# BGP in 120 minutes

**RIPE89**

**Wolfgang Tremmel**
**academy@de-cix.net**

# About me

➔Wolfgang Tremmel

➔studied Informatik (Uni Karlsruhe)

 ➔Degree: Diploma (1994)

➔Network Engineer at **XLINK**

 ➔Since 1996 Director NOC

 ➔Since 2000 Senior Network Planner DSL at **kpn Qwest**

➔2001 - 2005 Director Network Planning at VIA NET.WORKS

➔2006 - 2016 Manager Customer Support at **DE-CIX**
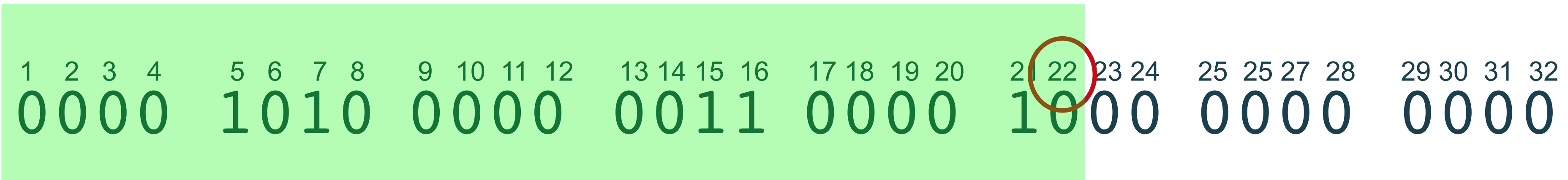
➔since 2016: Head of DE-CIX Academy

wolfgangtremmel1966

@wtremmel@hessen.social

**Where networks meet**

# What is BGP about?

# *IPv4 Prefixes*

# 10.3.8.0/22

| 1 | 2 | 3 | 4 | | 5 | 6 | 7 | 8 | | 9 | 10 | 11 | 12 | | 13 | 14 | 15 | 16 | | 17 | 18 | 19 | 20 | | 21 | 22 | 23 | 24 | | 25 | 25 | 27 | 28 | | 29 | 30 | 31 | 32 |

```
0000 1010 0000 0011 0000 1000 0000 0000
```

→ IPv4 and IPv6 addresses have a network and a host part

→ A **prefix** is just the network part

→ Important:

· **The boundary between network and host can be anywhere!**

# Characteristics of Prefixes: IPv4

# 10.3.8.0/22

**Prefix-Length: 0-32**

| 1 | 2 | 3 | 4 | | 5 | 6 | 7 | 8 | | 9 | 10 | 11 | 12 | | 13 | 14 | 15 | 16 | | 17 | 18 | 19 | 20 | | 21 | 22 | 23 | 24 | | 25 | 25 | 27 | 28 | | 29 | 30 | 31 | 32 |

0000 1010 0000 0011 0000 1000 0000 0000

**Notation:**
- 4 Numbers 0-255
- Separated by "."
- a "/", followed by

**Host-part all zero**

**32 Bits long**

# Characteristics of Prefixes: IPv6

**Prefix-Length: 0-128**

## 2003:de:274f:400::/64

0 01 02 03 04 05 06 07 08 09 0a 0b 0c 0d 0e 0f 10 11 12 13 14 15 16 17 18 19 1a 1b 1c 1d 1e 1f 20 21 22 23 24 25 26 27 28 29 2a 2b 2c 2d 2e 2f 30 31 32 33 34 35 36 37 38 39 3a 3b 3c 3d 3e 3f 40 41 42 43 44 45 46 47 48 49 4a 4b 4c 4d 4e 4f 50 51 52 53 54 55 56 57 58 59 5a 5b 5c 5d 5e 5f 60 61 62 63 64 65 66 67 68 69 6a 6b 6c 6d 6e 6f 70 71 72 73 74 75 76 77 78 79 7a 7b 7c 7d 7e 7f

**Notation:**
- 4 digit hex numbers (0-9,a-f)
- Separated by ":"
- "::" = fill up with zeros

**Host-part all zero**

**128 Bits long**

# How does BGP work?

# BGP is a protocol to announce prefixes
## Everybody has Neighbors



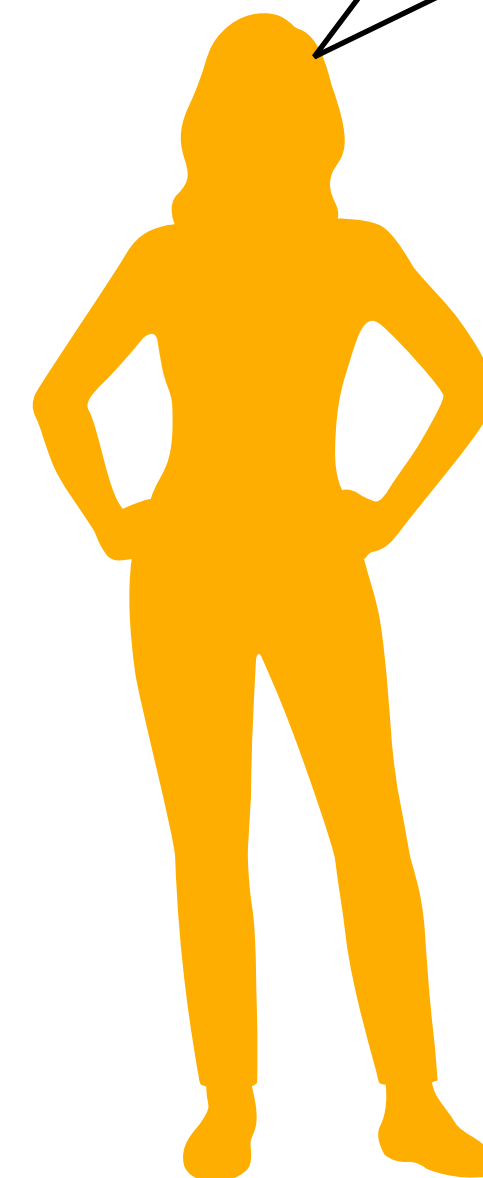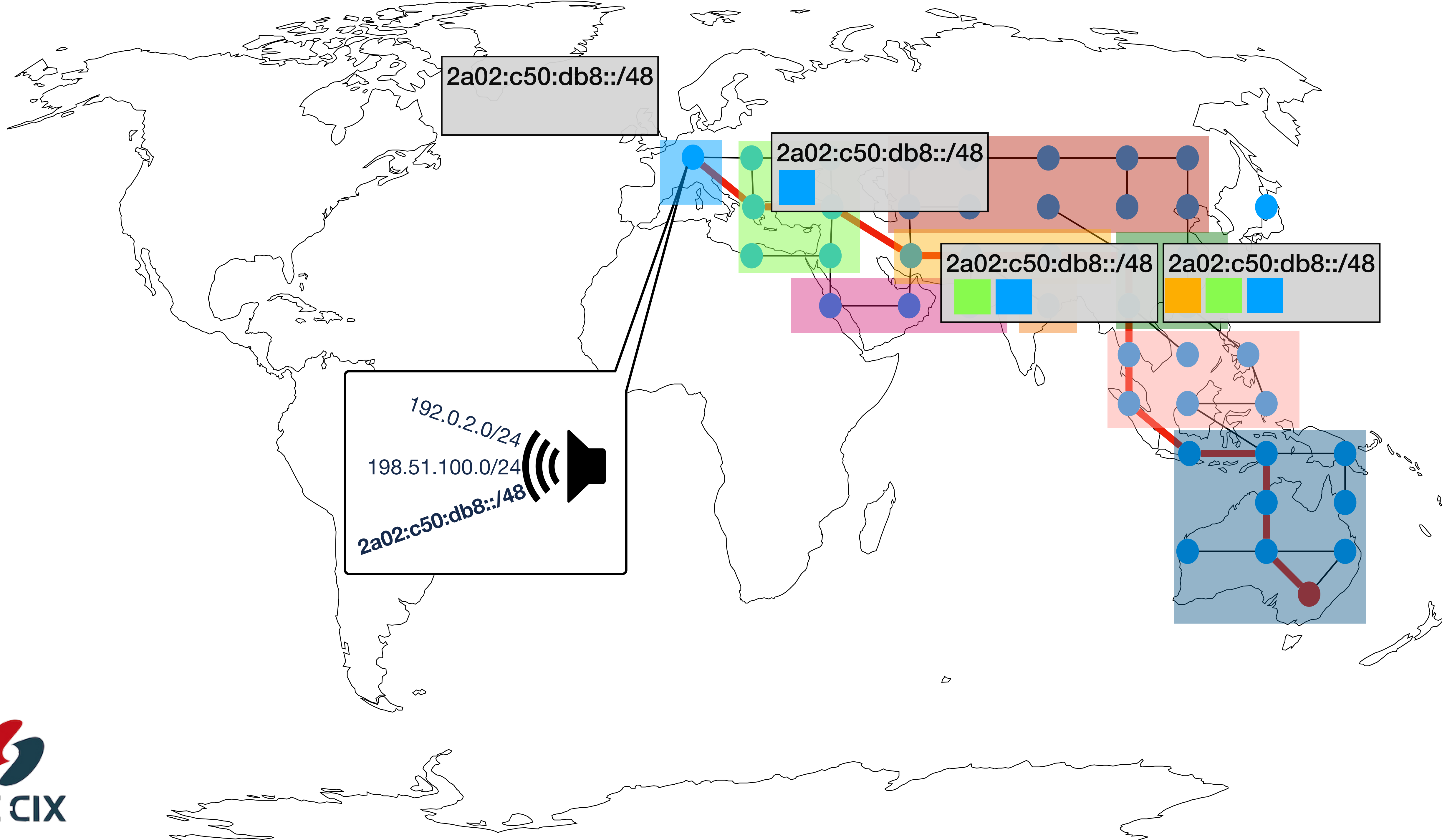I am **AS196610**, DE-CIX Academy, and I announce prefix **2a02:c50:db8::/48**

My neighbor AS196610 announces prefix 2a02:c50:db8::/48

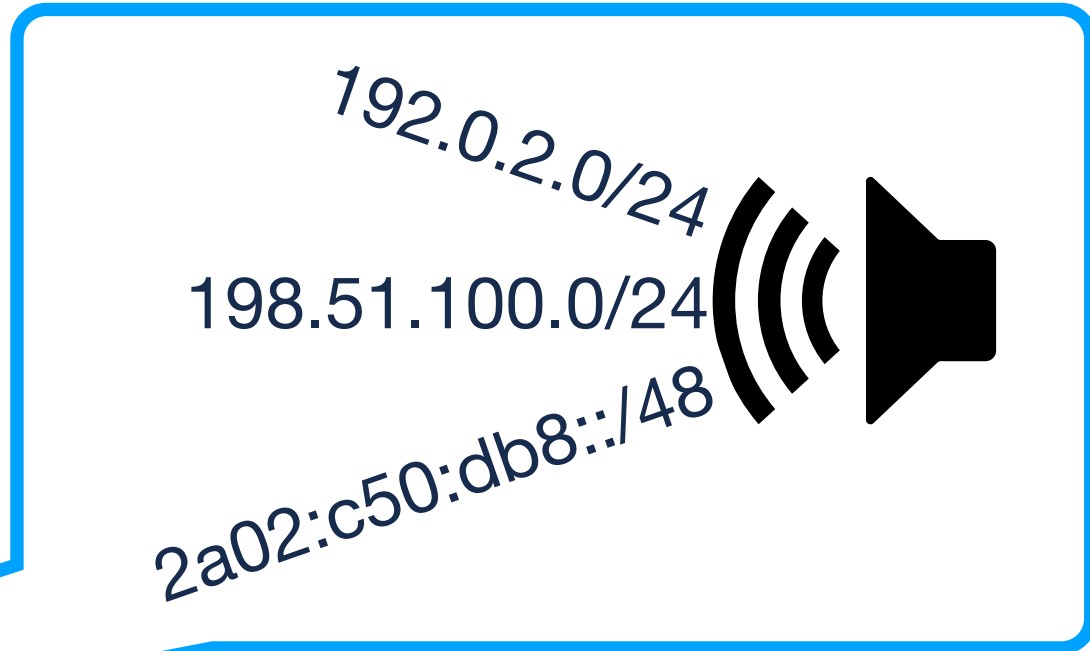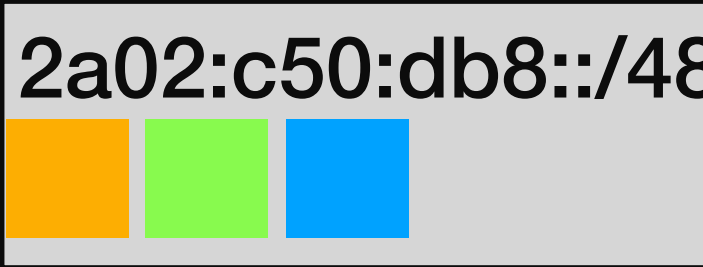My green neighbor told me, his neighbor AS196610 announces prefix 2a02:c50:db8::/48

