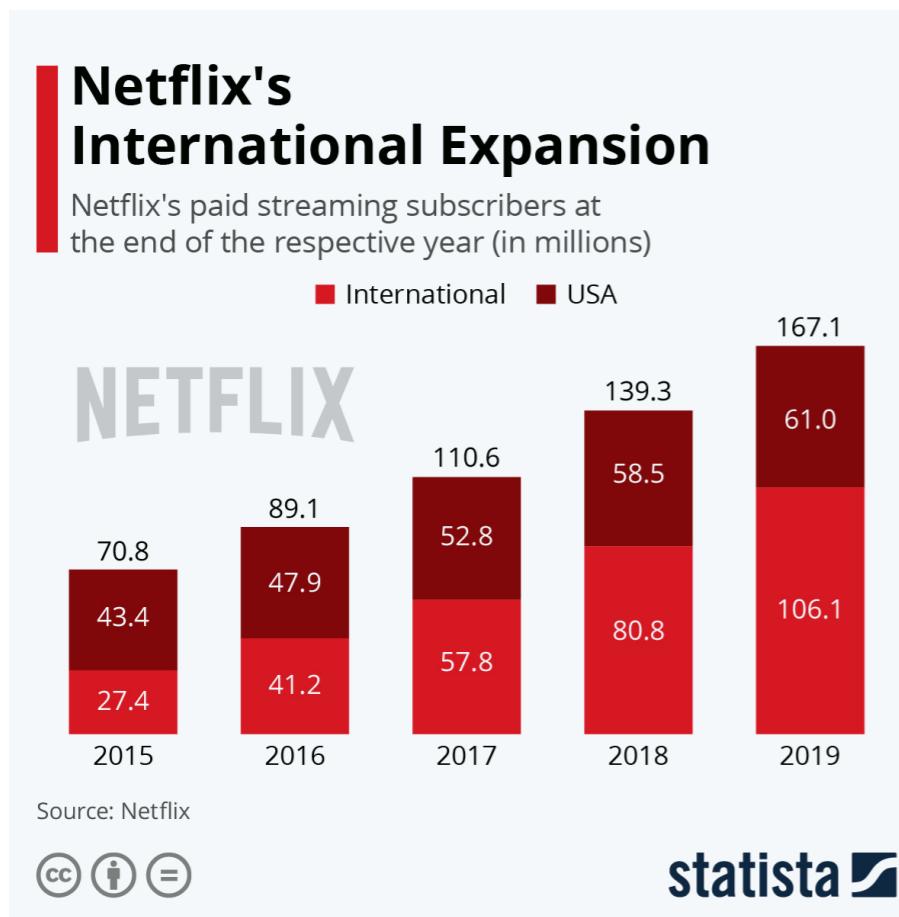
Internet growth is still going on strong



IP networks carry always more traffic



Most of the Internet traffic is video already and we keep watching more of more of them



Ever-increasing network traffic stresses
network infrastructures

One recent example

somehow arguable though

Netflix Lowers Video Quality to Aid Broadband ISP Congestion

Thursday, March 19th, 2020 (8:53 pm) - Score 5,555

[Email](#) | [Link News](#) [43 Comments](#)



In a somewhat unusual development Netflix, the global internet video streaming giant, has responded to some concerns about broadband congestion in other European countries by agreeing with the EU to lower the quality (bitrate) of its videos for the next 30 days (i.e. easing the strain on fixed line ISPs and mobile operators).

The move will result in the company reducing video bitrates across all of their streams, although they have yet to say if those paying extra for the highest quality are going to be compensated. Netflix have also yet to confirm if this will apply to the United Kingdom, where the majority of fixed line providers have yet to really struggle with managing increased demand ([we've covered this here](#)).

<https://www.ispreview.co.uk/index.php/2020/03/netflix-lowers-video-quality-to-aid-broadband-isp-congestion.html>

In this lecture, you'll learn
how to optimize the

- performance
- flexibility
- reliability

of large-scale network infrastructures

Techniques

Performance

Traffic Engineering

Load Balancing

Quality of Service

Multicast

Flexibility

Virtual Private Networks

Reliability

Fast Convergence

Techniques

Performance

Traffic Engineering

Load Balancing

Quality of Service

Multicast

Flexibility

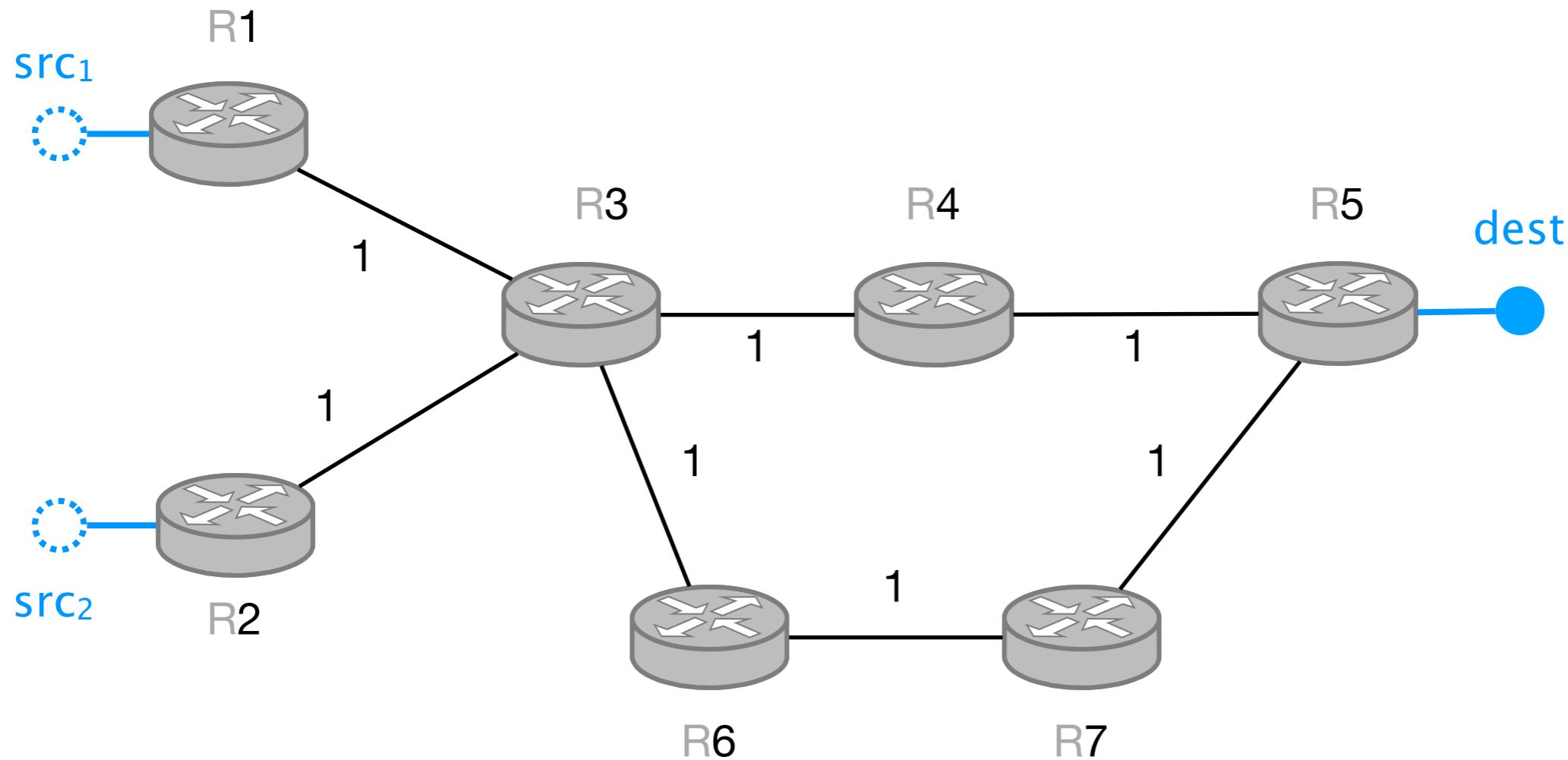
Virtual Private Networks

Reliability

Fast Convergence

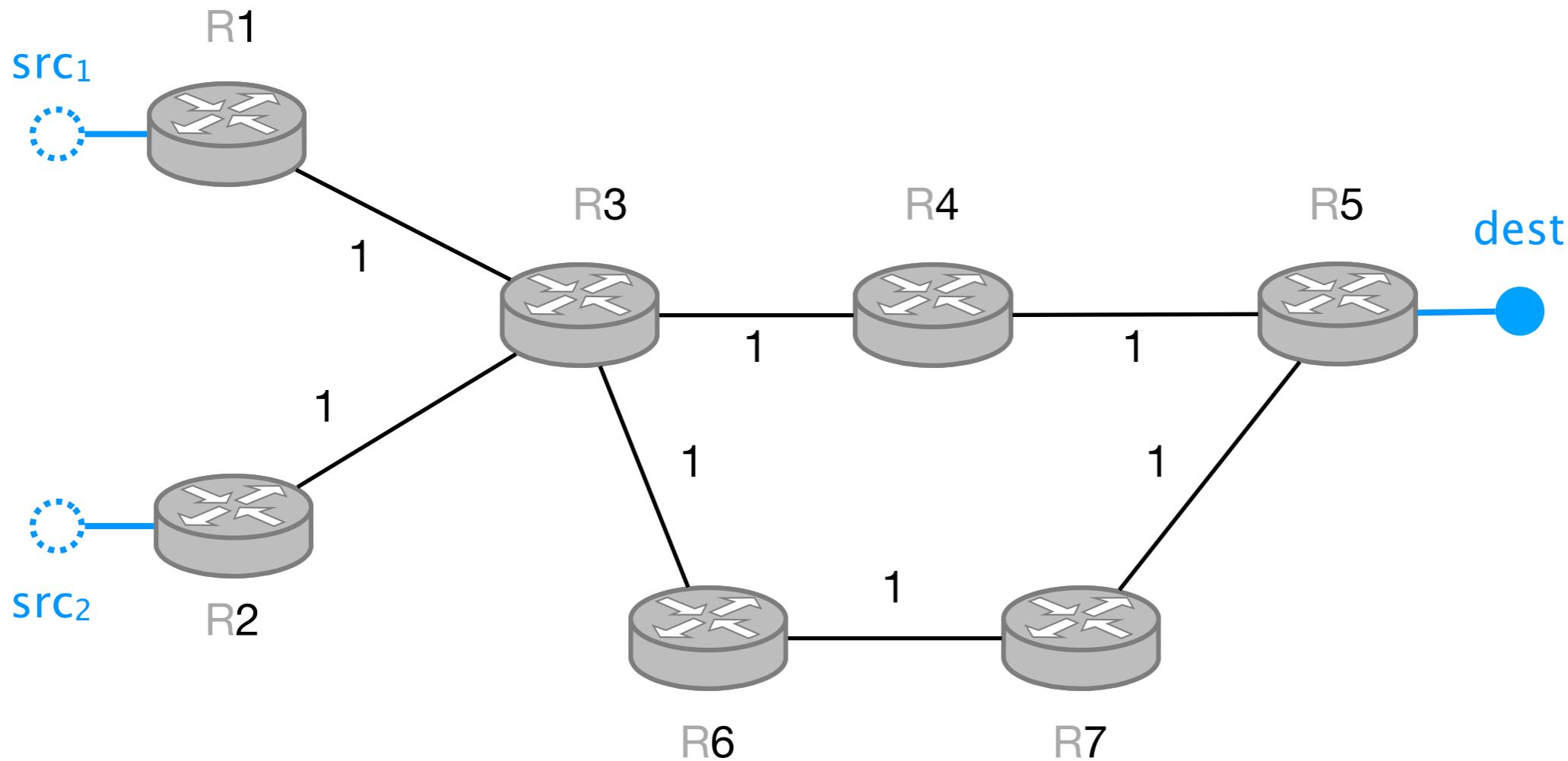
Consider this 10-Gbps network running vanilla OSPF

src₁ and *src₂* send traffic to *dest*, 1 TCP flow each

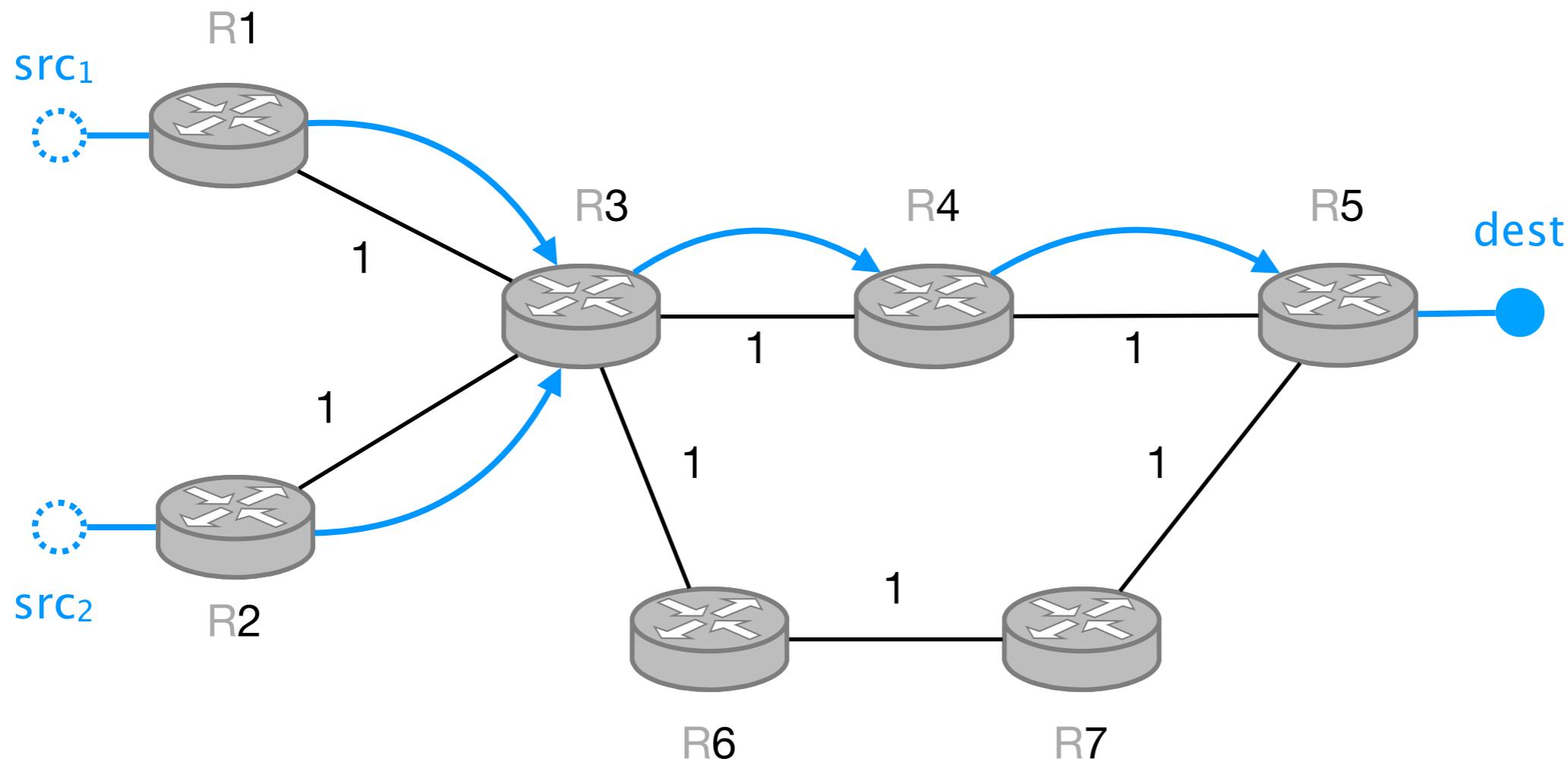


Consider this 10-Gbps network running vanilla OSPF

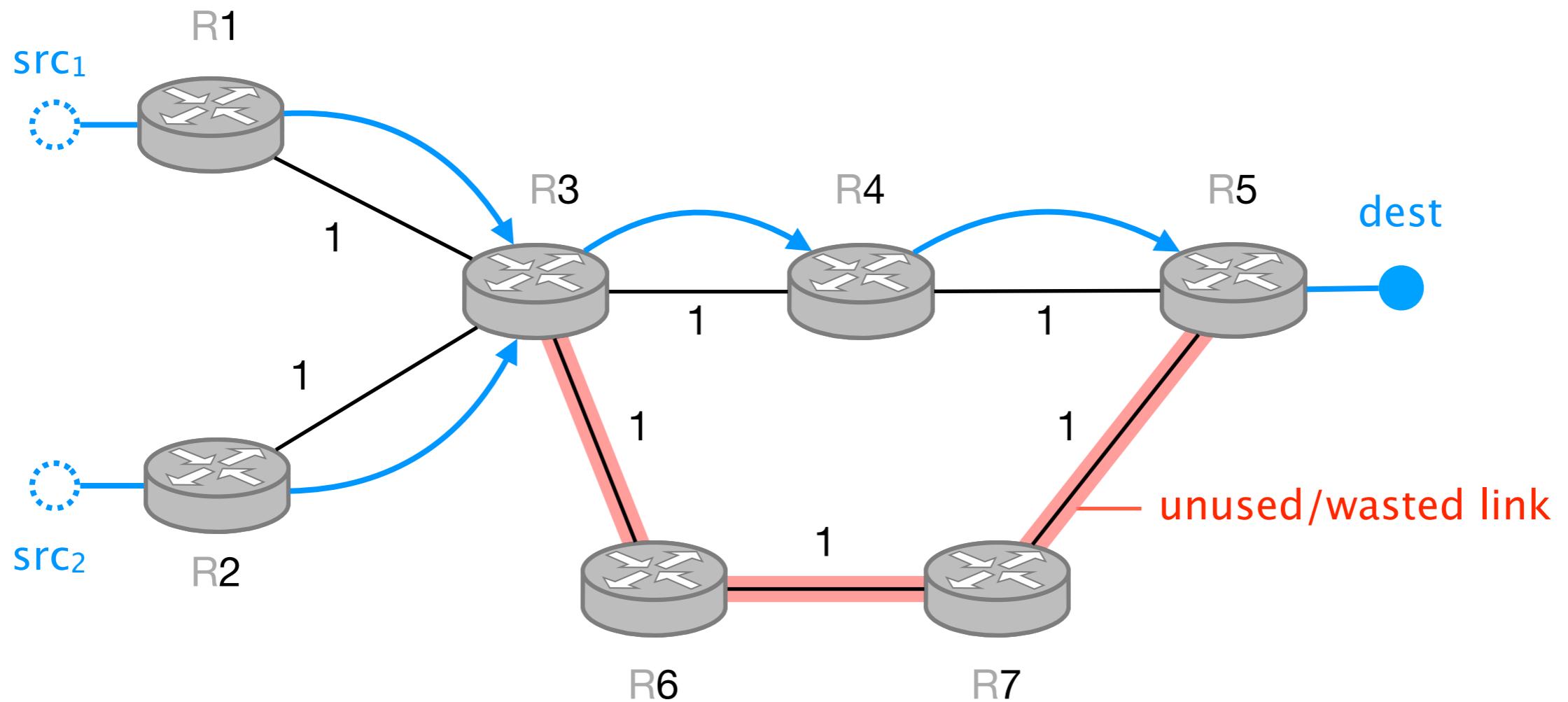
src_1 and src_2 send traffic to dest , 1 TCP flow each



What's the max. throughput $(\text{src}_1, \text{dest})$ and $(\text{src}_2, \text{dest})$ can achieve?

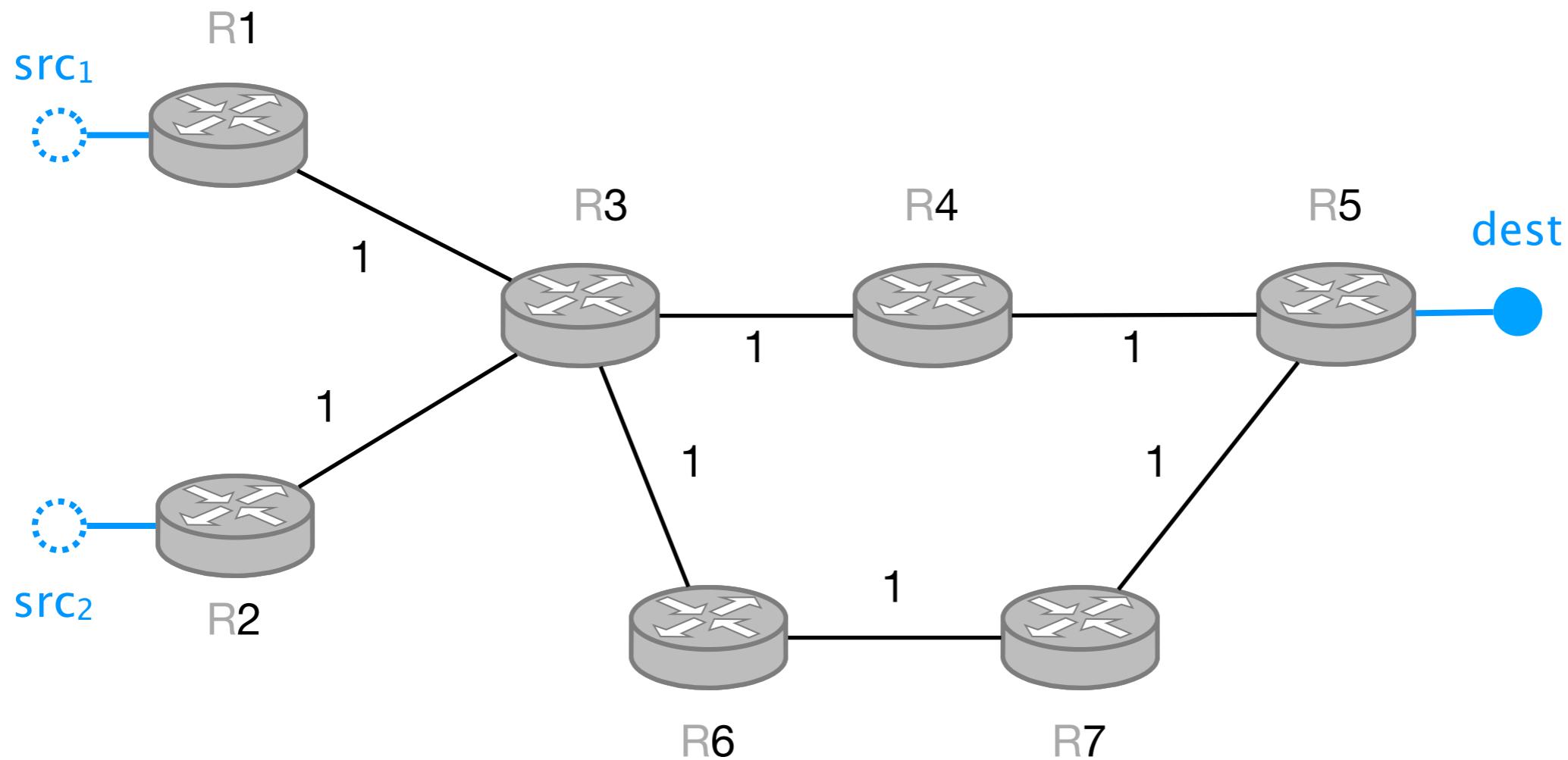


What's the max. throughput $(src_1, dest)$ and $(src_2, dest)$ can achieve? **5 Gbps**



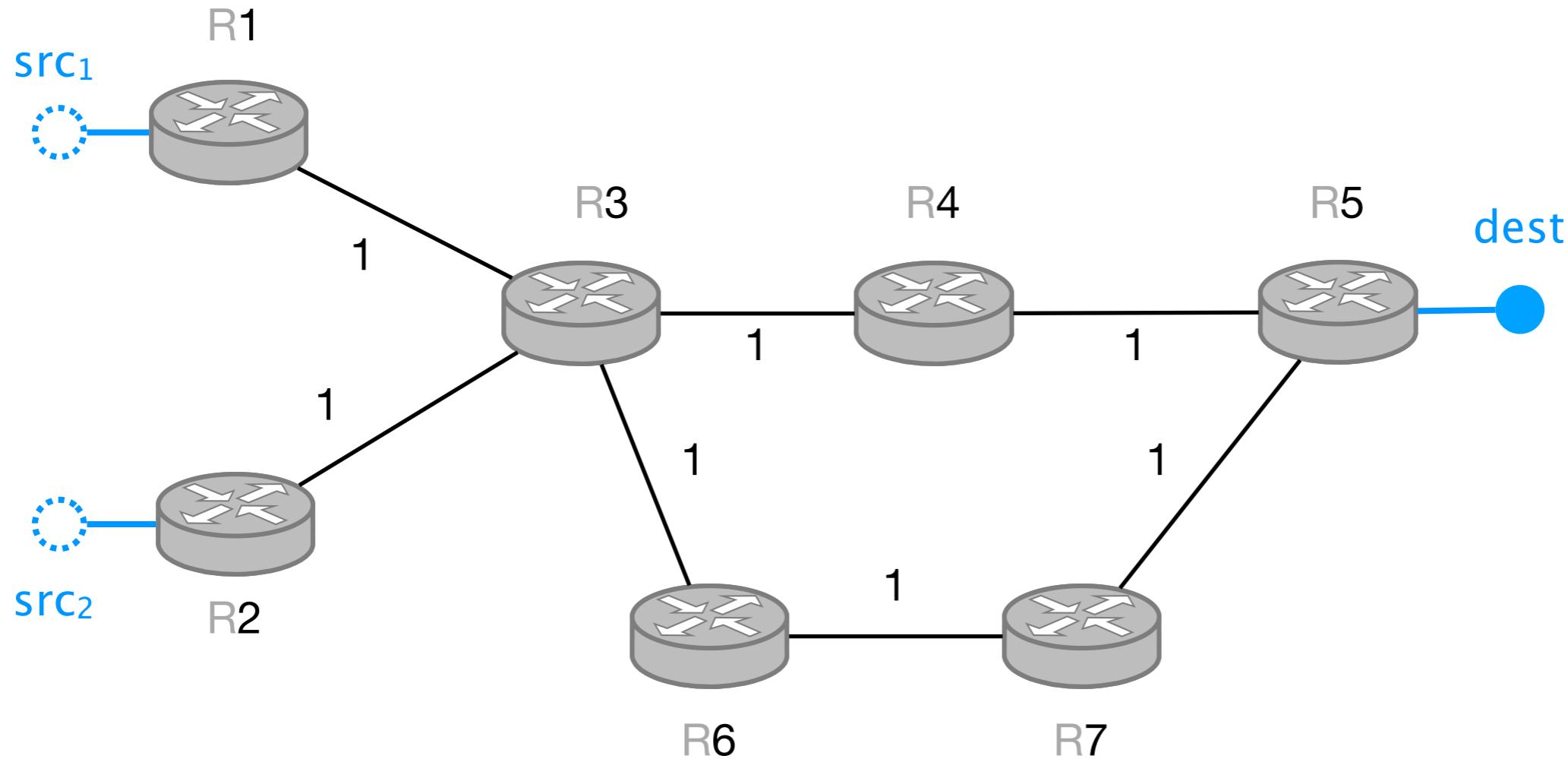
What's the max. throughput (**src₁**, **dest**) and (**src₂**, **dest**) can achieve? **5 Gbps**

What can we do better?



What can we do better?

What about better weights?



