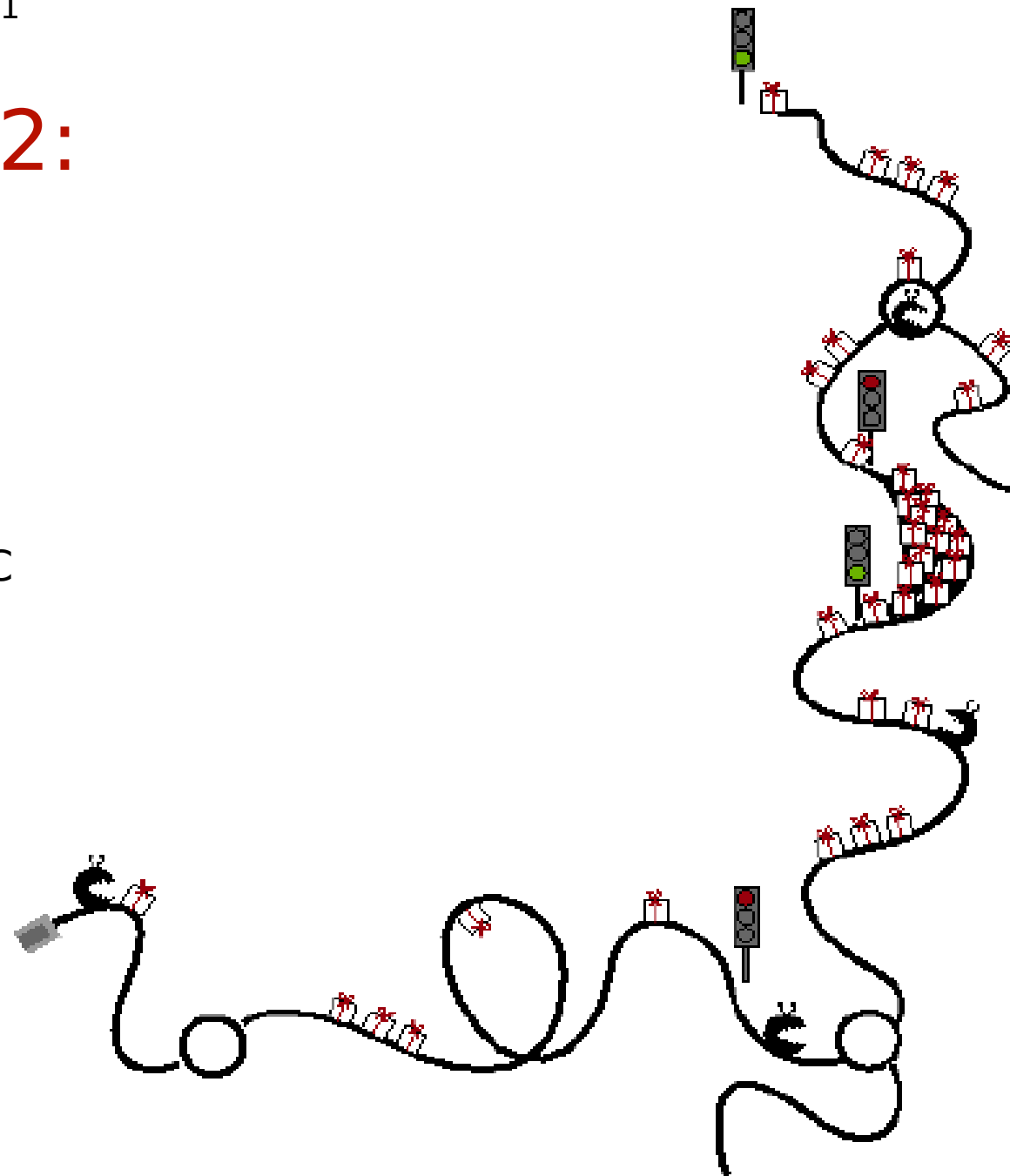# Lecture 11-12: Routing

Olof Hagsand KTH CSC

# Reading instructions

**Forouzan: TCP/IP Protocol Suite:**

Chapter 11: Unicast routing protocols

You need to complement with slides, especially if you do not make the routing lab

11.6 OSPF: Skip detailed packet descriptions
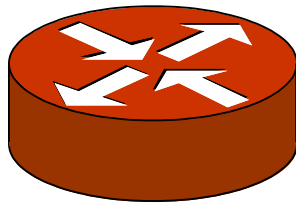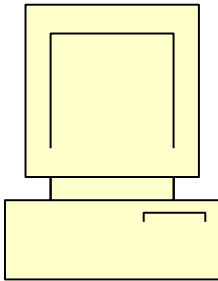
11.8 BGP: Skip detailed packet descriptions

