# BGP in 120 minutes

**RIPE89**

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

**Where networks meet**

www.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

# 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

| 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 |
|---|---|---|---|---|---|---|---|---|---|---|----|----|----|---|----|----|----|----|---|----|----|----|----|---|----|----|----|----|---|----|----|----|----|---|----|----|----|----|
| 0 | 0 | 0 | 0 | | 1 | 0 | 1 | 0 | | 0 | 0 | 0 | 0 | | 0 | 0 | 1 | 1 | | 0 | 0 | 0 | 0 | | 1 | 0 | 0 | 0 | | 0 | 0 | 0 | 0 | | 0 | 0 | 0 | 0 |

**Prefix-Length: 0-32**

**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  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

# BGP announces prefixes
## To neighbors

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

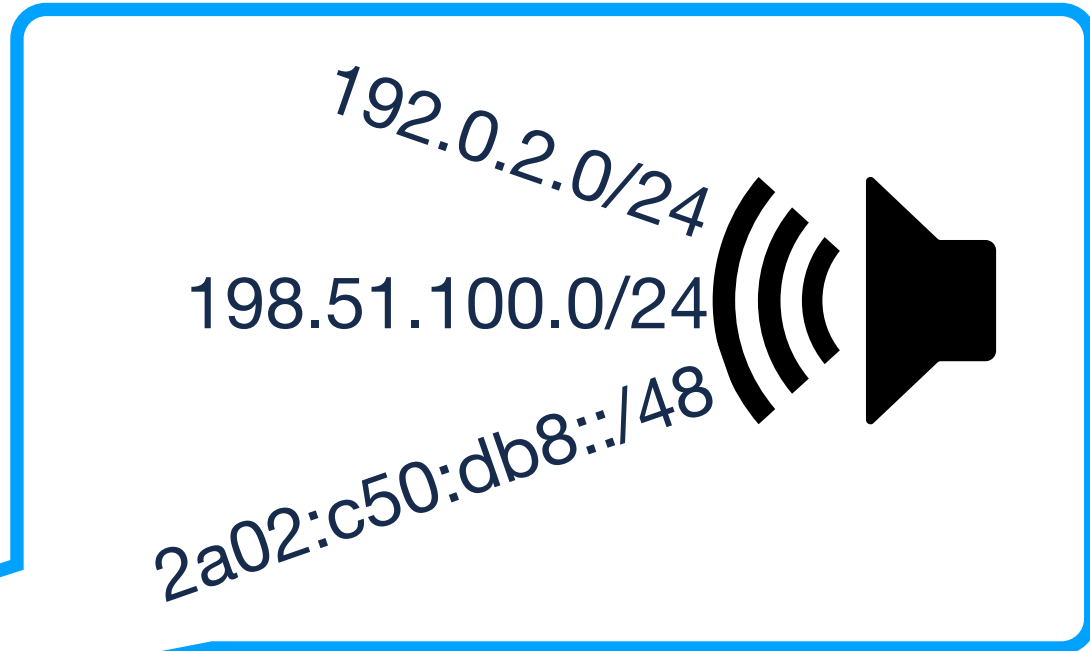I am **AS196610**, DE-CIX
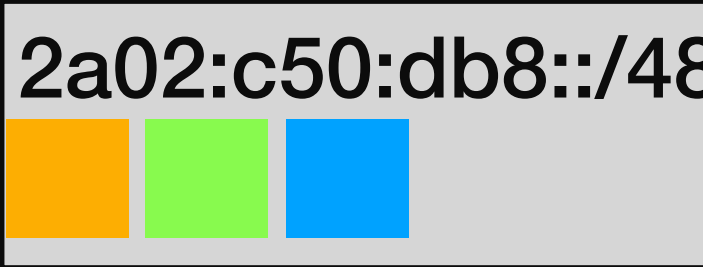Academy, and I announc
prefix
2a02:c50:db8::/48