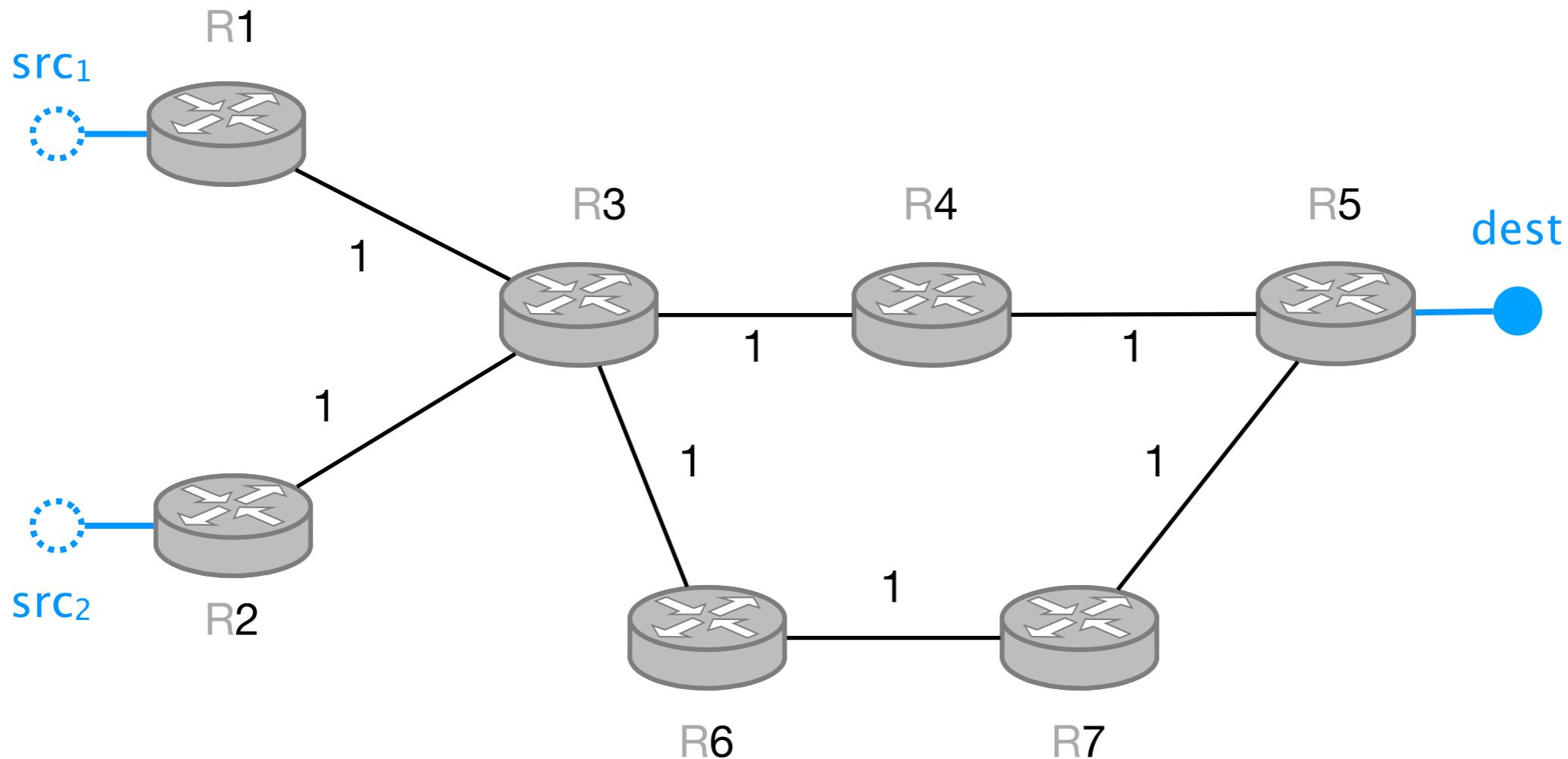
Equal-Cost Multi-Path (ECMP) is a routing strategy in which IP routers splits traffic over all best paths

Equal-Cost Multi-Path (ECMP) is a routing strategy
in which IP routers splits traffic over all **best** paths

in OSPF, best == shortest

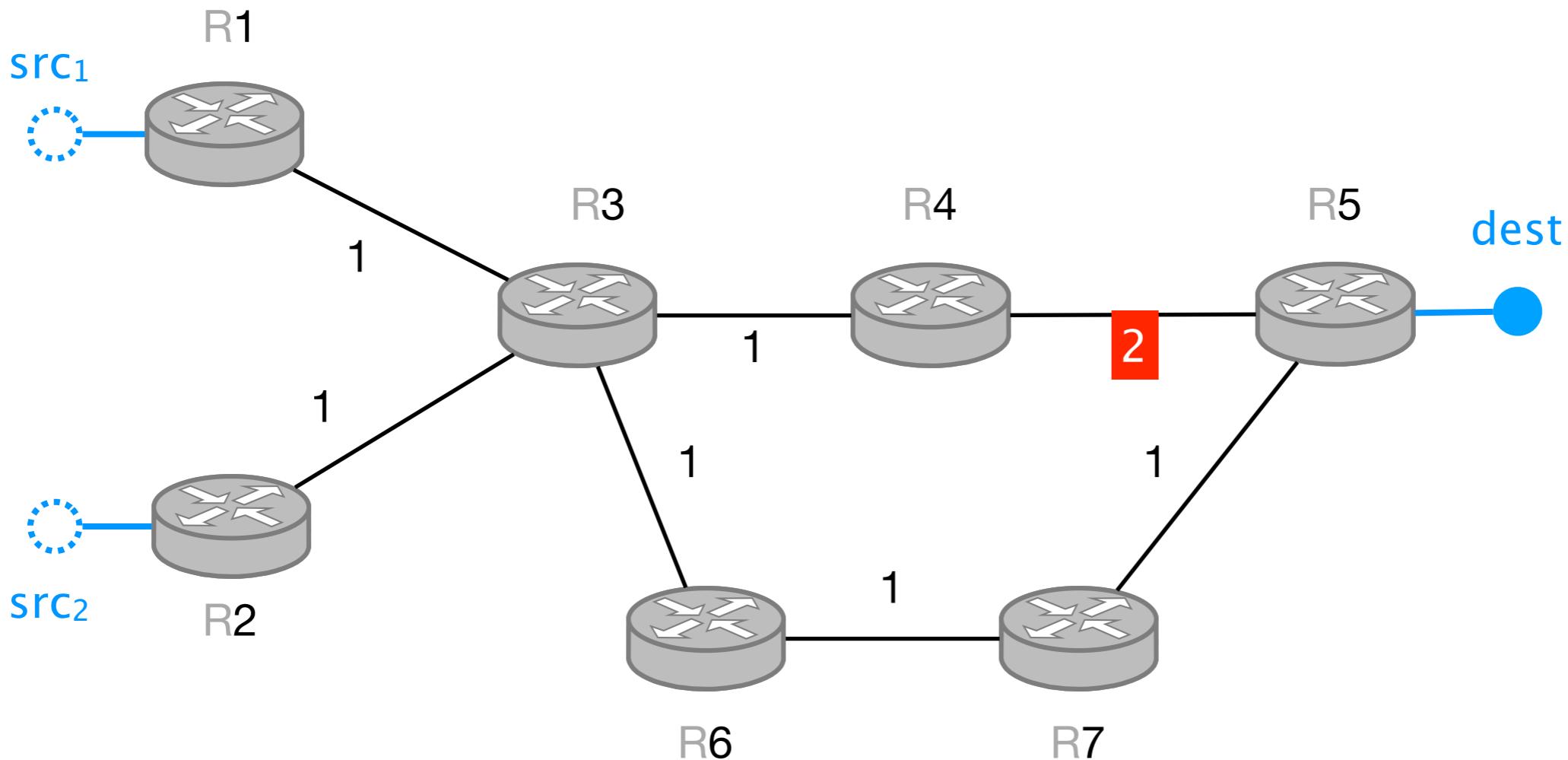
What can we do better?

What about better weights?

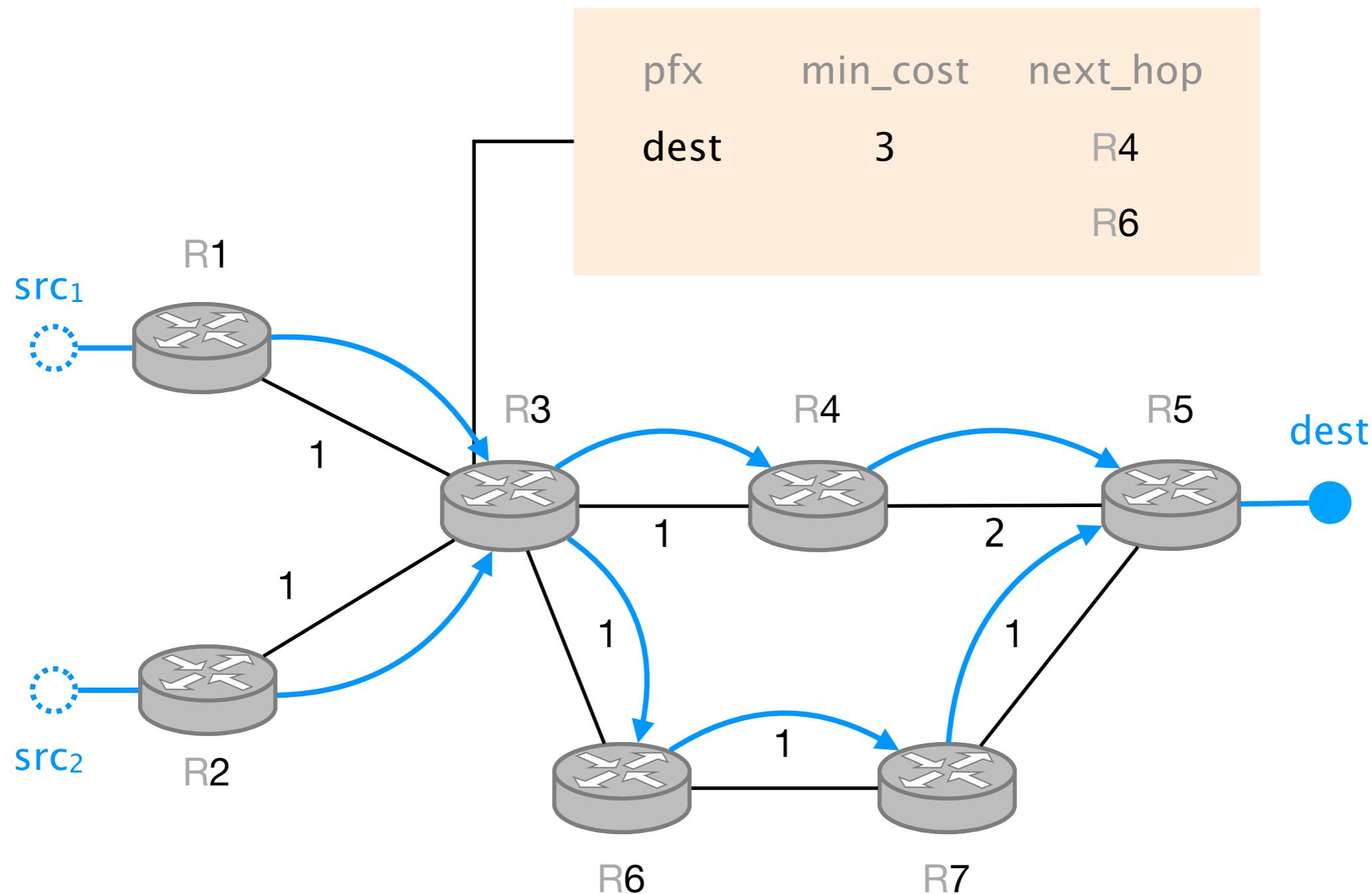


What can we do better?

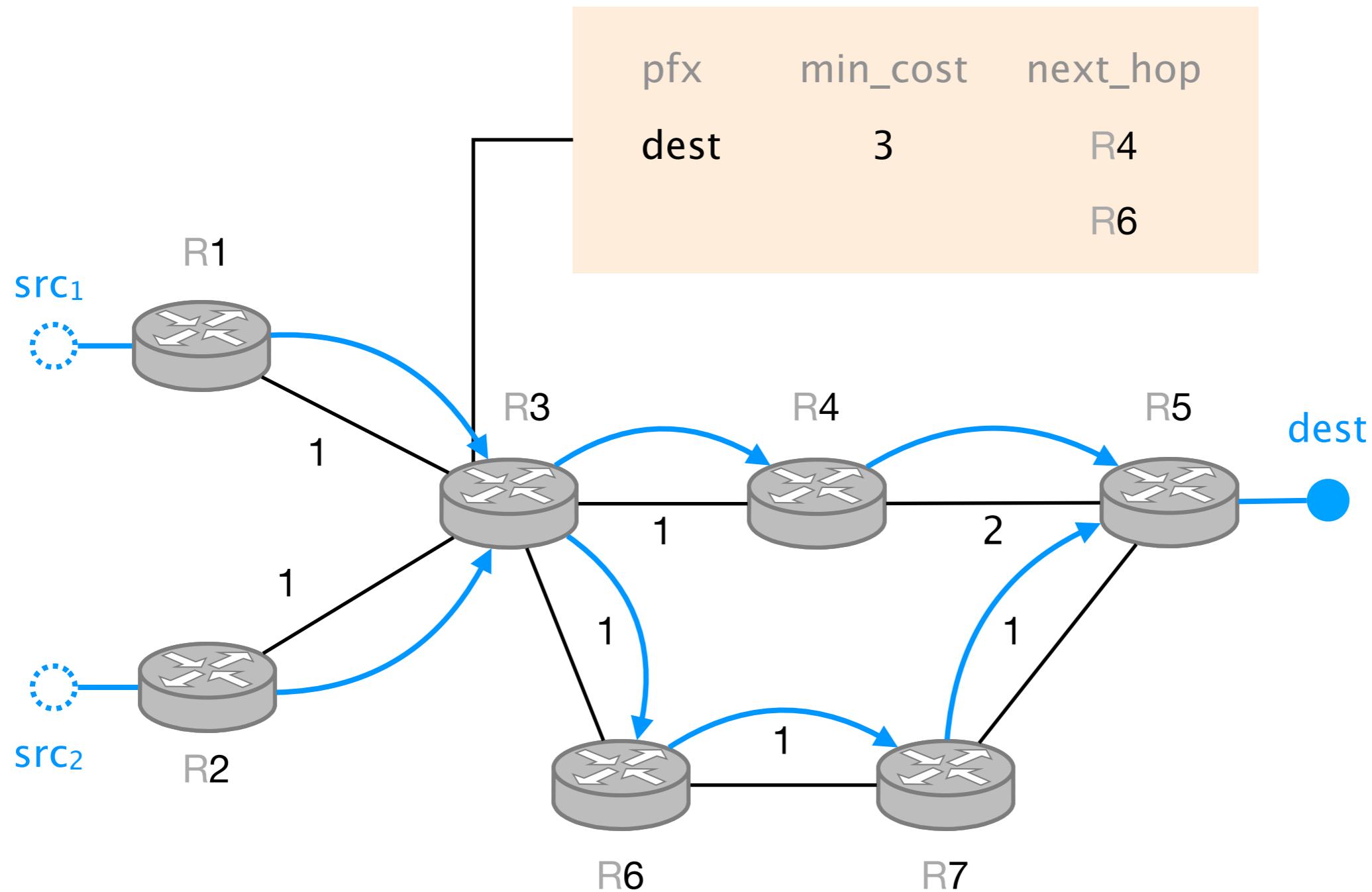
What about better weights?



R3's routing table



R3's routing table



Is it enough to guarantee we fully utilize the network?

Nope...

Equal-Cost Multi-Path (ECMP) is a routing strategy
in which IP routers **splits traffic** over all best paths

how?

Equal-Cost Multi-Path (ECMP) typically relies on stateless hash functions to split traffic

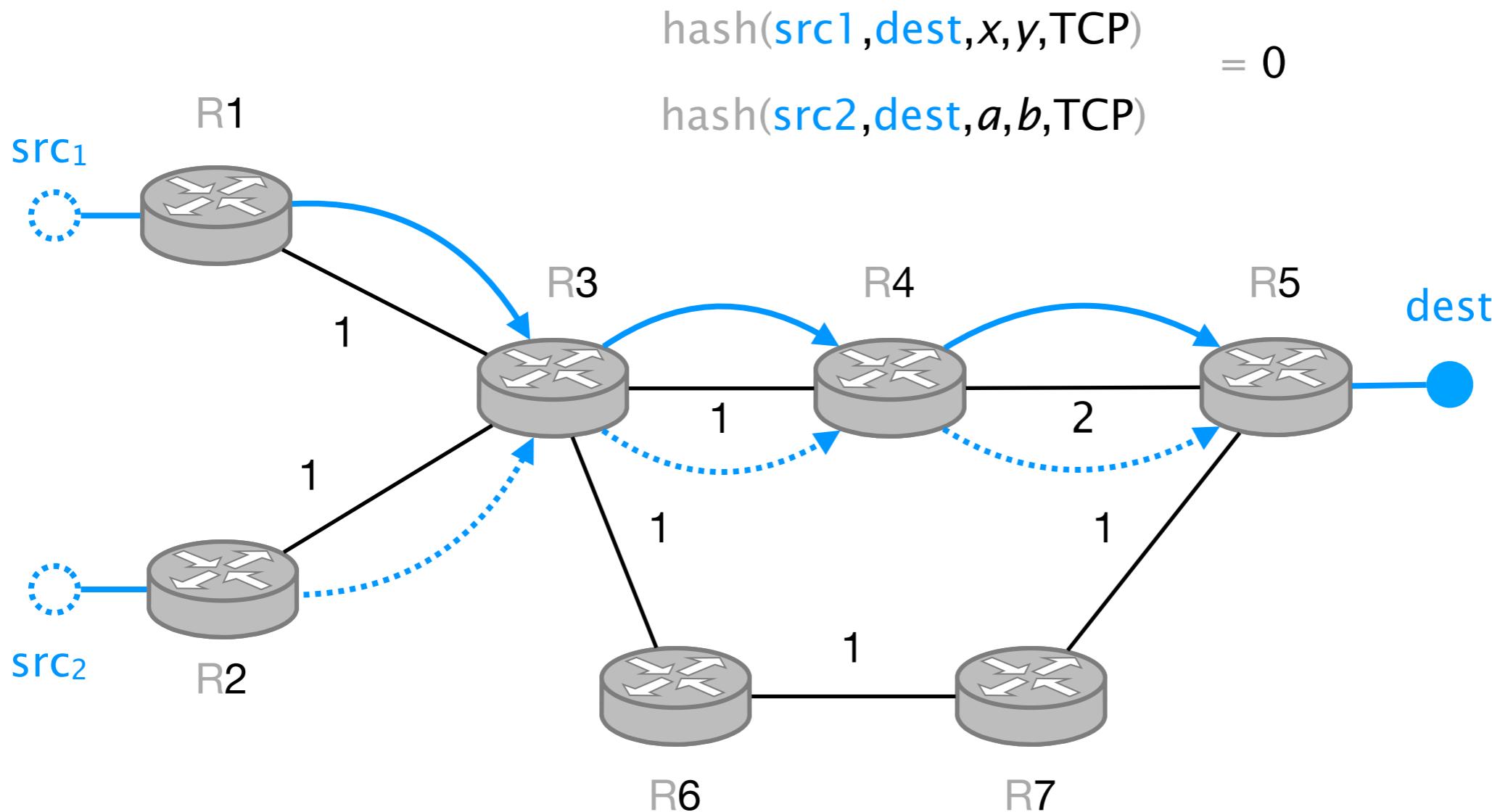
$$\text{ECMP next_hop} = \text{hash} \left[\begin{array}{l} \text{src_ip} \\ \text{dst_ip} \\ \text{src_port} \\ \text{dst_port} \\ \text{proto} \end{array} \right] \% \# \text{next_hops}$$

Property

All packets belonging to the same TCP flows go over the same path
(doing so avoids reordering)

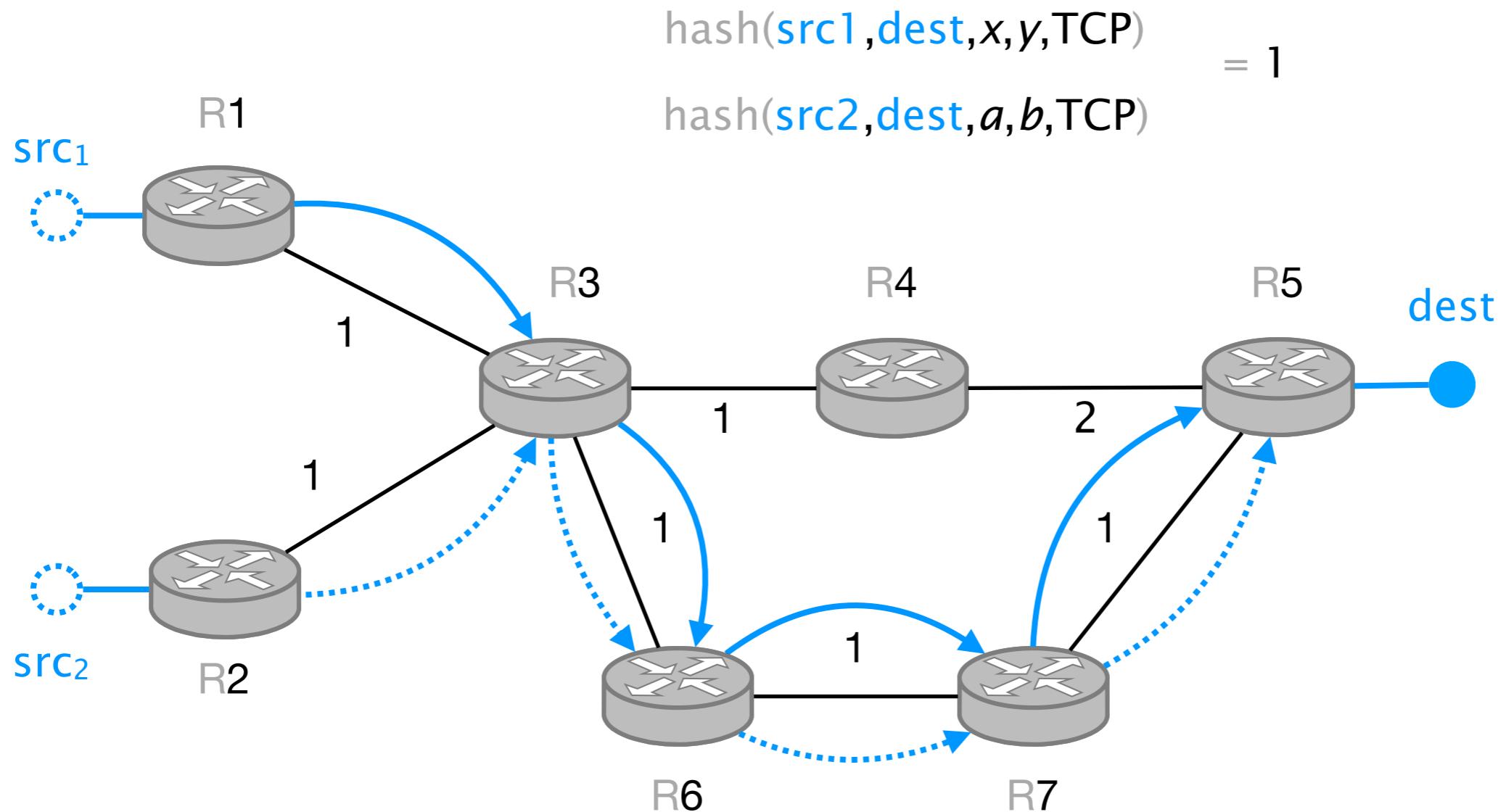
What can go wrong?

Hash collisions!

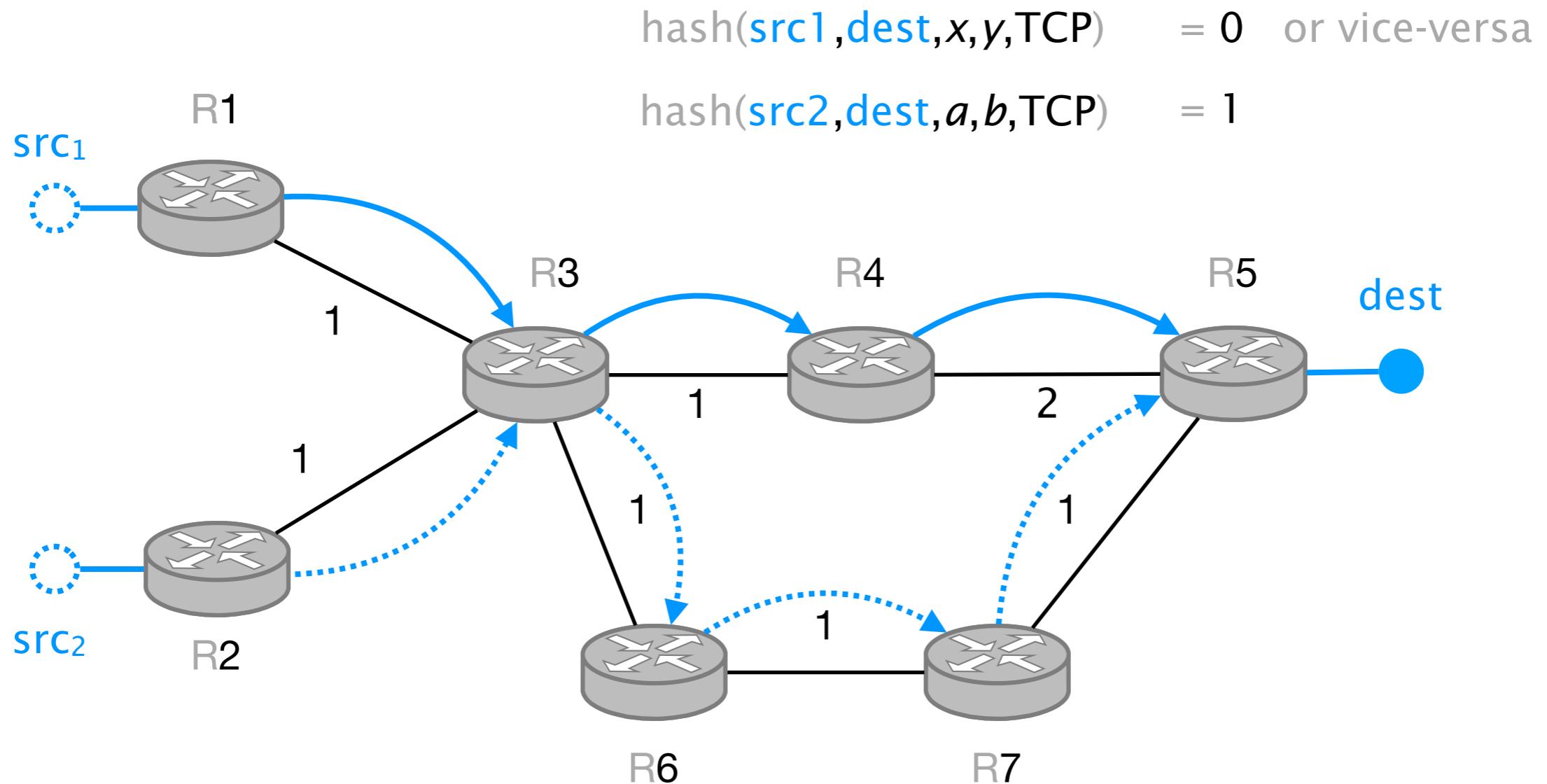


What can go wrong?

Hash collisions!

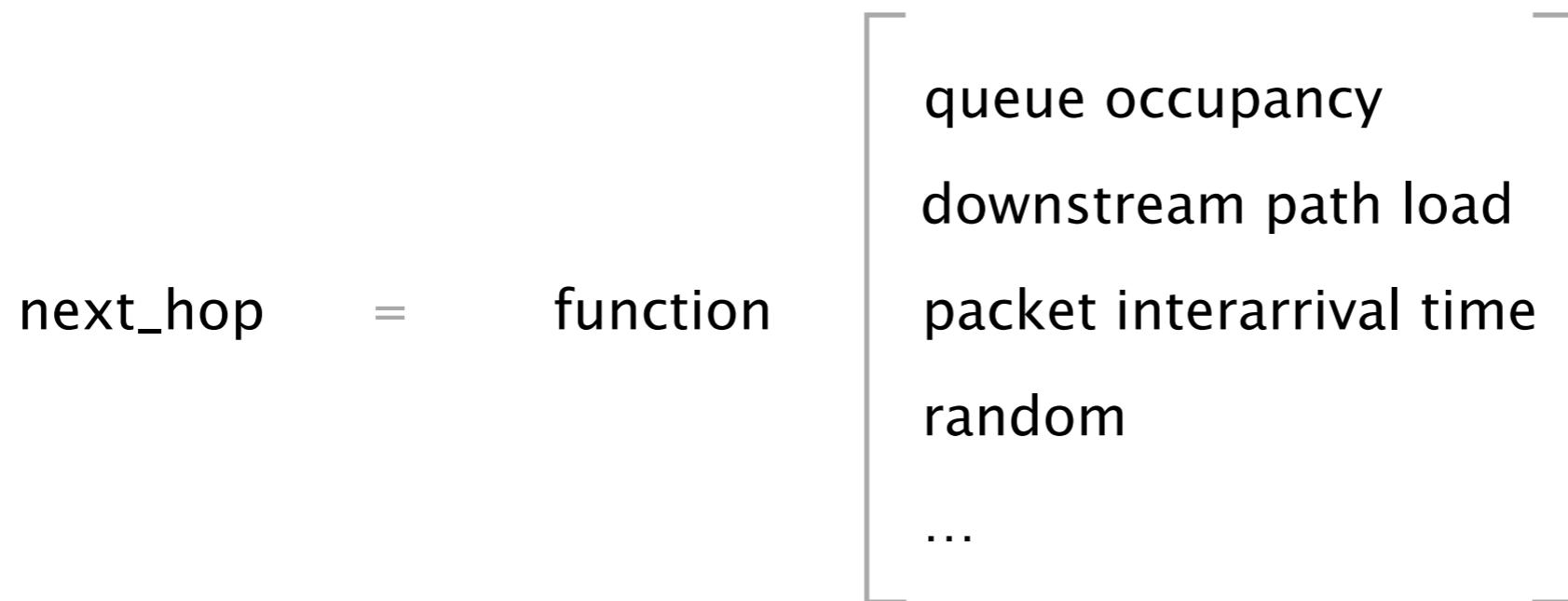


Traffic is effectively load balanced when the hash don't match

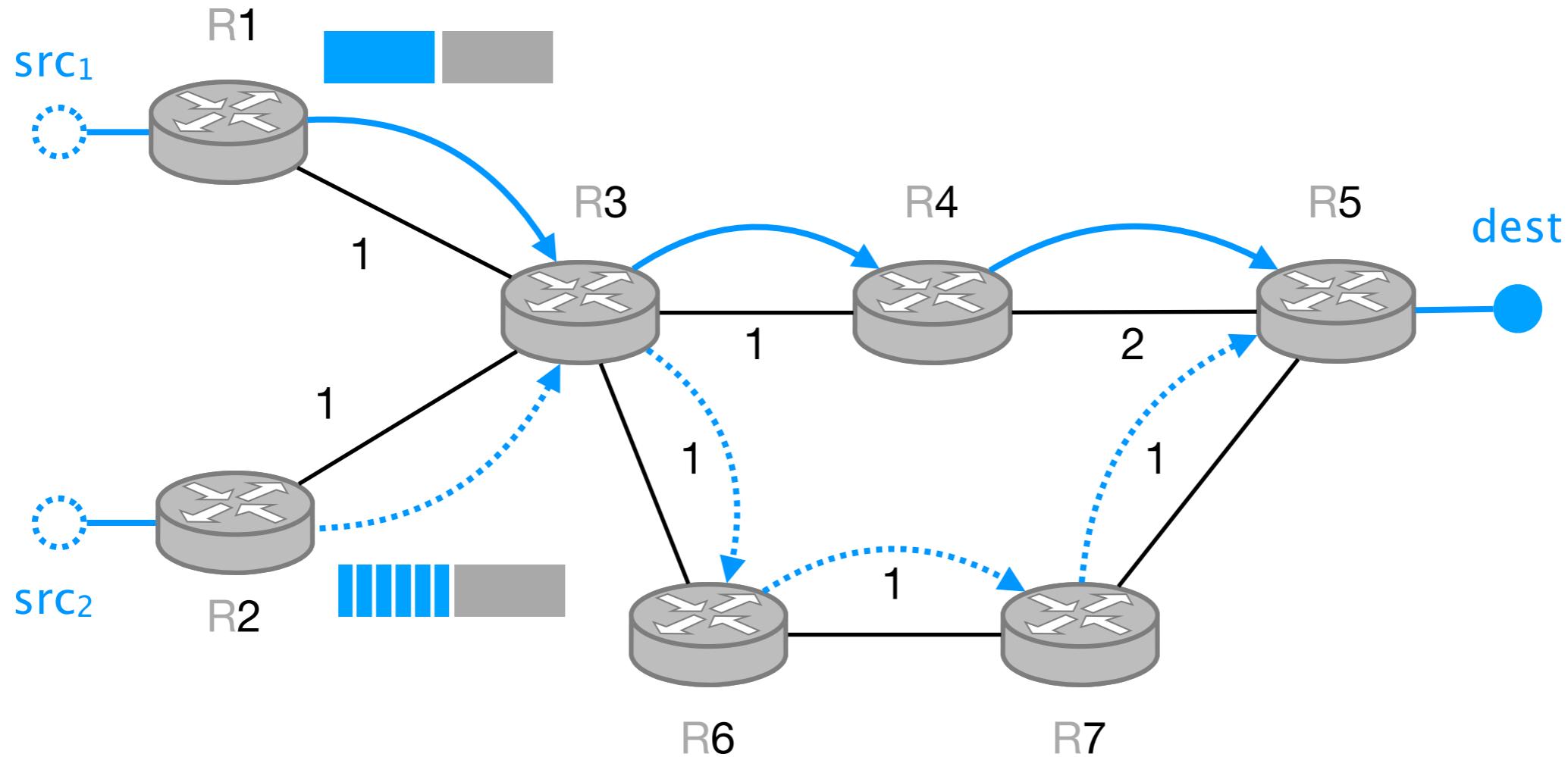


We'll see different ways to solve this problem including:
smarter load-balancing functions and label switching