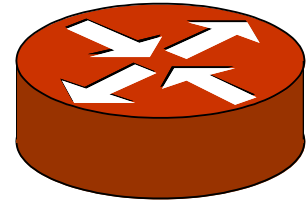**EP2120: Lab4 : Introduction to routing**

# Routers

# What is a router?

- Host (end-system)
  - One or many network interfaces
  - Cannot forward packets between them
- Router
  - Can forward packets between multiple interfaces
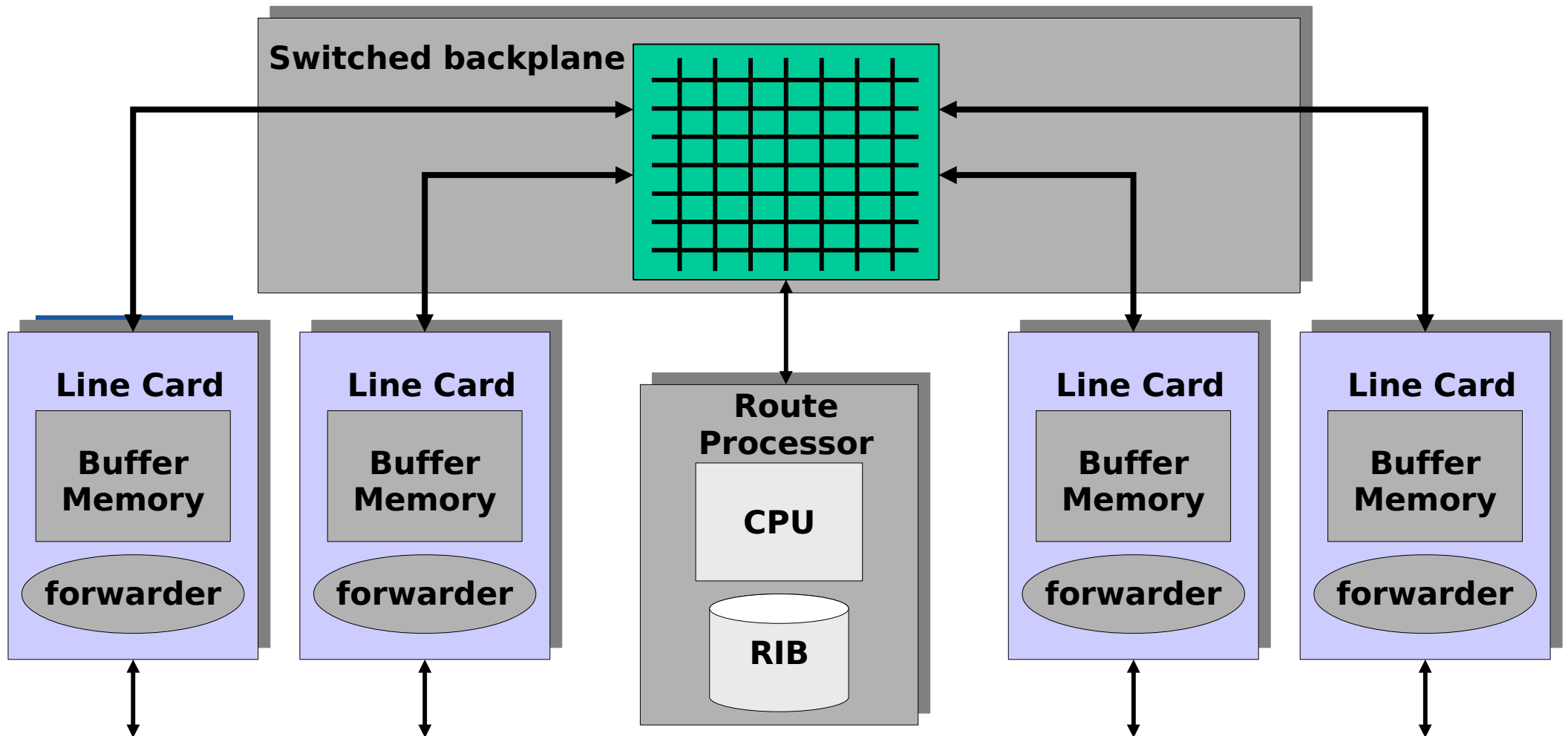  - Forwarding on Layer 3