- 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**

    2a02:c50:db8::/48

  - 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

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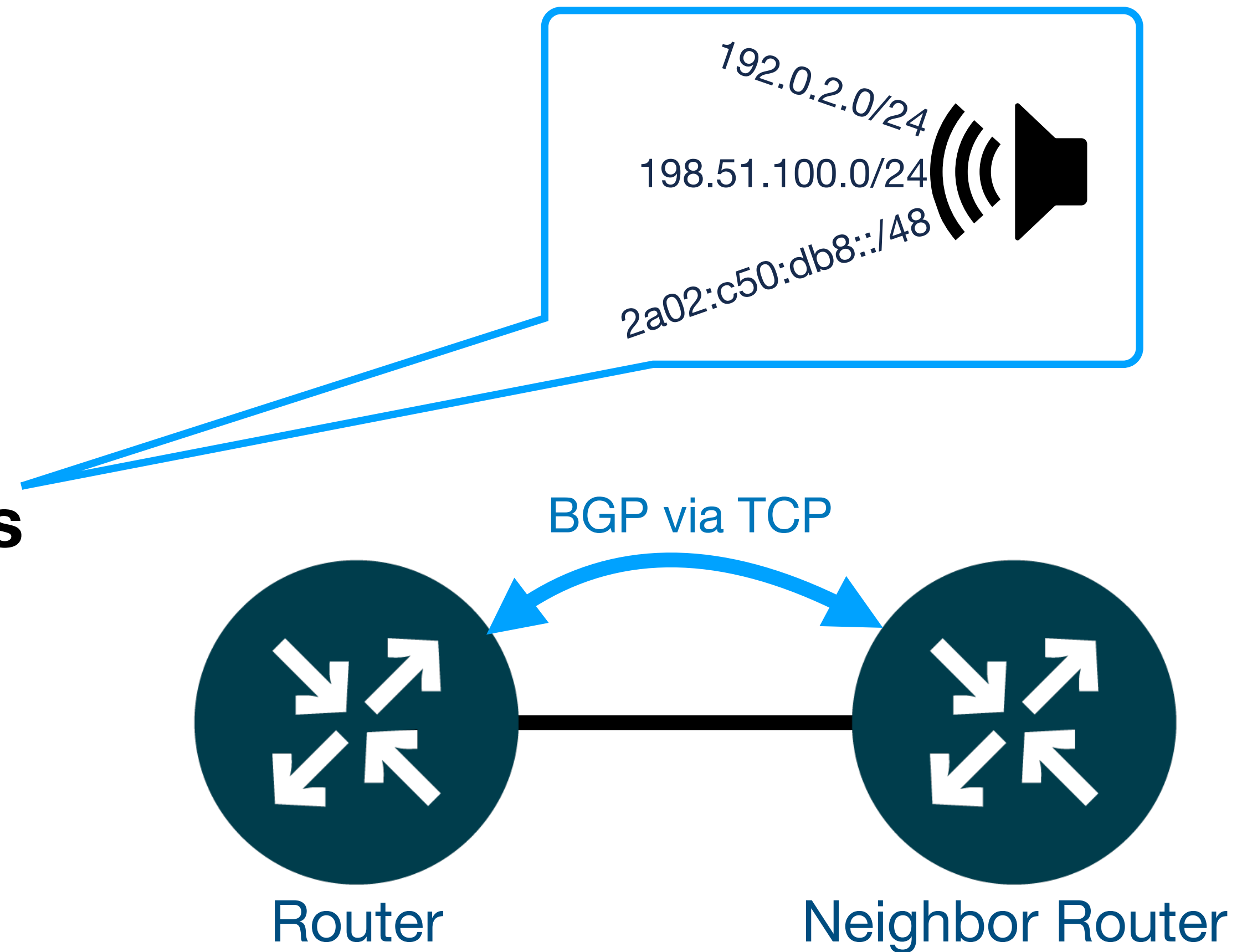
- Initially 16bit (0...65535) they are now 32bit long (0..."a lot")