DE-CIX Academy
AS196610

# BGP announces prefixes
## To neighbors

192.0.2.0/24
198.51.100.0/24
2a02:c50:db8::/48

I am **AS196610**, DE-CIX
Academy, and I announce
prefix
2a02:c50:db8::/48

2a02:c50:db8::/48

- BGP **announces** IP prefixes to **neighbors**

  - These neighbors have to be **configured**

  - Each BGP speaking device is part of an **Autonomous System**

  - The path these announcements take is recorded - this is called the **Autonomous System Path**

  - The AS Path shows which Autonomous Systems have forwarded the prefix announcement

  - The rightmost AS in the AS Path is called the "**Originator**"

DE-CIX

# What is an *Autonomous System*?

# What is an Autonomous System?

## Simple Definition

- A group of IP prefixes

  - But to route or announce them, you need hardware

  - A router (or multiple routers)

  - This router speaks BGP (to its neighbors)

  - And has an *Autonomous System Number* configured

- Another new term: **Autonomous System Number (ASN)**

Router

I am **AS196610**, DE-CIX Academy, and I announce prefix 2a02:c50:db8::/48

# Autonomous System Number

## or AS Number or ASN

- Initially 16bit (0...65535) they are now 32bit long (0..."a lot")

- AS numbers are globally unique

- Unique means, somebody has to administrate them

- This is the IANA (Internet Assiged Numbers Authority)

  - But they have delegated that task to the 5 RIRs (Regional Internet Registries)

  - So in Europe: Become a member of the RIPE NCC and request one

*"An AS has a **globally unique** number (sometimes referred to as an **ASN**, or Autonomous System Number) associated with it; this number is used in both the exchange of exterior routing information (between neighboring ASes), and as an **identifier of the AS** itself." ([RFC1930](RFC1930))*
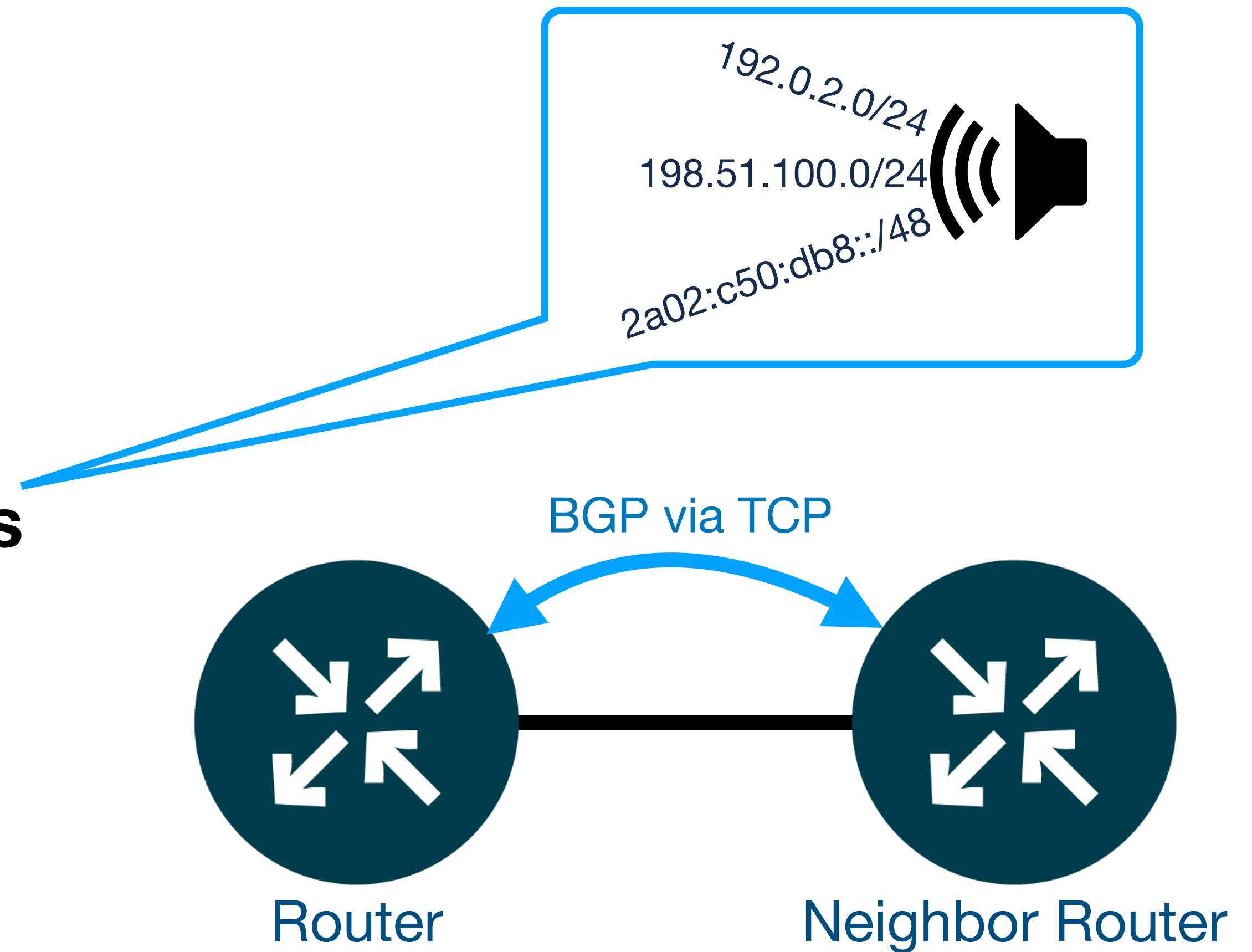
# BGP Announcing Prefixes

# BGP Neighbors
## Directly connected neighbors

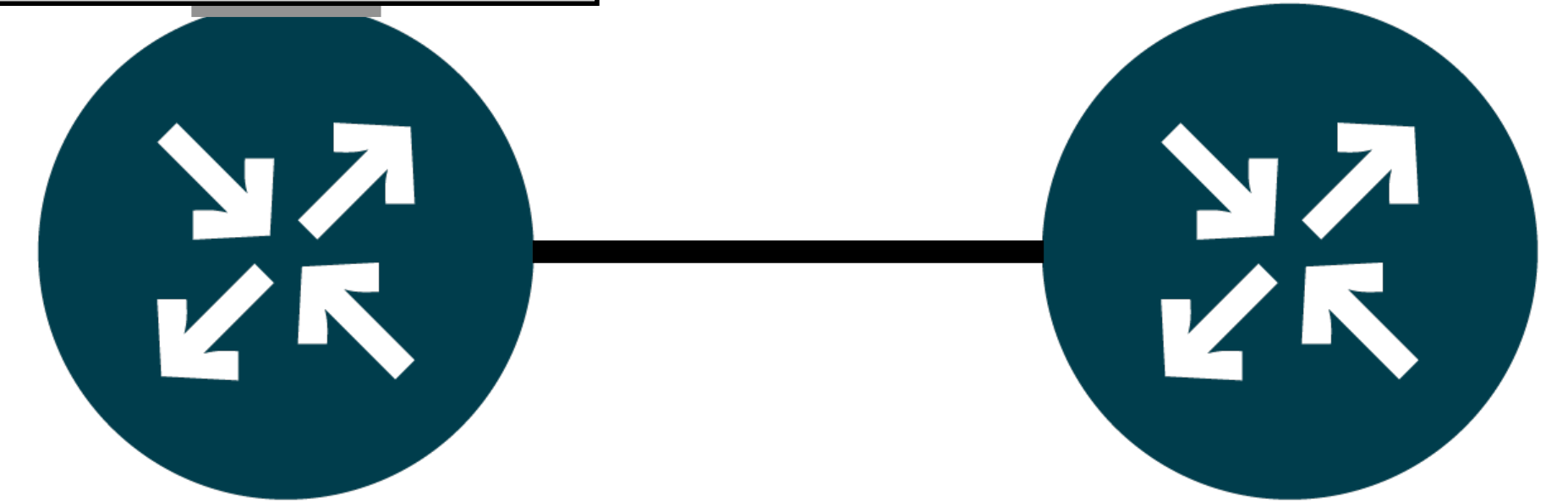192.0.2.0/24
198.51.100.0/24
2a02:c50:db8::/48

- BGP **announces** IP prefixes to **neighbors**

- These neighbors have to be **configured**

- BGP uses **TCP** to connect to a neighbor

- TCP brings already:

  - **Reliable transport** (sender knows that receiver got it)

  - **Flow control** (do not send faster than the receiver can receive)

  - **Framing** (putting BGP messages into packets)

BGP via TCP

Router

Neighbor Router

# BGP works incremental
## Using add- / withdraw- messages
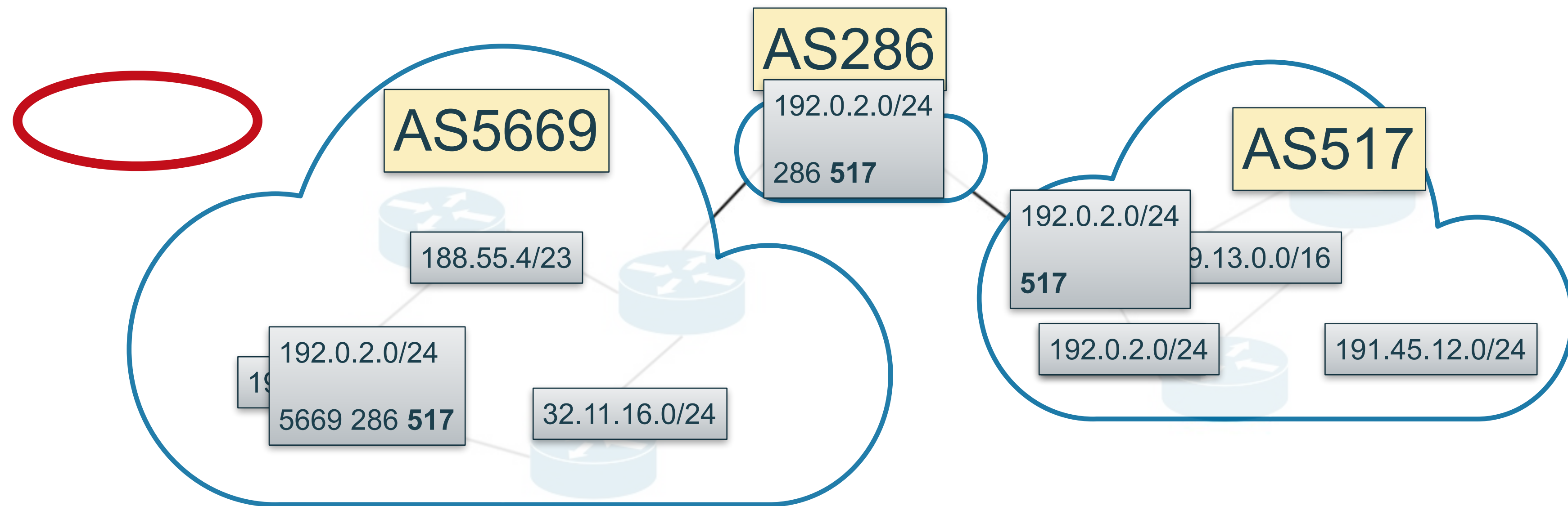
withdraw:
2a02:c50:db8::/48