# What does a router do?

- Packet forwarding
- Not only IPv4:
- IPv6, MPLS, Tunneling,...
    - (But never naming,..)
- Filter traffic
    - Access lists based on src/dst, etc.
- Metering/Shaping/Policing
    - Measuring, forming and dropping traffic
- Compute routes: build forwarding table
- In the "background": routing
- In "real-time": forwarding

# Inside a router, example

**Switched backplane**

**Line Card**

**Buffer Memory**

*forwarder*

**Line Card**

**Buffer Memory**

*forwarder*

**Route Processor**

**CPU**

**RIB**

**Line Card**

**Buffer Memory**

*forwarder*

**Line Card**

**Buffer Memory**

*forwarder*

- A router consists of linecards with ports interconnected by a backplane
- A route processor (RP) runs all routing protocols and management
- Modern routers do forwarding on the linecards (*data-plane*), often in hardware – the RP is only involved with *control-plane* processing.

# Routing algorithms
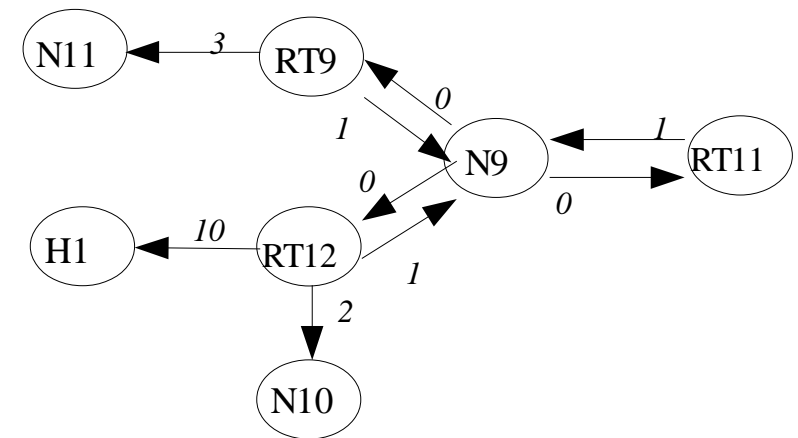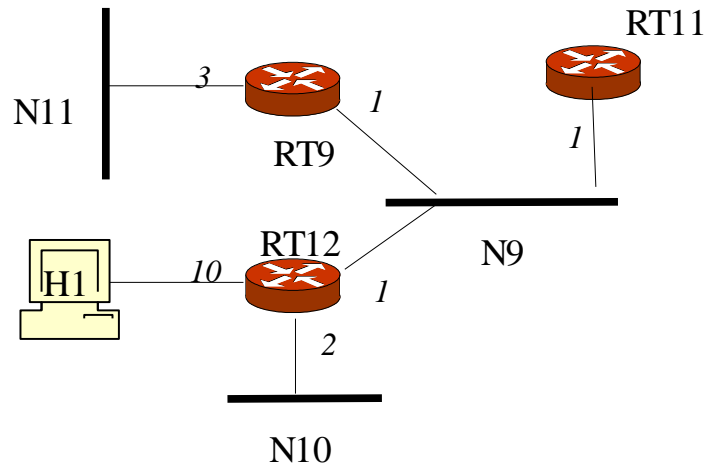
# Routing algorithms