- 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

# 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

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

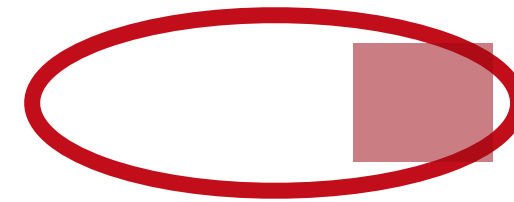- As long as the BGP session is "up", a router assumes its neighbors are "in sync" (= did not forget anything it sent)
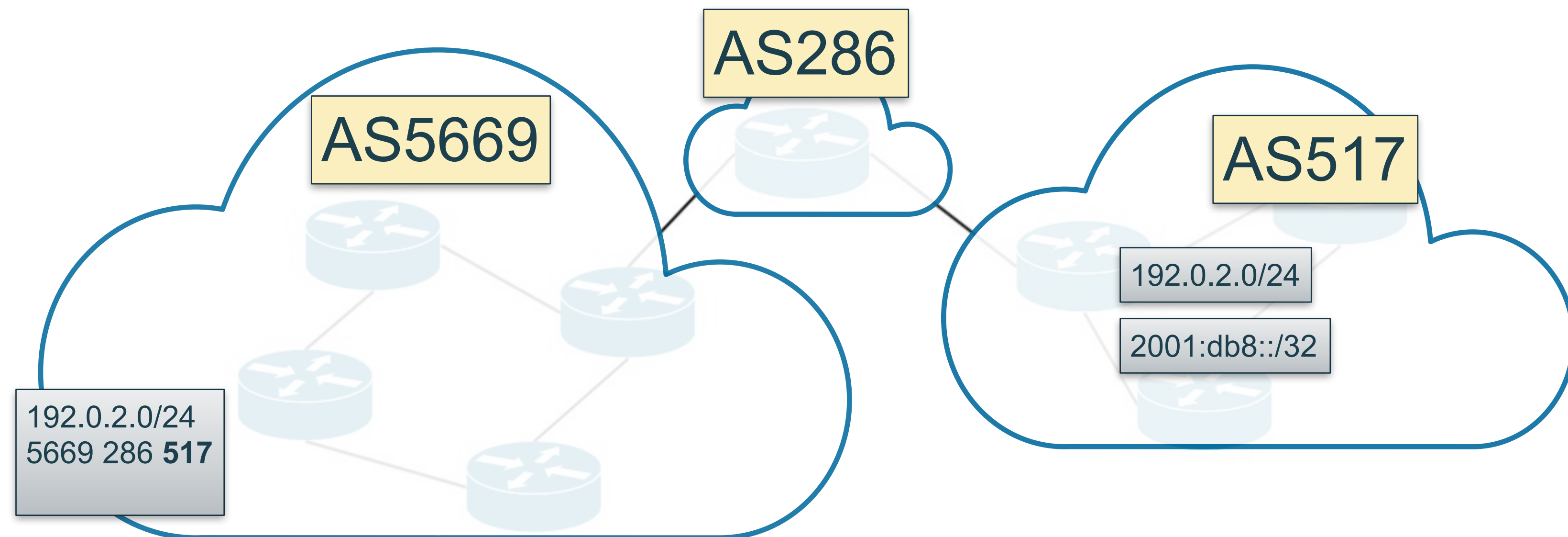
DE CIX