We'll see different ways to solve this problem including:
smarter load-balancing functions and label switching



We'll see different ways to solve this problem including:
smarter load-balancing functions and **label switching**



Techniques

Performance

Traffic Engineering

Load Balancing

Quality of Service

Multicast

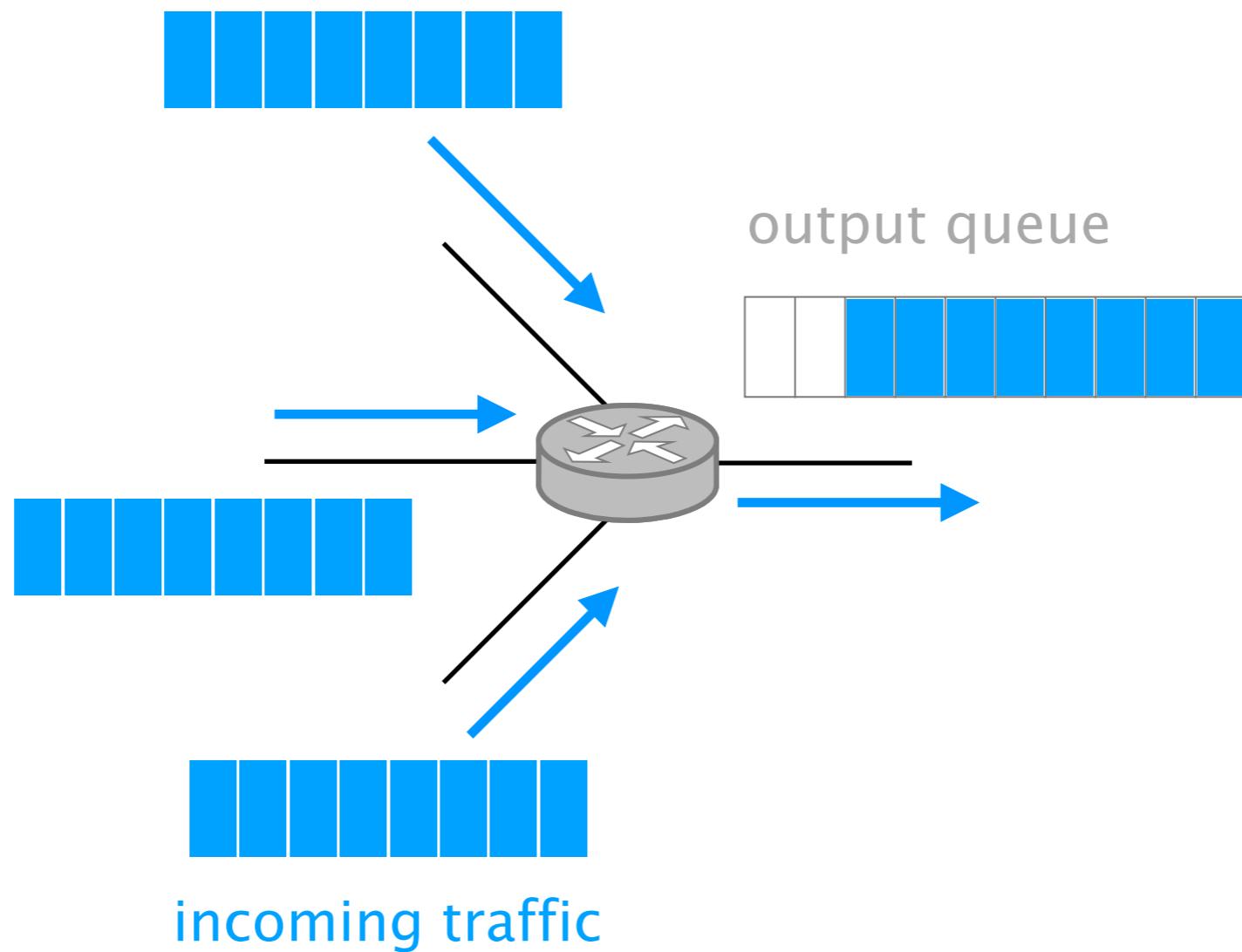
Flexibility

Virtual Private Networks

Reliability

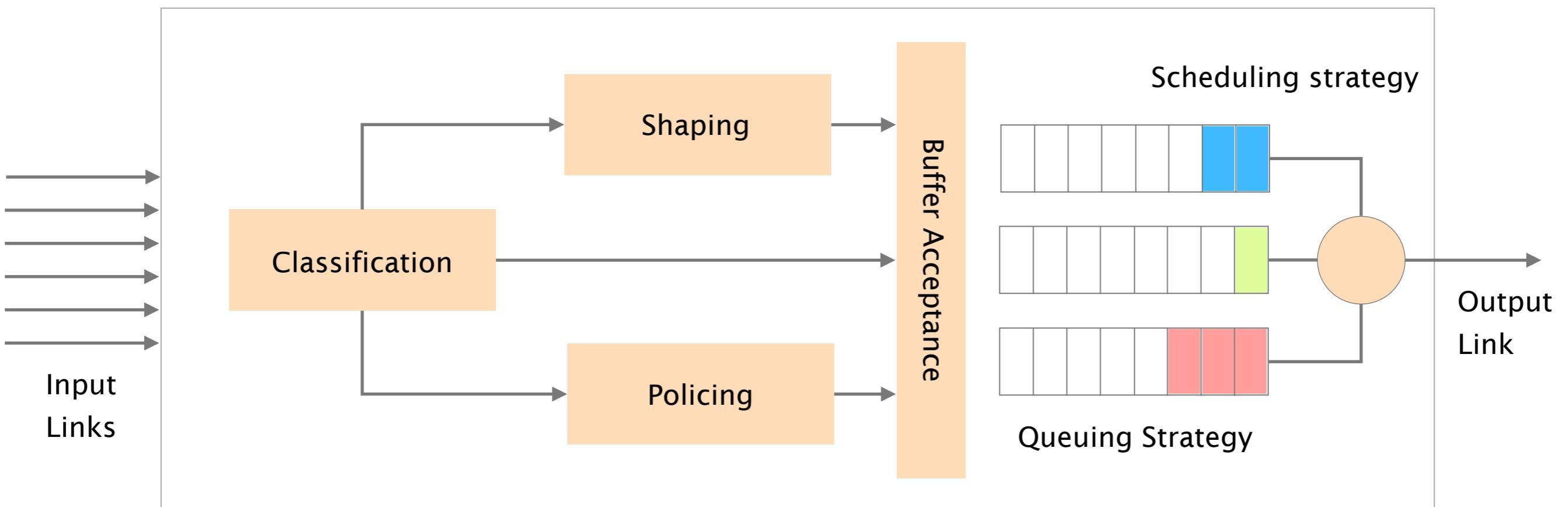
Fast Convergence

While traffic engineering helps,
it cannot always prevent congestion



We'll see different ways to manage congestion using Quality of Service (QoS)

QoS-enabled router



Techniques

Performance

Traffic Engineering

Load Balancing

Quality of Service

Multicast

Flexibility

Virtual Private Networks

Reliability

Fast Convergence

How do we deliver live videos to millions of clients?

Champions League final



vs



France

11.1

x

{8, 25}

=

{89, 279}

Germany

13.84

{110, 346}

million viewers

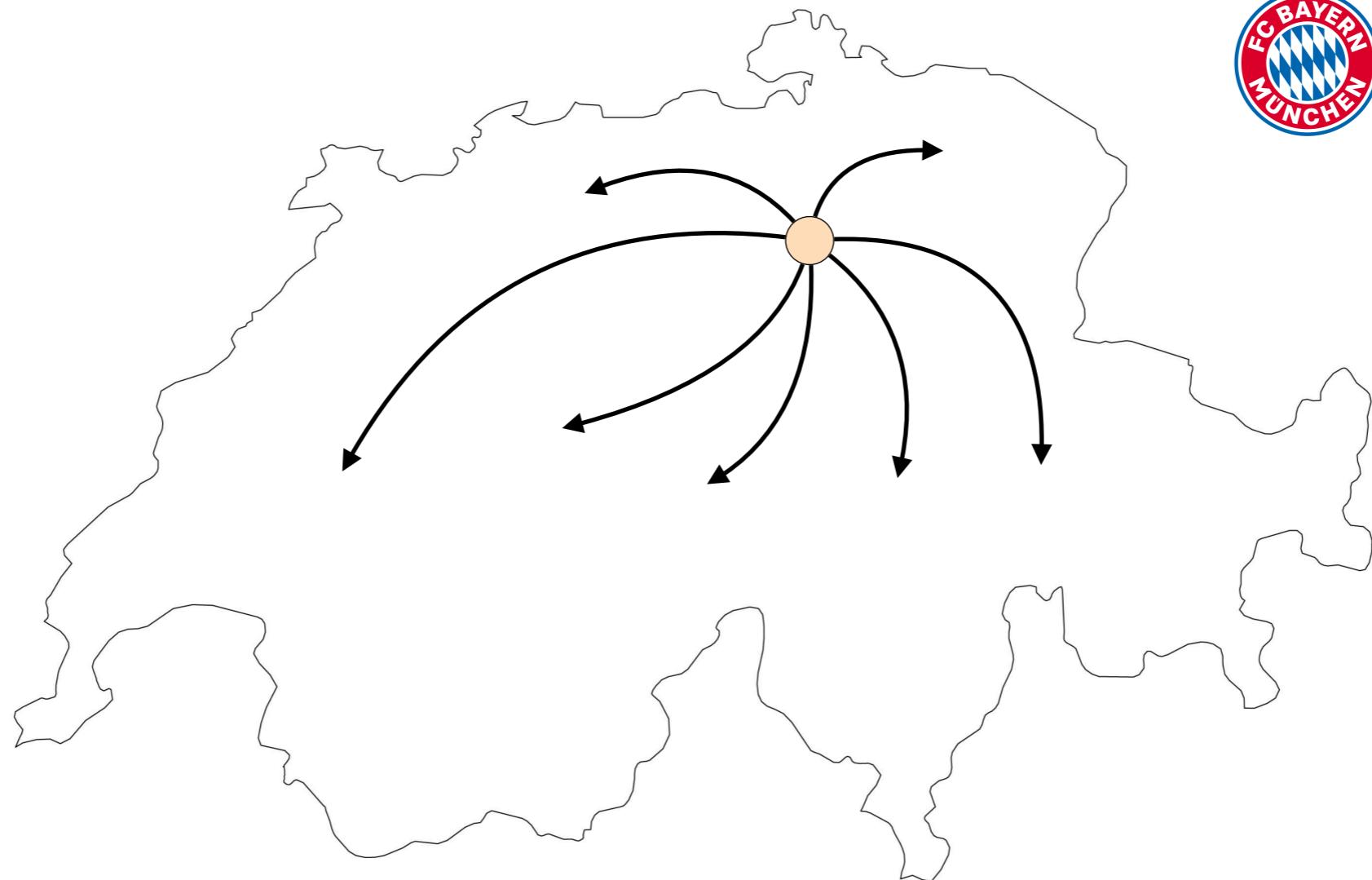
average

Mbps/feed

{HD, UHD}

Tbps

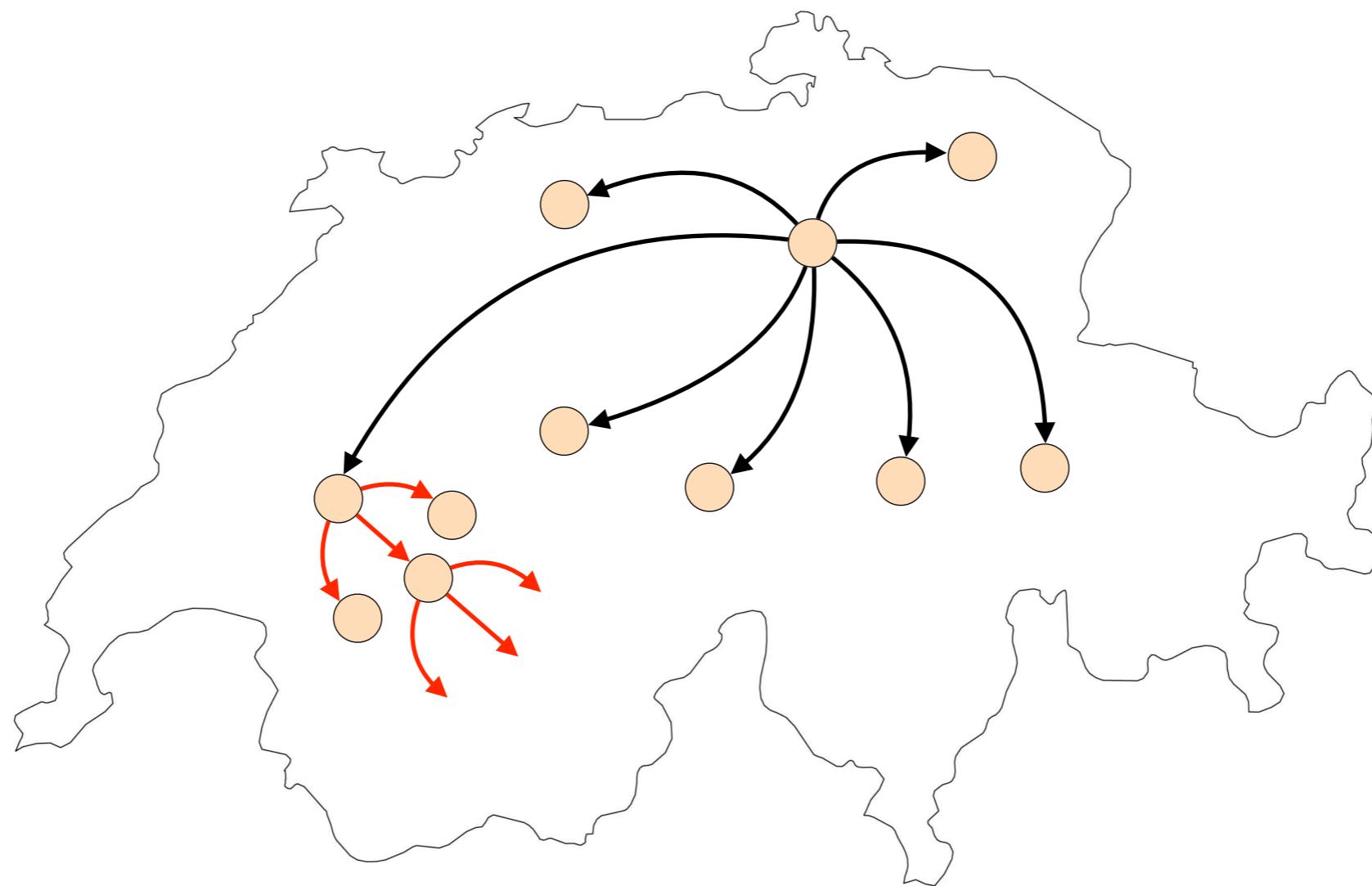
How do we deliver live videos to millions of clients?



UEFA
CHAMPIONS
LEAGUE



We'll see how the network itself can help doing that using IP multicast



Techniques

Performance

Traffic Engineering

Load Balancing

Quality of Service

Multicast

Flexibility

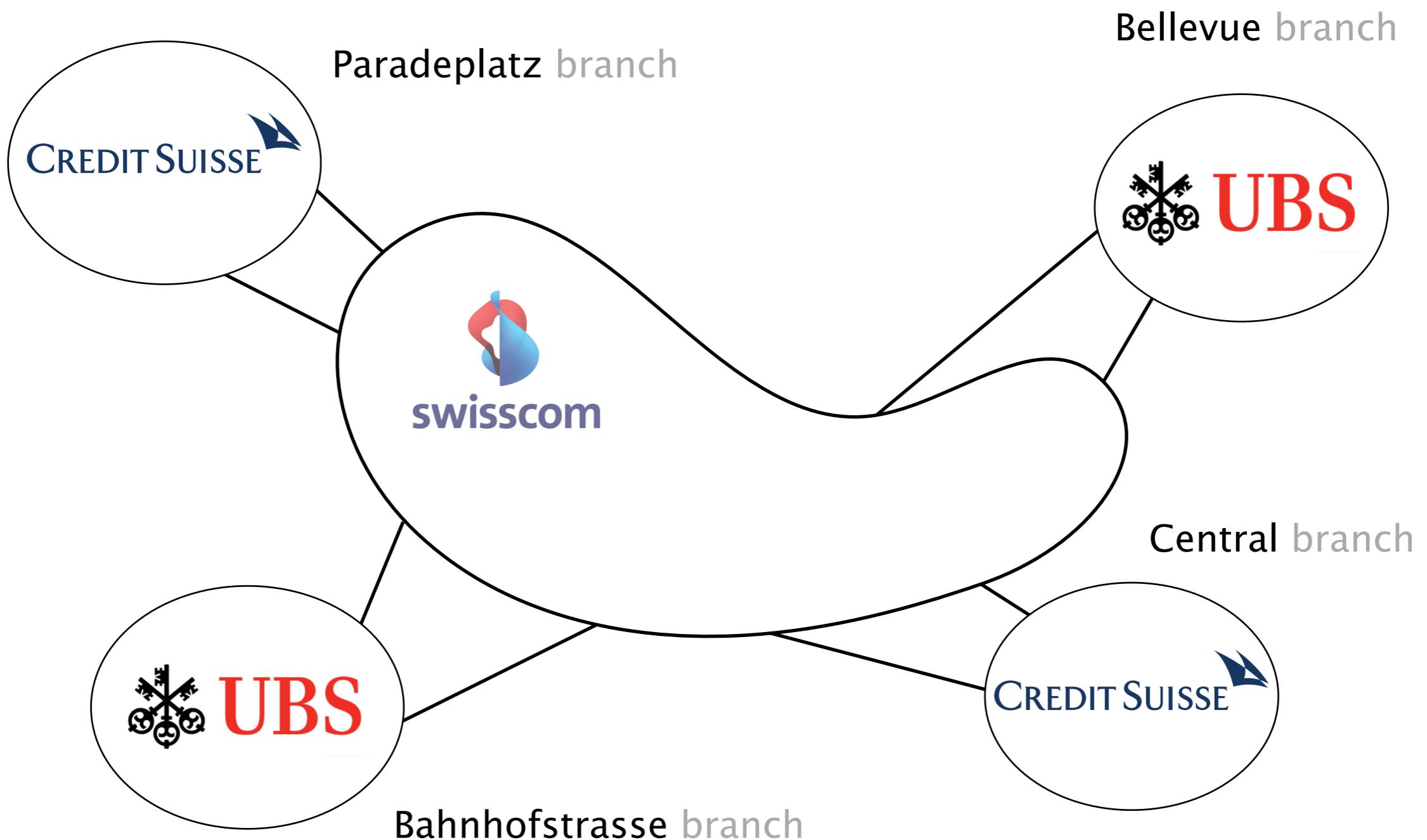
Virtual Private Networks

Reliability

Fast Convergence

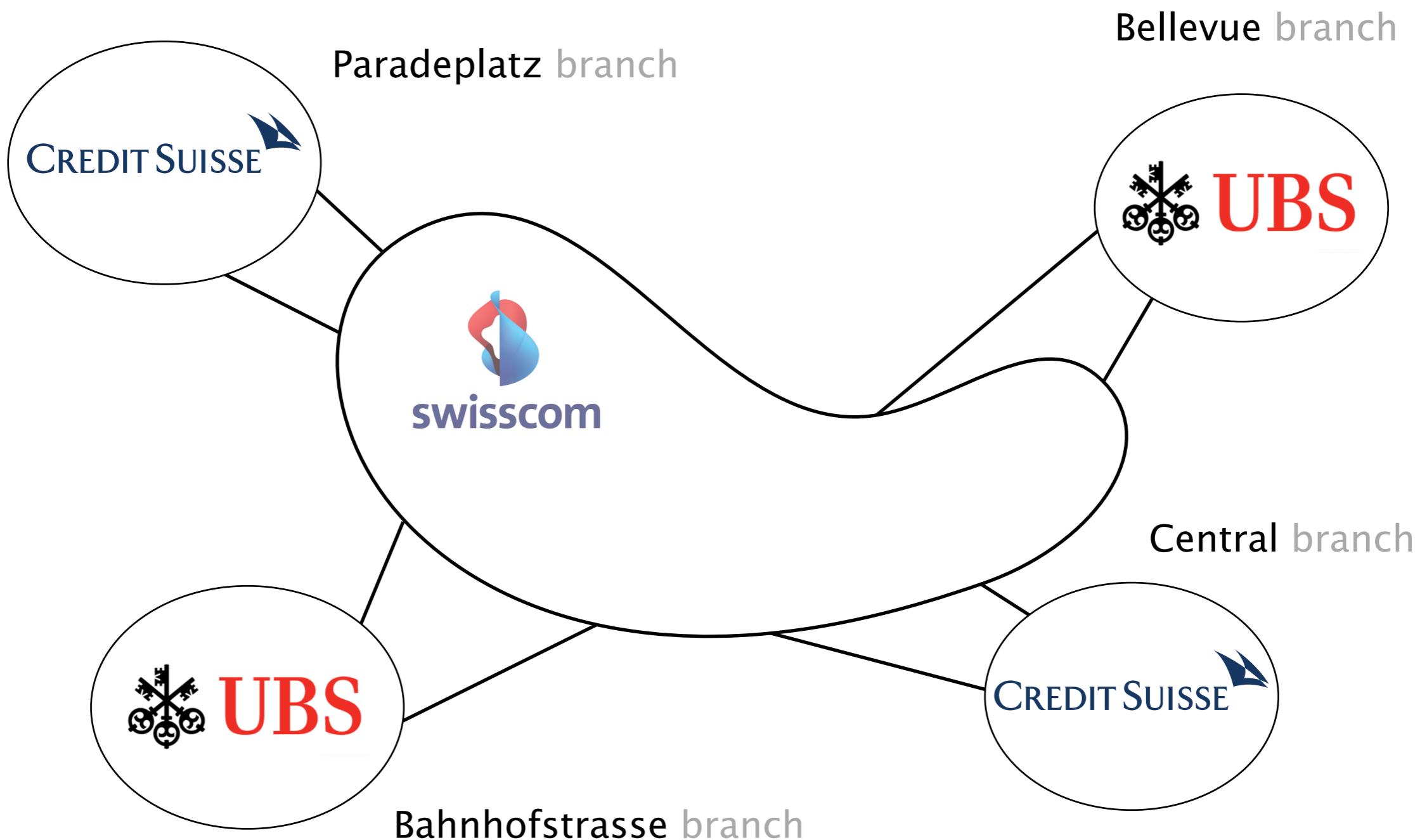
How do we virtualize a network infrastructure?

Consider this classical scenarios in which a provider (Swisscom) would like to interconnect multiple sites with each other



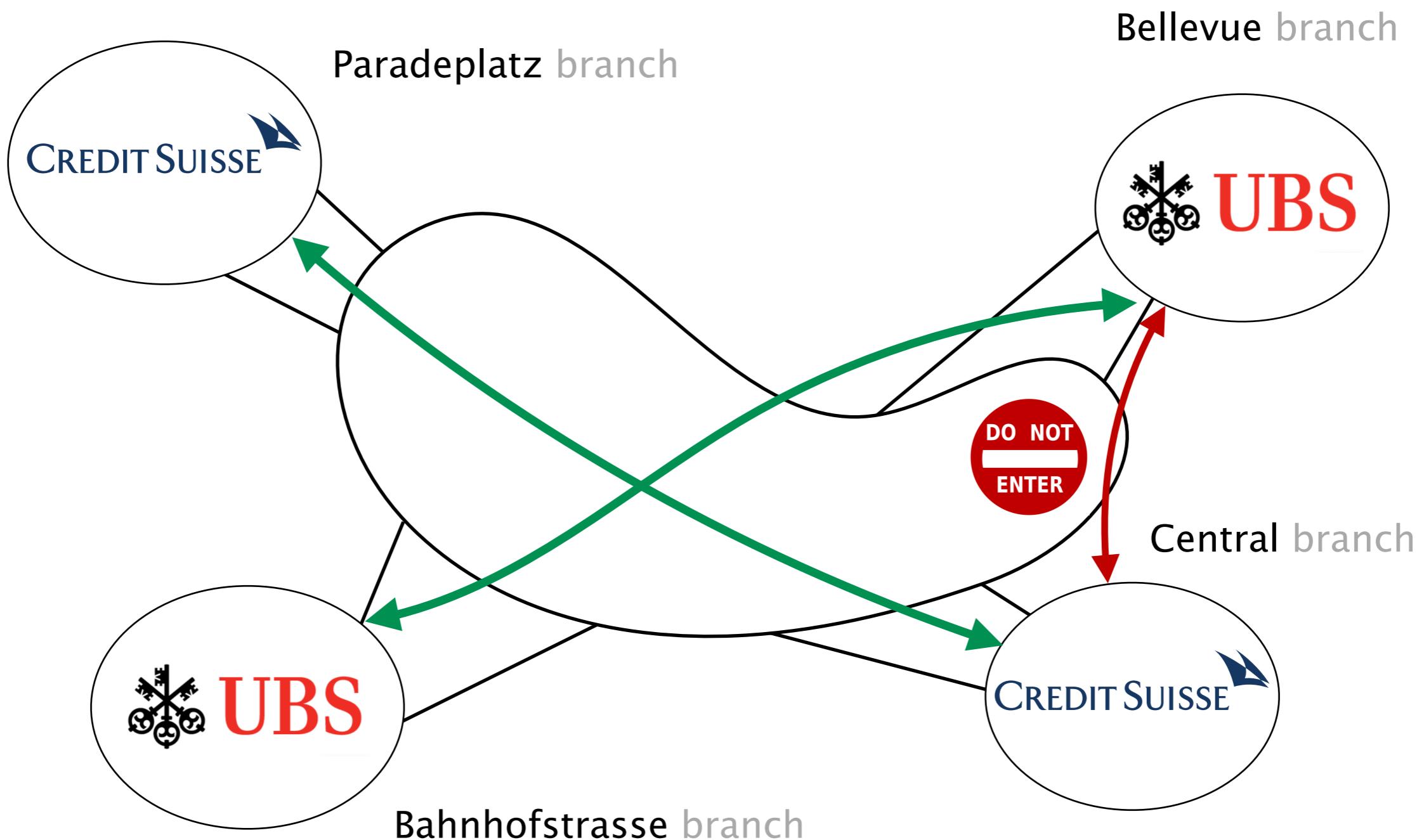
For obvious reasons...

Credit Suisse (resp. UBS) sites should *only* be able to talk to each other

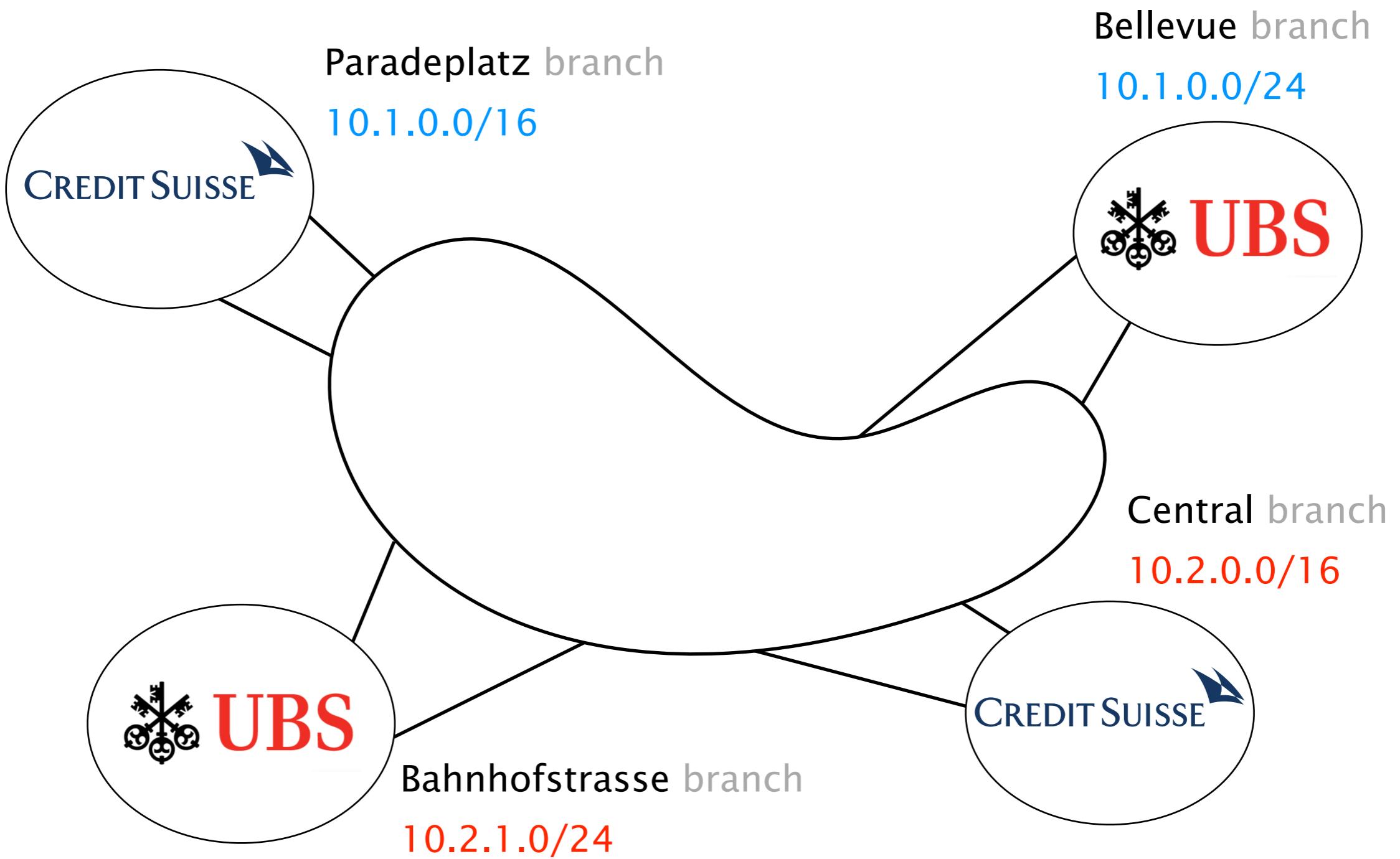


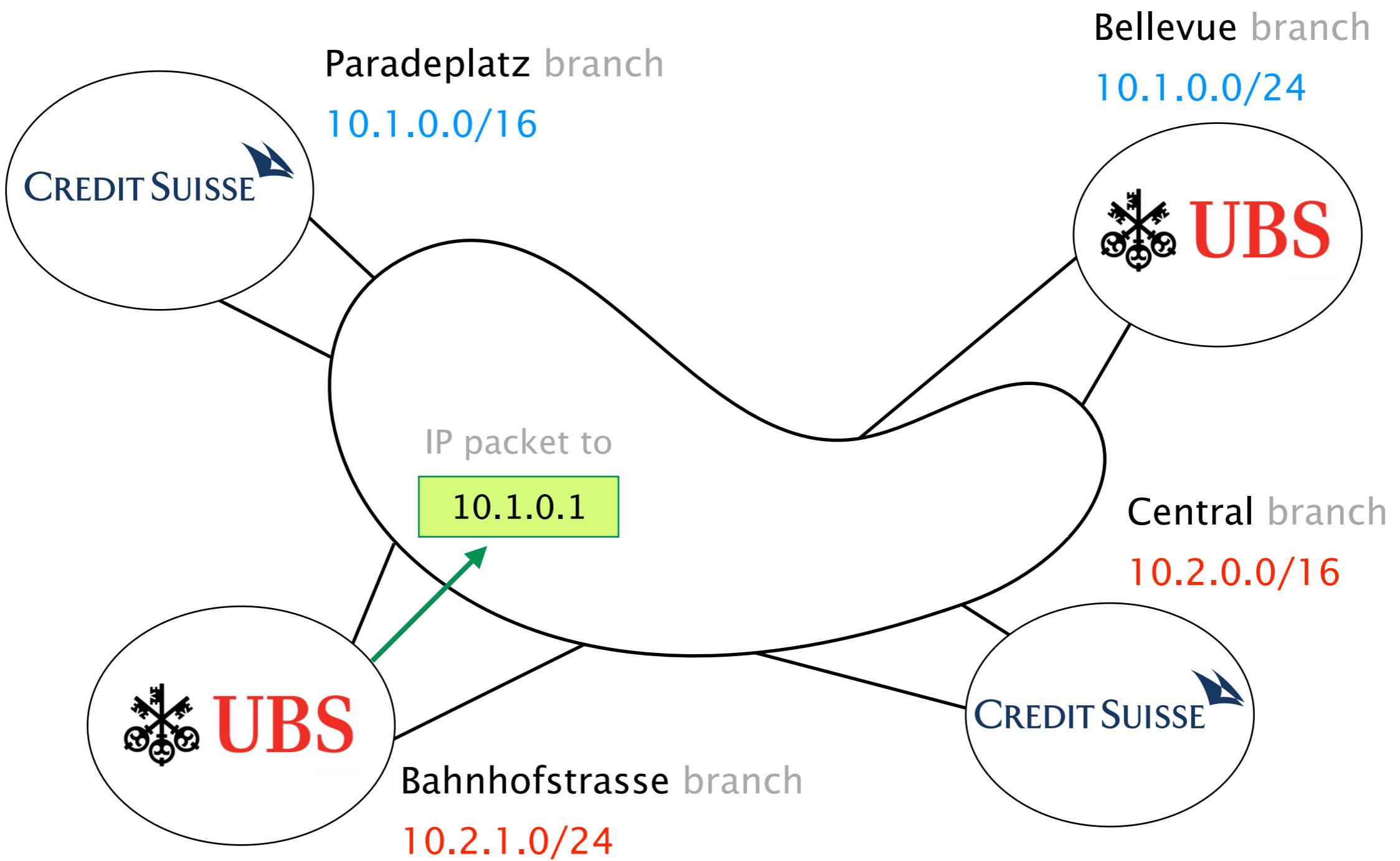
For obvious reasons...