- How does a router find a best path?
- Most solutions based on SPF (Shortest Path First) algorithms that are well known in graph theory.
  - Bellman-Ford
  - Dijkstra
- Link-State protocols (OSPF, IS-IS) use Dijkstra
- Distance-Vector protocols (RIP, IGRP, BGP) use Bellman-Ford
- Apart from that, there may also be other algorithms in
  - Multicast routing
  - Ad-hoc routing
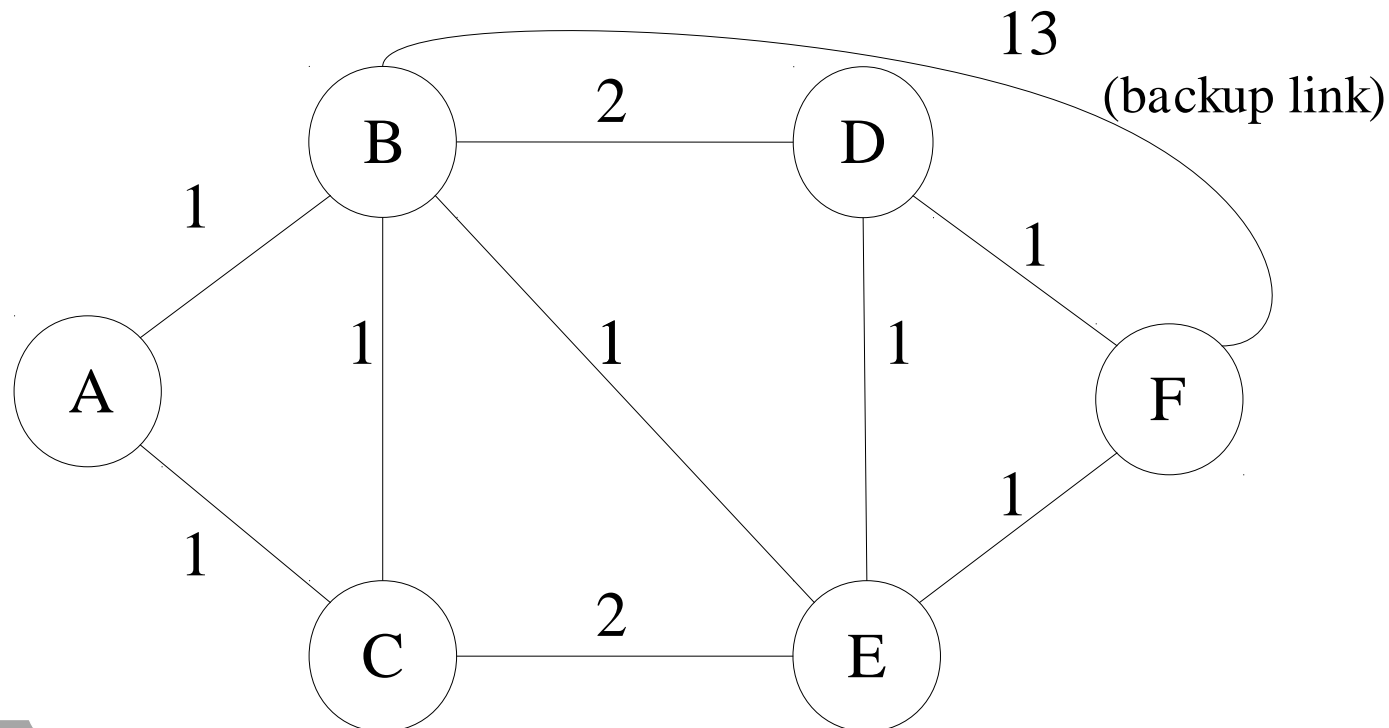  - Sensor networks
  - Delay-tolerant networks

# Graphs vs networks

- Algorithms are usually defined on graphs whereas protocols work on networks
- Graphs have nodes and edges whereas networks have interfaces, broadcast links, addresses, hierarchical layering, etc.
- Note the modelling of the broadcast link N9

# Shortest Path First (SPF)

- Given link metrics (weights) on each individual link
- Find the path (sequence of links) where the sum of the metrics of all links (cumulative cost) is lowest
- Equal cost multipath (ECMP): A *set* of paths with the least (same) cost
- What is the SPF from A to F?

# Alternative: Widest path first

- Numbers denote width: load or bandwidth
  - *Available* bandwidth
- It is easy to extend SPF algoritms with a widest-path computation rather than shortest path.
- What is the widest path from A -> E?

# Distance-Vector/Bellman-Ford

- Each router sends a list of distance-vectors (route with costs) to each neighbour periodically
- Every router selects the route with smallest metric (positive integer)
- The underlying algorithm is called Bellman-Ford.
- Protocols that use Bellman-Ford are called Distance-vector protocols