# 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

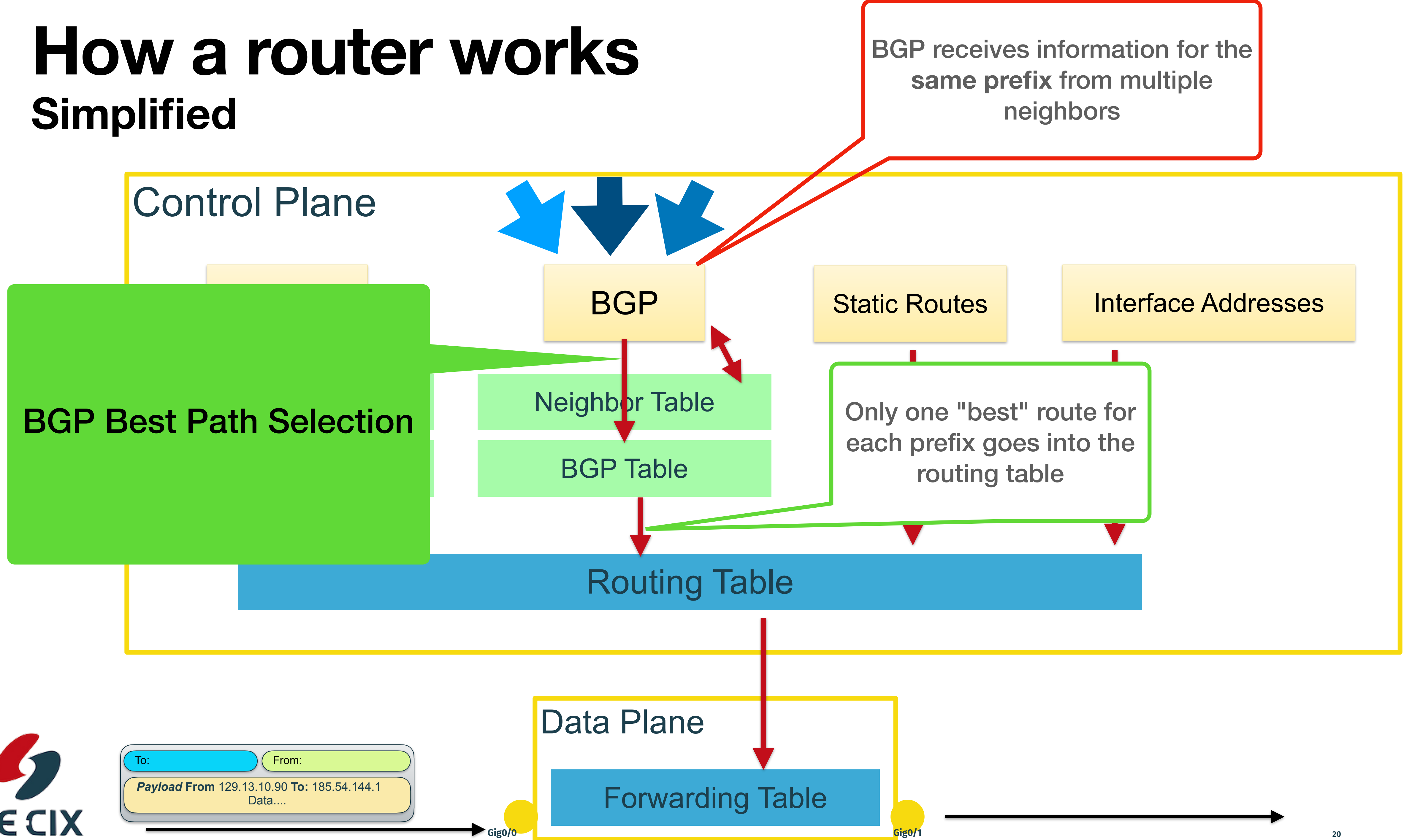2001:db8::/32

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
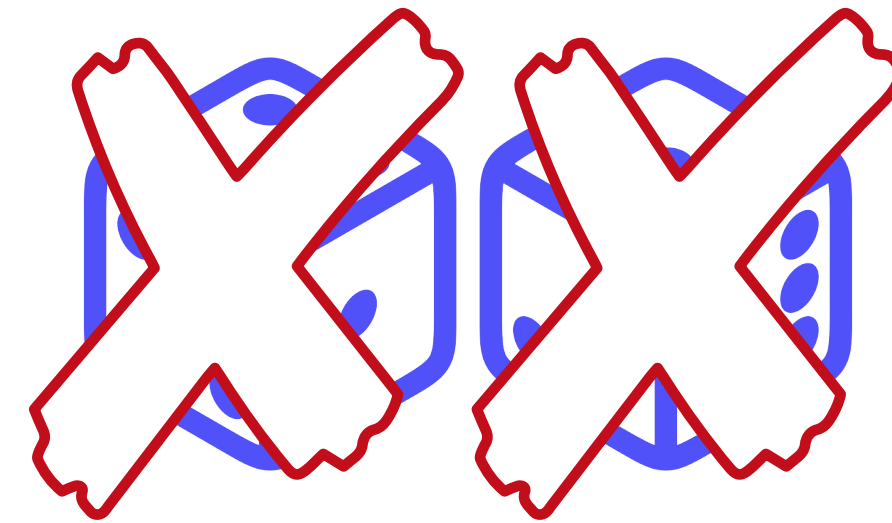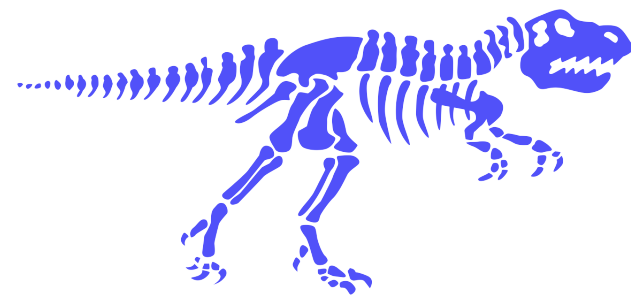
DE CIX