- At session setup, BGP announces "everything" to its neighbor

- After that, updates are **incremental**:

  - If BGP learns about a new prefix, it sends an **add**-message to neighbors

  - If a prefix goes away, it sends a **withdraw** message to neighbors

- As long as the BGP session is "up", a router assumes its neighbors are "in sync" (= did not forget anything it sent)
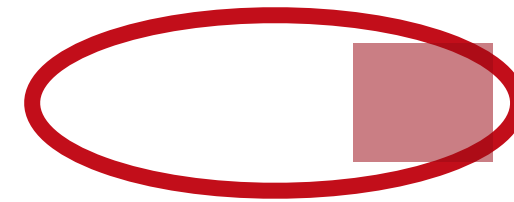
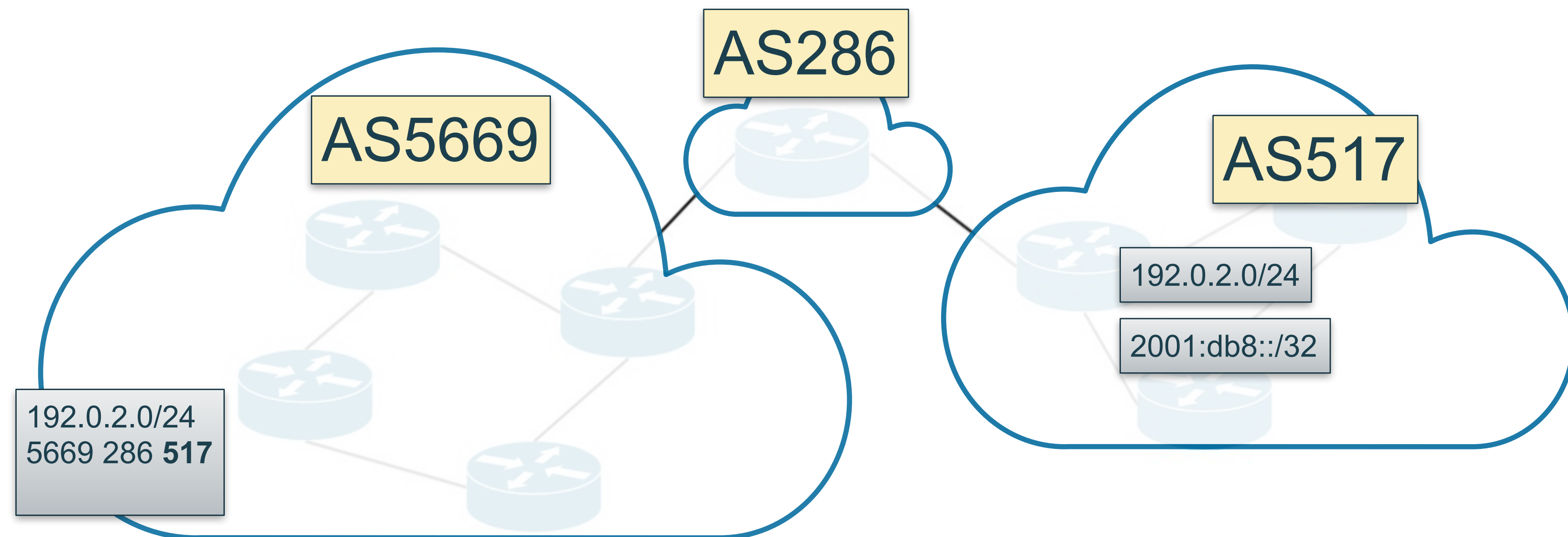# BGP Announcing Prefixes
## Building the AS path

# BGP Announcing Prefixes

→ Prefixes

→ AS Numbers

→ AS Path

Originator AS

AS286

AS5669

AS517

192.0.2.0/24

2001:db8::/32
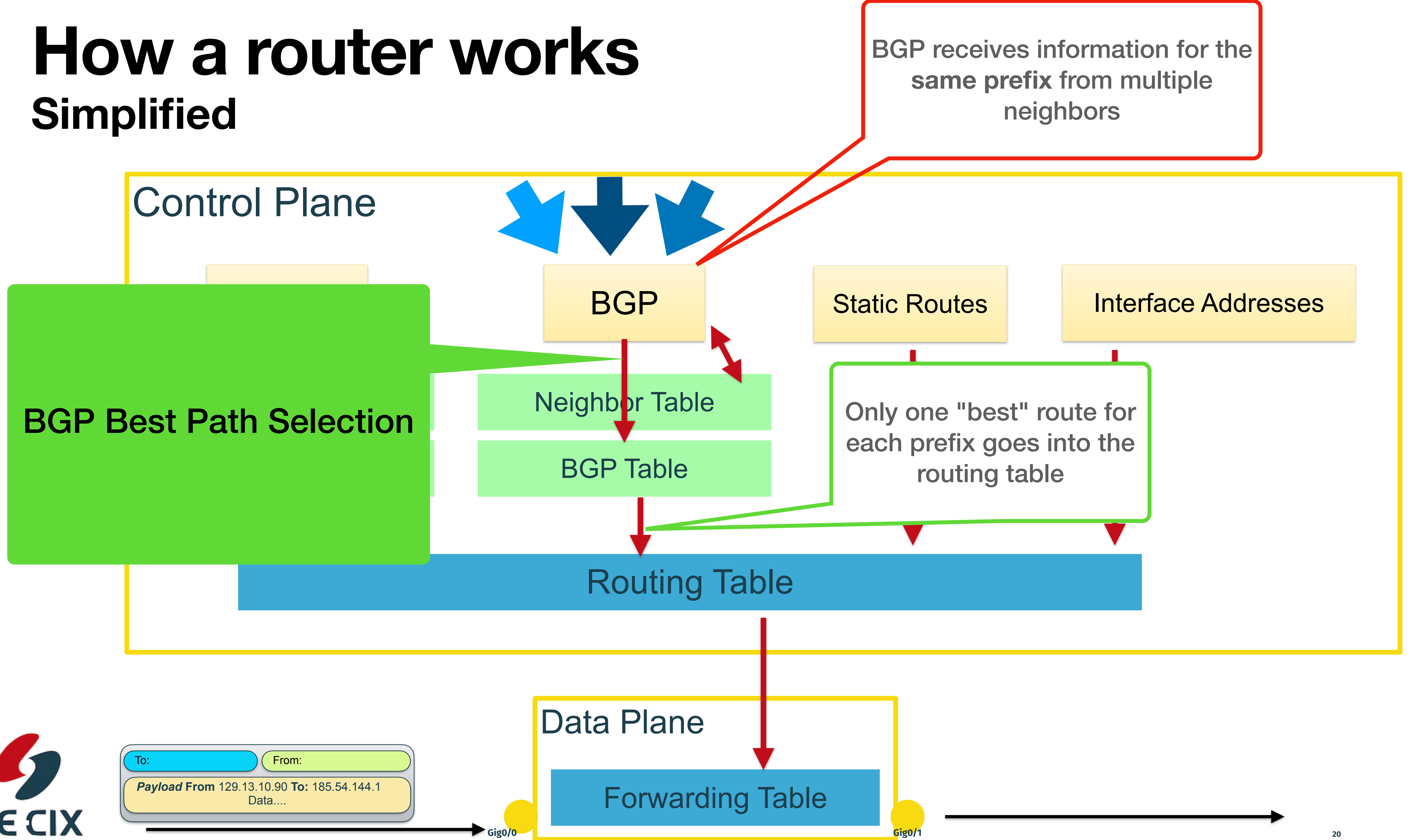
192.0.2.0/24
5669 286 **517**

# Attributes of BGP prefixes

**Not only the AS path**

- **Mandatory** attributes: have to be there

  - Example: AS-Path

- **Optional** attribute: are, well, optional

  - Example: MED


- **Transitive** attributes

  - are kept on the prefix and forwarded via BGP

- **Non-transitive** attributes

  - are added to a prefix and not forwarded by the receiver
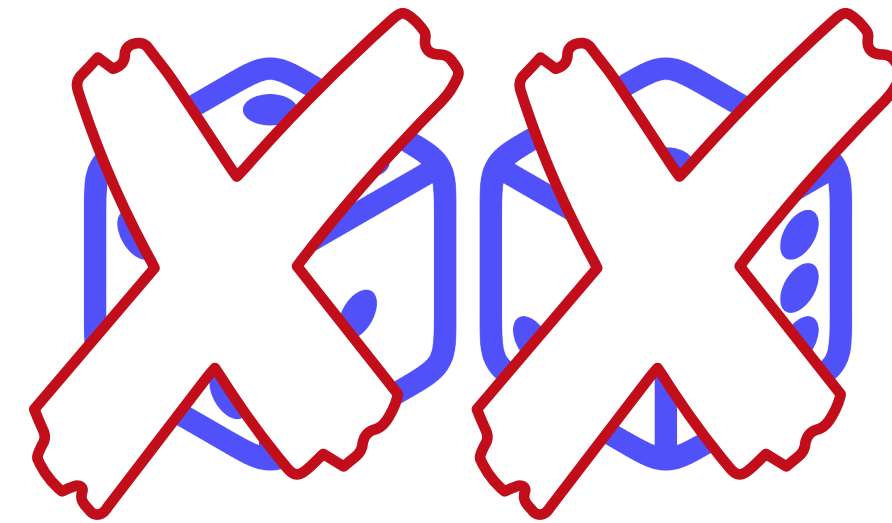
# How a router works
## Simplified



BGP receives information for the **same prefix** from multiple neighbors

Control Plane

BGP Best Path Selection

BGP

Static Routes

Interface Addresses

Neighbor Table

BGP Table

Only one "best" route for each prefix goes into the routing table

Routing Table

Data Plane

To:     From:

*Payload* **From** 129.13.10.90 **To:** 185.54.144.1 Data....