# Example: Distance-vector

**A:s initial state: (directly connected networks)**

| Dest | Cost | NextHop |
|------|------|---------|
| B | 1 | - |
| D | 3 | - |

**A distributes this DV to its neighbours (B and D)**

**A receives B:s (initial) distance vector**

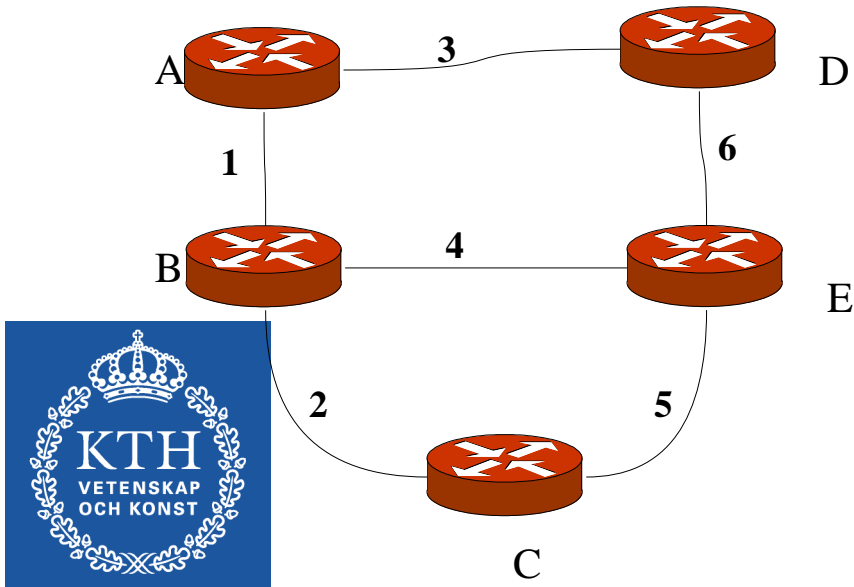| Dest | Cost |
|------|------|
| A | 1 |
| C | 2 |
| E | 4 |

**A:s state after merging B:s DV:**

| Dest | Cost | NextHop |
|------|------|---------|
| B | 1 | - |
| C | 3 | B |
| D | 3 | - |
| E | 5 | B |

**A distributes this DV to its neighbours (B and D)**

# Example: Complete and final state



**Link metric matrix**

|   | A | B | C | D | E |
|---|---|---|---|---|---|
| A |   | 1 |   | 3 |   |
| B | 1 |   | 2 |   | 4 |
| C |   | 2 |   |   | 5 |
| D | 3 |   |   |   | 6 |
| E |   | 4 | 5 | 6 |   |

A's Distance-Vector

|   | A | B | C | D | E |
|---|---|---|---|---|---|
| A | 0 |   |   |   |   |
| B |   | 0 |   |   |   |
| C |   |   | 0 |   |   |
| D |   |   |   | 0 |   |
| E |   |   |   |   | 0 |

**Initial state**

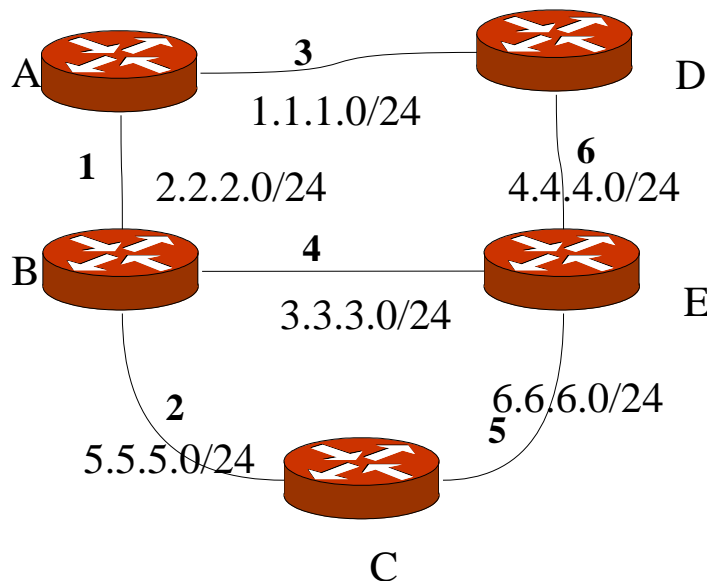|   | A | B | C | D | E |
|---|---|---|---|---|---|
| A | 0 | 1 | 3 | 3 | 5 |
| B | 1 | 0 | 2 | 4 | 4 |
| C | 3 | 2 | 0 | 6 | 5 |
| D | 3 | 4 | 6 | 0 | 6 |
| E | 5 | 4 | 5 | 6 | 0 |

**Final state**

# The algorithm

- Keep a table with an entry for each destination D in the network.
- Store the metric M (distance) and next-hop N for each D in the table.
- Periodically, send the table to all neighbors (the distance-vector).
- For each update that comes in from neighbor N' (to D with a new metric):
  - Add the cost of the link to N' to the new metric to get M'.
  - Replace the route if M' < M.
  - If N = N', always replace the route.
- In most protocols, M is bounded, typically to 16. This upper bound is defined as unreachable(infinity).