# How a router works
## Simplified



Control Plane

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

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

DE-CIX

20

# BGP Best Path Selection

# BGP Best Path Selection Algorithm
**Motivation**

- 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

10.3.8.0/22

286 517

AS64496

10.3.8.0/22

517

AS64500

eBGP

10.3.8.0/22

64497 517

AS64497

DE-CIX

# 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**

eBGP

AS64497

# Let's add peering

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**

eBGP

AS64496

AS286

AS517

AS64500

DE CIX

eBGP

AS64497

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 | | |

# 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 | | |

# Consider the following network



Receives Prefixes via eBGP

Prefixes

AS64496

eBGP

Traffic

iBGP

Provides Transit Service

iBGP

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?

# *BGP Route Selection Algorithm:*

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

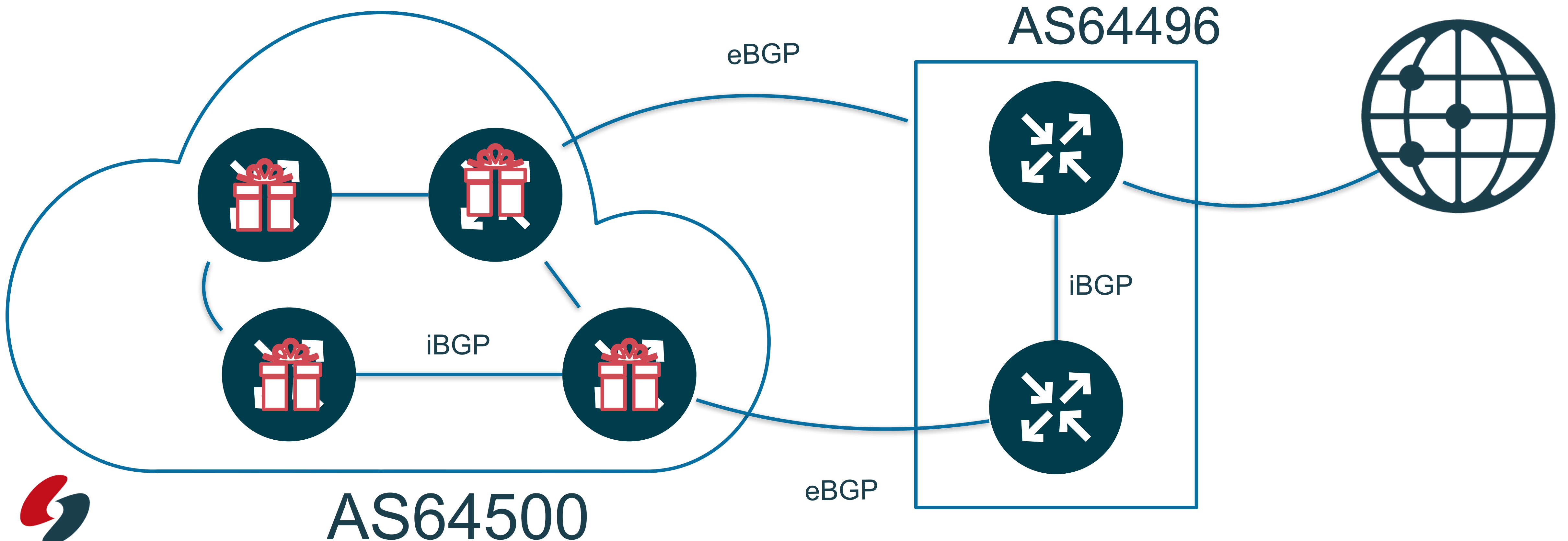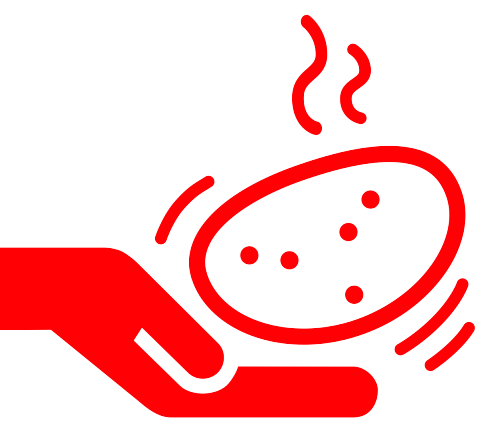| | | |
|---|---|---|
| 1 | NextHop reachable? | Continue if "yes" |
| 2 | Local Preference | higher wins |
| 3 | AS Path Length | shorter wins |
| 4 | Origin Type | IGP over EGP over Incomplete |
| 5 | | |
| 6 | | |
| 7 | | |
| 8 | | |
| 9 | | |
| 10 | | |