Forwarding Table

Gig0/0     Gig0/1

20

# BGP Best Path Selection

# BGP Best Path Selection Algorithm
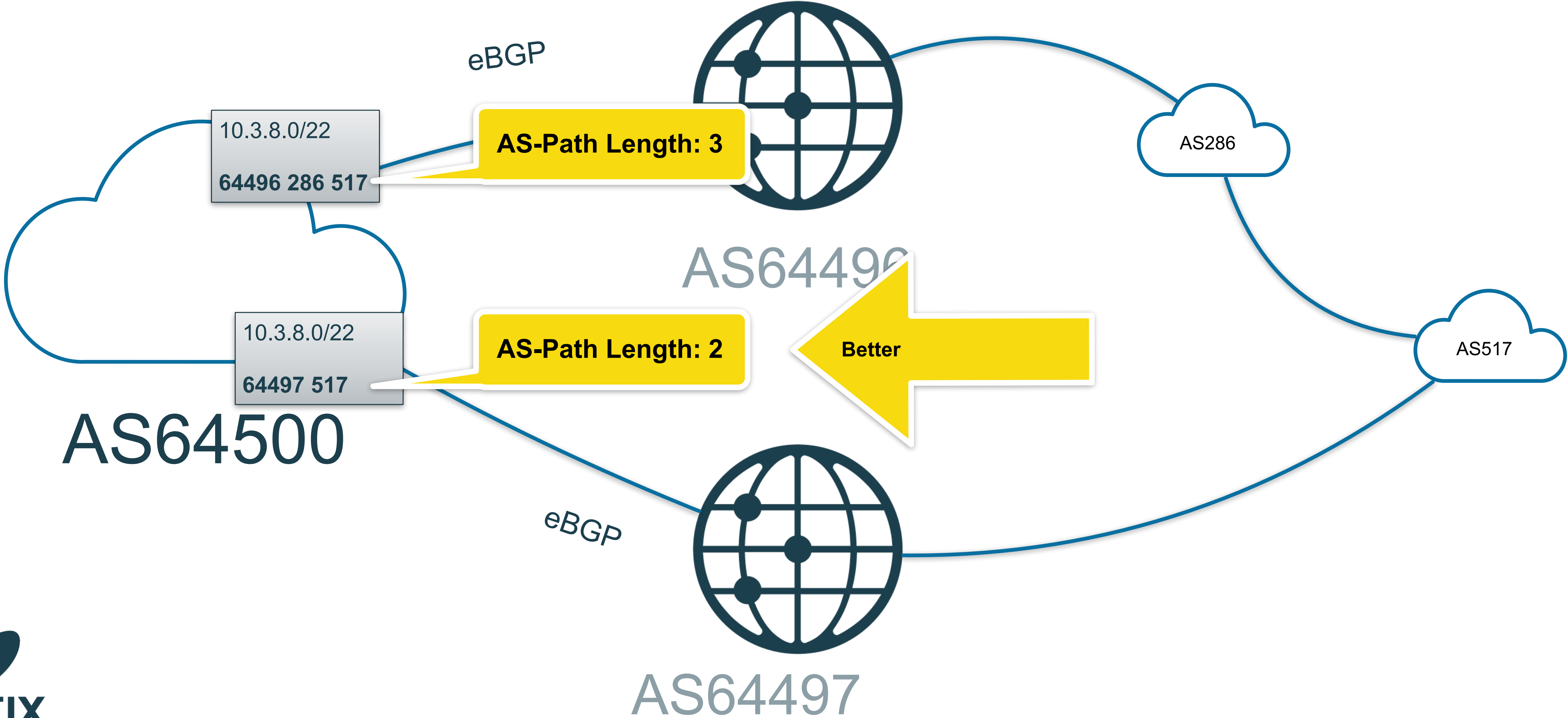## Motivation

**1**

- Only one single path for each destination is needed (and wanted)

- Decision must be based on attributes

- And must not be random, but deterministic

- Some of the criteria will sound strange

- Some are really outdated

- So lets have a look how this works...

# Let's get started.... with two upstreams



eBGP

10.3.8.0/22

64496 286 517

AS64496

10.3.8.0/22

286 517

10.3.8.0/22

517

AS64500

eBGP

10.3.8.0/22

64497 517

AS64497

*Where networks meet*

www.de-cix.net

DE-CIX

23

# Let's get started.... with two upstreams



eBGP

10.3.8.0/22

**64496 286 517**

**AS-Path Length: 3**

AS64496

AS286

10.3.8.0/22

**64497 517**

**AS-Path Length: 2**

**Better**

AS517

AS64500

eBGP

AS64497

# BGP Best Path Selection

| 1 | NextHop reachable? | Continue if "yes" |
|---|---|---|
| 2 | | |
| 3 | | |
| 4 | | |
| 5 | | |
| 6 | | |
| 7 | | |
| 8 | | |
| 9 | | |
| 10 | | |

AS-Path Length: 3

AS-Path Length: 2

Better

# Let's add peering



eBGP

10.3.8.0/22

**64496 286 517**

AS64496

10.3.8.0/22

**286 517**

eBGP

10.3.8.0/22

**64497 517**

AS64500

10.3.8.0/22

**517**

DE CIX

eBGP

AS64497

DE CIX

# Let's add peering



eBGP

10.3.8.0/22
**64496 286 517**

10.3.8.0/22
**286 517**

AS-Path Length: 2

10.3.8.0/22
**64497 517**

AS-Path Length: 2

AS64500

AS64496

AS286

AS517

AS64497

eBGP

DE CIX

**Where networks meet**

www.de-cix.net

27

# BGP Best Path Selection

| | | |
|---|---|---|
| 1 | NextHop reachable? | Continue if "yes" |
| 2 | | |
| 3 | AS Path Length | shorter wins |
| 4 | | |
| 5 | | |
| 6 | | |
| 7 | | |
| 8 | | |
| 9 | | |
| 10 | | |

AS-Path Length: 2

AS-Path Length: 2

# *Local Preference*

➔ Higher wins

➔ Integer value (32bit, 0-4294967295)

➔ Propagated via iBGP inside an Autonomous System

➔ Usually set using rules when receiving prefixes

➔ Typical  values:

· Customer prefixes:   10000

· Peering prefixes:     1000

· Upstream prefixes:   10

Why am I not using "100" here?

| 1 | NextHop reachable? | Continue if "yes" |
|---|---|---|
| 2 | Local Preference | higher wins |
| 3 | AS Path Length | shorter wins |
| 4 | | |
| 5 | | |
| 6 | | |
| 7 | | |
| 8 | | |
| 9 | | |
| 10 | | |

DE-CIX

*Where networks meet*

www.de-cix.net

29

# BGP Route Selection: Origin Type

→ Origin Type is a "historical" attribute

→ Three possible values:

→ IGP - route is generated by BGP network statement - "i"

→ EGP - route is received from EGP - **"e"**

→ incomplete - redistributed from another protocol -
   **"?"** as the "real source" is unknown

→ *This rule is not really important*

→ Fun fact: There are prefixed in the global routing table marked "e"

**E**xterior **G**ateway **P**rotocol

Predecessor of BGP which is no longer used

| 1 | NextHop reachable? | Continue if "yes" |
|---|---|---|
| 2 | Local Preference | higher wins |
| 3 | AS Path Length | shorter wins |
| 4 | | |
| 5 | | |
| 6 | | |
| 7 | | |
| 8 | | |
| 9 | | |
| 10 | | |

Where networks meet

# Consider the following network

Receives Prefixes via eBGP

Prefixes

AS64496

eBGP

Traffic

iBGP

iBGP

Provides Transit Service

AS64500

eBGP

DE-CIX

*Where networks meet*

www.de-cix.net

# Consider the following network

→ There are two circuits

→ AS64496 wants one of them preferred

→ How to tell AS64500?

?

Prefixes

AS64496

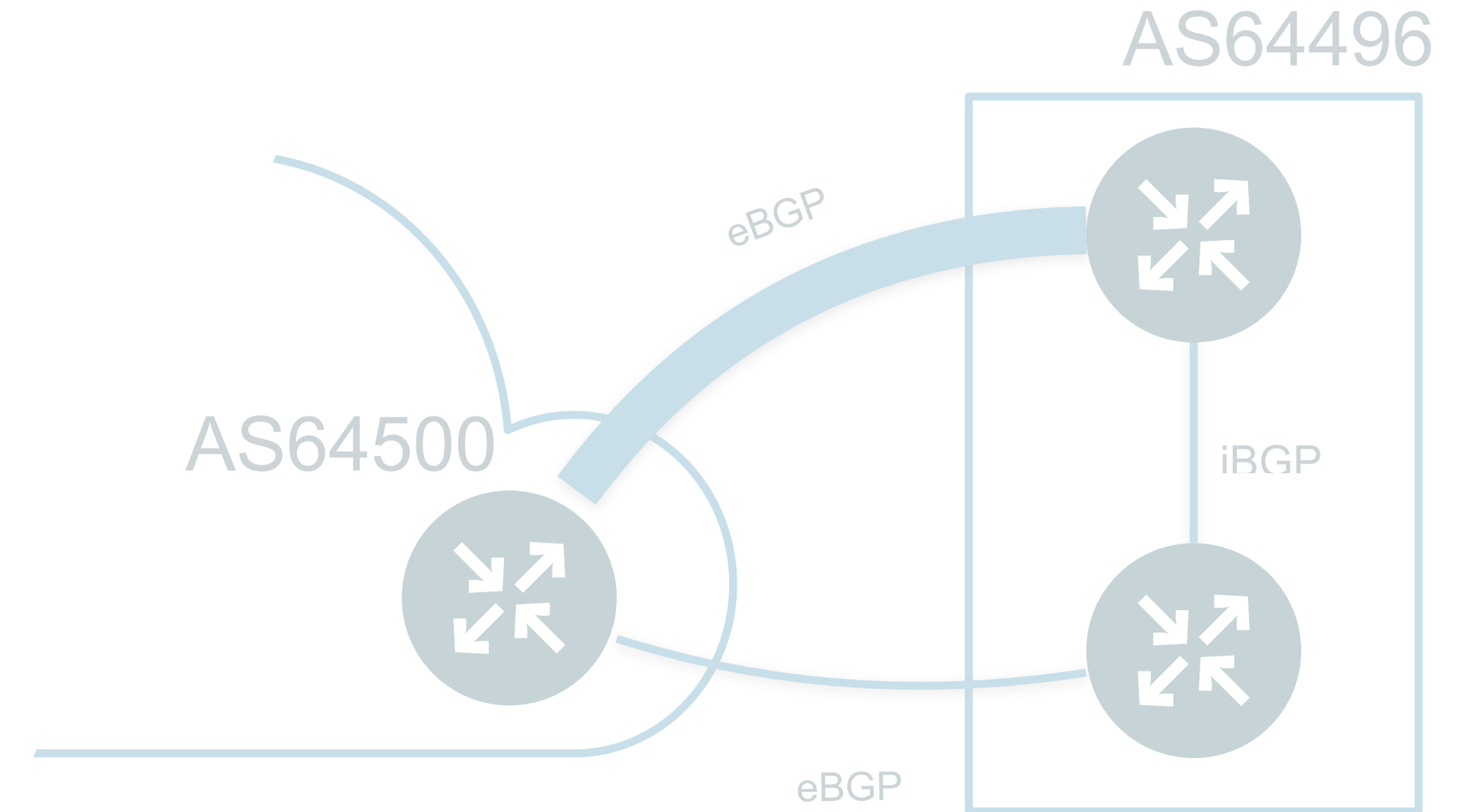eBGP

Traffic

iBG

iBGP

AS64500

eBGP

Where networks meet

www.de-cix.net

# BGP Route Selection Algorithm:

## How to tell your neighbor where you prefer traffic?

| 1 | NextHop reachable? | Continue if "yes" |
|---|---|---|
| 2 | Local Preference | higher wins |
| 3 | AS Path Length | shorter wins |
| 4 | Origin Type | IGP over EGP over Incomplete |
| 5 | | |
| 6 | | |
| 7 | | |
| 8 | | |
| 9 | | |
| 10 | | |

**DE-CIX**

**Where networks meet**

# BGP Route Selection Algorithm: MED

AS64496

AS64500

eBGP

iBGP

eBGP

➔ MED = **M**ulti-**E**xit **D**iscriminator

➔ Only compared if next-hop AS is the same

➔ 32bit value (0..4294967294)

➔ Lower wins

➔ Optional (does not have to be there),
non-transitive (does not get forwarded)

➔ A missing MED can be treated as "best" (=0, default)
or "worst" (=4294967294)