Credit Suisse (resp. UBS) sites should *only* be able to talk to each other

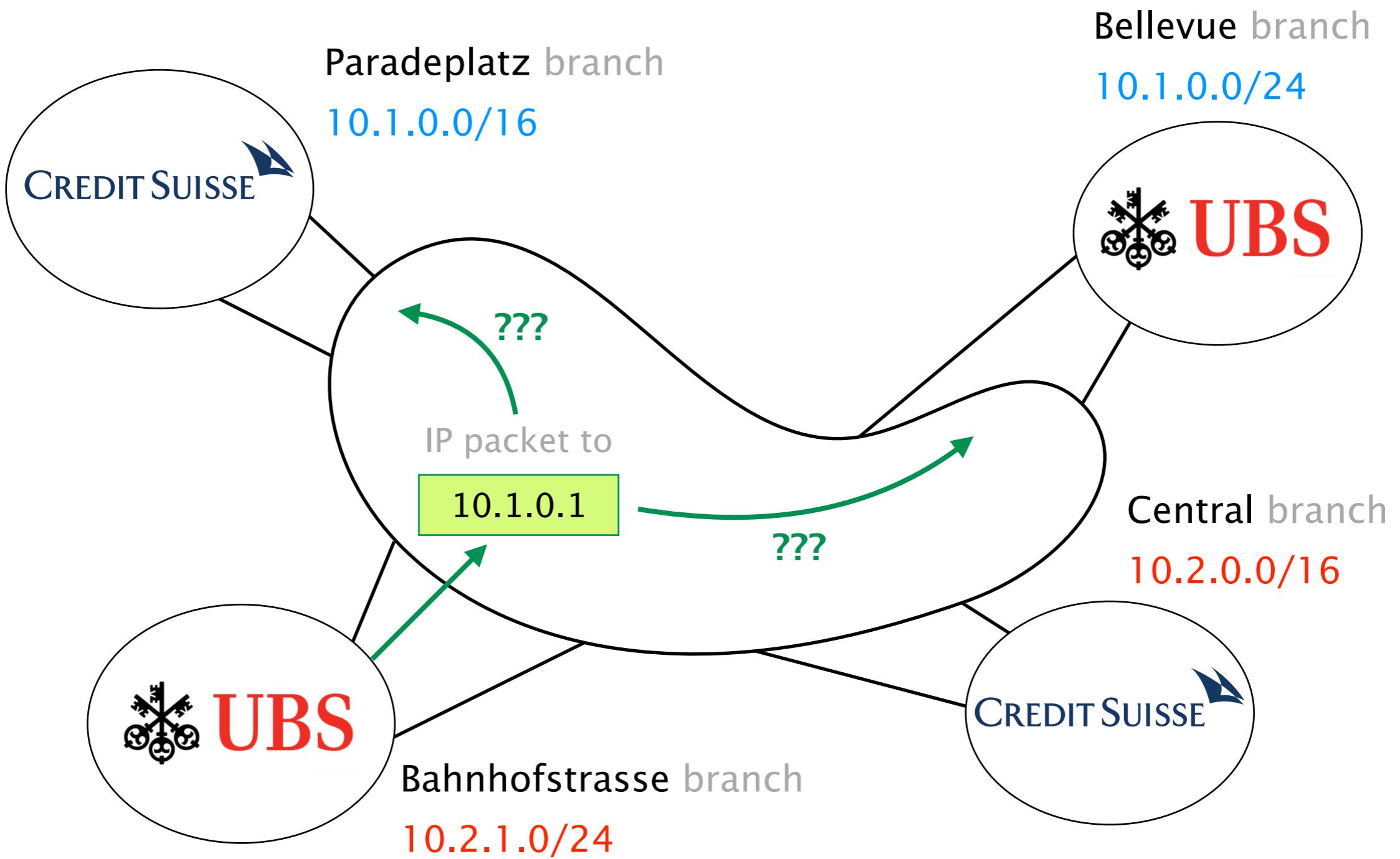


How does Swisscom enable such service given that Credit Suisse and UBS address space can overlap?

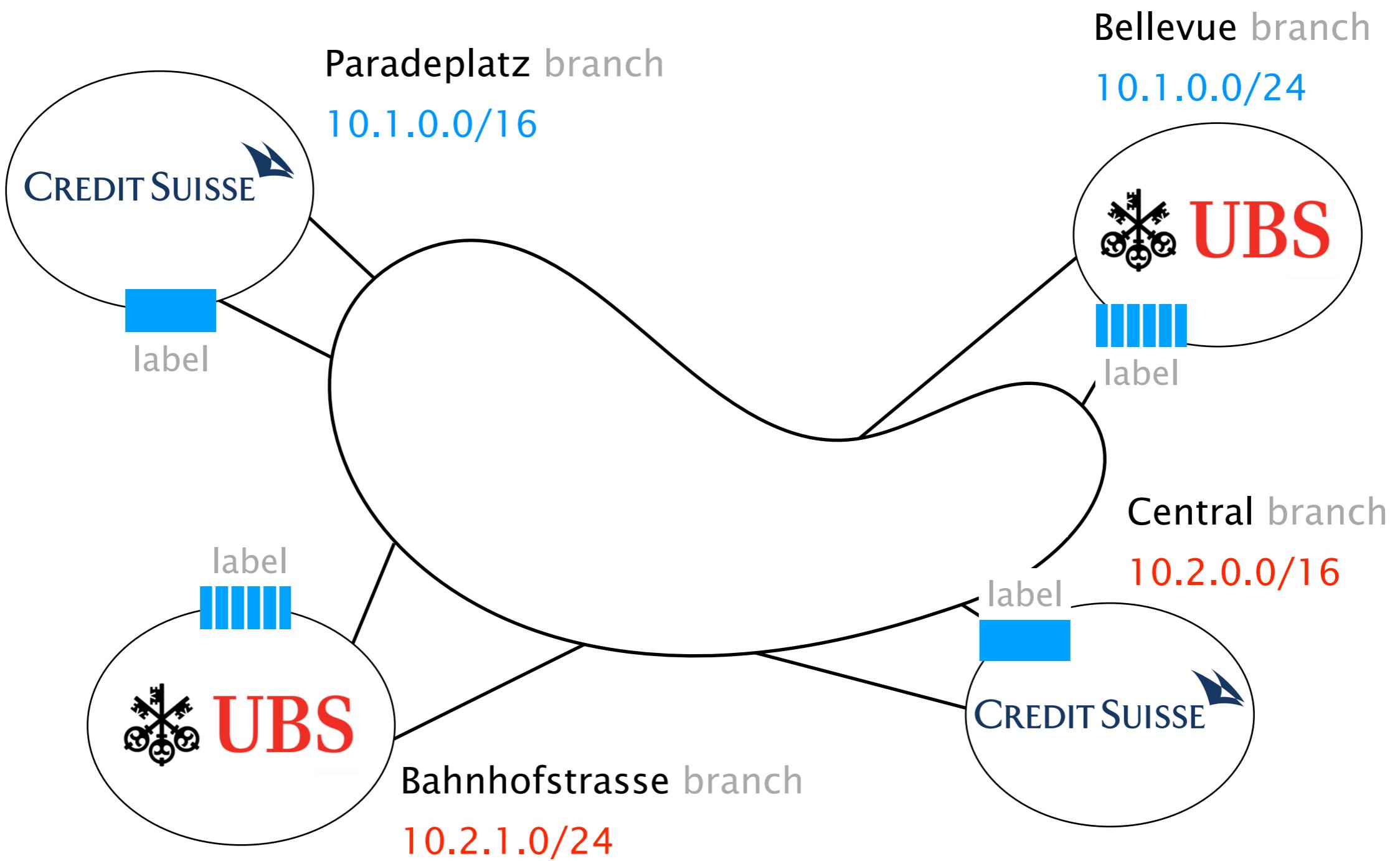


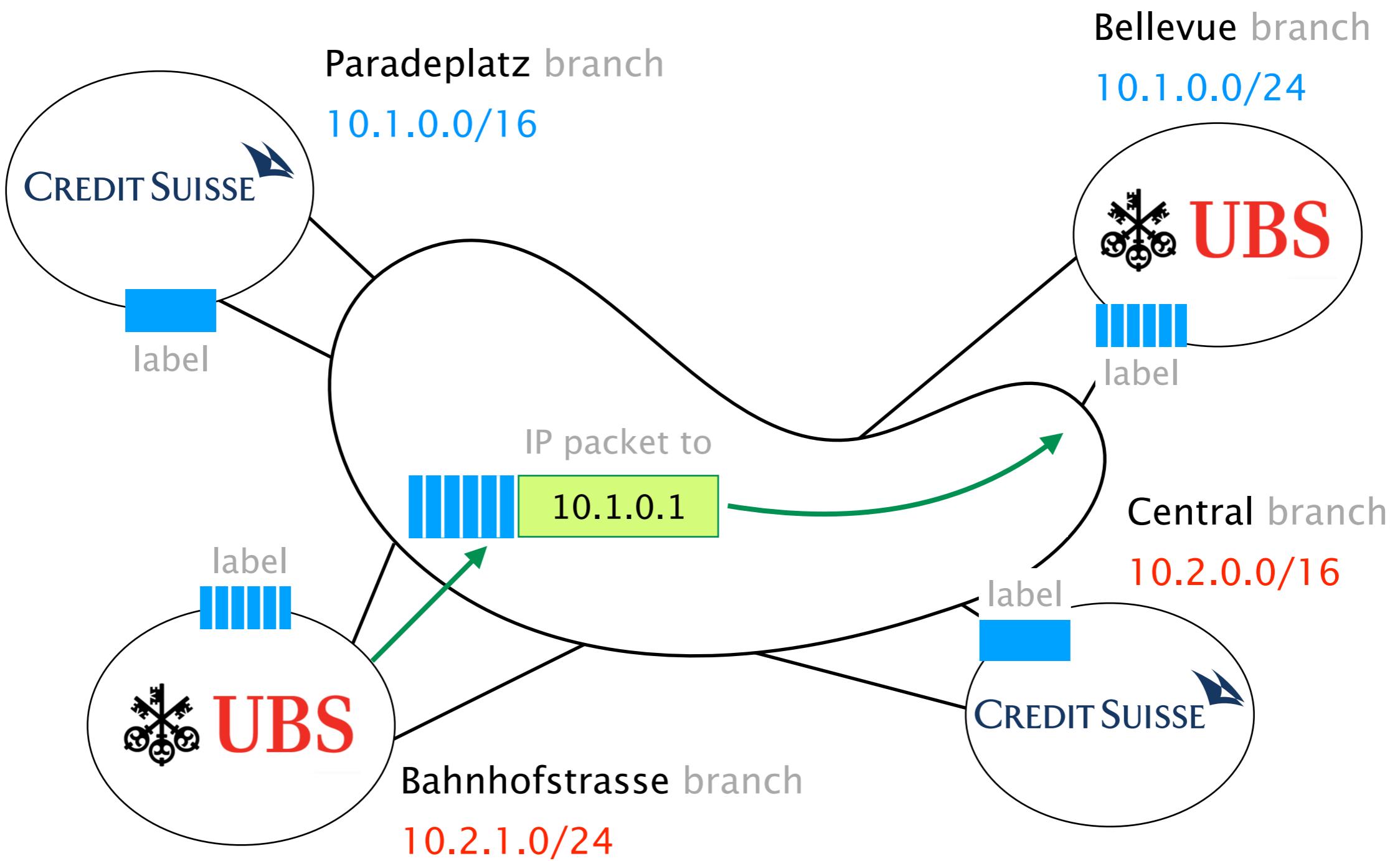


How does Swisscom enable such service given that Credit Suisse and UBS address space can overlap?



We'll see different how ISPs achieve this
using label switching again





Techniques

Performance

Traffic Engineering

Load Balancing

Quality of Service

Multicast

Flexibility

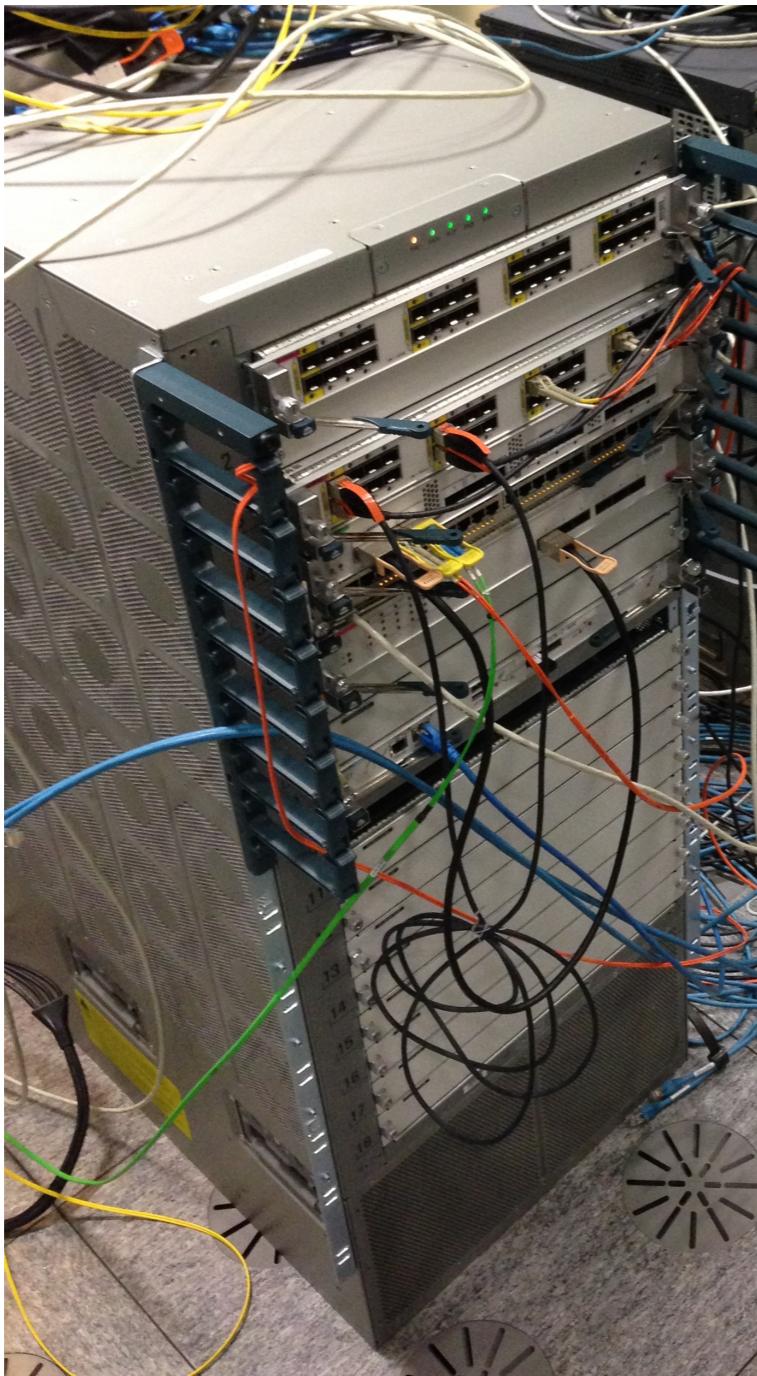
Virtual Private Networks

Reliability

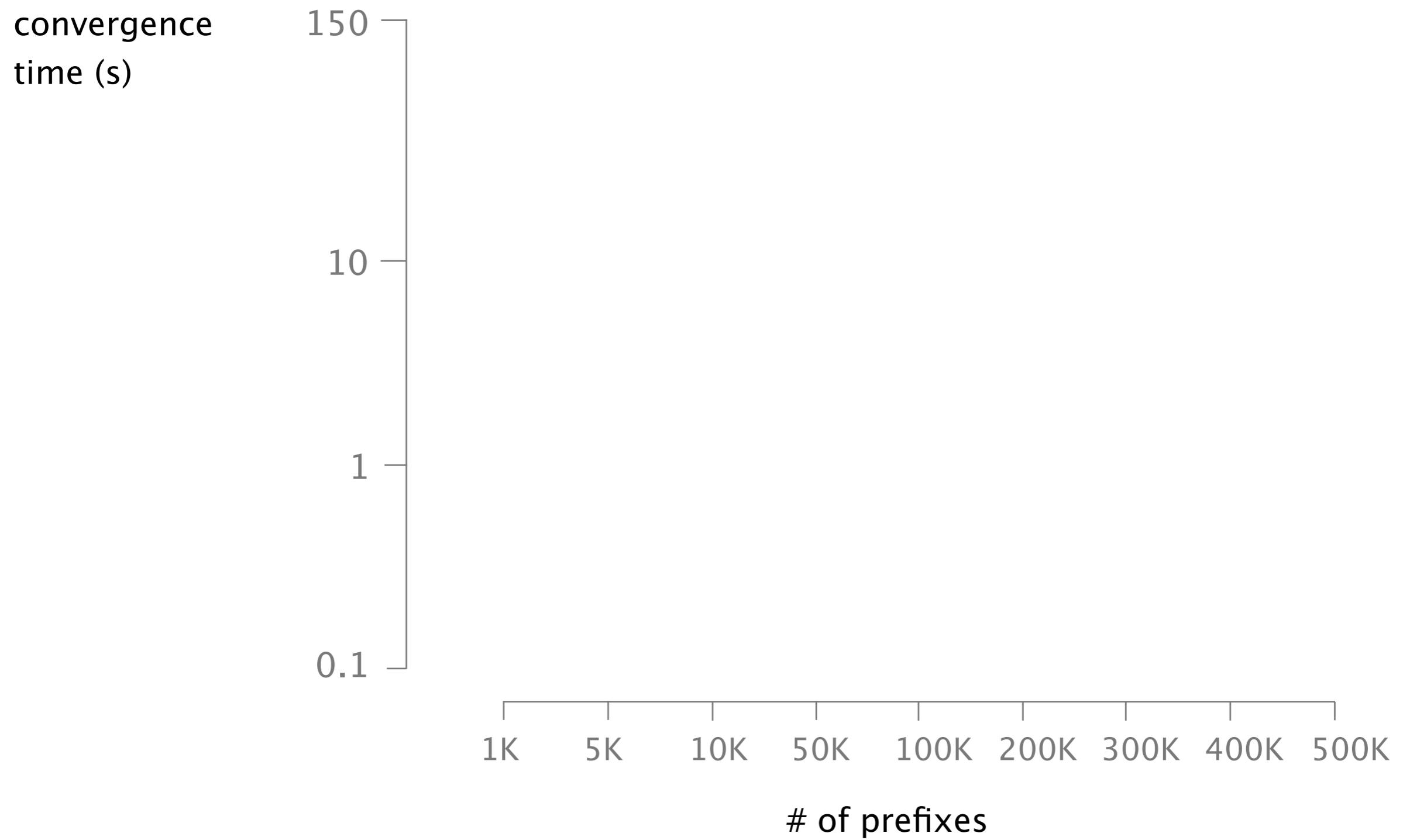
Fast Convergence

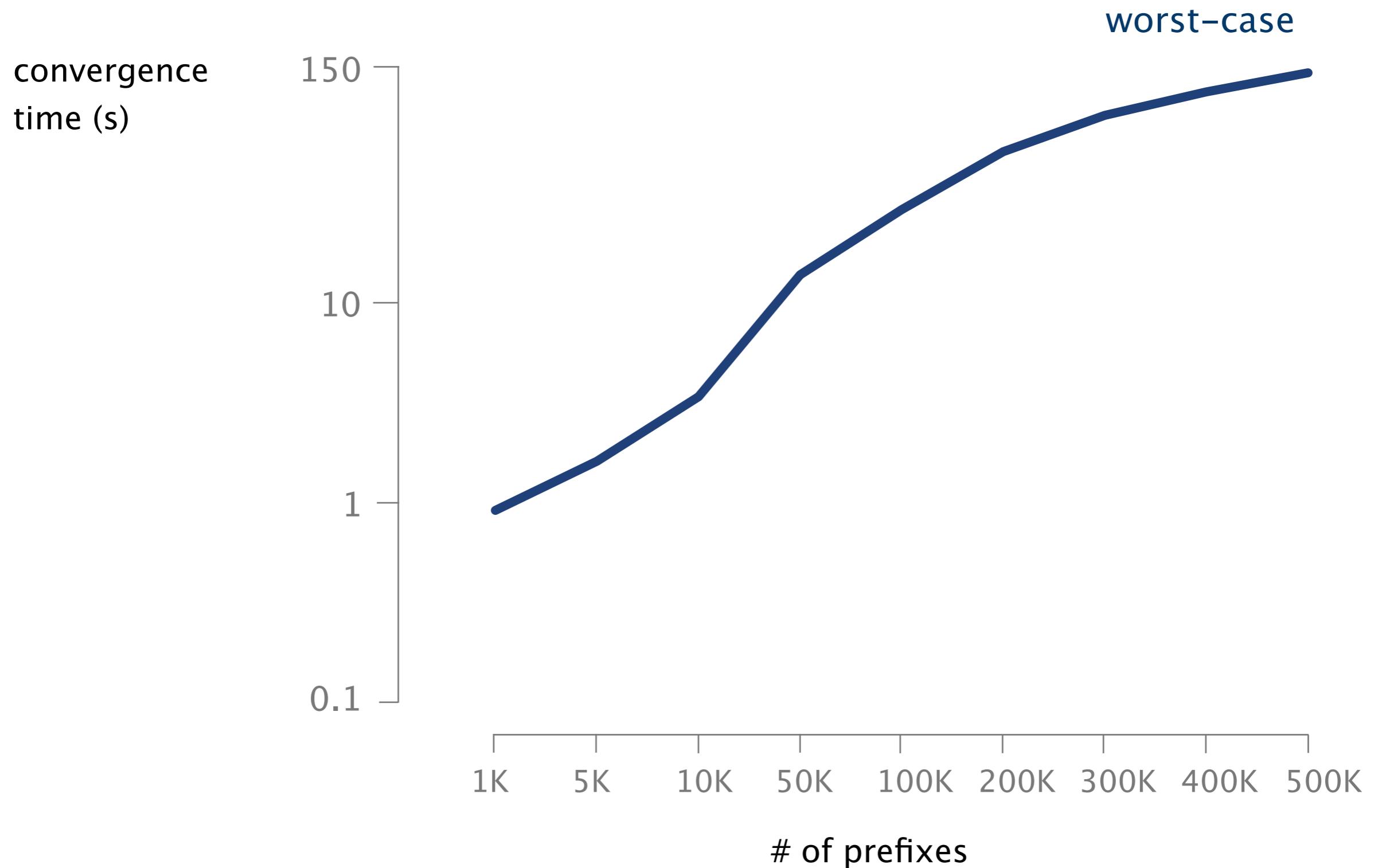
IP networks converge slowly upon failures
at least, by default

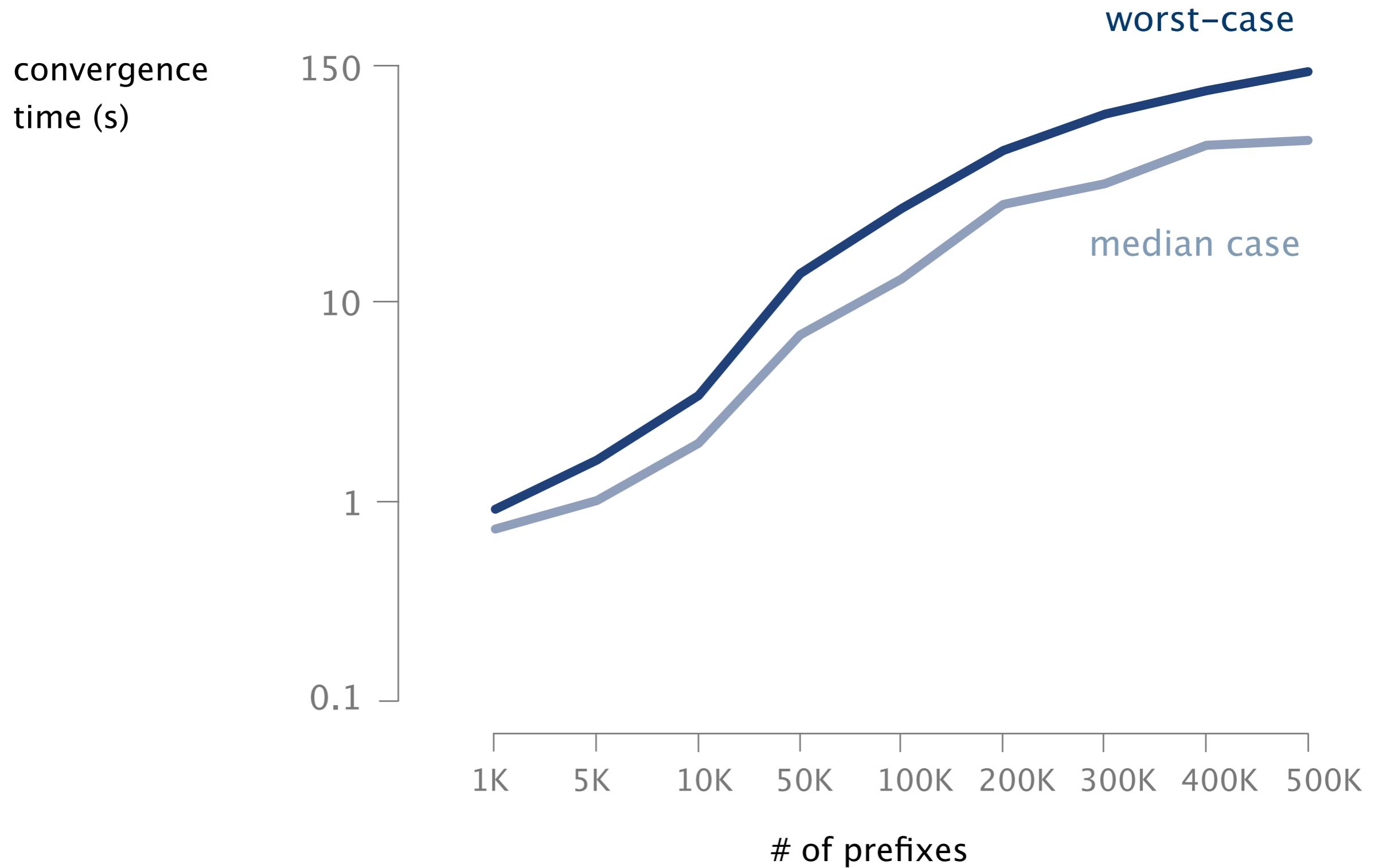
We measured how long it takes for
an ETH router to converge



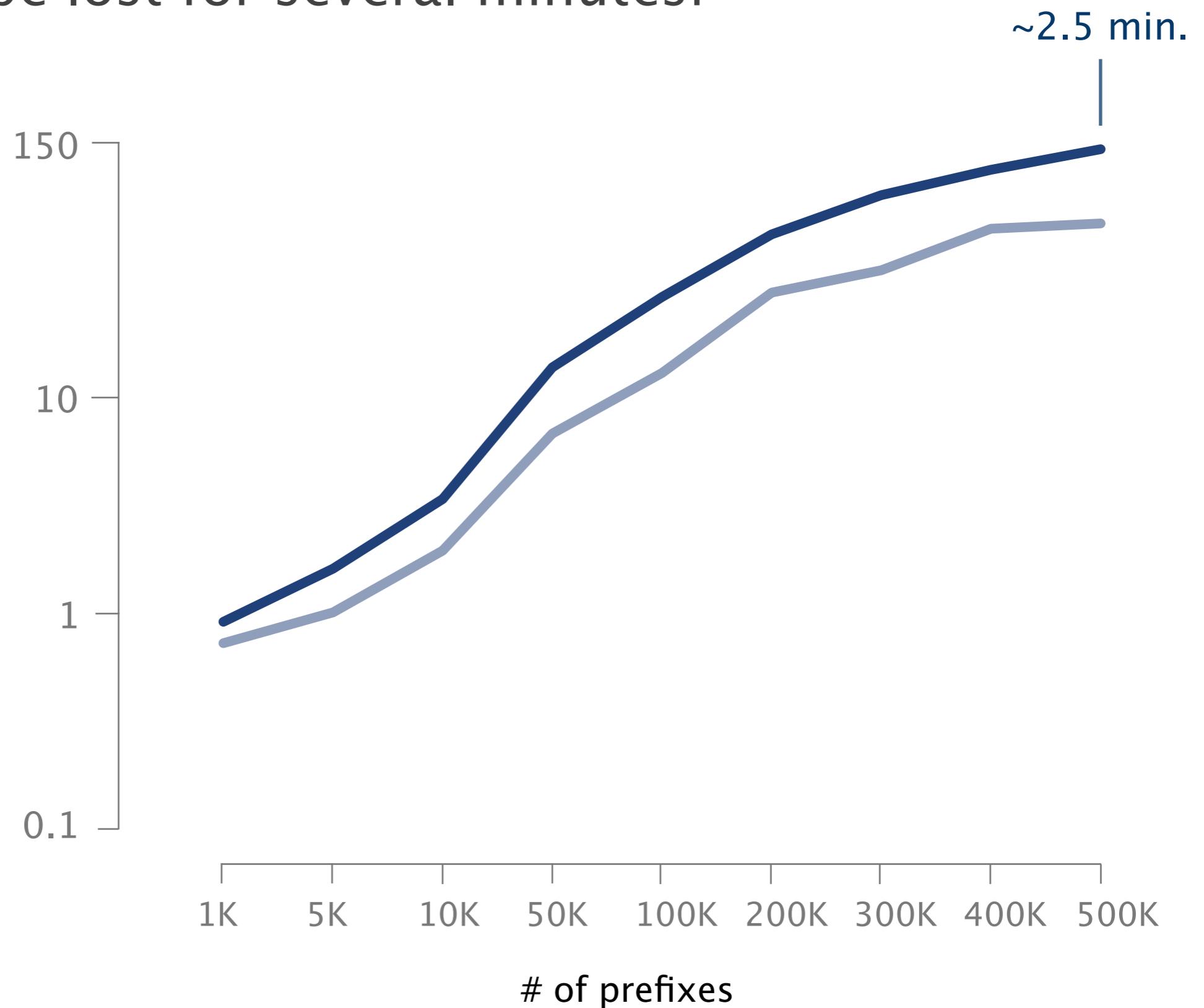
Cisco Nexus 9k
ETH's **recent routers**







Traffic can be lost for several minutes!