**DE-CIX**

**Where networks meet**

# BGP Route Selection Algorithm: MED

➜ 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

AS64496

eBGP

AS64500

iBGP

eBGP

# 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**

AS64500

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

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

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

eBGP

eBGP wins

eBGP

iBGP

iBGP

eBGP wins

# BGP Route Selection : nearest exit wins

eBGP

AS64496

iBGP

iBGP

eBGP

AS64500

DE-CIX

Where networks meet

www.de-cix.net

# 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 | | |

# 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 | | |

**Where networks meet**

# 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 |

# 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

# Other versions of this presentation

# BGP in 120 minutes
## What we did today

120

BGP!

- 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

# 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

# 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

# Experiment time!

# Lets play with a BGP router
**You just need your browser**

https://bgplab.as196610.net:7000/

DE CIX

# https://bgplab.as196610.net:7000/



Things to try:

- show bgp summary
- show bgp ipv4
- show bgp ipv6

DE-CIX

# DE-CIX Academy BGP lab

- The lab is open source and available for download

- Get it here:

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

DE CIX

# Network relationships

# The Internet
## A typical Internet Service Provider

# The Internet
## Adding "Upstream"

# The Internet
## Adding a 2nd ISP

# The Internet
## Data transport via upstreams

# The Internet
## More direct via "peering"

# The Internet
## Peering on multiple levels

DE-CIX

# The Internet
## Peering on multiple levels



DE CIX

# 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

# 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)

# 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



DE CIX

# 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

# What is the RIPE database?

## Documenting our resources

- A public resource database

- It documents:

  - AS numbers, their owners and their use

  - IP resources, their owners and their use

  - AS-sets, lists of ASes

- To access it, you can use the "whois" command

- Or go to the RIPE database [website](#)

```
aut-num:         AS196610
as-name:         DECIX-Academy
descr:           DE-CIX Academy Educational Network
org:             ORG-DtGI1-RIPE
adinet6num:          2a02:c50::/32
as-set:          AS-DECIX-HAM-RS-V6
descr:           ASN of DE-CIX Hamburg custo
descr:           DE-CIX Hamburg
admin-c:         DXSU6695-RIPE
tech-c:          DXSU6695-RIPE
mnt-by:          DECIX-MNT
remarks:         look at AS-DECIX-HAM for DE
remarks:         look at AS-DECIX-HAM-CONNEC
remarks:         Visit http://ham.de-cix.net
members:         AS42
members:         AS112
members:         AS250
members:         AS680
members:         AS1680
members:         AS1820
```

# More Information?
## RIPE Database Training

- The training is free

- The training is online

- Just go the ripe.net website



https://academy.ripe.net/enrol/index.php?id=9

# 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?

DE CIX

# 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

# 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

# BGP Communities

# BGP Communities

→A transitive, optional BGP attribute

→**Transitive**: Once attached, it stays until removed

→**Optional**: it does not have to be there

→"BGP Communities are like a sticker on a suitcase"

**DE-CIX**

*Where networks meet*

*www.de-cix.net*

# *"Original" BGP Communities*

➔ Definition:

"*A community is a group of destinations which share some common property*"

➔ Introduced in RFC1997 in year 1996

➔ A community is expressed by a 32Bit-Number

➔ High 16 bit are the AS defining the low 16 bits

   ➔ Notation: "6695:1000", "5669:32000"

➔ You can attach as many communities as you like (within reason)

   ➔ BGP max message size is 4096 Bytes

# *What are they useful for? Information!*

**198.51.100.0/24**          **80.81.192.15**          **from 80.81.192.15**

**Path: 1301 286 517**

**Origin IGP, metric 0, localpref 100, valid, external**

Frankfurt

DE-CIX

**Where networks meet**

# Informational Communities

**198.51.100.0/24**          **80.81.192.15**          **from 80.81.192.15**

**Path: 1301 286 517**

**Origin IGP, metric 0, localpref 100, valid, external**

**Received from:**          Upstream

# *Example: Encode geographical information*

# **65010:1**

Example: "1" here means geographical community

ISO-Country-Codes here …
**250** - France
**276** - Germany
**840** - USA

You may encode the continent here (if you are global) like:
• 1 = Europe
• 2 = North America
• 3 = Asia …

Just an Example!