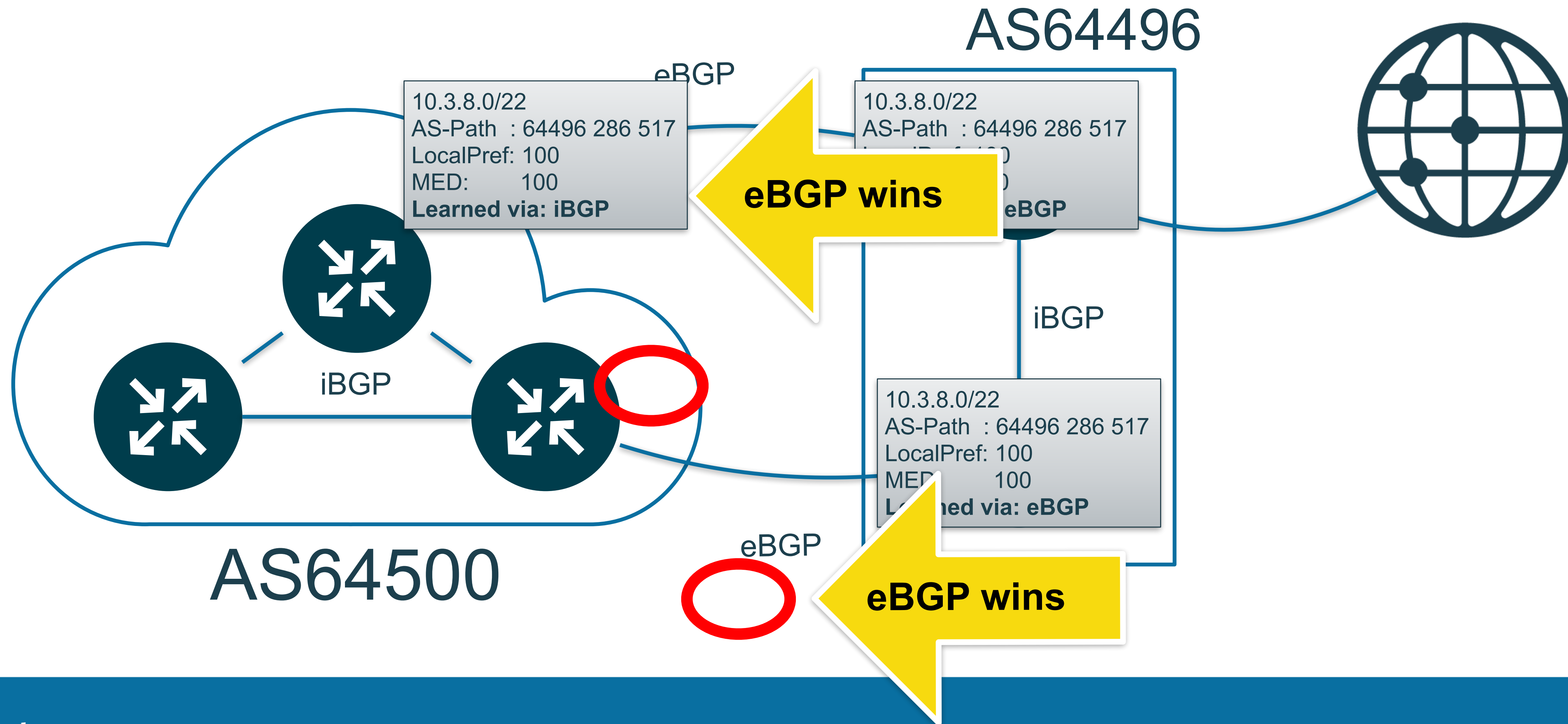➔ And of course you can override whatever you receive

# BGP Route Selection : Hot Potato Rules

| | | | |
|---|---|---|---|
| 1 | NextHop reachable? | Continue if "yes" | |
| 2 | Local Preference | higher wins | |
| 3 | AS Path Length | shorter wins | |
| 4 | Origin Type | IGP over EGP over Incomplete | |
| 5 | MED | lower wins | |
| 6 | | | |
| 7 | | | |
| 8 | | | |
| 9 | | | |
| 10 | | | |

# BGP Route Selection : eBGP wins

AS64496

```
10.3.8.0/22
AS-Path  : 64496 286 517
LocalPref: 100
MED:       100
Learned via: iBGP
```

eBGP

eBGP wins

```
10.3.8.0/22
AS-Path  : 64496 286 517
         100
eBGP
```

iBGP

iBGP

AS64500

```
10.3.8.0/22
AS-Path  : 64496 286 517
LocalPref: 100
MED        100
Learned via: eBGP
```

eBGP

eBGP wins

DE-CIX

*Where networks meet*

www.de-cix.net

# BGP Route Selection : nearest exit wins

AS64496

eBGP

iBGP

iBGP

AS64500

eBGP

DE-CIX

# BGP Route Selection : Age / Stability

| | | |
|---|---|---|
| 1 | NextHop reachable? | Continue if "yes" |
| 2 | Local Preference | higher wins |
| 3 | AS Path Length | shorter wins |
| 4 | Origin Type | IGP over EGP over Incomplete |
| 5 | MED | lower wins |
| 6 | eBGP, iBGP | eBGP wins |
| 7 | Exit | nearest wins |
| 8 | | |
| 9 | | |
| 10 | | |

**Where networks meet**

# BGP Route Selection : Age / Stability

➔ Exact phrasing is (Cisco):
"When both paths are external, prefer the path that was received first"

➔ So this applies only if a router has two (or more) eBGP sessions

➔ Which happens quite often when connecting to Internet Exchanges

# BGP Route Selection : Last Resort

| | | |
|---|---|---|
| 1 | NextHop reachable? | Continue if "yes" |
| 2 | Local Preference | higher wins |
| 3 | AS Path Length | shorter wins |
| 4 | Origin Type | IGP over EGP over Incomplete |
| 5 | MED | lower wins |
| 6 | eBGP, iBGP | eBGP wins |
| 7 | Exit | nearest wins |
| 8 | Age of route | older wins |
| 9 | | |
| 10 | | |

# BGP Route Selection : Last Resort

→ Router ID: lower wins

→ Neighbor IP: lower wins

→ Rules of last resort

→ ...because at the end one and only one best path has to be selected

→ Usually path selection stops before it gets to these two rules.

**BGP**
**Last Exit** ↗

| | | |
|---|---|---|
| 1 | NextHop reachable? | Continue if "yes" |
| 2 | Local Preference | higher wins |
| 3 | AS Path Length | shorter wins |
| 4 | Origin Type | IGP over EGP over Incomplete |
| 5 | MED | lower wins |
| 6 | eBGP, iBGP | eBGP wins |
| 7 | Exit | nearest wins |
| 8 | Age of route | older wins |
| 9 | Router ID | lower wins |
| 10 | Neighbor IP | lower wins |

*Where networks meet*

# *BGP Route Selection : Summary*

| 1 | NextHop reachable? | Continue if "yes" |
|----|--------------------|--------------------|
| 2 | Local Preference | higher wins |
| 3 | AS Path Length | shorter wins |
| 4 | Origin Type | IGP over EGP over Incomplete |
| 5 | MED | lower wins |
| 6 | eBGP, iBGP | eBGP wins |
| 7 | Exit | nearest wins |
| 8 | Age of route | older wins |
| 9 | Router ID | lower wins |
| 10 | Neighbor IP | lower wins |

*Where networks meet*

www.de-cix.net

# Network relationships

# The Internet
## A typical Internet Service Provider

# The Internet
## Adding "Upstream"



DE CIX

# The Internet
## Adding a 2nd ISP



DE CIX

# The Internet
**Data transport via upstreams**

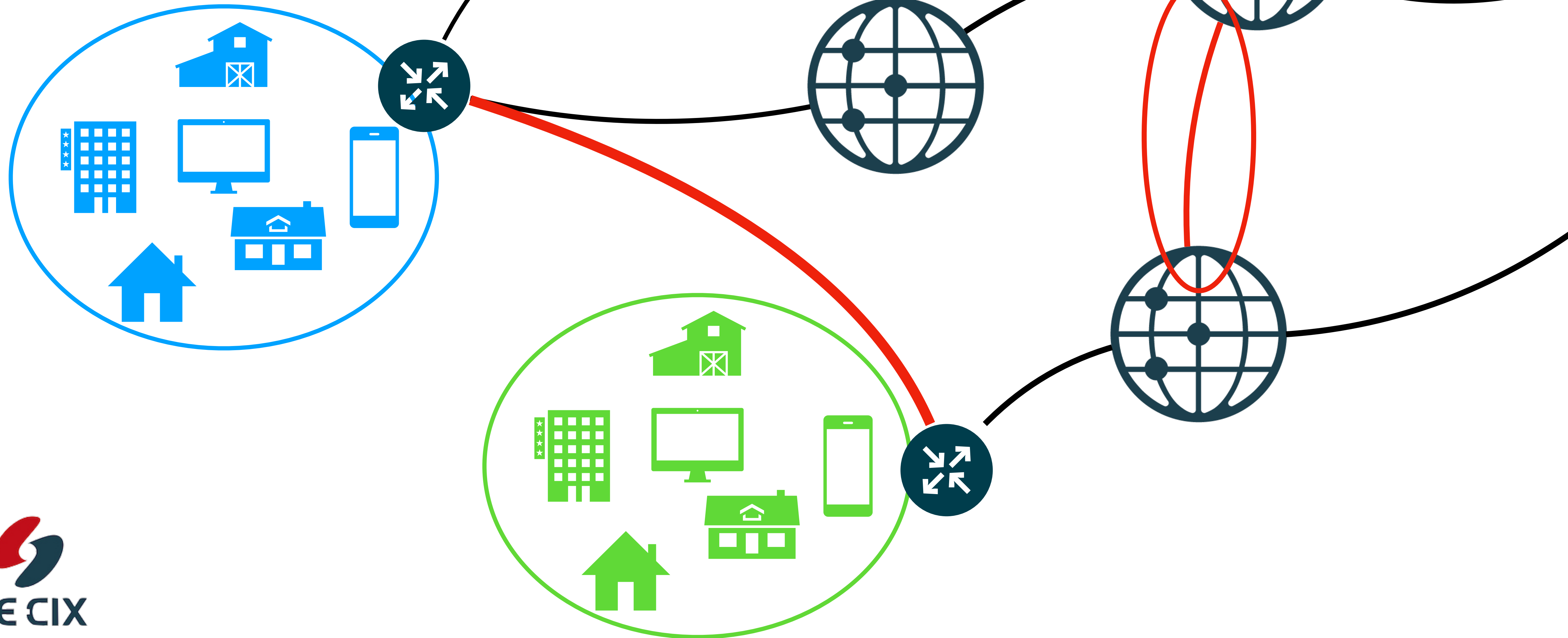# The Internet
## More direct via "peering"
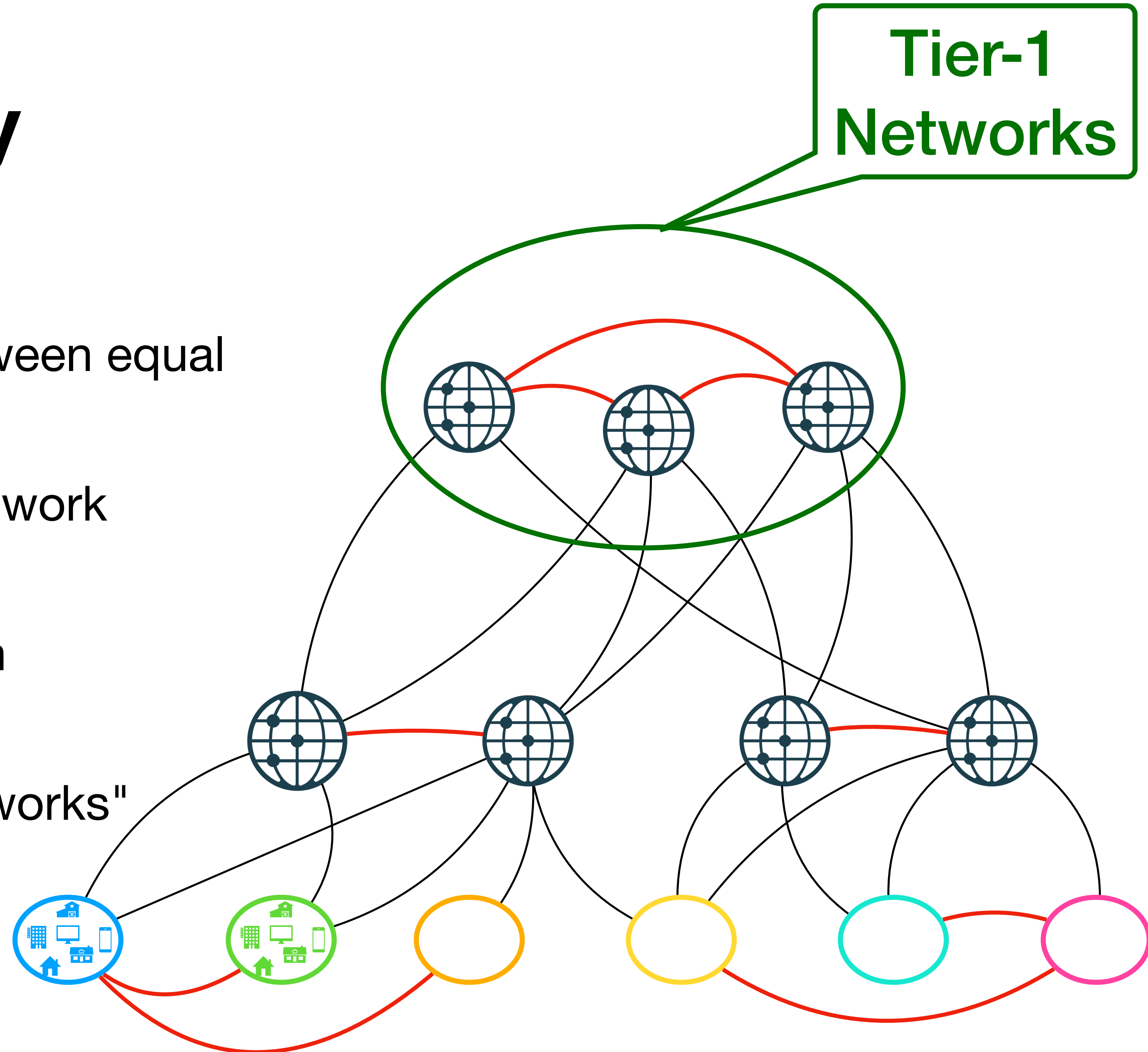
# The Internet
## Peering on multiple levels