Why does it take so long?

	detect the failure	ms
	report the failure	ms
	generate/flood routing messages	10s of ms
	recompute routing paths	10s of ms
bottleneck	communicate next-hops to the line cards	100s of ms
	install new next-hops	10s of ms

We'll see different Fast Reroute techniques to speed up convergence time

Fast Reroute
technologies

detect the failure	hw acceleration
report the failure	
generate/flood routing messages	pre-computation
recompute routing paths	
communicate next-hops to the line cards	pre-provisionning
install new next-hops	fast activation
max. convergence time for <i>any</i> failure	
<1 sec	

Techniques

Performance

Traffic Engineering

Load Balancing

Quality of Service

Multicast

Flexibility

Virtual Private Networks

Reliability

Fast Convergence

Besides learning about the techniques per-se,
you'll learn how to operate and implement them

Besides learning about the techniques per-se,
you'll learn how to operate and implement them



operate
your own IP network
in virtual labs

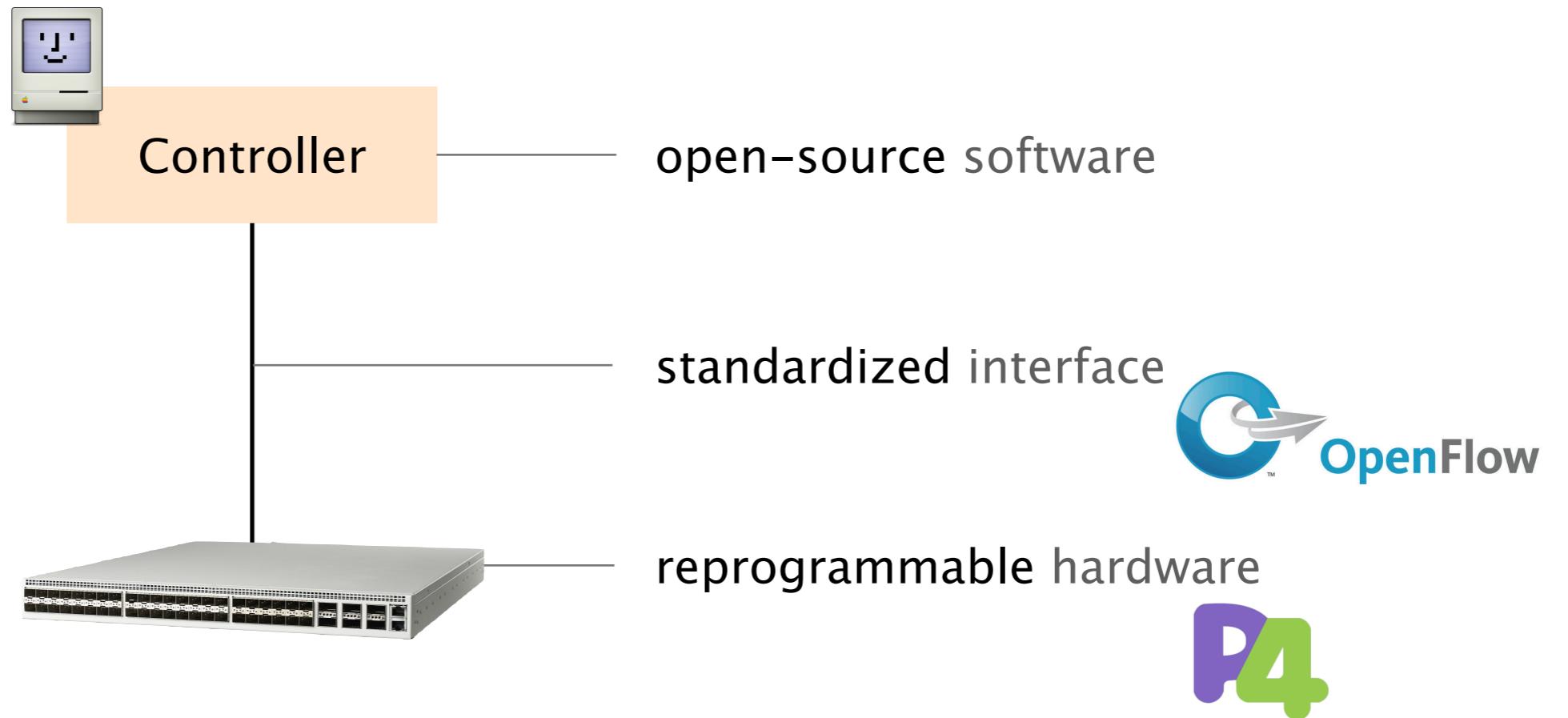


implement
your own forwarding logic
in programmable networks

Until recently, network devices tended to be completely locked down



Things are changing though as
networks are becoming programmable



Network programmability is attracting tremendous industry interest...

VMware Acquires Once-Secretive Start-Up Nicira for \$1.26 Billion

JULY 23, 2012 AT 1:25 PM PT

[Tweet](#) | [Share](#) | [+1](#) | [Share](#) | [Print](#)

VMware, the software company best known for its virtualization technology that forms the backbones of so-called cloud computing today, said it will pay \$1.26 billion for Nicira, a networking start-up that has sought to do to networks what VMware has done to computers.

The news comes on the same day that VMware was to report quarterly earnings. And while I don't usually cover VMware's earnings, I may as well mention the results: The company reported revenue for the quarter ended June rose to \$1.12 billion, while earnings on a per-share basis were 68 cents. Analysts had been expecting sales of \$1.12 billion and earnings of 66 cents.

Nicira had been running in stealth mode for quite awhile; [I got to reveal](#) its plans to the world last February.

The deal amounts to a nice payoff for Nicira's investors including Andreessen Horowitz, Lightspeed Venture Partners and NEA, as well as VMware founder Diane Greene and venture capitalist Andy Rachleff.



With \$600M Invested in SDN Startups, the Ecosystem Builds



Scott Raynovich, June 10, 2014

[Tweet](#) [in](#) [f](#) [g+](#) [d](#)



More than \$600 million has been invested in at least two dozen [software-defined networking \(SDN\)](#) startups so far, according to Rayno Report research. You can expect that to continue to climb. With the SDN ecosystem starting to take hold with a broad range of alliances and distribution partnerships, we're just getting started.

The [Arista IPO will help build visibility](#) for next-generation, software-driven networking. But Arista is selling its own hardware and is not an SDN pure-play. A new line of [SDN startups](#), with a more radical approach to software-based networking, is building momentum. These newer SDN startups are just getting their gear into customers' hands and starting to build sales channels, so you can expect a long revenue ramp.

This excitement is boosting startup valuations, according to [Rayno Report research](#). There are now at least ten [SDN startups](#) with valuations over \$100 million. As I reported in April, a recent investment in [Cumulus Networks](#) pushed up the valuation of the private company north of \$300 million, according to industry sources. [Big Switch](#), which did a deal in 2012 valuing it near \$170 million, took money from [Intel](#) in 2013, most likely boosting its valuation to over \$200 million, according to several sources.

Related Articles

[How to Effectively Embed SDN in the Enterprise](#)

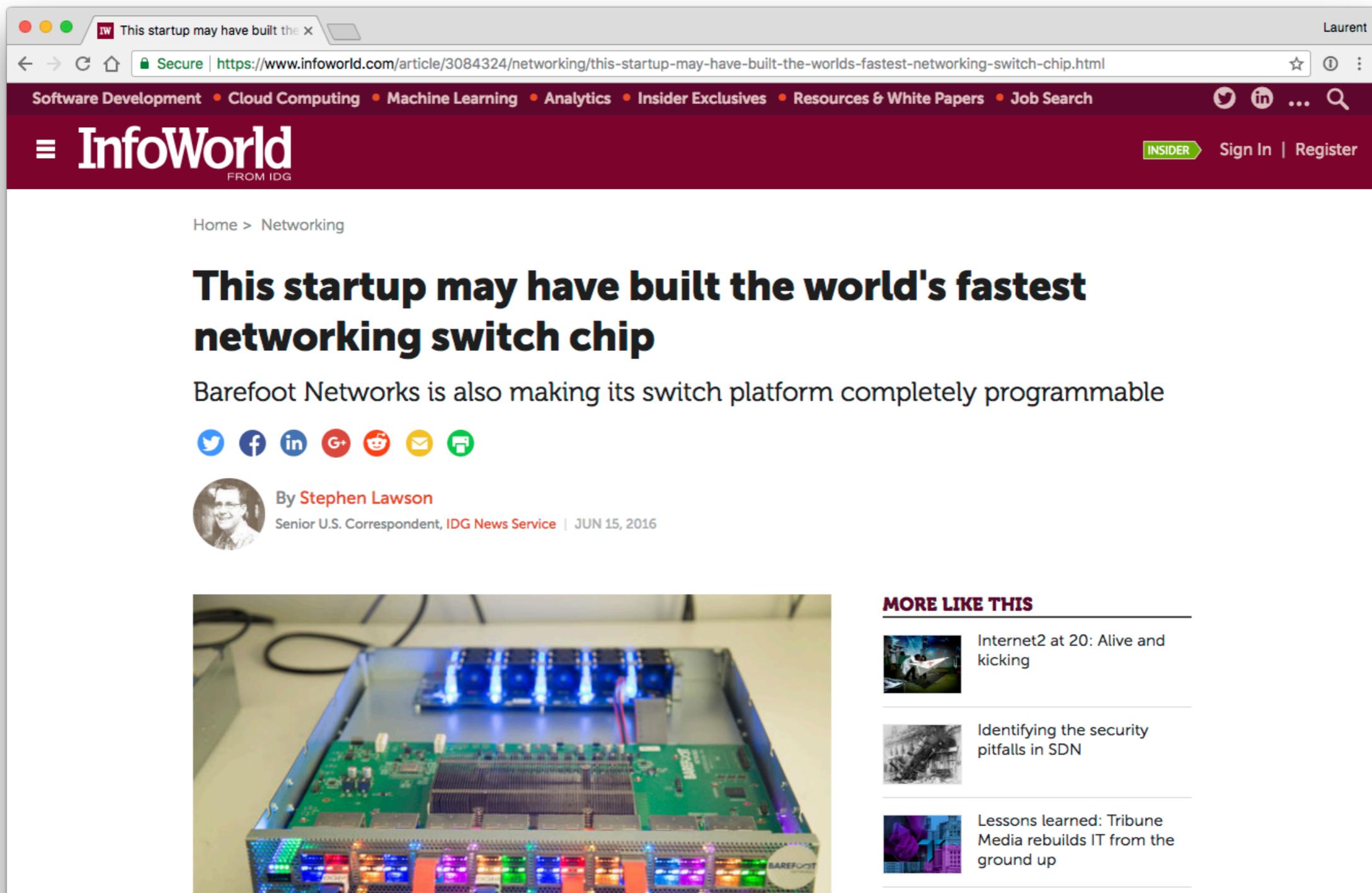
[NFV and SDN: What's the Difference Two Years Later?](#)

[sFlow Creator Peter Phaal On Taming The Wilds Of SDN & Virtual Networking](#)

[Featured Article: Bringing Data-Driven SDN to the Network Edge](#)

[NFV Delivers Pervasive Intelligence for MNOs](#)

Barefoot Networks (Stanford startup) started to produce re-programmable network hardware in 2013



A screenshot of a web browser displaying an InfoWorld article. The title of the article is "This startup may have built the world's fastest networking switch chip". The article is by Stephen Lawson, a Senior U.S. Correspondent for IDG News Service, dated June 15, 2016. The page includes social sharing icons for Twitter, Facebook, LinkedIn, Google+, Reddit, Email, and Print. Below the article is a photograph of a Barefoot Networks network switch hardware, showing its green circuit boards and blue LED status lights. To the right of the article is a sidebar titled "MORE LIKE THIS" featuring three other news articles with small thumbnail images.

This startup may have built the world's fastest networking switch chip

By [Stephen Lawson](#)
Senior U.S. Correspondent, IDG News Service | JUN 15, 2016

[Twitter](#) [Facebook](#) [LinkedIn](#) [Google+](#) [Reddit](#) [Email](#) [Print](#)

MORE LIKE THIS

-  Internet2 at 20: Alive and kicking
-  Identifying the security pitfalls in SDN
-  Lessons learned: Tribune Media rebuilds IT from the ground up

In June 2019,
Barefoot Networks was acquired by Intel

THE WALL STREET JOURNAL.

Europe Edition ▾ | September 22, 2019 | Print Edition | Video

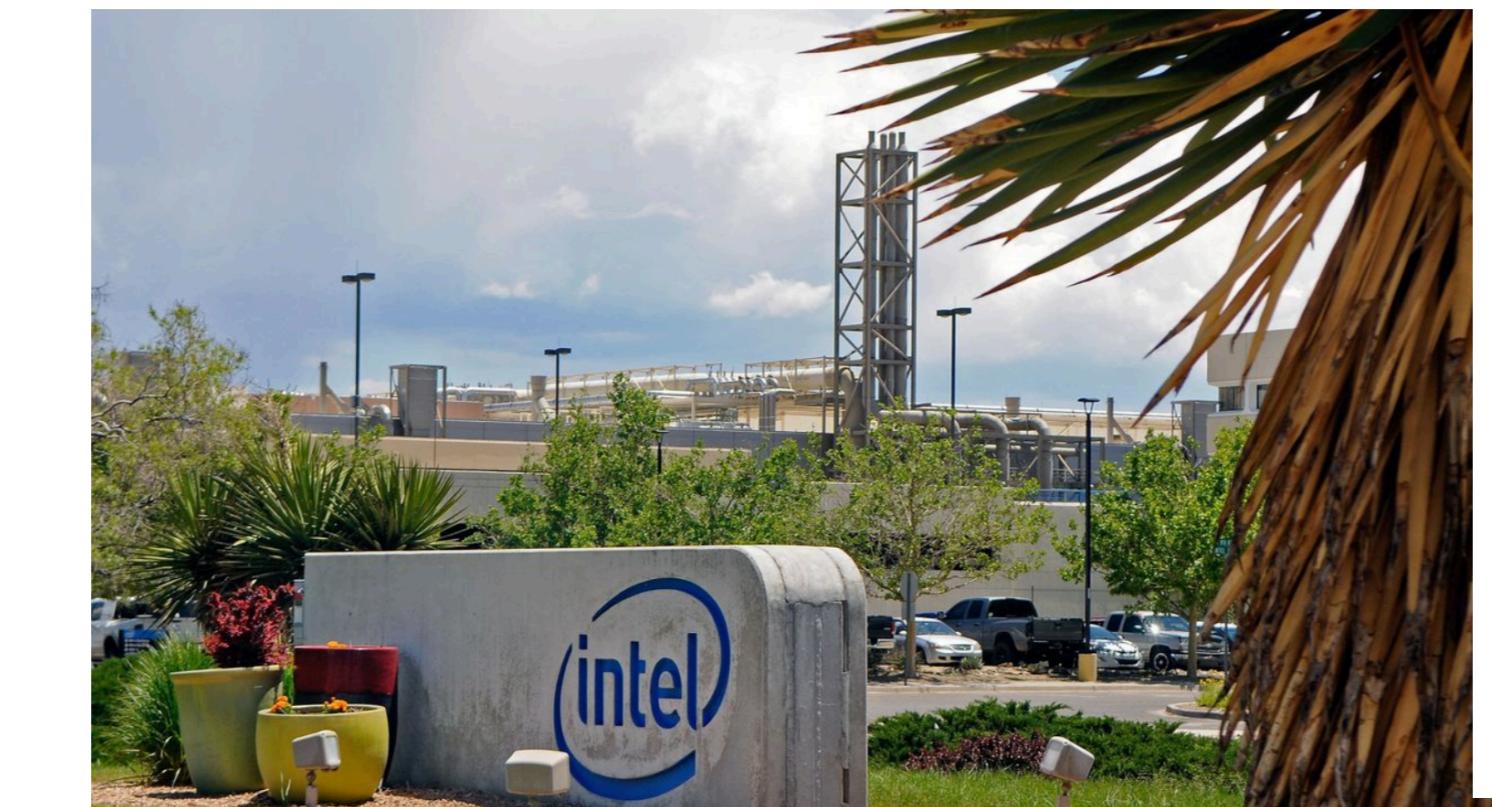
Home World U.S. Politics Economy Business Tech Markets Opinion Life & Arts Real Estate WSJ. Magazine

SHARE

TECH

Intel Agrees to Acquire Networking Startup Barefoot Networks

Barefoot Networks is backed by Google, Alibaba, Tencent and Goldman Sachs



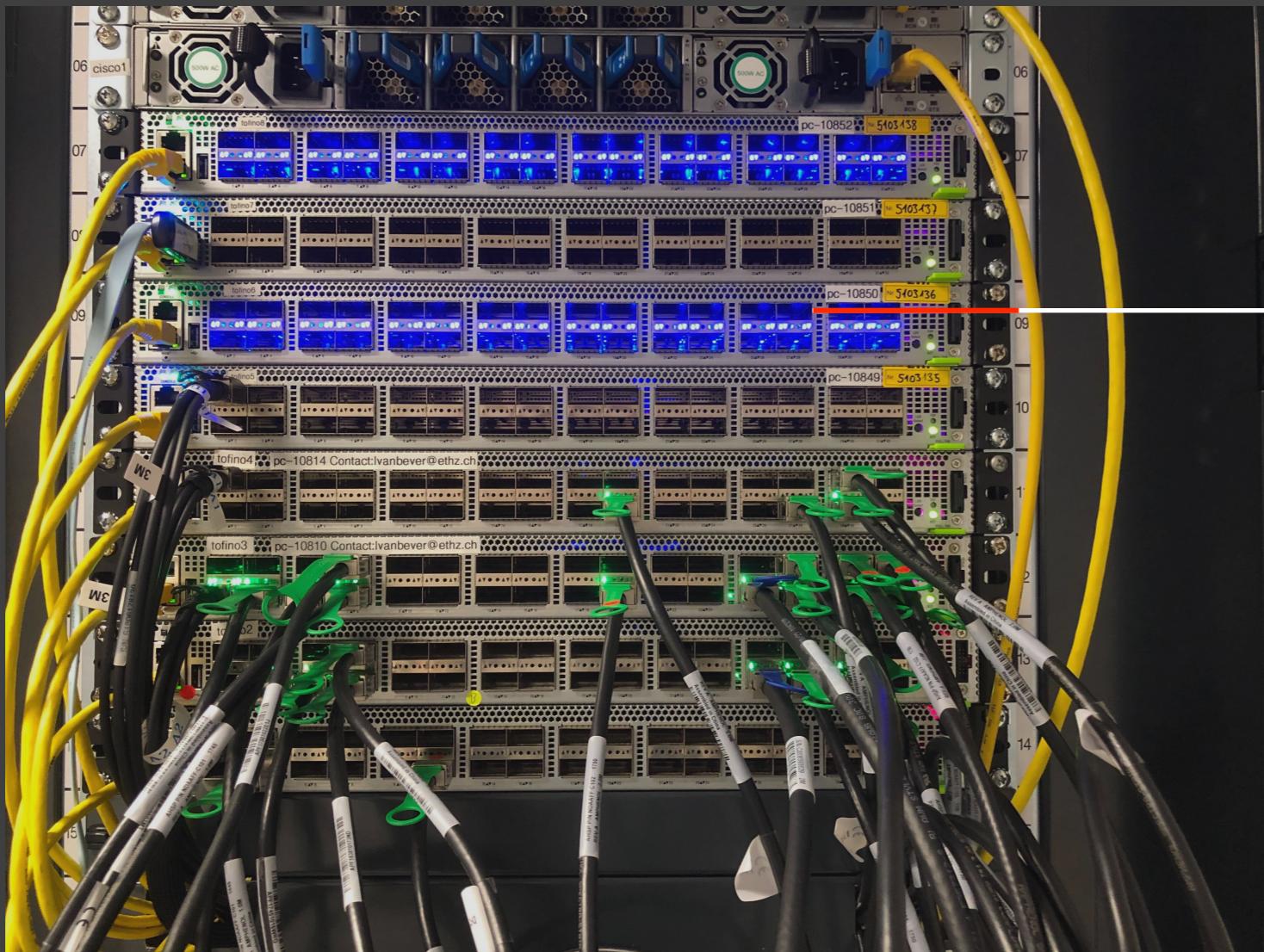
P4 is a domain-specific language which describes how a switch should process packets



<https://p4.org>

A sneak peek at our own networking lab

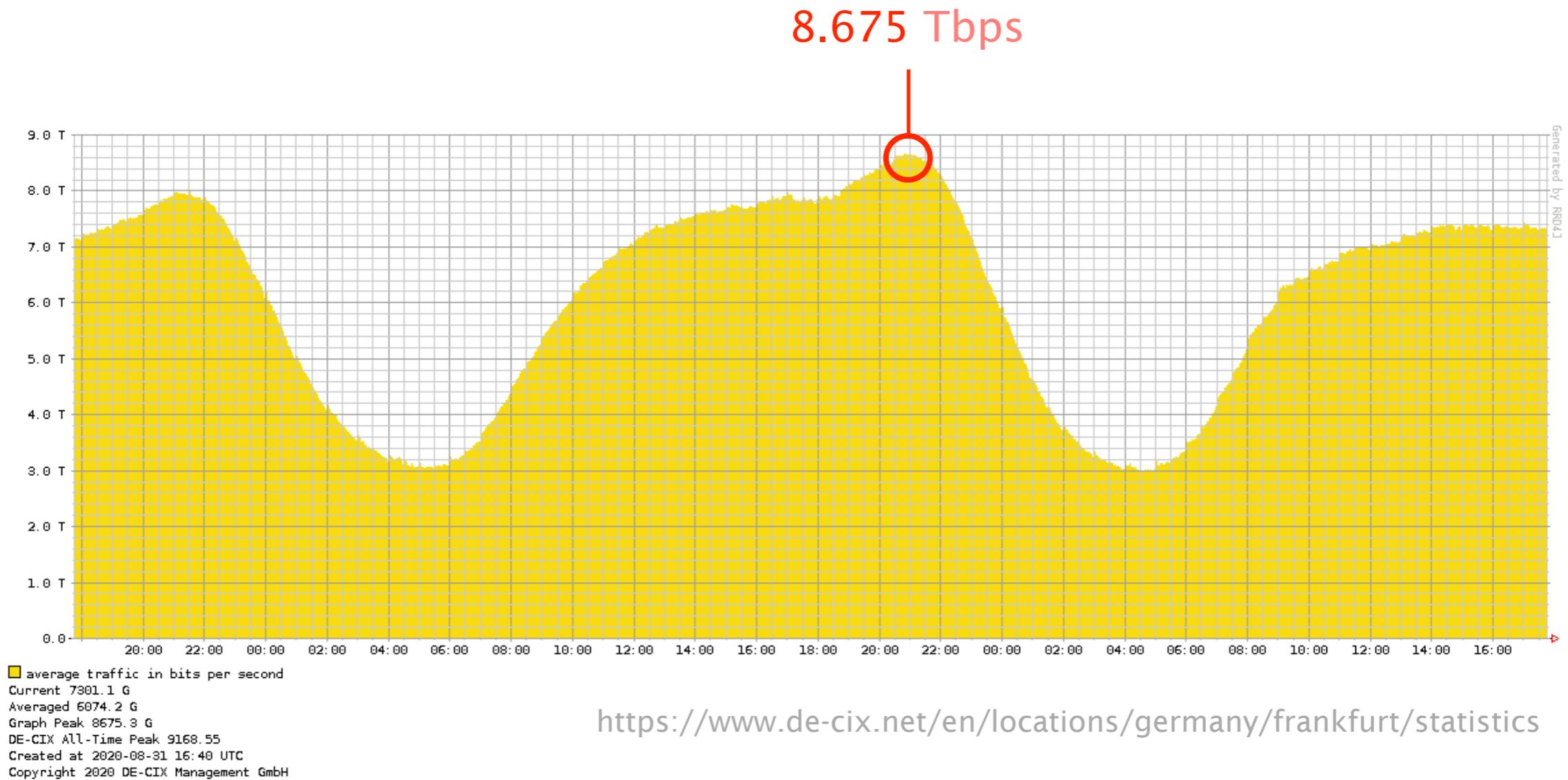
8x Wedge100BF-32X. Total capacity: 25.6 Tbps



32x QSFP28 ports
25/40/50/100 GbE

8x Wedge100BF-32X. Total capacity: **25.6 Tbps**

~3x what DE-CIX sees at peak time!



Course Organization

The course is gonna be divided in two blocks

Lectures/Exercices

~10 weeks

Group project

~4 weeks

in teams of 3

The course is gonna be divided in two blocks

Lectures/Exercices

~10 weeks

Group project

~4 weeks

in teams of 3

There will be 2h of lectures & 2h of exercises

Tue 14-16 Lecture

Tue 16-18 Practical exercises

Exercises are *not* graded *but* will help for the project & exam

Both will take place **online**

The course is gonna be divided in two blocks

Lectures/Exercices

~10 weeks

Group project

~4 weeks

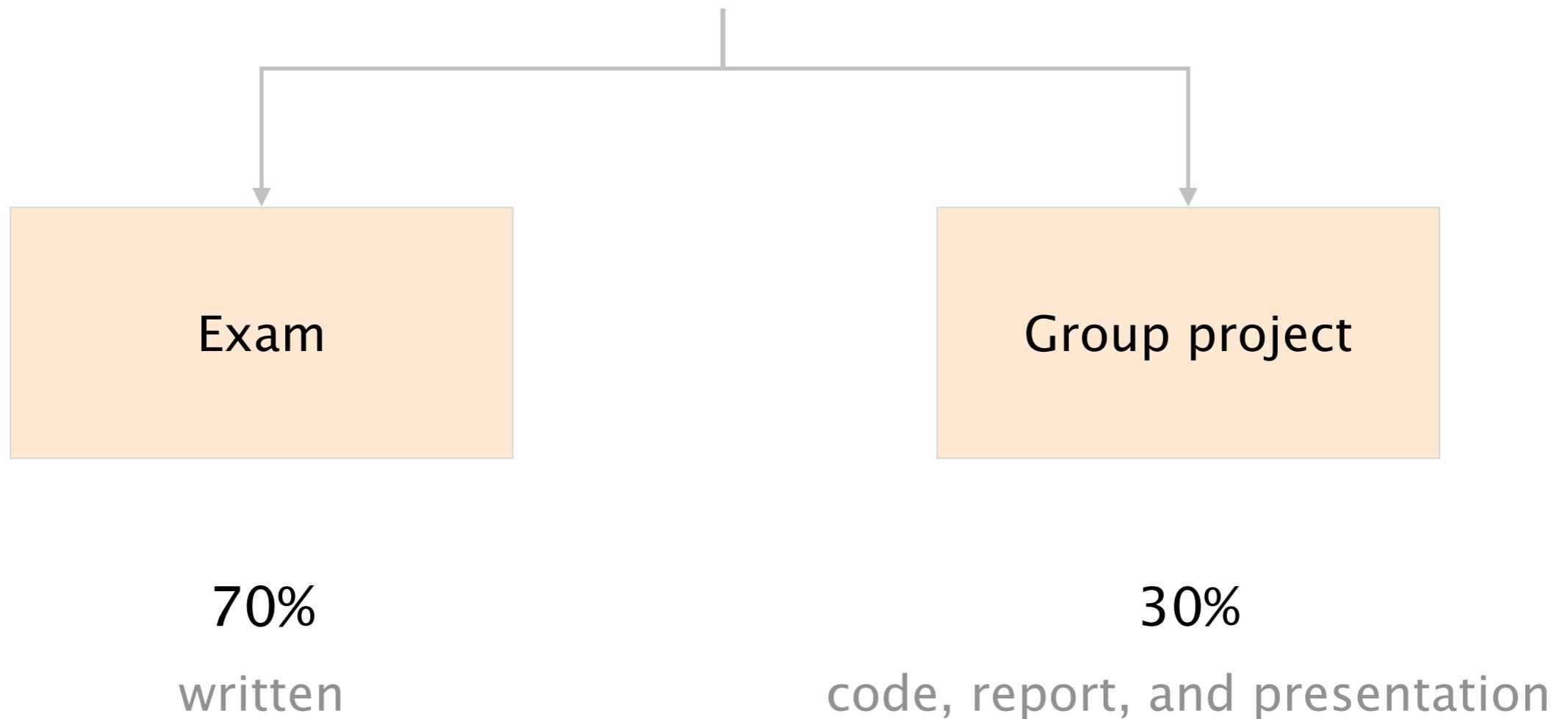
in teams of 3

In the project,
you'll partake in a class-wide challenge

Make *your* network the *fastest one*

We'll provide more information in the coming weeks
and, yes, *there will be prizes* for the best teams!