DE-CIX

# *Example: Encode logical information*

## 65010:2____

Example: "2" here means logical source

Upstream? Peering? Customer?
1 = Upstream
2 = Private Peer
3 = Peer at an IXP
4 = Customer

More details here, like:
- Customer ID
- Upstream location
- up to you!

*Where networks meet*

*www.de-cix.net*

85 85

# *What are they useful for? Action!*

**198.51.100.0/24**          80.81.192.15          from 80.81.192.15

**Path: 65010**

**Origin IGP, metric 0, localpref 100, valid, external**

Announce to customers/peers

**DE-CIX**

Encoding up to you!

# Action Communities: Encoding

→ Again - you only have two 16bit numbers ... (with original BGP Communities)

→ Some Ideas ...

- If you want your customers to send you "actions"

  - You really should have them put your AS number into the first 16bit number

  - You **must scrub** everything they should not send on incoming

- Possible actions:

  - (not) announce to upstream, peers, customers

  - fine granular announcement control (geographically, by IXP, ...)

  - announce with longer AS path

  - change *local preference*

  - Blackhole

# Action Communities: Well-Known

→ A couple of communities are pre-defined by RFCs

→ NO-EXPORT

- Do not send the prefix to eBGP neighbours (other ASes)

→ NO-ADVERTISE

- Do not send the prefix to anyone (not even internal via iBGP)

→ NO-PEER

- Do not send to any peers

→ BLACKHOLE

- Sink all traffic to prefixes tagged with this community

- Most commonly used with host routes

- Implies NO-EXPORT

# 32Bit AS? No luck with original communities

# 65010:12345

➜ Two 16-bit numbers

➜ No way to encode a 32Bit AS number and something else ...

- RFC4360 -  Extended Communities

➜ Extended Communities - Lots of new features

- In total 2*32Bits

- Introducing a "type" field

- Possible to encode 16Bit Type, 32Bit AS, 16Bit Data

**DE·CIX**

*Where networks meet*

# Extended Communities

| 00 | 01 | 02 | 03 | 04 | 05 | 06 | 07 | 08 | 09 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |

I  T  Type high  Type low

Value

Value

→ **I** = Type is IANA assigned (= well known) or private

→ **T = 0**: Transitive across AS borders

→ **T = 1**: Non-Transitive - should be removed before forwarding to another AS

→ **Type**:  Types are either IANA-assigned or experimental. For a list of assigned types see the RFC

→ **Value**: 48 Bits, meaning is dependent on type

→ Standardized in 2006

# Extended Communities and 32Bit ASes

| 00 | 01 | 02 | 03 | 04 | 05 | 06 | 07 | 08 | 09 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| I | T | 0x02 | | | | | | 0x02 or 0x03 | | | | | | | | 32Bit-AS Number | | | | | | | | | | | | | | | |
| 32Bit-AS Number (continued) | | | | | | | | | | | | | | | | Value | | | | | | | | | | | | | | | |

→ **You can encode a 32Bit AS-Number**

  → **and a 16 Bit value**

# Extended Communities and 32Bit ASes

| 00 | 01 | 02 | 03 | 04 | 05 | 06 | 07 | 08 | 09 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| I | T | 0x02 | | | | | | 0x02 or 0x03 | | | | | | | | 16Bit-AS Number | | | | | | | | | | | | | | | |
| 32Bit Value | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |

→ You can encode a 32Bit AS-Number

  → and a 16 Bit value

→ or a 16Bit AS-Number

  → and a 32 Bit value


→ 32Bit AS and 32Bit Value?