# The Internet
## Peering on multiple levels

# Peering Hierarchy
## Peering on multiple levels

- Peering happens usually between equal size networks

- Peering takes place on all network levels

- The "top ones" only peer with each other

  - They are called "Tier-1 networks"



Tier-1 Networks

DE CIX

# Public tools for BGP

# Public tools for BGP

## RIPE Stat

- Operated by the RIPE NCC (same entity handing out AS numbers in this region)

- Details about prefixes, ASes and more

- just check it out at https://stat.ripe.net

# Public tools for BGP

## bgp.tools

- Private initiative

- Free, offer premium monitoring service for a fee

- just check it out at https://bgp.tools

# Public tools for BGP

## bgp.he.net

- Operated by Hurricane Electric ([he.net](he.net))

- Free, but shows only HEs point of view

- just check it out at [https://bgp.he.net](https://bgp.he.net)



DE-CIX

# Public tools for BGP

## BGP Alerter

- Open source tool running locally

- Using data from public datasets

  - like ris.ripe.net

- Get the source or a precompiled binary from https://github.com/nttgin/BGPalerter

```
Wolfgangs-MacBook-Pro-273:Downloads wtremmel$ ./bgpalerter-macos-x64
Loaded config: /Users/wtremmel/Downloads/config.yml
Impossible to load config.yml. A default configuration file has been generated.
BGPalerter, version: 1.32.0 environment: production
? The file prefixes.yml cannot be loaded. Do you want to auto-configure BGPalerter? Yes
? Which Autonomous System(s) you want to monitor? (comma-separated, e.g., 2914,3333) 196610
? Do you want to be notified when your AS is announcing a new prefix? Yes
? Do you want to be notified when a new upstream AS appears in a BGP path? Yes
? Do you want to be notified when a new downstream AS appears in a BGP path? Yes
Getting announced prefixes of AS196610
Total prefixes detected: 2
Generating monitoring rule for 2a02:c50:db8::/48
Generating monitoring rule for 91.214.253.0/24
Detected upstreams for 196610: 1239, 13786, 15704, 15830, 20485, 24889, 25091, 29075, 30781, 31133, 321
4, 34019, 34549, 34927, 35280, 35710, 37468, 39351, 41327, 4230, 43350, 43727, 4455, 47605, 47734, 4836
2, 49697, 50629, 51531, 6939, 8447, 8758, 8932, 8966, 9002
Detected downstreams for 196610: 10122, 10310, 10466, 11284, 11403, 12297, 12335, 12389, 12418, 12430,
12479, 12540, 12578, 12668, 12714, 12741, 13094, 13213, 13287, 13335, 13414, 13536, 136907, 137409, 137
86, 138915, 14061, 14537, 14593, 14928, 15133, 15599, 15672, 15682, 15699, 15704, 15754, 15757, 15930,
15954, 16164, 16552, 17378, 18001, 1820, 1828, 18966, 19318, 19551, 196709, 19689, 197204, 197267, 1975
18, 197826, 198367, 199226, 199290, 199434, 199524, 199599, 199610, 199952, 199976, 200030, 200350, 200
380, 200845, 201359, 201746, 201776, 202054, 202087, 202173, 202207, 202334, 202486, 20253, 202766, 202
813, 202829, 202844, 202984, 203099, 203724, 203936, 20473, 204773, 204805, 204861, 205022, 205627, 205
675, 205697, 20655, 206810, 20710, 20764, 207785, 207923, 209141, 20940, 209674, 209835, 210123, 210756
, 211157, 211227, 211826, 21719, 21859, 21949, 22356, 22418, 22697, 22742, 23393, 23470, 23764, 24429,
24482, 24663, 24768, 25292, 25532, 25549, 262589, 263444, 2635, 266925, 267613, 2683, 27257, 27611, 280
07, 28189, 2860, 28761, 28891, 28917, 2906, 29117, 29119, 29124, 29226, 29303, 29337, 29470, 29479, 296
32, 29802, 29838, 29852, 30081, 30833, 31214, 31500, 31514, 31769, 31950, 32035, 3218, 32217, 3223, 324
25, 3267, 32787, 32934, 3316, 3327, 33353, 33438, 33570, 34123, 34352, 34879, 35168, 35280, 35394, 3552
2, 35539, 35598, 35699, 36236, 36351, 36591, 36891, 37468, 38040, 39020, 39063, 39134, 39328, 39337, 39
386, 394102, 39684, 39691, 396986, 396998, 398465, 398930, 399100, 40545, 40676, 40805, 4134, 4136, 414
46, 41617, 41690, 41721, 41731, 41798, 42, 4230, 42325, 42473, 42511, 42518, 4258, 42632, 42649, 42947,
 43160, 43298, 43727, 43832, 43996, 44020, 44128, 44391, 44670, 44814, 47321, 47541, 47542, 47569, 4776
4, 47775, 47787, 48084, 48249, 48287, 48293, 48348, 48366, 48524, 48719, 48739, 48846, 48848, 49403, 49
544, 49697, 49724, 49776, 49779, 49813, 50060, 50304, 50509, 50646, 50923, 51531, 51681, 51764, 51865,
52091, 52320, 52468, 53766, 53828, 53991, 54113, 5467, 54994, 5505, 5518, 55256, 55805, 55818, 56630, 5
6814, 56958, 57073, 57363, 57365, 57463, 57624, 57724, 57877, 57910, 57976, 58310, 59865, 60068, 60280,
 60488, 60767, 6079, 60840, 60917, 61031, 61090, 61461, 61832, 62044, 62240, 62668, 62904, 63399, 63949
, 64049, 6507, 6774, 6789, 6866, 6939, 7195, 7713, 8002, 8242, 8301, 8331, 8359, 8400, 8629, 8764, 8966
, 9009, 9049, 9110, 9304, 9498
Generating generic monitoring rule for AS196610
Done!
Monitoring 91.214.253.0/24
Monitoring 2a02:c50:db8::/48
Monitoring AS196610
```

# Public tools for BGP

## ExaBGP

- Open source tool to "talk" BGP

- Use cases:

  - for testing or even in production

  - announce prefixes

  - with any attributes you want

- https://github.com/Exa-Networks/exabgp

```
ubuntu@bgplab:~/BGPLab/experiment-02$ exabgp exabgp.conf
14:04:55 | 1493    | welcome      | Thank you for using ExaBGP
14:04:55 | 1493    | version      | 4.2.17
14:04:55 | 1493    | interpreter  | 3.10.6 (main, May 29 2023, 11:10:38) [GCC 11.3
14:04:55 | 1493    | os           | Linux bgplab 5.15.0-76-generic #83-Ubuntu SMP
TC 2023 x86_64
14:04:55 | 1493    | installation |
14:04:55 | 1493    | cli control  | named pipes for the cli are:
14:04:55 | 1493    | cli control  | to send commands  /run/exabgp.in
14:04:55 | 1493    | cli control  | to read responses /run/exabgp.out
14:04:55 | 1493    | configuration | performing reload of exabgp 4.2.17
14:04:55 | 1493    | reactor      | loaded new configuration successfully
```

DE CIX

# Public tools for BGP

## DE-CIX Academy BGP lab

- For teaching a BGP seminar

- Based on FRRouting

- Runs (multiple) routers in Docker containers

- Just needs a linux server as host

- Get it at https://gitlab.com/de-cix-public/team-academy/bgp/BGPLab

# Managing BGP relationships

# The lazy Network Manager

**How to keep record of your peers**

# Setting up BGP sessions
## Standard procedure

- Contact your neighbor

- Exchange a few emails

- Configure BGP

# Years later...

# You need to contact your neighbor

**But where did I put the contact information**

- I might have my original emails somewhere

- Or I put the contact information into an Excel sheet

- Or I configured it as a comment on my router

- Or....

# But then you notice...

# But then you notice...

**Surprise, surprise...**

- The contact you emailed with works no longer there

- The company name of your peer has changed

- The email address you have (peering@...) is no longer valid

- What now?

# There is a solution

# Why not have a common database?

**For networks who peer...**

- Put contact information into a central database

- Make it accessible for all networks who peer

- Everybody maintains their own information (hopefully)

- If you need some information, simply look it up

DE CIX

# PeeringDB

## A database for networks who peer

- Free for users

- Financed by sponsoring

- Some public information

- Contact data is private

- Check it out at https://peeringdb.com

# Other versions of this presentation

# BGP in 120 minutes
## What we did today

- Length: 90-120 minutes

- Features:

  - me talking

  - you asking questions

- Covers:

  - The very basics of BGP

  - Up and including BGP best path selection / more depending on time

120

BGP!

# BGP 4-5 hour workshop
## Not just the basics...

- Length: 4-5 hours, including at least one break

- Happened a number of times at workshop Sunday at DENOG

- Features:

  - Me talking

  - You asking questions

  - Limited number of **lab experiments** using FRRouting

- Covers:

  - The very basics of BGP

  - Up and including BGP best path selection

  - BGP Communities if time permits

4-5h

BGP!

?