Your final grade



Exam

70%
written

Examples (more in the lecture)

How would you solve
problem <X>?

How would you optimize
this network design for <X>?

Is the P4 program <X> correct?

... important to do the exercises

Your dream team of teaching assistants



Romain
head TA



Edgar



Roland



Thomas



Maria



Albert



Alexander



Ege

+ Eric, Yannick
followed lecture last year

We'll use Moodle as course platform check it out regularly

Advanced Topics in Communication Networks

Dashboard / My courses / Advanced Topics in Communication Networks



In previous courses, you have learned about the general Internet architecture and the main principles that make it work. But the reality is somewhat more complex. There are many ways to configure, use, and optimize networks; this is what we are going to cover in this course. This year, we will focus on the following questions.

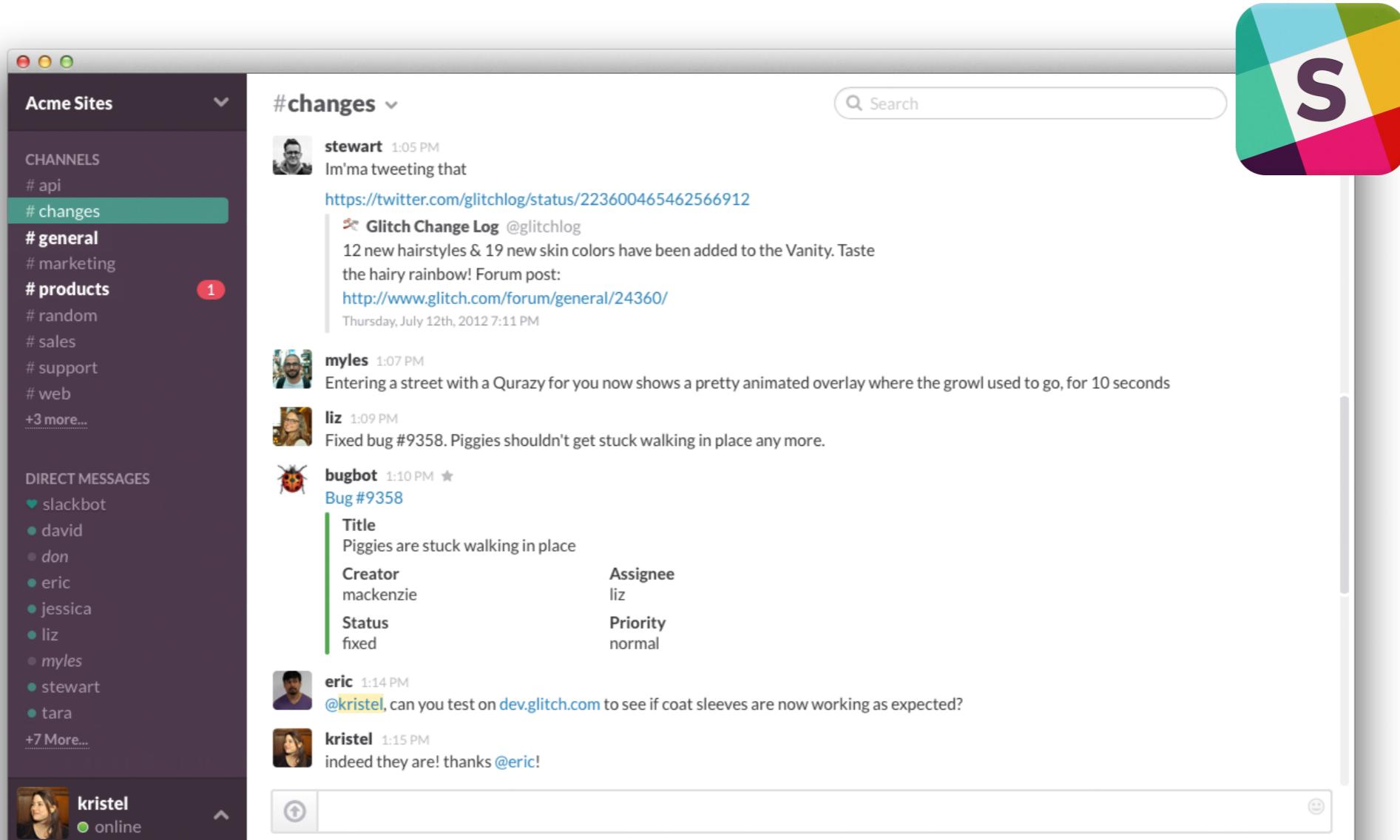
How do we optimize the

- performance,
- flexibility, and
- reliability

of (large) network infrastructures?

By the end of this course, you will know some of the most common methods used to address these questions and you will be able to apply them to configure an actual network, as if you were an actual network operator.

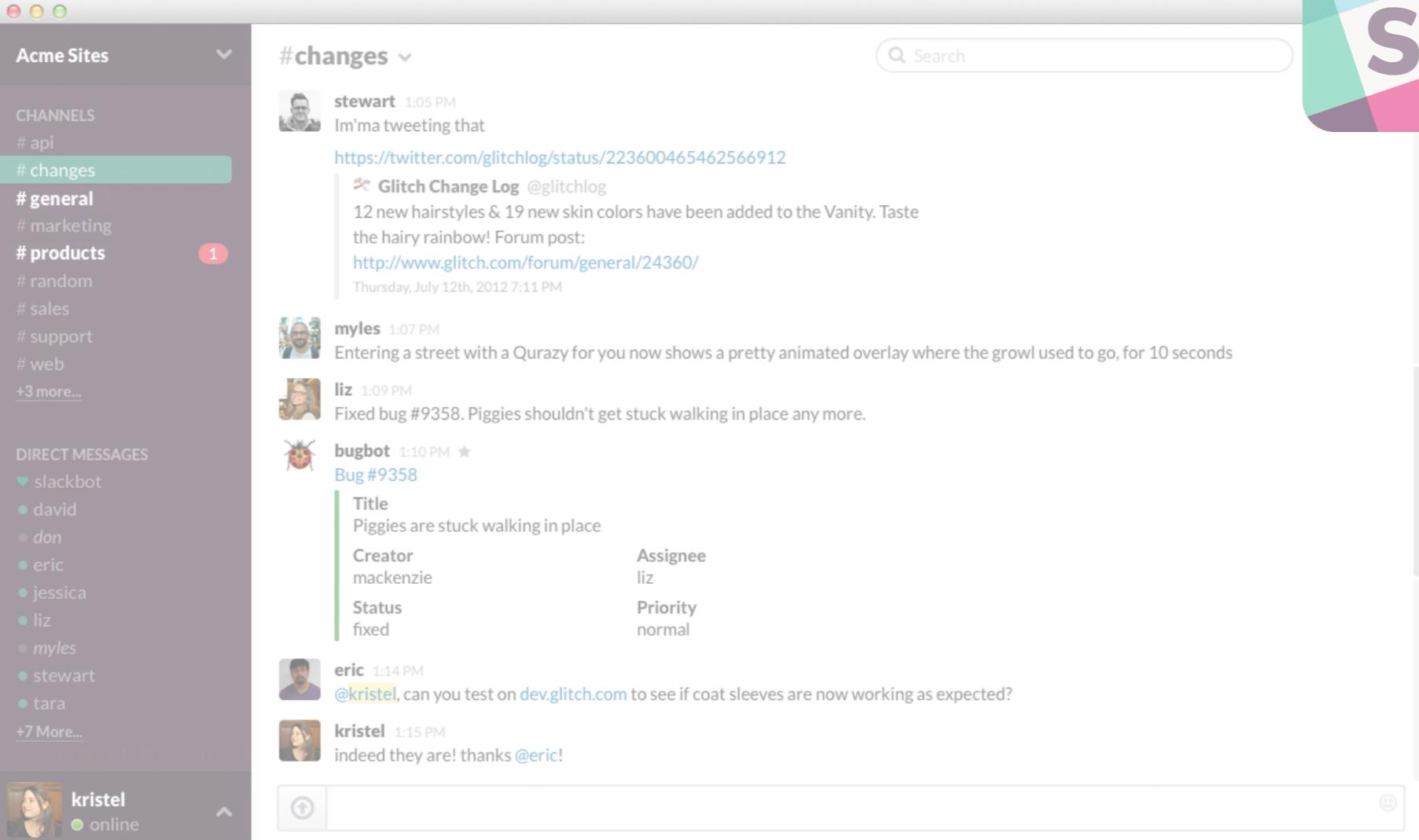
We'll use Slack to discuss about the course, exercises, and projects



Web, smartphone and desktop clients available

Register today using your real name

<https://adv-net20.slack.com/signup>



The screenshot shows the Slack desktop application interface. On the left, there's a sidebar with a team dropdown set to "Acme Sites". Under "CHANNELS", the "# changes" channel is selected, highlighted with a green bar. Other channels listed include # api, # general, # marketing, # products (with 1 unread message), # random, # sales, # support, and # web. Under "DIRECT MESSAGES", several users are listed: slackbot, david, don, eric, jessica, liz, myles, stewart, tara, and +7 More... At the bottom of the sidebar, a user named kristel is shown as online.

The main window displays the "# changes" channel feed. A message from stewart at 1:05 PM says "Im'ma tweeting that" with a link to <https://twitter.com/glitchlog/status/223600465462566912>. Below it, a message from @Glitch Change Log at 1:09 PM discusses new hairstyles and skin colors added to the game. A message from myles at 1:07 PM mentions an animated overlay for Qurazy. liz fixed bug #9358 regarding piggies getting stuck. A bug report from bugbot at 1:10 PM details a piggie stuck walking issue, assigned to liz with priority normal. Eric and kristel then discuss testing coat sleeves on dev.glitch.com.

A colorful icon with a large letter "S" is positioned in the top right corner of the main window area.

Web, smartphone and desktop clients available

Should ***you*** take this course?

If you like computer networks, ***yes!***

that said...

You shouldn't take the course if

- you *hate* programming
- you can't work during the semester
- you expect 10+ years of exam history

We've almost completely revamped the lecture,
some things will def. need adjusting. ***Please say so!***



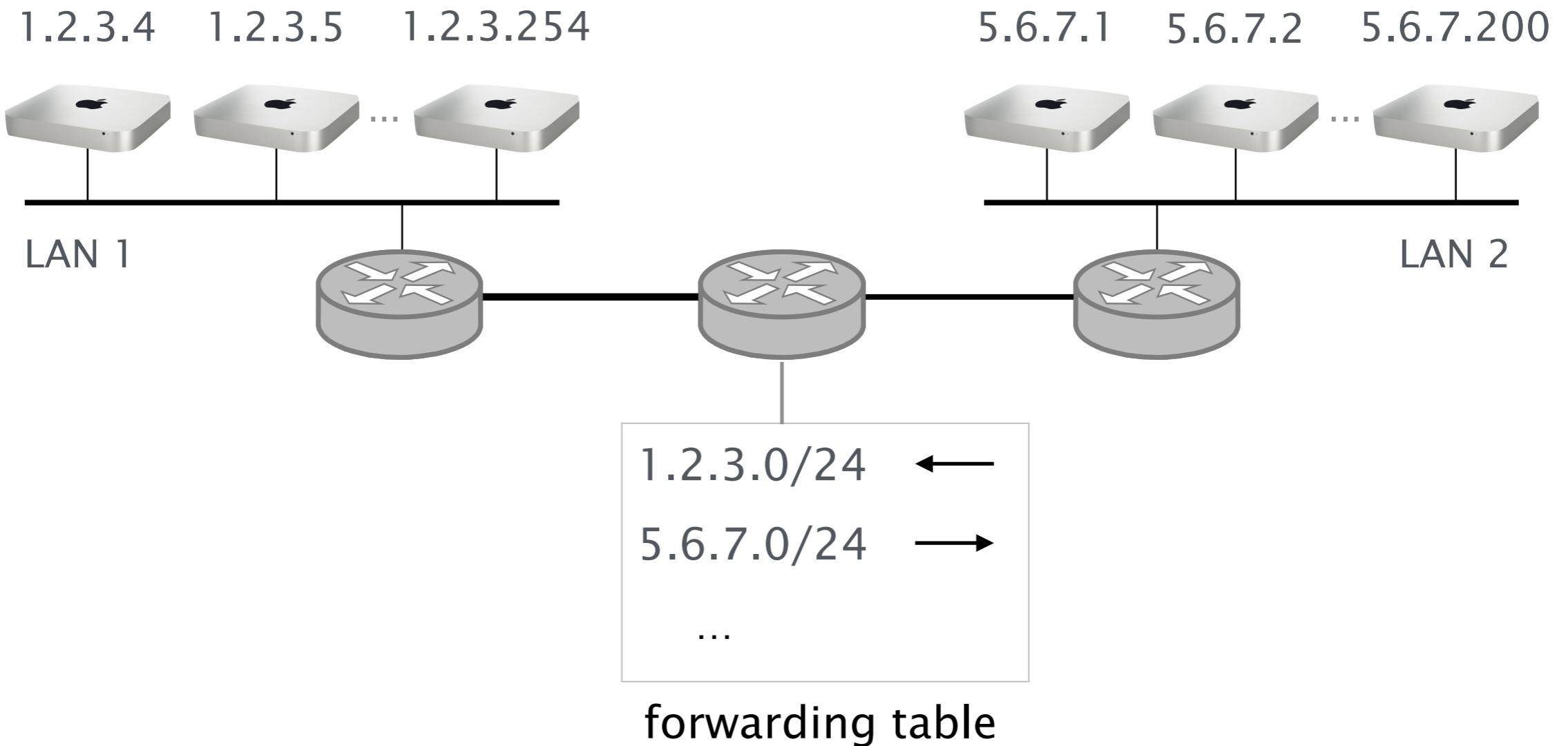
We do our best to take your feedback into account... so shoot!

Your first



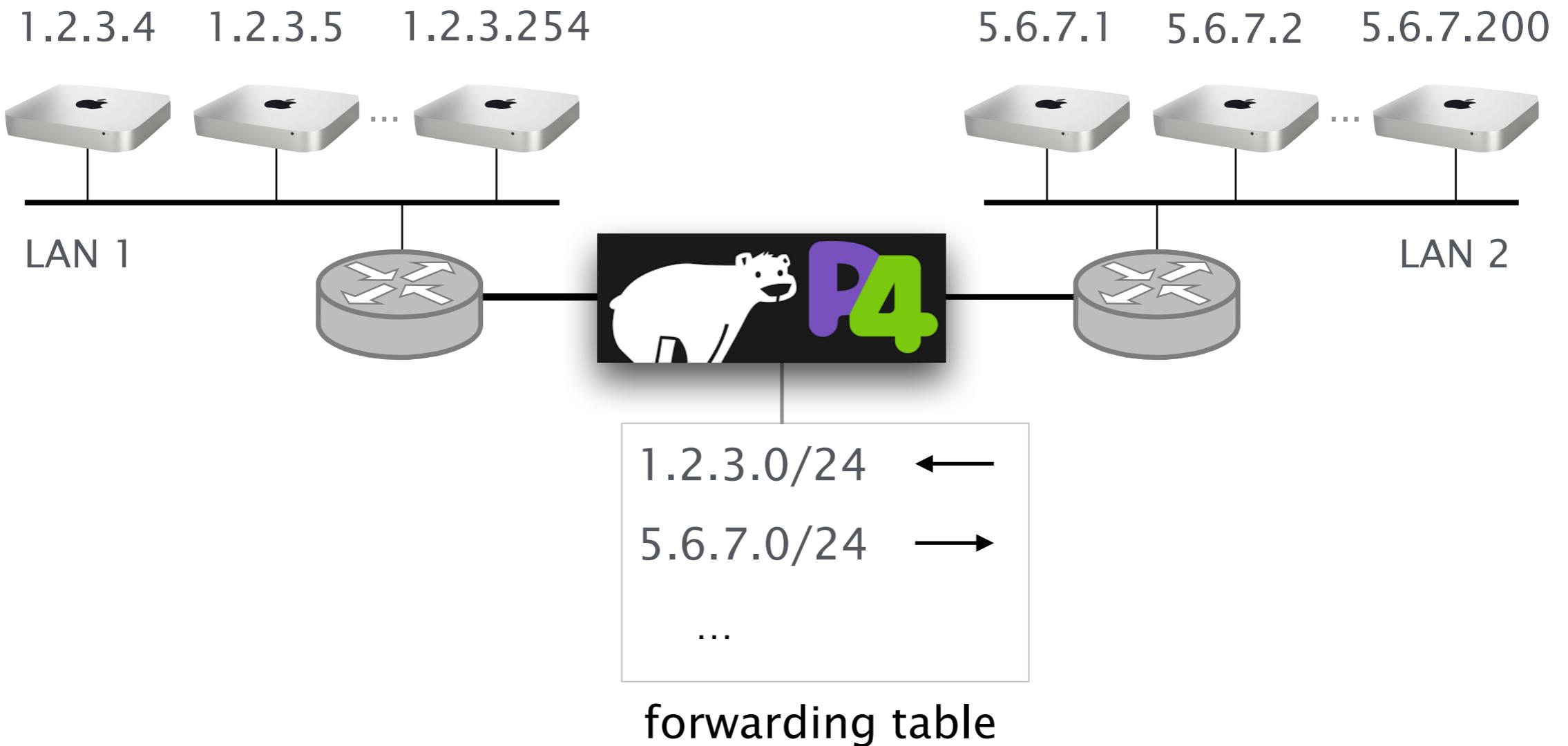
program

IP forwarding in a traditional router



IP forwarding

in P4?

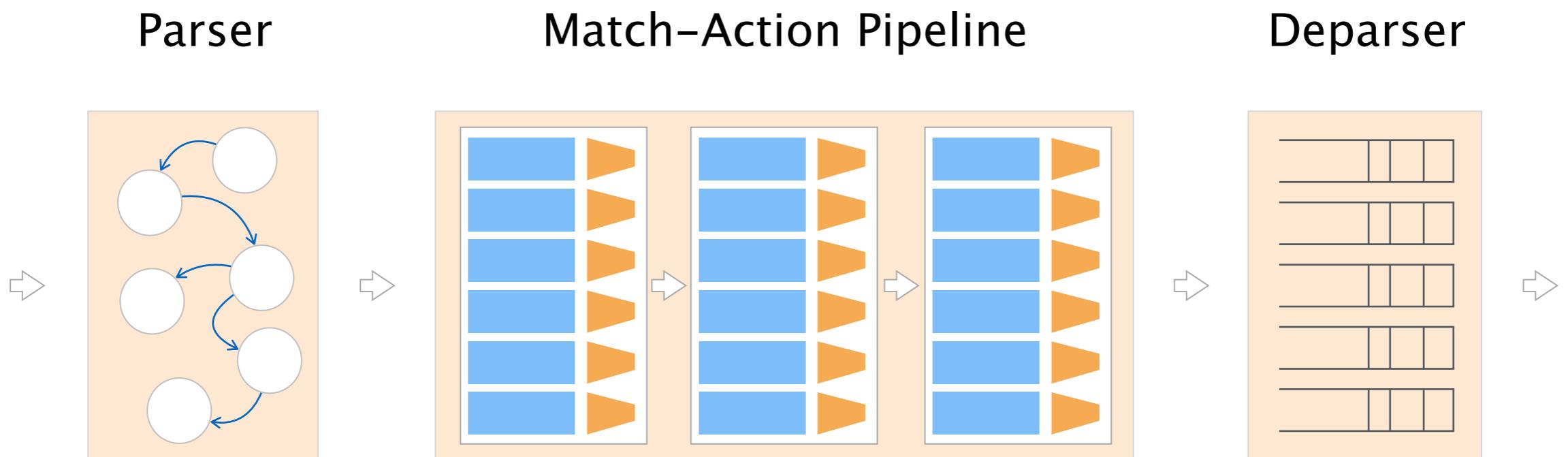


When forwarding an IP packet,
an IP router performs four actions:

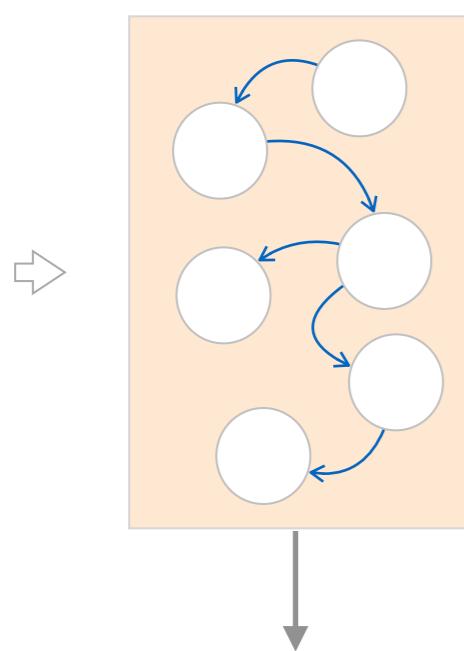
- Lookup the next hop(s) to use
- Update the MAC addresses
- Decrement TTL
- Forward packets to output port(s)

Each of them should be implemented in P4

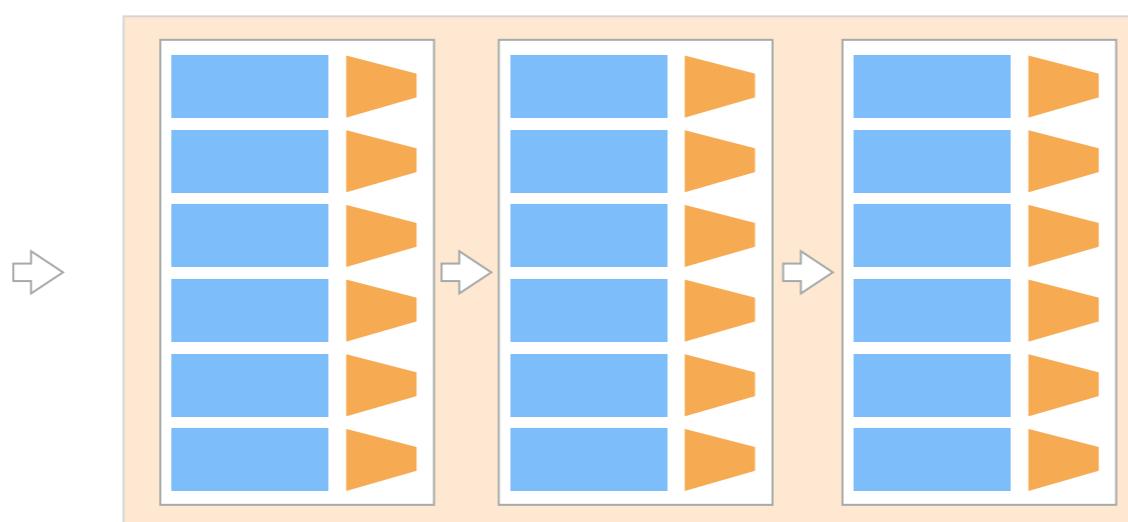
A P4 program consists of three basic parts



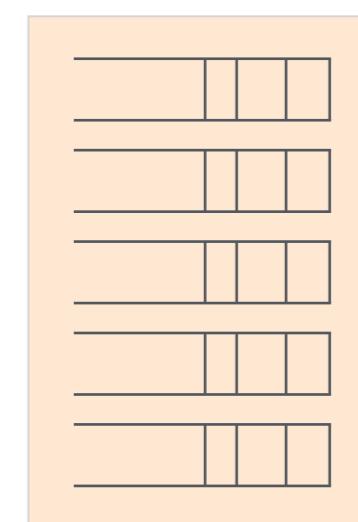
Parser



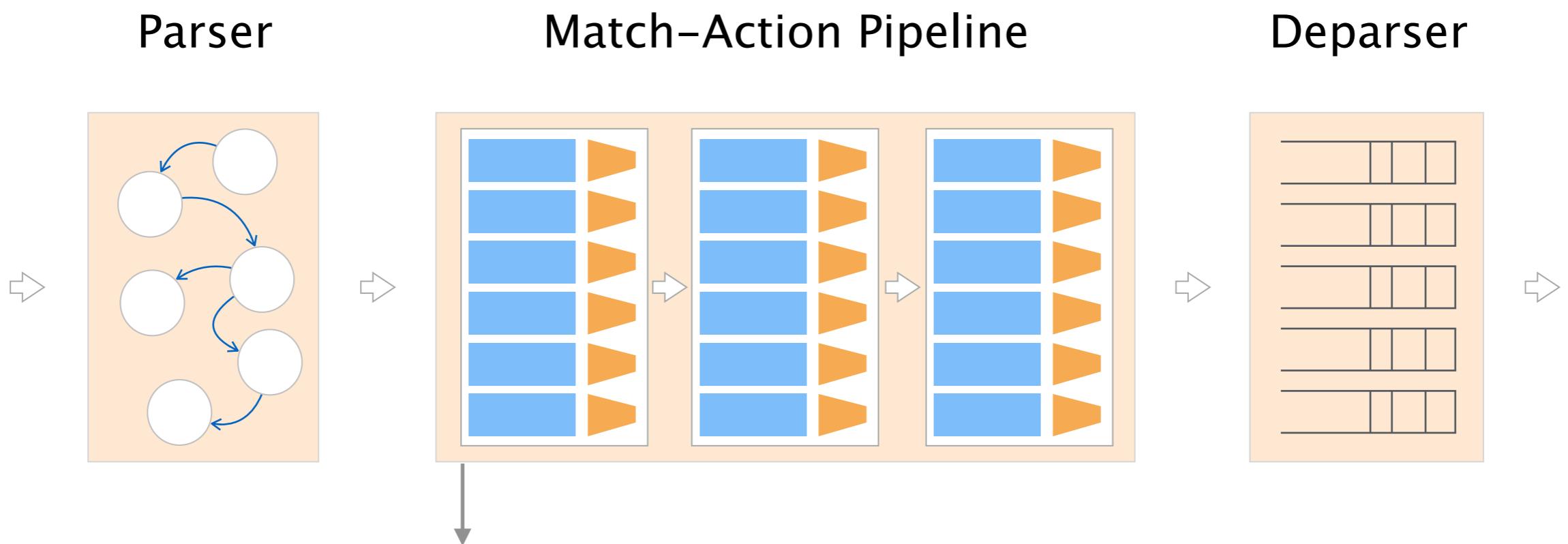
Match–Action Pipeline



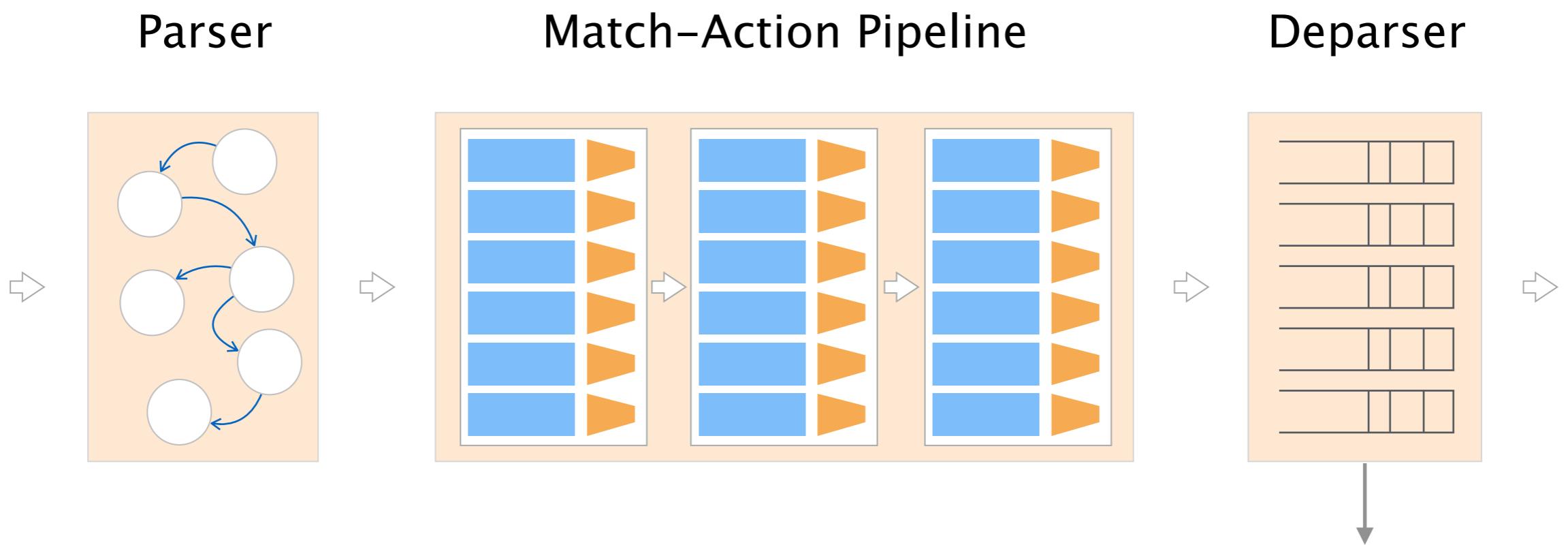
Deparser



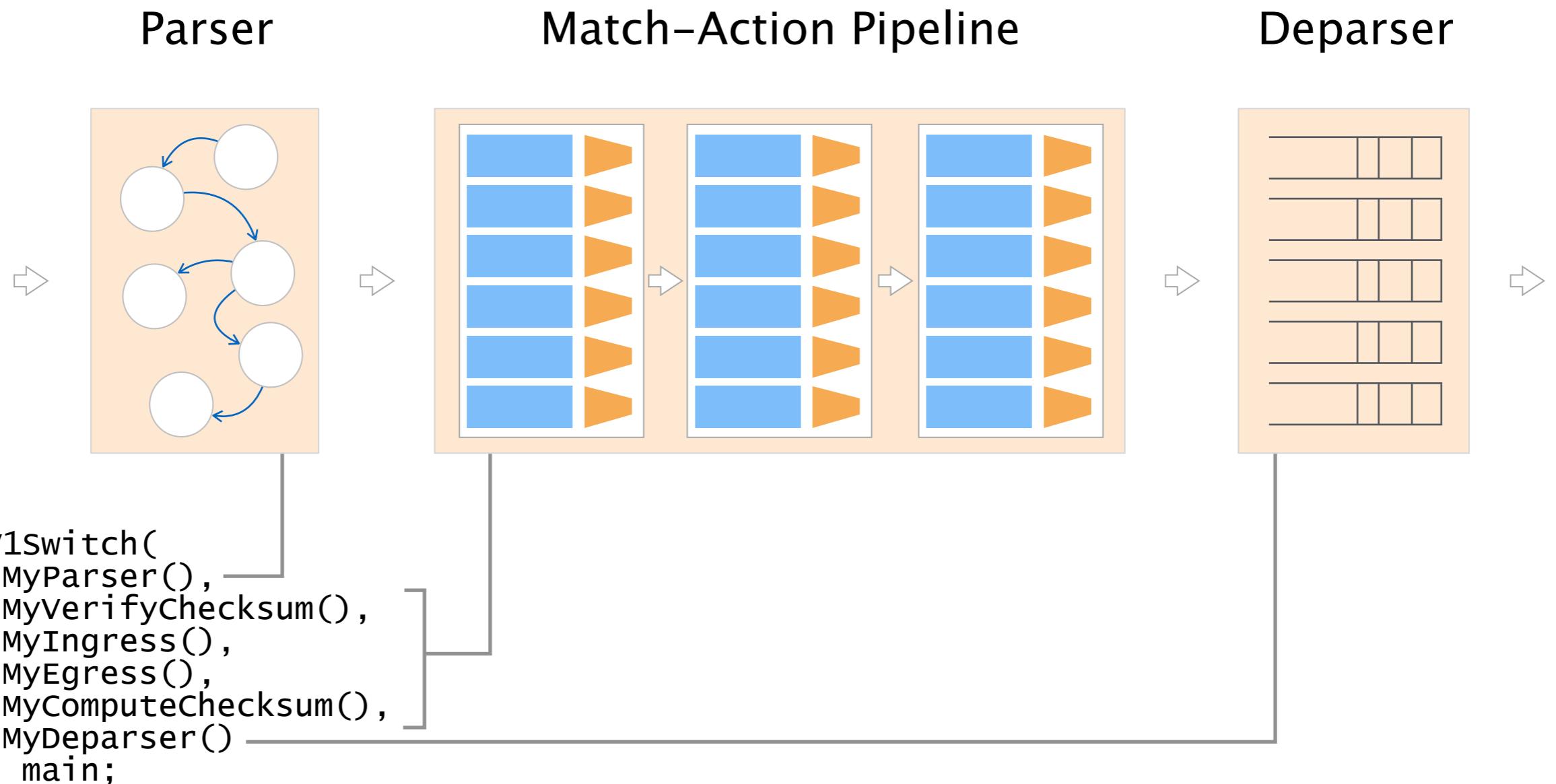
Declares the headers that
should be extracted from the packet