  → **not possible!**

# Extended communities use cases

→ Notation:

- Similar to original communities: **RT:6500000:1234** or **RT:1234:6500000**

→ Disadvantages

- Only 48bits in total

- Only one 32Bit value is possible (and one 16Bit value)

- RT, RO and other types confusing to many operators

→ Conclusion

- Another community version was needed

- It took the IETF a while to realize that (11 years)

# Introducing: Large Communities

| 00 | 01 | 02 | 03 | 04 | 05 | 06 | 07 | 08 | 09 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| Global Administrator (32Bit AS) |||||||||||||||||||||||||||||||
| Local Data Part 1 (Function) |||||||||||||||||||||||||||||||
| Local Data Part 2 (Parameter) |||||||||||||||||||||||||||||||

➡ Very simple - three 32Bit values (finally something useful)

➡ Global Administrator:

- An AS number (in 32Bit notation)

- Has defined meaning of two other fields

- May have published that meaning

➡ Local Data

- Can be seen as "just two 32Bit numbers"

- Or as "Function" / "Parameter"

# Large BGP Communities

| 00 | 01 | 02 | 03 | 04 | 05 | 06 | 07 | 08 | 09 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Global Administrator (32Bit AS) | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Local Data Part 1 (Function) | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Local Data Part 2 (Parameter) | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |

→ Notation:

    → Similar to Original Communities: **196610:100:65000010**

→ Defined in two RFCs:

    → RFC8092: BGP Large Communities Attribute

    → RFC8195: Use of BGP Large Communities

→ A dedicated website exists: http://largebgpcommunities.net

    → Keeping track of Implementations, News etc.

# BGP Communities and the DE-CIX Route Servers
## Default Behaviour

route server
## AS6695

192.168.1.0/24
AS-Path: 64500

.113.99/24

203.0.113.1/24

203.0.113.2/24

203.0.113.3/24

203.0.113.4/24

192.168.1.0/24
AS-Path: 64500

**Router of ISP1**
**AS64500**

**Router of ISP2**
**AS65501**

**Router of ISP3**
**AS65502**

**Router of ISP4**
**AS65503**

# *Do not announce to any AS*

**route server**
**AS6695**

```
0:6695
6695:0:0
```

203.0.113.99/24

203.0.113.1/24

203.0.113.2/24

203.0.113.3/24

203.0.113.4/24

192.168.1.0/24
AS-Path: 64500
**0:6695**

**Router of ISP1**
**AS64500**

**Router of ISP2**
**AS65501**

**Router of ISP3**
**AS65502**

**Router of ISP4**
**AS65503**

# Do not announce to any AS, but announce to AS65501

**route server**
**AS6695**

192.168.1.0/24
AS-Path: 64500

.113.99/24

```
0:6695  6695:65501
6695:0:0  6695:1:65501
```

203.0.113.1/24

192.168.1.0/24
AS-Path: 64500
**0:6695**

203.0.113.2/24          203.0.113.3/24          203.0.113.4/24

**Router of ISP1**          **Router of ISP2**          **Router of ISP3**          **Router of ISP4**
**AS64500**          **AS65501**          **AS65502**          **AS65503**

# Do not announce to AS65501

**route server**
**AS6695**

192.168.1.0/24
AS-Path: 64500

.113.99/24

```
0:65501
6695:0:65501
```

203.0.113.1/24        203.0.113.2/24    203.0.113.3/24    203.0.113.4/24

192.168.1.0/24
AS-Path: 64500

**0:65501**

**Router of ISP1**    **Router of ISP2**   **Router of ISP3**   **Router of ISP4**

**AS64500**           **AS65501**          **AS65502**          **AS65503**

# Prepend 1 time to AS65503

**route server**

**AS6695**

192.168.1.0/24

AS-Path: 64500

13.99/24

```
65001:65503
6695:101:65503
```

203.0.113.1/24

203.0.113.2/24

203.0.113.3/24

203.0.113.4/24

192.168.1.0/24

AS-Path: 64500

**65001:65503**

**Router of ISP1**

**AS64500**

**Router of ISP2**

**AS65501**

**Router of ISP3**

**AS65502**

**Router of ISP4**

**AS65503**

# Add NO-EXPORT to AS65502

**route server**

**AS6695**

6695:901:65502

192.168.1.0/24
AS-Path: 64500

13.99/24

203.0.113.1/24

192.168.1.0/24
AS-Path: 64500
**6695:901:65502**

203.0.113.2/24

203.0.113.3/24

203.0.113.4/24

**Router of ISP1**
**AS64500**

**Router of ISP2**
**AS65501**

**Router of ISP3**
**AS65502**

**Router of ISP4**
**AS65503**

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

DE-CIX

# BGP: Autonomous Systems
**RFCs**

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

- RFC6793: 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)

-