# 3.5 Day BGP Seminar
## All and everything

- Length: 3.5 days, starting Monday noon, finishing Thursday late afternoon,

- Classroom seminar, max. 14 attendees

- Features:

  - Me talking

  - You asking questions

  - Extensive number of lab experiments using FRRouting

- Covers:

  - All of BGP

  - Including BGP Security, Traffic Engineering, Peering Relationships

  - Tools useful for BGP and peering

# BGP and Security

# *Protect your routing infrastructure*

# Reference Document on BGP Security

# RFC 7454

**Where networks meet**

www.de-cix.net

# RFC 7454

**BGP Operations and Security**

Abstract

The Border Gateway Protocol (BGP) is the protocol almost exclusively used in the Internet to exchange routing information between network domains.  Due to this central nature, it is important to understand the security measures that can and should be deployed to prevent accidental or intentional routing disturbances.

This document describes measures to protect the BGP sessions itself such as Time to Live (TTL), the TCP Authentication Option (TCP-AO), and control-plane filtering.  It also describes measures to better control the flow of routing information, using prefix filtering and automation of prefix filters, max-prefix filtering, Autonomous System (AS) path filtering, route flap dampening, and BGP community scrubbing.

# Simple Measures

➔ Easy to implement

➔ Easy to maintain

➔ ...but only of limited use

➔ ...still should be implemented

➔ List of measures:

1. **Maximum Prefix**

# Maximum Prefix



Sending:
**10** prefixes

AS64500

Receiving:
usually about **10** prefixes

iBGP

Where networks meet

www.de-cix.net

# Maximum Prefix

Sending:
**345000** prefixes

AS64500

ALARM

Receiving:
usually about **10** prefixes

iBGP

# Maximum Prefix

Sending:
**345000** prefixes

AS64500

ALARM

Receiving:
usually about **10** prefixes

iBGP

**Maximum Prefix**:
- Define a *threshold*
- Define an action if threshold is hit
  - Usually tear down the BGP session

# *Maximum Prefix*

➔Good counter-measure against misconfigured peers

➔Possible actions:

➔Tear down session (until manual intervention)

➔Tear down and restart (after $n$ minutes)

➔Warning only

➔Best practices:

➔Set threshold high enough (like 10* usual size)

➔Configure a warning at 90% - so you **see** it!

ALARM

# Simple Measures

➔Easy to implement

➔Easy to maintain

➔...but only of limited use

➔...still should be implemented

➔List of measures:

    1. Maximum Prefix

    2. **MD5 Session Password / TCP AO**

# MD5 Session Password / TCP AO

AS64500

iBGP

BGP Session via TCP

TC

# MD5 Session Password / TCP AO

➔Set the same password on each side

➔Password is used to MD5 sign **each** TCP packet by the sender

➔Receiver checks the signature

➔If it does not match, packet is silently discarded

➔Still used, even MD5 no longer state of the art

➔More modern approach: TCP-AO (authentication option) with stronger hashes

➔Recommendation: Use this for iBGP, but not for eBGP

➔Important: **You need some password management**!

Where networks meet

# Simple Measures

➔ Easy to implement

➔ Easy to maintain

➔ ...but only of limited use

➔ ...still should be implemented

➔ List of measures:

    1. Maximum Prefix

    2. MD5 Password

    3. **IP Time-to-live security**

# IP Time-to-live security

BGP TCP Packet
**IP TTL = 1**

AS64500

iBGP

BGP Session via TCP

# IP Time-to-live security

AS64500

iBGP

BGP Session via TCP

# IP Time-to-live security



BGP TCP Packet
**IP TTL = 255**

AS64500

iBGP

BGP Session via TCP

**Where networks meet**

www.de-cix.net

88

# *IP Time-to-live security*

→Send IP packets with initial TTL of 255

→Receiver checks if value is really 255

→If not, packet is silently discarded

→Very easy to implement (just enable it)

→But **must be configured on both sides**

→Defined in RFC5082

**Where networks meet**

*www.de-cix.net*

# *BGP Filtering*

# BGP Filtering



**Blocklist**

**The Good Stuff***

**Raw Input**

**Allowlist**

* Diagram according to Job Snijders

# *BGP Filtering*

**Blocklist**

➔Filtering received prefixes

➔ Prefix filtering

➔AS Path filtering

➔RPKI

* Diagram according to Job Snijders

**DE·CIX**

**Where networks meet**

**www.de-cix.net**

# *Prefix filtering*

**Martians**

➔Block non-routable IPv4 prefixes like:

➔Private IPv4 space

　➔10.0.0.0/8, 172.16.0.0/12, 192.168.0.0/16

　➔IPv4 networks reserved for documentation purposes

　➔IPv4 multicast address space - 224.0.0.0/4

　➔IPv4 reserved for "future use" - 240.0.0.0/4

➔For IPv6

　➔Allow only 2000::/3

　➔Block everything else

```
ip prefix-list ipv4-unwanted permit 192.168.0.0/16 le 32
ip prefix-list ipv4-unwanted permit 172.16.0.0/12 le 32
ip prefix-list ipv4-unwanted permit 10.0.0.0/8 le 32
!
route-map upstream-in deny 100
   match ip address prefix-list ipv4-unwanted
```

**DE-CIX**

*Where networks meet*

www.de-cix.net

93

# *Prefix filtering*

➔Filter against too small and too large prefixes

➔IPv4:

 ➔Prefix sizes are /8 - /24

 ➔Block everything smaller or larger (Exception: Blackholing)

➔IPv6:

 ➔Prefix sizes are /19 - /48

➔You might allow a default
 route from your upstream
 providers

```
ip prefix-list ipv4-unwanted permit 0.0.0.0/0 ge 25
ip prefix-list ipv4-unwanted permit 0.0.0.0/0 ge 1 le 7
!
ipv6 prefix-list ipv6-unwanted permit ::/0 ge 49
ipv6 prefix-list ipv6-unwanted permit ::/0 le 18
!
route-map upstream-in deny 100
   match ip address prefix-list ipv4-unwanted
   match ipv6 address prefix-list ipv6-unwanted
```

Prefix Size

**DE-CIX**

*Where networks meet*

www.de-cix.net

# *More Prefix filtering*

→IXP Lan Prefixes (and their more specifics)

→Why? Have a look....

80.81.196.61/21

80.81.192.1          .2          .3          IXP Lan: 80.81.192.0/21

# *More Prefix filtering*

→IXP Lan Prefixes (and their more specifics)

→Why? Have a look....

BGP Prefix:
80.81.192.0/24

80.81.196.61/21

80.81.192.1   .2   .3   IXP Lan: 80.81.192.0/21

DE CIX

# *More Prefix filtering*

→IXP Lan Prefixes (and their more specifics)

→Your own prefixes

→Your customers prefixes (for the same reasons)

```
ip prefix-list ipv4-unwanted permit 80.81.192.0/21 le 32
!
ipv6 prefix-list ipv6-unwanted permit 2001:7f8::/64 le 128
!
route-map upstream-in deny 100
  match ip address prefix-list ipv4-unwanted
  match ipv6 address prefix-list ipv6-unwanted
```

# AS path filtering

→even if the prefix is totally legit, the AS path might be bad

→if your own AS is in the path, prefixes are filtered automatically

→but you need to filter against...

→private ASes (64512-65534 + 4200000000-4294967294)

→reserved ASes - see IANA Special-Purpose AS Numbers

→**anywhere in the AS path!**

```
203.0.113.0/24     192.0.2.1  517 48854 65101 65102 203453 203453 203453 i
```

# AS path filtering

➔private ASes (64512-65534 + 4200000000-4294967294)

➔reserved ASes - see IANA Special-Purpose AS Numbers

➔regular expressions can be used

➔but do not overdo it!

  ➔ _(6451[2-9]|645[2-9][0-9]|64[6-9][0-9]{2}|65[0-4][0-9]{2}|655[0-2][0-9]|6553[0-5])_

➔completely valid, but unreadable

➔better split it up

# *Filtering from Customers*

→Here we need the "Allowlist"

→remember?

→From customers allow only

→Customers prefixes

→Customers ASes (anywhere in the path)

→Use this to create an Allowlist **per customer**

Blocklist

Raw
Input

Allowlist

**The Good
Stuff**

DE-CIX

# *Control your Announcements*

## *have good*



DE CIX

# MANRS

→Mutually

→Agreed

→Norms for

→Routing

→Security

# MANRS

➔Prevent propagation of incorrect routing information

 ➔Filter incoming - what you do not let in, you cannot announce

 ➔Do not announce anything outgoing you should not

➔Prevent traffic with spoofed source IP addresses

➔Facilitate global operational communication and coordination between
 network operators

 ➔= "talk to each other"

➔Facilitate validation of routing information on a global scale

# *Conclusion*

# *Conclusion*

→Protect your BGP routers and sessions

→Filter incoming

   →Unwanted IP Prefixes

   →Bogus ASes in the path

   →Allow from customers using an **Allowlist**

→Filter outgoing

   →Make sure you announce only valid prefixes

**The Good Stuff**

Blocklist

Raw Input

Allowlist

**DE·CIX**

# More security: RPKI

# Part 1

# RPKI - What is it?

# Certificate - based proof of address assignment

→There are only five entities handing out IP resources

→These are the five RIRs



Source: https://www.ripe.net/about-us/what-we-do/ripe-ncc-service-region

## Certificate - based proof of address assignment

➔There are only five entities handing out IP resources

 ➔These are the five RIRs

➔These are the trust-anchors in this model

➔They have contractual proof who they gave which resource

➔And sign you a certificate for it

*Certificate*

```
91.214.253.0/24
was assigned to
DE-CIX Academy
```

Signed by
RIPE NCC

DE-CIX

**Where networks meet**

# What can you do with this certificate?

➔ You can create a ROA - Route Origin Authorization

➔ ROAs contains three values:

  ➔ The IP prefix it is for

  ➔ An Autonomous System Number you allow to originate that prefix

  ➔ A range of allowed netmasks for this prefix

➔ And of course its digitally signed

**Where networks meet**

www.de-cix.net

# RPKI - Resource Public Key Infrastructure

➔Digitally signed route objects

➔Resource holders can get their resources signed

➔And can define how they are announced

   ➔Define an Origin-AS

   ➔Define a maximum length of a prefix

   ➔This is called a "ROA" (Route Origin Authorisation)

➔Routers can use this to validate BGP announcements

192.0.2.0/24

*My Prefix*

# RPKI - Resource Public Key Infrastructure

➔What problem does it want to solve?