Defines the tables and
the processing logic



Declares how
the output packet
will look on the wire



```
#include <core.p4>
#include <v1model.p4>
```

Libraries

```
const bit<16> TYPE_IPV4 = 0x800;
typedef bit<32> ip4Addr_t;
header ipv4_t {...}
struct headers {...}
```

Declarations

```
parser MyParser(...) {
    state start {...}
    state parse_etherent {...}
    state parse_ip4 {...}
}
```

Parse packet headers

```
control MyIngress(...) {
    action ipv4_forward(...) {...}
    table ipv4_lpm {...}
    apply {
        if (...) {...}
    }
}
```

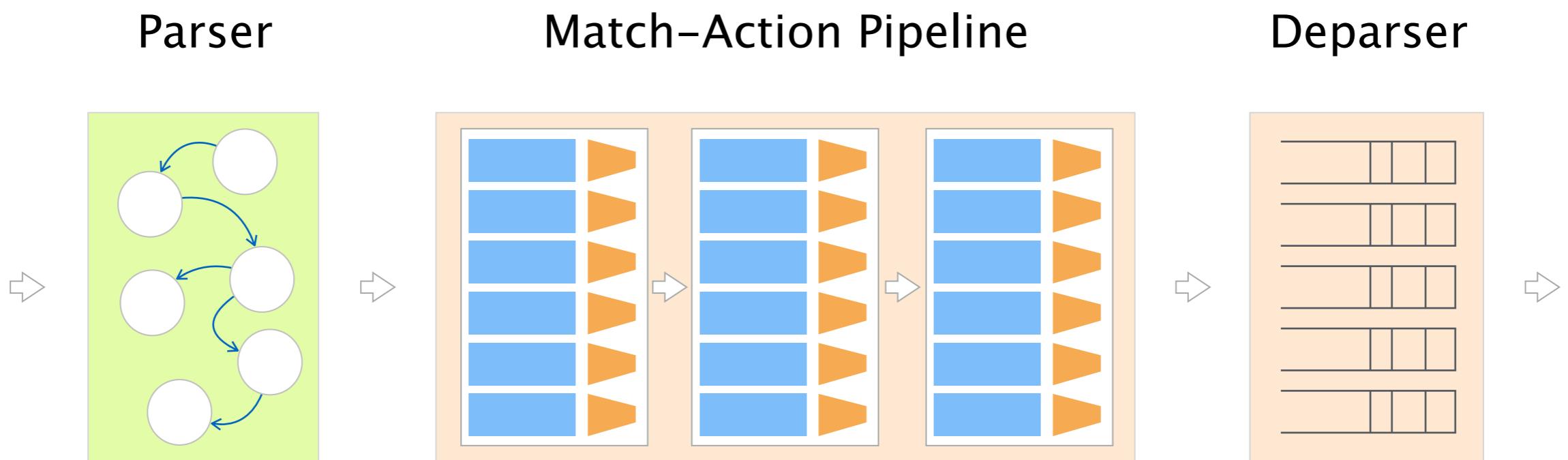
Control flow
to modify packet

```
control MyDeparser(...)
```

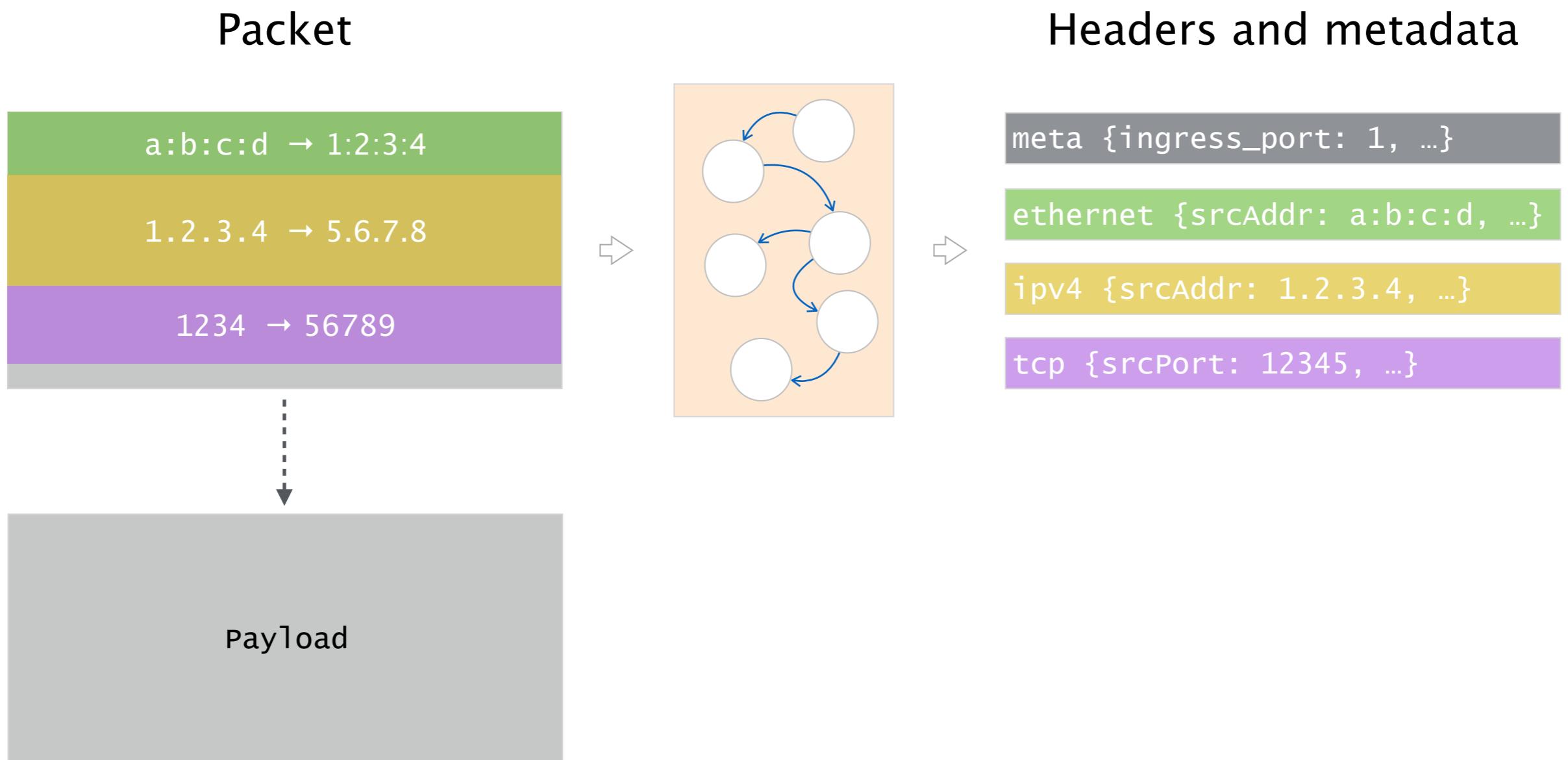
Assemble
modified packet

```
v1Switch(
    MyParser(),
    MyVerifyChecksum(),
    MyIngress(),
    MyEgress(),
    MyComputeChecksum(),
    MyDeparser()
) main;
```

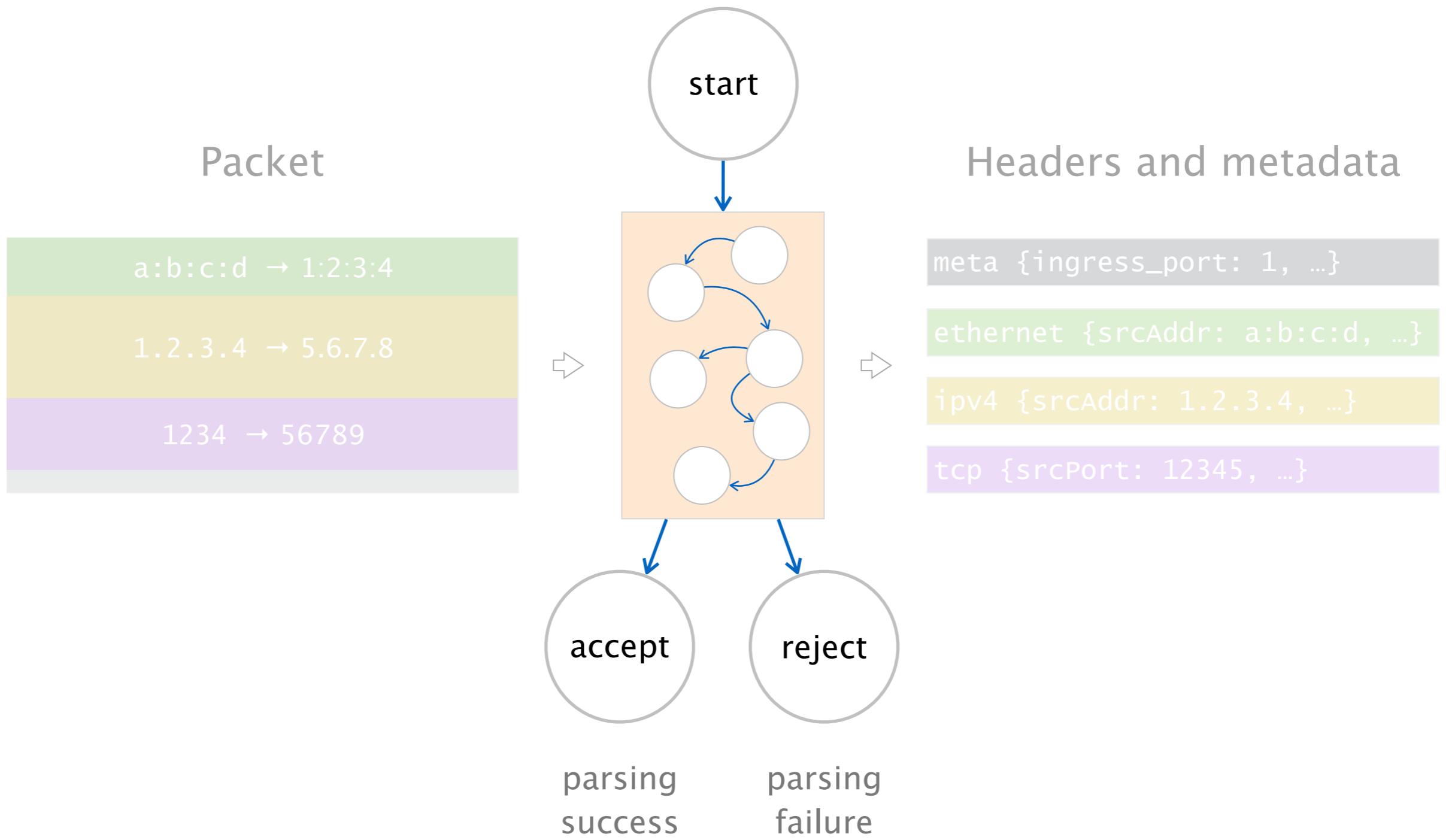
“main()”

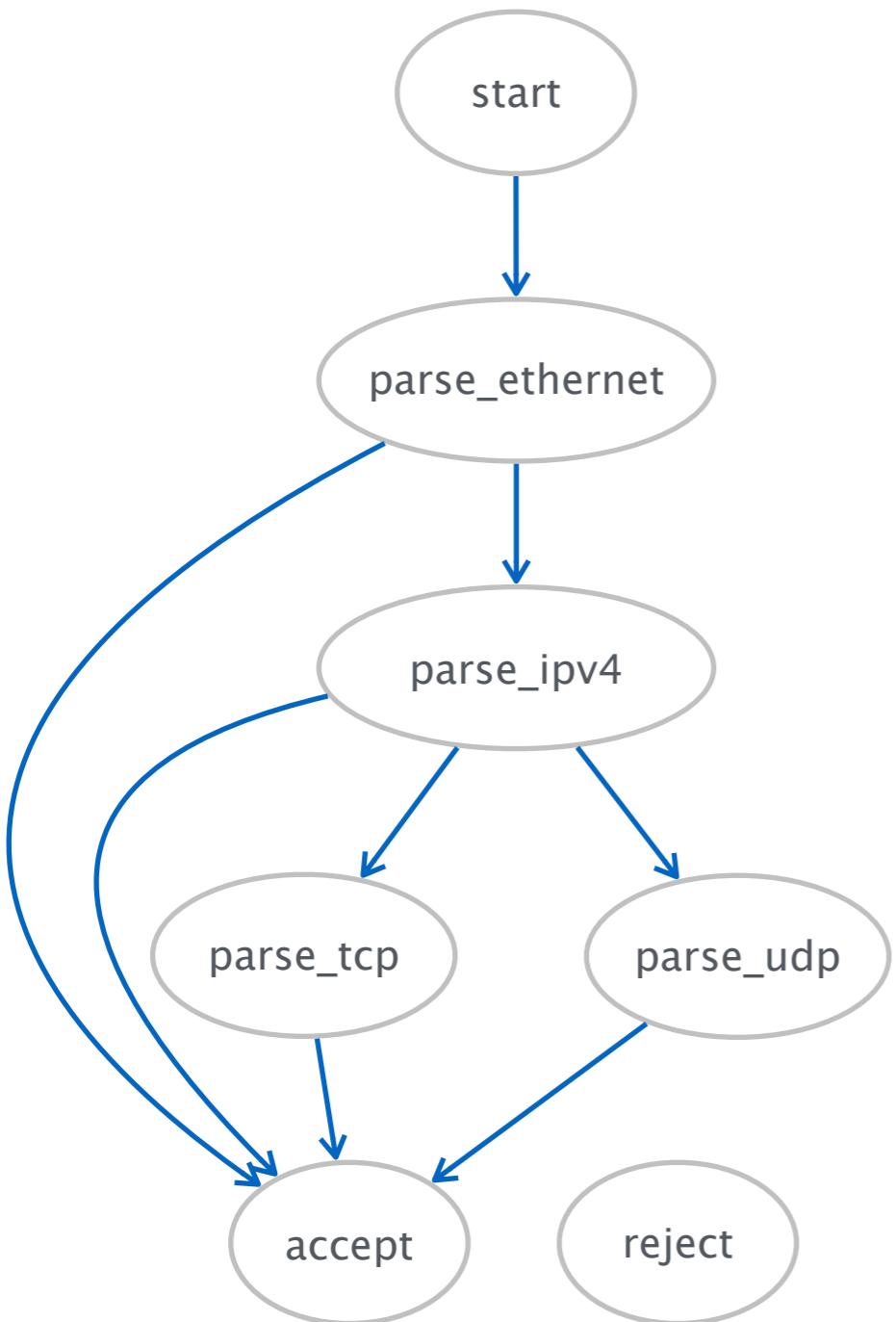


The parser uses a state machine
to map packets into headers and metadata



The parser has three predefined states:
start, accept and reject





```

parser MyParser(...) {

    state start {
        transition parse_ethernet;
    }

    state parse_ethernet {
        packet.extract(hdr.ethernet);
        transition select(hdr.ethernet.etherType) {
            0x800: parse_ipv4;
            default: accept;
        }
    }

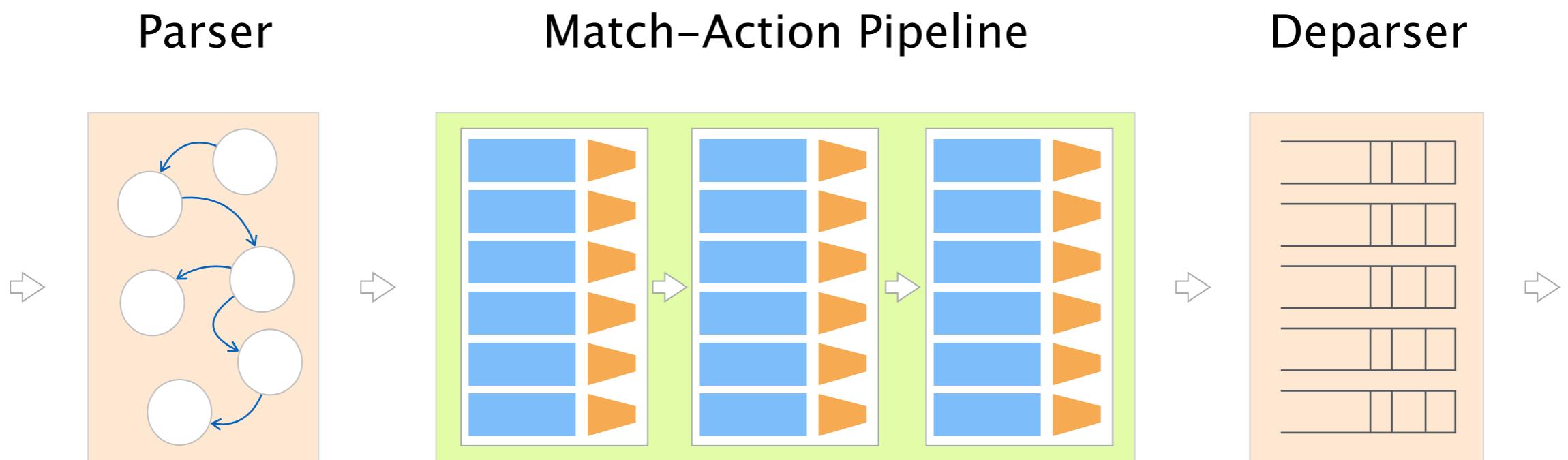
    state parse_ipv4 {
        packet.extract(hdr.ipv4);
        transition select(hdr.ipv4.protocol) {
            6: parse_tcp;
            17: parse_udp;
            default: accept;
        }
    }

    state parse_tcp {
        packet.extract(hdr.tcp);
        transition accept;
    }

    state parse_udp {
        packet.extract(hdr.udp);
        transition accept;
    }

}

```



Basic building blocks of P4 programs

Control

Control flow

describes how headers should be processed

Actions

fragments manipulating headers/metadata

Tables

map user-defined keys with actions

Control

Control flow

describes how headers should be processed

Actions

fragments manipulating headers/metadata

Tables

map user-defined keys with actions

Control flow expresses an imperative program which describes how packets are processed

```
control MyIngress(inout headers hdr,  
                  inout metadata meta,  
                  inout standard_metadata_t std_meta) {  
  
    bit<9> port; → Variable declaration  
  
    apply {  
        port = 1  
        std_meta.egress_spec = port;  
        hdr.ethernet.srcAddr = hdr.ethernet.dstAddr;  
        hdr.ethernet.dstAddr = 0x2;  
        hdr.ipv4.ttl = hdr.ipv4.ttl - 1;  
    }  
}
```

↓
Control flow

Control

Control flow

describes how headers should be processed

Actions

fragments manipulating headers/metadata

Tables

map user-defined keys with actions

Actions allow to re-use code

similar to functions in C

```
control MyIngress(inout headers hdr,
                   inout metadata meta,
                   inout standard_metadata_t std_meta) {

    action ipv4_forward(macAddr_t dstAddr,
                        egressSpec_t port) {
        std_meta.egress_spec = port;
        hdr.ethernet.srcAddr = hdr.ethernet.dstAddr;
        hdr.ethernet.dstAddr = dstAddr;
        hdr.ipv4.ttl = hdr.ipv4.ttl - 1;
    }

    apply {
        ipv4_forward(0x123, 1);
    }
}
```



Control

Control flow

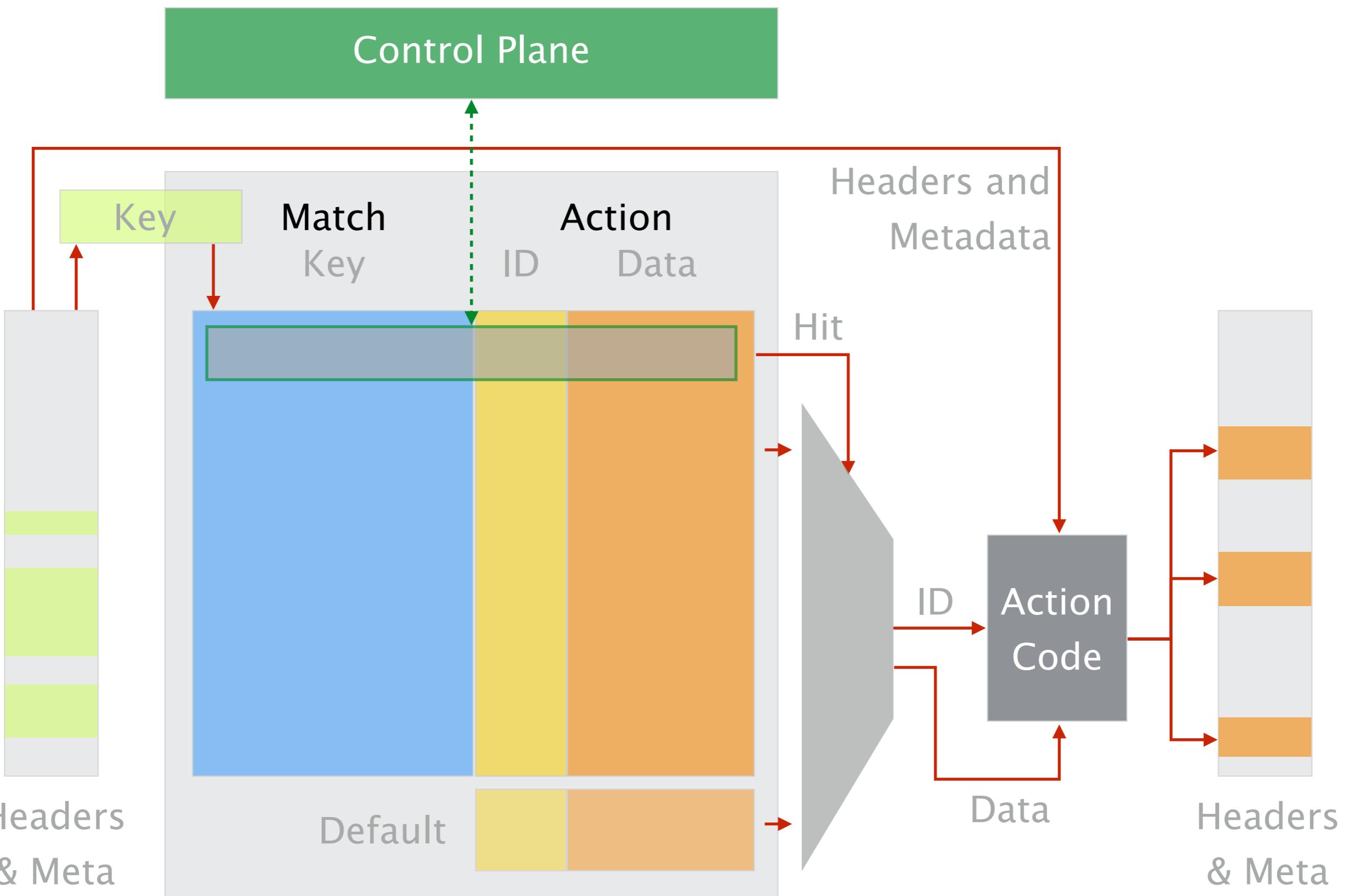
describes how headers should be processed

Actions

fragments manipulating headers/metadata

Tables

map user-defined keys with actions



```
table [ ] {  
    key = { [ ] : [ ]; } [ ] ;  
    actions = { [ ] } [ ] ;  
    size = [ ];  
    default_action = [ ]; } [ ]
```

Table name

Field(s) to match

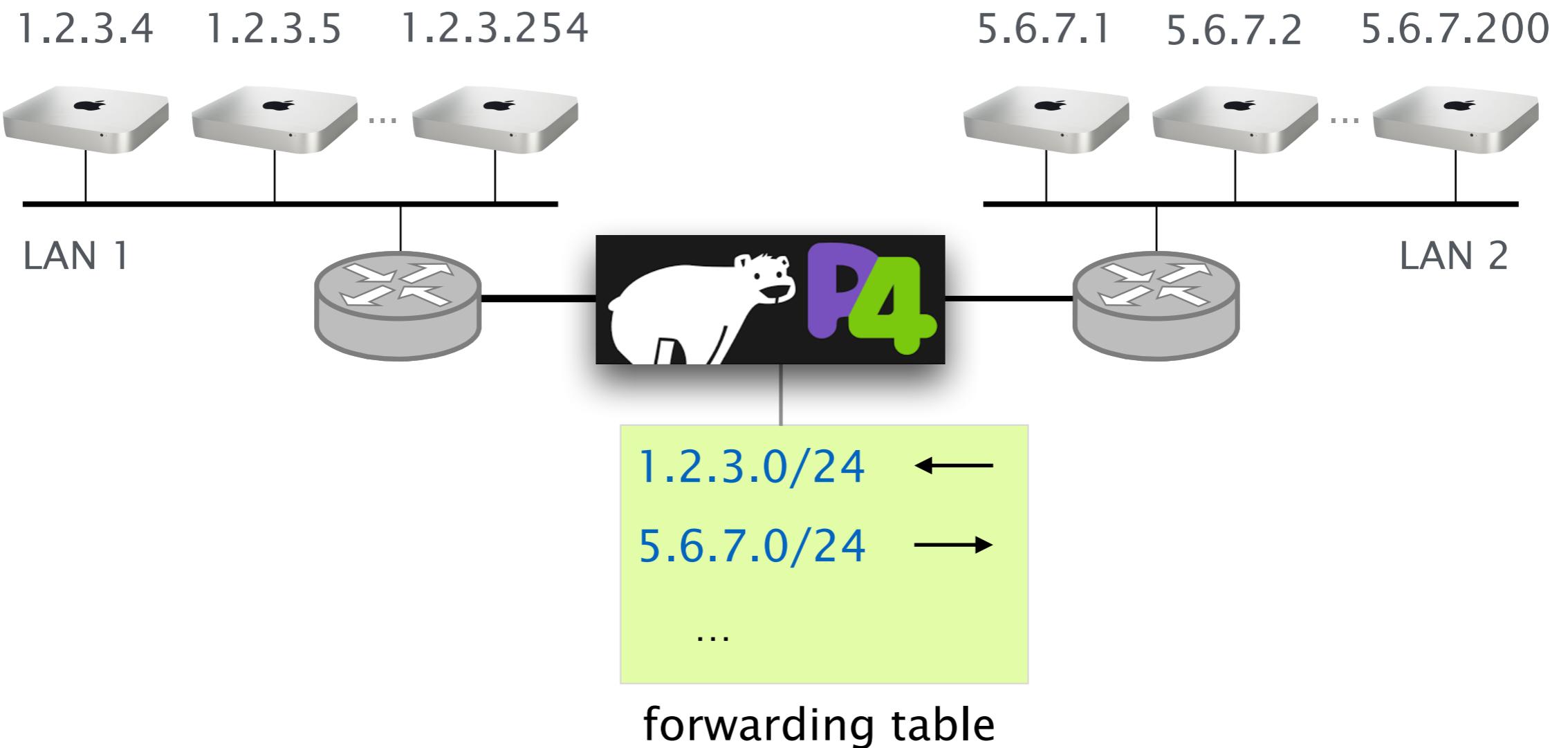
Match type

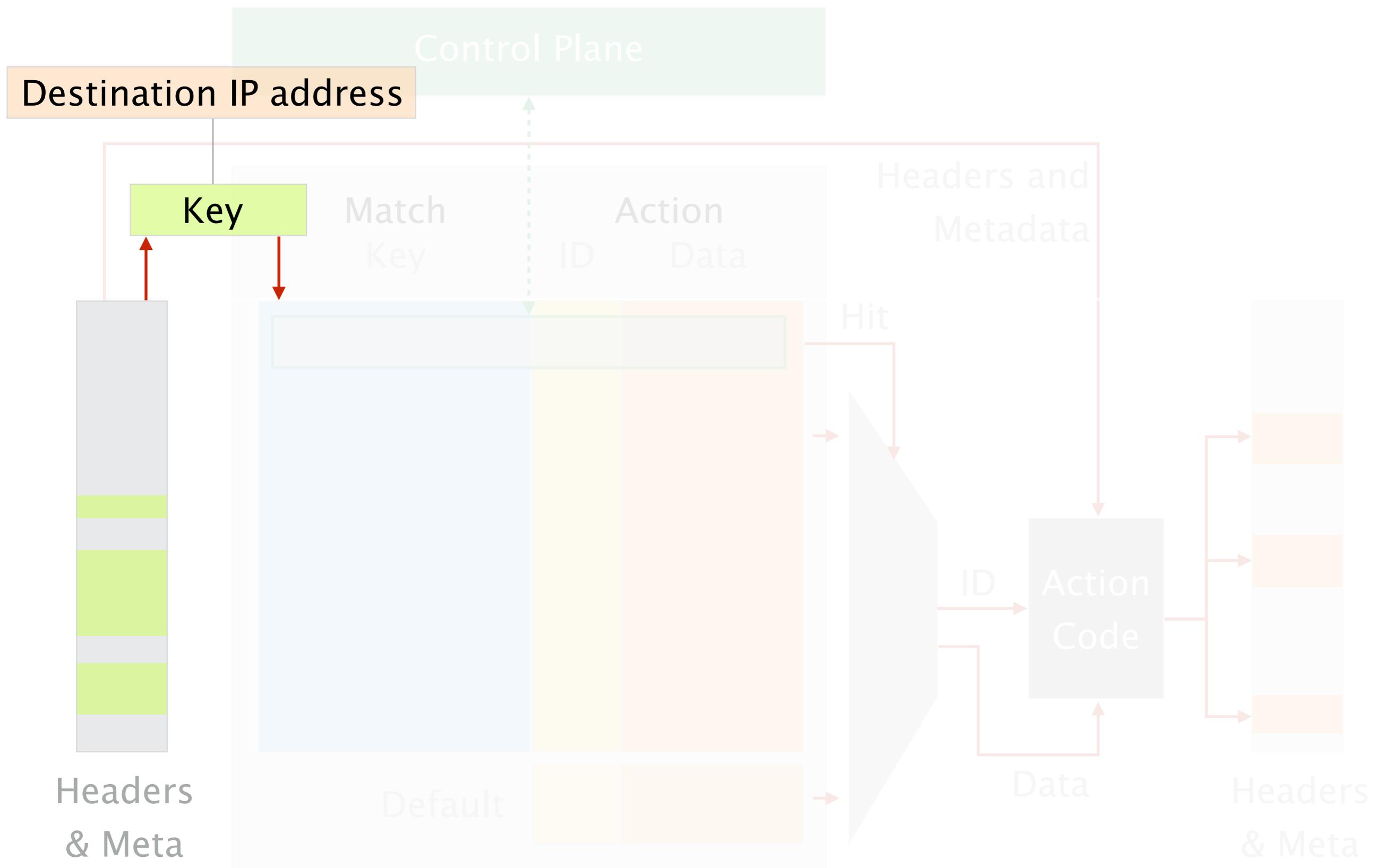
Possible actions

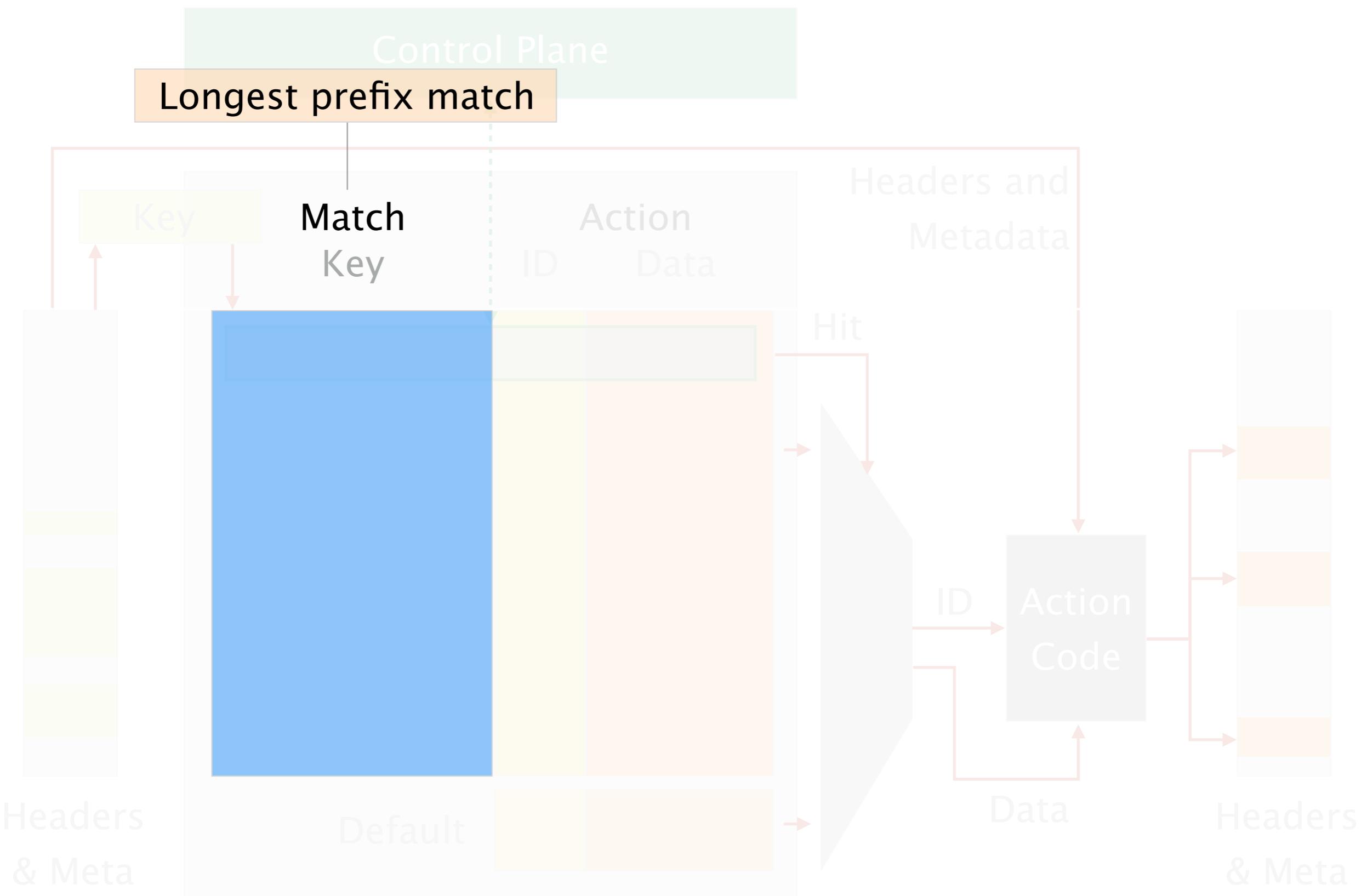
Max. # entries in table

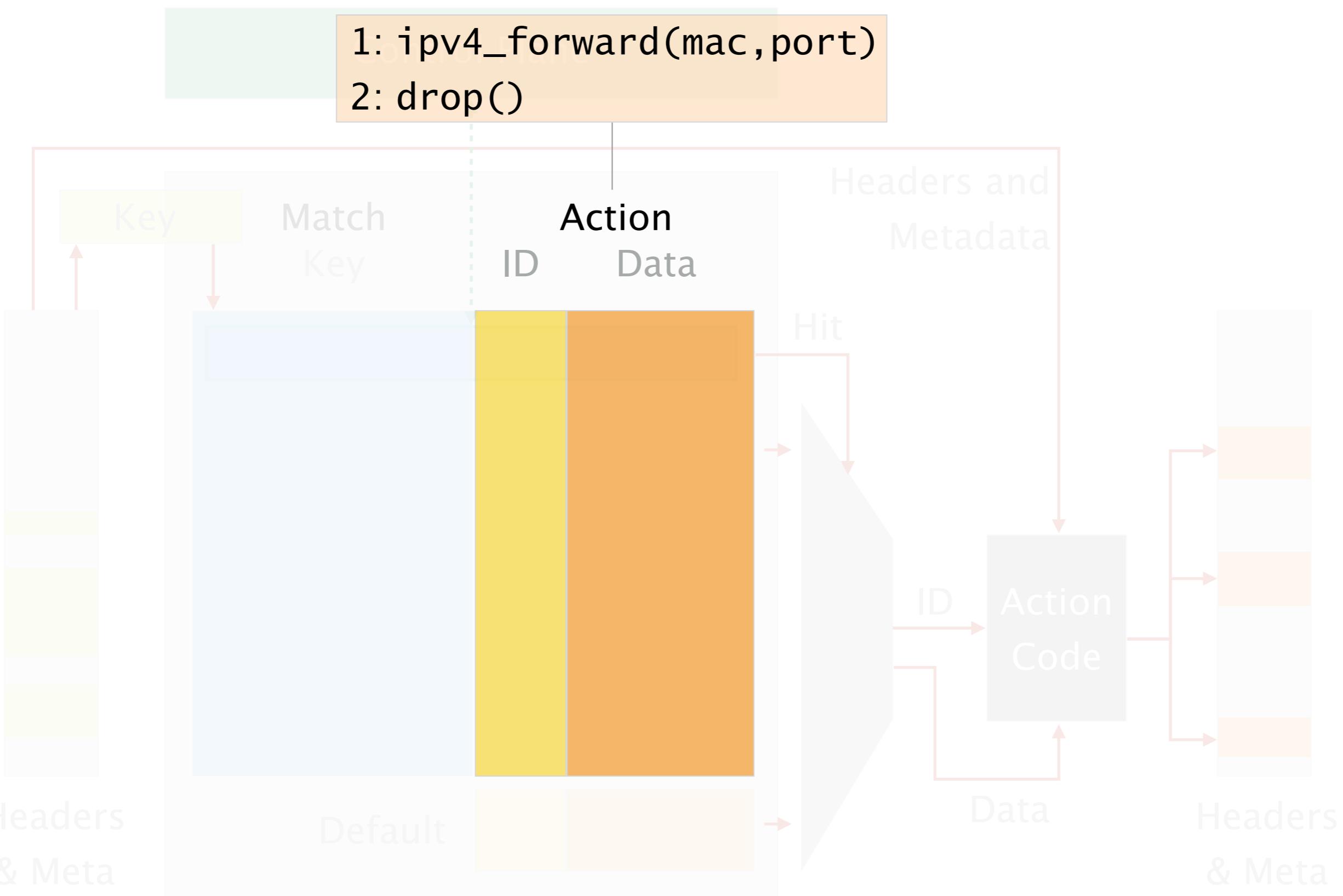
Default action

Example: IP forwarding table









```
table ipv4_1pm {  
    key = {  
        hdr.ipv4.dstAddr: 1pm;  
    }  
    actions = {  
        ipv4_forward;  
        drop;  
    }  
  
    size = 1024;  
    default_action = drop();
```

Table name

Destination IP address

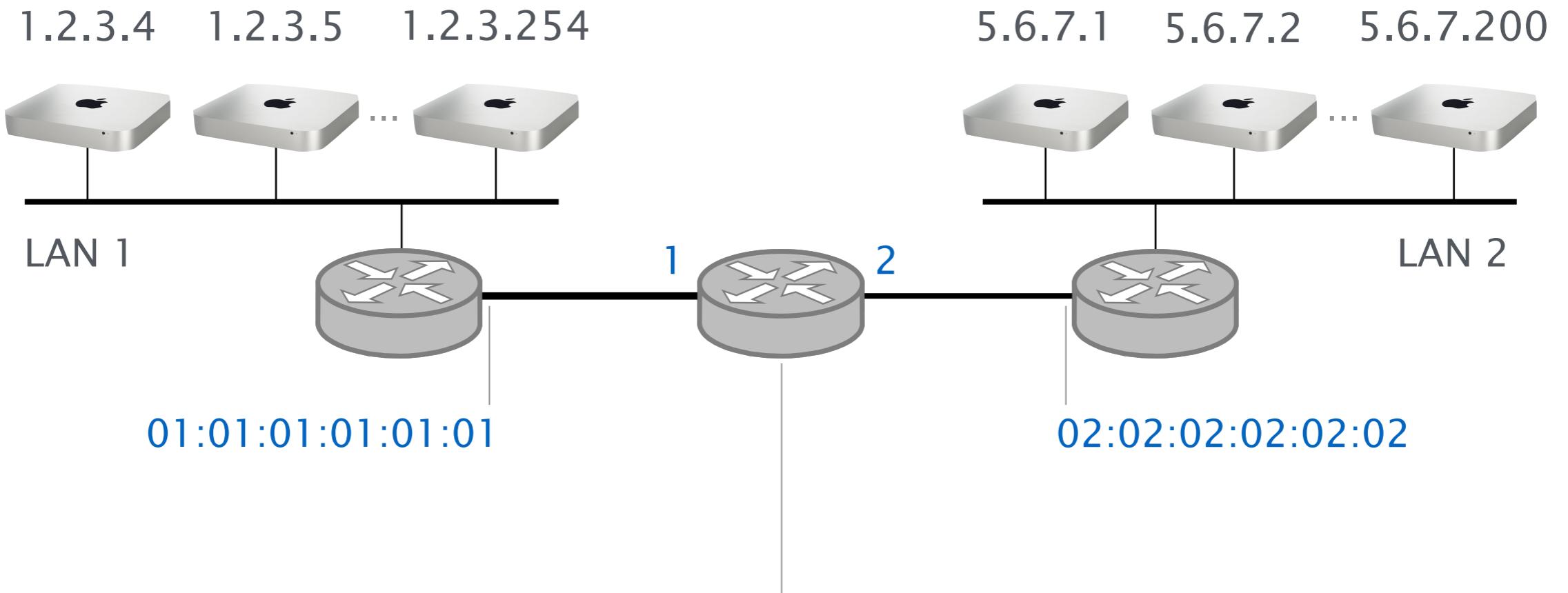
Longest prefix match

Possible actions

Max. # entries in table

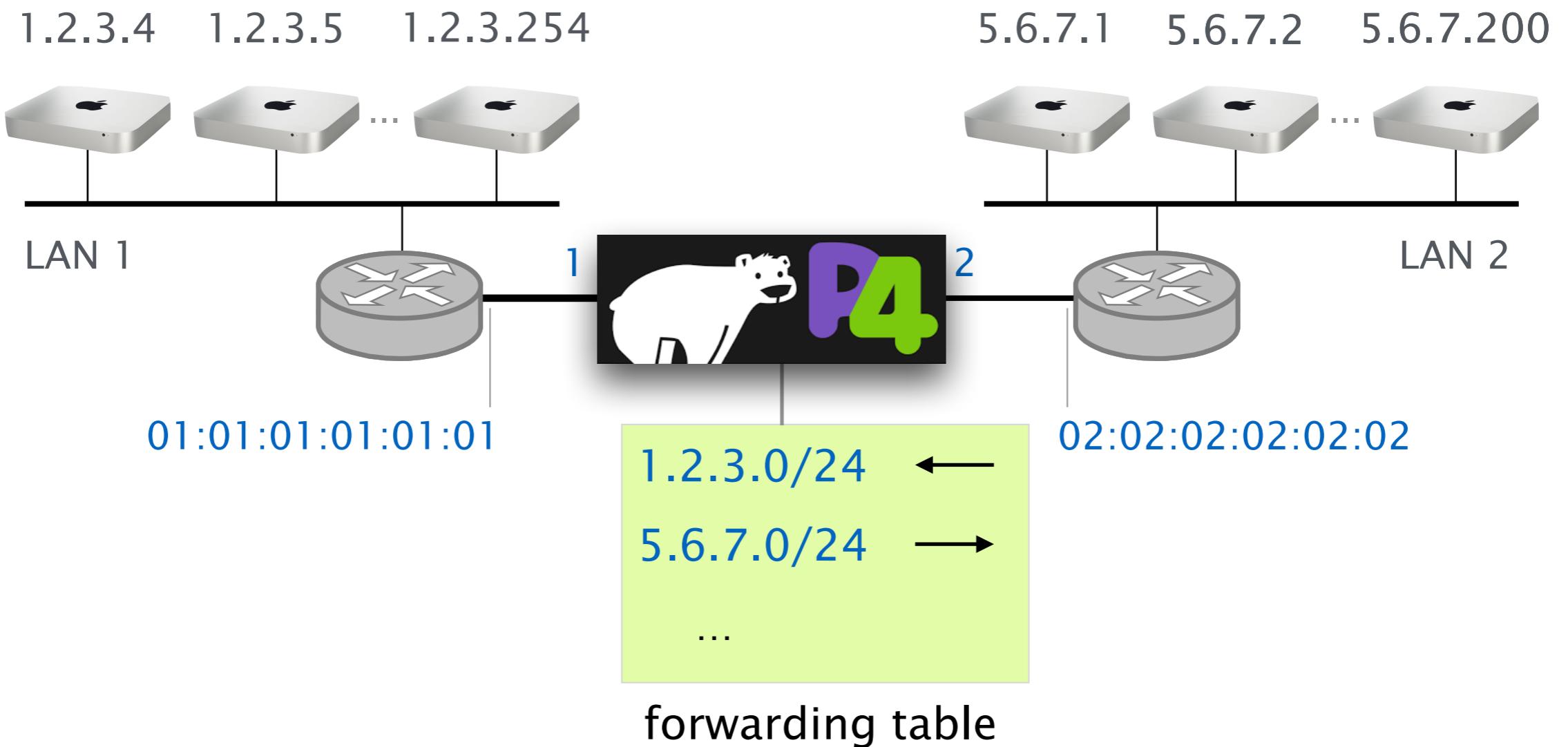
Default action

Example: IP forwarding table



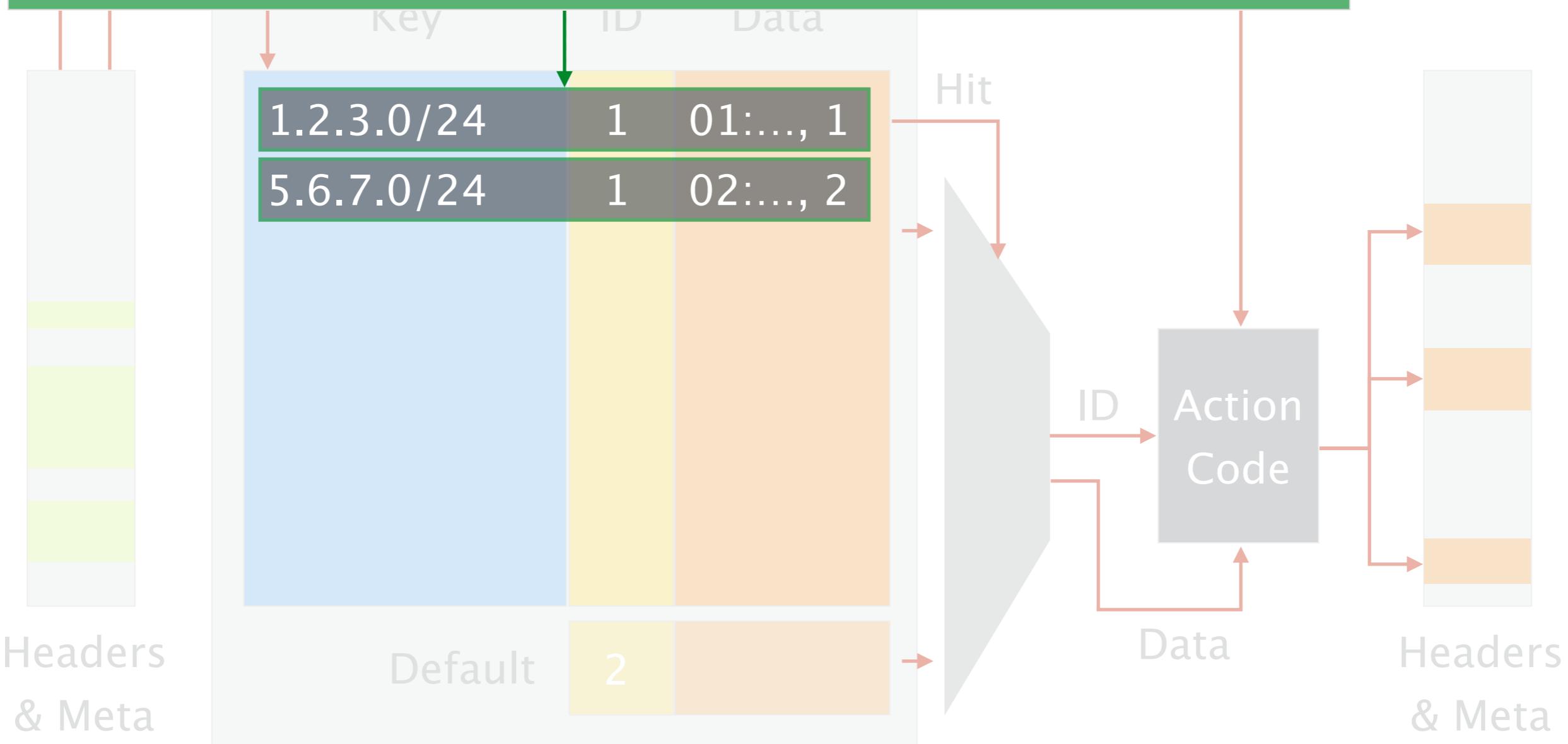
```
action ipv4_forward(macAddr_t dstAddr, egressSpec_t port) {
    standard_metadata.egress_spec = port;
    hdr.ethernet.srcAddr = hdr.ethernet.dstAddr;
    hdr.ethernet.dstAddr = dstAddr;
    hdr.ipv4.ttl = hdr.ipv4.ttl - 1;
}
```

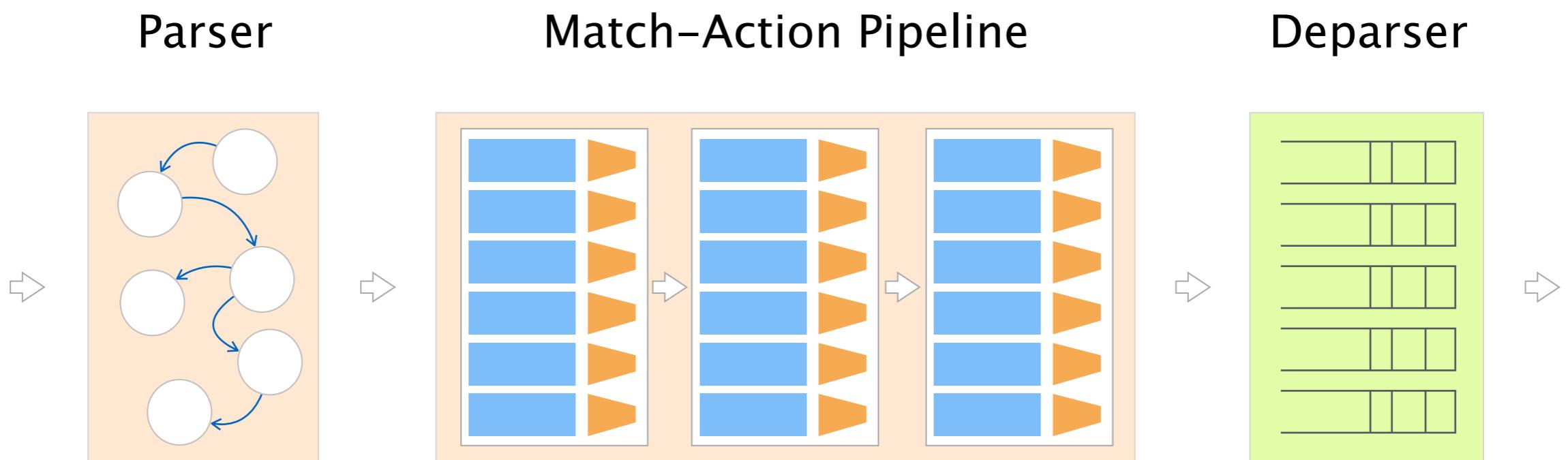
Example: IP forwarding table



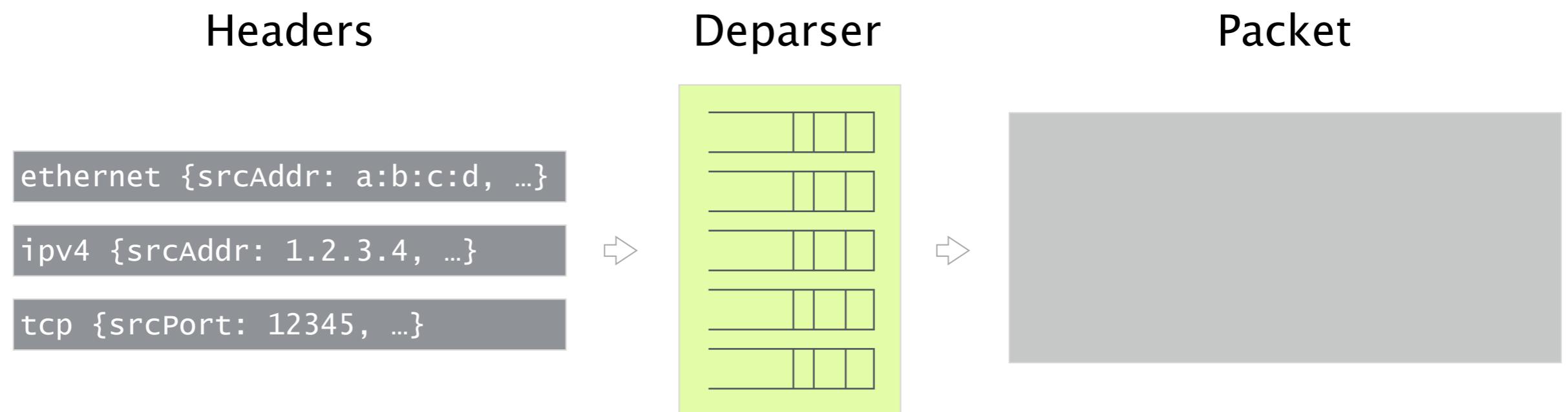
Control Plane

```
table_add ipv4_1pm ipv4_forward 1.2.3.0/24 => 01:01:01:01:01:01 1  
table_add ipv4_1pm ipv4_forward 5.6.7.0/24 => 02:02:02:02:02:02 2
```





The Deparser assembles the headers back into a well-formed packet



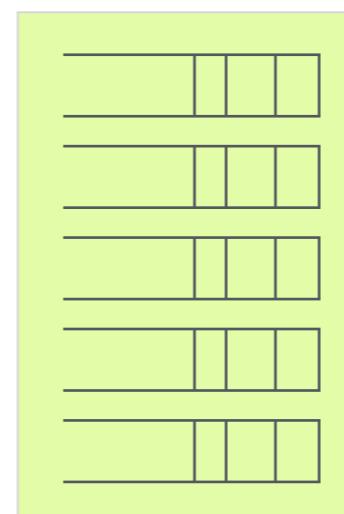
Headers

```
ethernet {srcAddr: a:b:c:d, ...}
```

```
ipv4 {srcAddr: 1.2.3.4, ...}
```

```
tcp {srcPort: 12345, ...}
```

Deparser



Packet

```
a:b:c:d → 1:2:3:4
```

```
control MyDeparser(packet_out packet, in headers hdr) {
    apply {
        packet.emit(hdr.ethernet);
    }
}
```

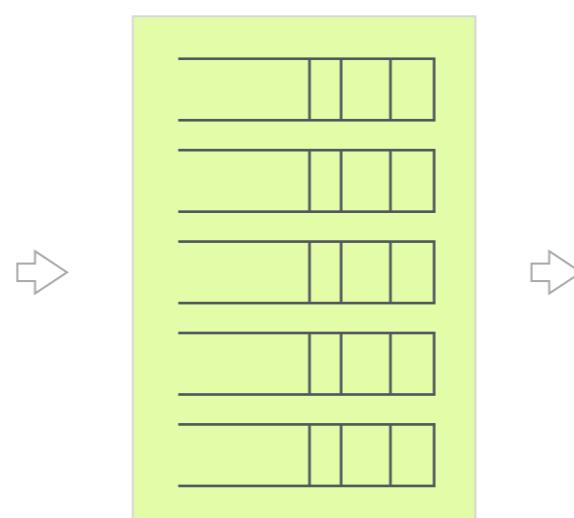
Headers

```
ethernet {srcAddr: a:b:c:d, ...}
```

```
ipv4 {srcAddr: 1.2.3.4, ...}
```

```
tcp {srcPort: 12345, ...}
```

Deparser



Packet

```
a:b:c:d → 1:2:3:4
```

```
1.2.3.4 → 5.6.7.8
```

```
control MyDeparser(packet_out packet, in headers hdr) {
    apply {
        packet.emit(hdr.ethernet);
        packet.emit(hdr.ipv4);
    }
}
```

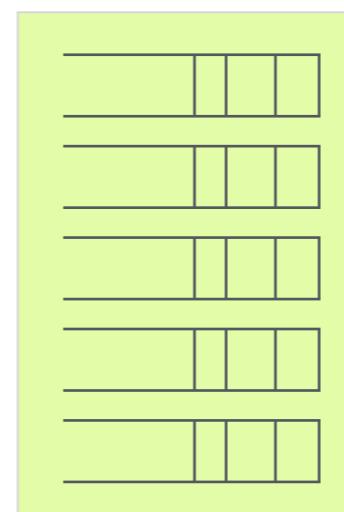
Headers

```
ethernet {srcAddr: a:b:c:d, ...}
```

```
ipv4 {srcAddr: 1.2.3.4, ...}
```

```
tcp {srcPort: 12345, ...}
```

Deparser



Packet

```
a:b:c:d → 1:2:3:4
```

```
1.2.3.4 → 5.6.7.8
```

```
1234 → 56789
```

```
control MyDeparser(packet_out packet, in headers hdr) {
    apply {
        packet.emit(hdr.ethernet);
        packet.emit(hdr.ipv4);
        packet.emit(hdr.tcp);
    }
}
```

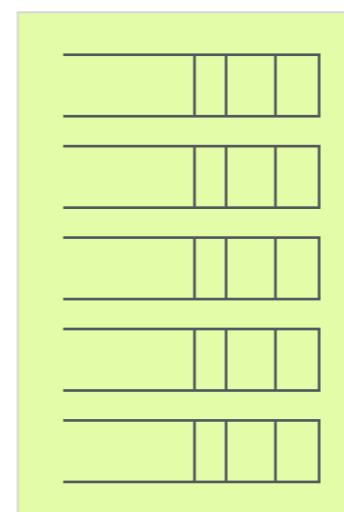
Headers

ethernet {srcAddr: a:b:c:d, ...}

ipv4 {srcAddr: 1.2.3.4, ...}

tcp {srcPort: 12345, ...}

Deparser



Packet

a:b:c:d → 1:2:3:4

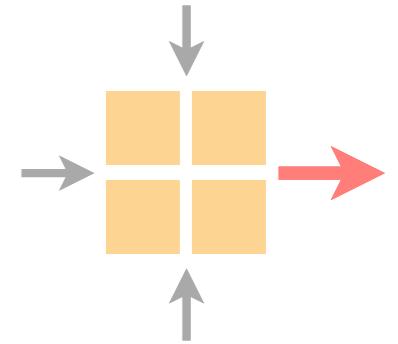
1.2.3.4 → 5.6.7.8

1234 → 56789

Payload

Advanced Topics in Communication Networks

Internet Routing and Forwarding



Laurent Vanbever
nsg.ee.ethz.ch

ETH Zürich (D-ITET)
15 Sep 2020