➔Certificate - based proof of resource assignment

➔Resources are IP prefixes and AS numbers

➔Verifiable originator AS for each prefix

➔Only allow certain prefix lengths

# *Example of a ROA*

➔ So, if your prefix is 91.214.253.0/24

➔ You might allow AS196610 to originate

➔ The prefix and also a /25 more-specific

➔ So the ROA looks like:

➔ 91.214.253/24        AS196610        /24-/25

➔ Or in detail:

```
{
  "filename": "b2zDxaYsNBGNNz0Iu93sUJUQ27I.roa",
  "asn": "AS196610",
  "validity_period": "2019-01-01T01:20:09.000Z -
2020-07-01T00:00:00.000Z",
  "signing_time": "2019-01-01T01:20:09.000Z",
  "prefixes": [
    {
      "prefix": "91.214.253.0/24",
      "maxLength": 25
    }
  ],
  "validation_result": {
    "isValid": true,
    "error": [],
    "warning": []
  }
}
```

**DE CIX**

*Where networks meet*

*Part 2*

# RPKI - setting it up

# *Your certificates and ROAS*

➔ Hosted RPKI: Easiest deployment

  ➔ Your RIR hosts your certificates and ROAs for you

  ➔ takes care of signing and key roll over

  ➔ Most RIRs have a nice web interface for that

➔ Non-Hosted RPKI: Run everything on your own

  ➔ If you heard about RPKI here the first time,
    this is probably not what you want to do

# Part 3

# Validating your ROAs

# So what does a *validator* do?

→Fetch resource certificates and ROAs from RIRs (via rsync)

→Validates the "chain of trust"

 →Check signatures of certificates

 →Check signatures of ROAs

→Supplies a *validated cache* for your routers

RIPE NCC

ARIN

APNIC

lacnic

AFRINIC

Fetch

Certificates and ROAs

Server

Caching Validator

Validated Cache

RPKI-RTR

Router

*Where networks meet*

*www.de-cix.net*

# *Example Validator:* ROUTINATOR

→by NLNetlabs

→https://nlnetlabs.nl/projects/rpki/routinator/

→written in RUST

→runs either directly

→or inside a Docker container

→Open Source

→Small footprint

→Very easy to install

# Part 4

# RPKI at DE-CIX

# DE-CIX Route Servers are using RPKI

Server

Caching
Validator

Validated
Cache

Certificates and ROAs

**ROA**
198.51.100/24 - /24
AS64499

**route server**

**AS6695**

203.0.113.99/24

invalid

203.0.113.1/24

203.0.113.2/24

203.0.113.3/24

203.0.113.4/24

198.51.100.0/24
AS-Path: 64500

DE-CIX

# DE-CIX Route Servers are using RPKI



Server

Caching
Validator

Validated
Cache

Certificates and ROAs

no ROA

route server
AS6695

192.0.2.0/24
AS-Path: 64500

.113.99/24

**?** Unknown

203.0.113.1/24

192.0.2.0/24
AS-Path: 64500

203.0.113.2/24

203.0.113.3/24

203.0.113.4/24

DE CIX

*Where networks meet*

*www.de-cix.net*

1

# BGP Error Handling

# Motivation: Blog entry

**Aug 29 2023**

## Grave flaws in BGP Error handling



Border Gateway Protocol is the de facto protocol that directs routing decisions between different ISP networks, and is generally known as the "glue" that holds the internet together. It's safe to say that the internet we currently know would not function without working BGP implementations.

However, the software on those networks' routers (I will refer to these as edge devices from now on) that implements BGP has not had a flawless track record. Flaws and problems do exist in commercial and open source implementations of the world's most critical routing protocol.

https://blog.benjojo.co.uk/post/bgp-path-attributes-grave-error-handling

# What happened?

- 2023-06-02: a small network announced one of their prefixes with a "bad" attribute

- This attribute was not understood by their immediate neighbors, and so the announcement was re-announced with the "bad" attribute unchanged

- Further away, Juniper routers "kind of" understood the attribute, saw it was bad, and **dropped the session they received it from**.

- So many routers, seemingly unrelated to the originator of the prefix, suddenly dropped sessions.

# Reminder: How BGP works

# BGP Neighbors
## Directly connected neighbors

192.0.2.0/24
198.51.100.0/24
2a02:c50:db8::/48

- BGP **announces** IP prefixes to **neighbors**

- These neighbors have to be **configured**

- BGP uses **TCP** to connect to a neighbor

- TCP brings already:

  - **Reliable transport** (sender knows that receiver got it)

  - **Flow control** (do not send faster than the receiver can receive)

  - **Framing** (putting BGP messages into packets)

BGP via TCP

Router

Neighbor Router

DE-CIX

# BGP works incremental
## Using add- / withdraw- messages

withdraw:
2a02:c50:db8::/48

- At session setup, BGP announces "everything" to its neighbor

- After that, updates are **incremental**:

  - If BGP learns about a new prefix, it sends an **add**-message to neighbors

  - If a prefix goes away, it sends a **withdraw** message to neighbors

- As long as the BGP session is "up", a router assumes its neighbors are "in sync" (= did not forget anything it sent)

# BGP Message Types
## TCP containing BGP messages

- BGP has the following message types:

  - **OPEN** - initial message for setting up a session

  - **UPDATE** - incremental routing updates: adds and withdraws

  - **KEEPALIVE** - send this if you have nothing to send

  - **NOTIFICATION** - to tell the other side there was an error, and then close the BGP session.

DE CIX

# Focus today: UPDATE message

## Incremental routing updates

- Update messages can contain multiple things:

  - A list of "adds"

    - Named "Network Layer Reachability Information"

    - With common attributes

  - A list of "withdraws"

    - Withdraws do not have attributes

| BGP Message |
|---|
| Marker (16 octets, all "1"s) |
| Total length (2 octets) |
| Type (Update = 2) |
| Withdrawn route length (in octets) |
| List of withdrawn routes (variable length) |
| Path attribute length (in octets) |
| List of path attributes (variable length) |
| "Network Layer Reachability Information" (= List of added prefixes) (variable length) |

# Attributes of BGP prefixes

**Update message details**

- **Mandatory** attributes: have to be there

  - Example: AS-Path

- **Optional** attribute: are, well, optional

  - Example: MED

- **Transitive** attributes

  - are kept on the prefix and forwarded via BGP

  - Even **(!)** when not understood by the forwarding device

- **Non-transitive** attributes

  - are added to a prefix and not forwarded by the receiver

# How is this realized in the protocol?

# More about BGP attributes

## Flags

- First byte of any attribute is "Flags"

  - Optional (1) or Well-Known (0)

  - Transitive (1) or Non-Transitive (0) (well-known is always transitive)

  - Partial (1) or Complete (0) ("Partial" only for optional transitive)

  - Extended Length Bit (0 = one length octet, 1 = two length octets)

  - Rest of the flags are unused

| Attribute Flags | | | | | Attribute Type Code | Length | |
|---|---|---|---|---|---|---|---|
| Opt | Tran sitiv e | Part ial | Ext Len | Unused | 1 Octet | 1 or 2 Octets Depending on ExtLen Flag | |

# More about BGP attributes

**Attribute Type Origin**

- Well-known, Transitive, Mandatory, (and quite simple)

- Length is one octet

- Possible values:

  - 0 - IGP

  - 1 - EGP

  - 2 - Incomplete

| Attribute Flags | | | | | Attribute Type Code | Length | Value |
|---|---|---|---|---|---|---|---|
| Opt **0** | Transi tive **1** | Parti al **0** | Ext Len **0** | Unused **0000** | 1 Octets 0x1 | 1 Octet 0x1 | 0x0 = IGP 0x1 = EGP 0x2 = Incomplete |

DE CIX

# More about BGP attributes

**Attribute Type AS Path**

- Well-known, Mandatory

- Realized as "sequence of segments"

- Segment: (type, length, value)

  - Type = 1: Unordered set of ASes traversed

  - Type = 2: Ordered Set of ASes traversed

- Length: Number of ASes in value part

| Segment | | |
|---|---|---|
| **Type** | **Length** | **Value** |
| 1= Unordered Set<br>2 = Ordered Set | 1 Octet | List of AS numbers |

# But what happens if an attribute is malformed?

# More about BGP attributes

**Attribute Type Origin**

- Well-known, Transitive, Mandatory, (and quite simple)

- Length is one octet

- Possible values:

  - 0 - IGP

  - 1 - EGP

  - 2 - Incomplete

> 0,1,2 are valid values
>
> What happens if there is a "3" in the value field?

| Attribute Flags | | | | | Attribute Type Code | Length | Value |
|---|---|---|---|---|---|---|---|
| Opt **0** | Transi tive **1** | Parti al **0** | Ext Len **0** | Unused **0000** | 1 Octets 0x1 | 1 Octet 0x1 | 0x0 = IGP 0x1 = EGP 0x2 = Incomplete |

# Invalid values in attributes

## RFC4271 - Section 6.3 "Update Message Error Handling"

All errors detected while processing the UPDATE message MUST be indicated by sending the NOTIFICATION message with the Error Code UPDATE Message Error.  The error subcode elaborates on the specific nature of the error.

If the ORIGIN attribute has an undefined value, then the Error Sub-code MUST be set to Invalid Origin Attribute.  The Data field MUST contain the unrecognized attribute (type, length, and value).

- What was a NOTIFICATION message again?

- **NOTIFICATION** - to tell the other side there was an error, and then close the BGP session.

# Error in Update Messages
## Shut down the BGP session

- So, if any of the well-known attributes contain an error

  - A notifcation is sent back

  - And the BGP session is closed

- This is a problem and a possible attack vector

- This was addressed in RFC7606:
  "Revised Error Handling for BGP UPDATE Messages"

DE CIX

# Revised Error Handling for BGP UPDATE Messages

**Treat the UPDATE like a WITHDRAW**

- RFC7606 addresses the problem of dropping BGP sessions because of errors

  - A number of RFCs who address error handling are updated

  - In most cases now session is no longer dropped, but the malformed UPDATE is now treated like a WITHDRAW

  - Read the RFC for details.

# *Conclusion*

# *Conclusion*

→ BGP Error handling has been improved over the years

→ In case of malformed attributes, BGP today handles an announcement like a withdrawal

→ Implementation bugs may cause major disruptions

→ The quality of a BGP implementation is also affected by how quickly critical bugs are fixed

→ See the original blog entry about vendor reaction times

**Where networks meet**

www.de-cix.net

https://de-cix.net/academy

# Links and further reading

# DE-CIX Academy Resources
## Lab and documentation

- DE-CIX Academy BGP Lab:
  https://gitlab.com/de-cix-public/team-academy/bgp/BGPLab

- Book: "BGP for networks who peer"
  https://github.com/wtremmel/BGP-for-networks-who-peer

- DE-CIX YouTube Channel:  https://www.youtube.com/@DE-CIX

DE CIX

# AS - Numbers
## How to request an AS number

- Giving AS numbers to the RIRs: iana.org

- Requesting an AS number, links for:

  - ARIN

  - Lacnic

  - APNIC

  - RIPE NCC

  - Afrinic

# BGP: Autonomous Systems
## RFCs

- <u>RFC1930</u>: Guidelines for creation, selection, and registration of an Autonomous System (AS)

- <u>RFC6793</u>: BGP Support for Four-Octet Autonomous System (AS) Number Space

# Routing
## Relevant RFCs

- **RFC4632**: Classless Inter-domain routing (CIDR)

# IPv6
## Relevant RFCs

- RFC4291: IPv6 addressing architecture

# BGP - Best Path Selection
## RFCs and Implementations

- RFC4271 - A Border Gateway Protocol 4 (BGP-4)

  - *Next Hop* is defined in Section 5.1.3

  - *AS Path* is defined in Section 5.1.2

  - *Local Preference*: Section 5.1.5

  - *Origin*: Section 5.1.1

  - *Multi Exit Discriminator (MED)*: Section 5.1.4

  - see 9.1 for the BGP best path selection algorithm

- BGP Best Path Selection by vendor

  - Cisco

  - Juniper

  - Mikrotik

  - Nokia

  - BIRD

  - FRRouting

| 1 | NextHop reachable? | Continue if "yes" |
|---|---|---|
| 2 | Local Preference | higher wins |
| 3 | AS Path | shorter wins |
| **4** | **Origin Type** | **IGP over EGP over Incomplete** |
| **5** | **MED** | **lower wins** |
| **6** | **eBGP, iBGP** | **eBGP wins** |
| **7** | **Exit** | **nearest wins** |
| **8** | **Age of route** | **older wins** |
| **9** | **Router ID** | **lower wins** |
| **10** | **Neighbor IP** | **lower wins** |

RFCs are Internet standards issued by the Internet Engineering Task Force (IETF)

# BGP Attributes
## Relevant RFCs

- BGP attribute types:

  - Registering new types: RFC2042

  - Published in BGP Parameters database at IANA

# BGP Security
## Relevant RFCs

- RFC7454 - BGP Operations and Security

- Password protect BGP sessions

  - RFC2385 (obsolete) - Protection of BGP Sessions via the TCP MD5 Signature Option

  - RFC5925 - The TCP Authentication Option

- RFC5082 - The Generalized TTL Security Mechanism (GTSM)

RFCs are Internet standards issued by the Internet Engineering Task Force (IETF)

# ~~Relevant~~ RFCs
## Historical (obsolete)

- [RFC827](): Exterior Gateway Architecture (EGP) (historical, obsolete)

-