# Going to real networks

- IP networks require destinations and nexthops (not just nodes)
    - Destinations are networks eg 192.16.32.0/24
    - Next-hops are IP addresses, eg 192.16.32.1
- Suppose the topology changes , eg routers, links crash?
    - Use timers (counters) and age the entries
    - Send updates every (e.g.) 30s
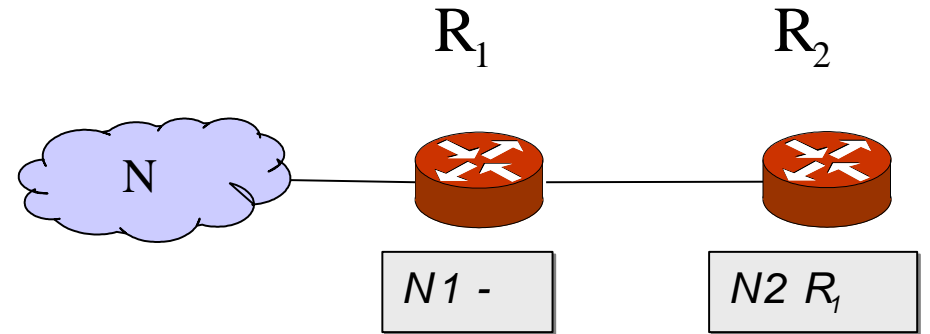    - If you do not hear from a router in (e.g.) 180s, mark it as invalid

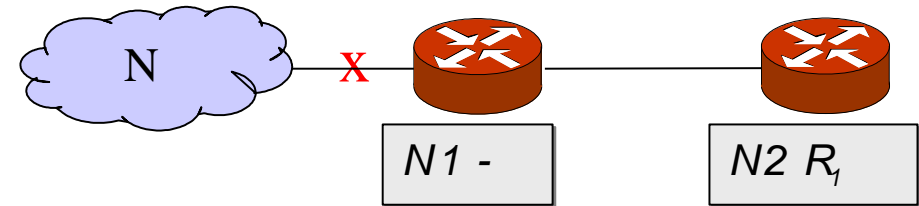| Dest | Cost | NextHop |
|---|---|---|
| 1.1.1.0/24 | 3 | - |
| 2.2.2.0/24 | 1 | - |
| 3.3.3.0/24 | 5 | 2.2.2.2 |
| 4.4.4.0/24 | 9 | 1.1.1.2 |
| 5.5.5.0/24 | 3 | 2.2.2.2 |
| 6.6.6.0/24 | 8 | 2.2.2.2 |

**Converged routing state of A**

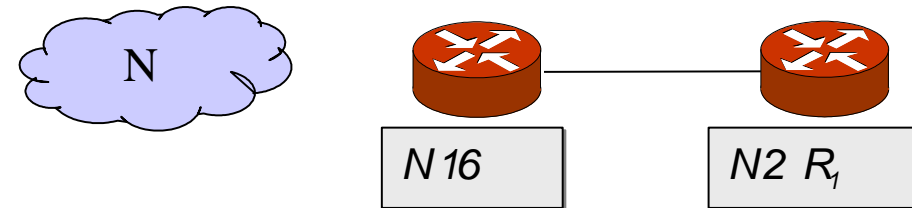# D.V. Problem: Count to Infinity (Two-node instability)

$R_1$           $R_2$

Initially, $R_1$ and $R_2$ both have a route to N with metric 1 and 2, respectively.

N

| N 1 - | N 2 $R_1$ |

The link between $R_1$ and N fails.

N   X

| N 1 - | N 2 $R_1$ |

Now $R_1$ removes its route to N, by setting its metric to 16 (infinity).

N

| N 16 | N 2 $R_1$ |

Now two things can happen: Either $R_1$ reports its route to $R_2$. Everything is fine.

N      N 16

| N 16 | N 16 |

# D.V. Problem: Count to Infinity

$R_1$  $R_2$

The other alternative is that $R_2$, which still has a route to N, advertises it to $R_1$. Now things start to go wrong: packets to N are looped until their TTL expires!

N

N2

Loop!

N3 $R_2$  N2 $R_1$

Eventually (~10-20s), $R_1$ sends an update to $R_2$. The cost to N increases, but the loop remains.

N

N3

Loop!

N3 $R_2$  N4 $R_1$

Yet some time later, $R_2$ sends an update to $R_1$.

...

N

N4

Loop!

N5 $R_2$  N4 $R_1$

Finally, the cost reaches infinity at 16, and N is unreachable. The loop is broken!
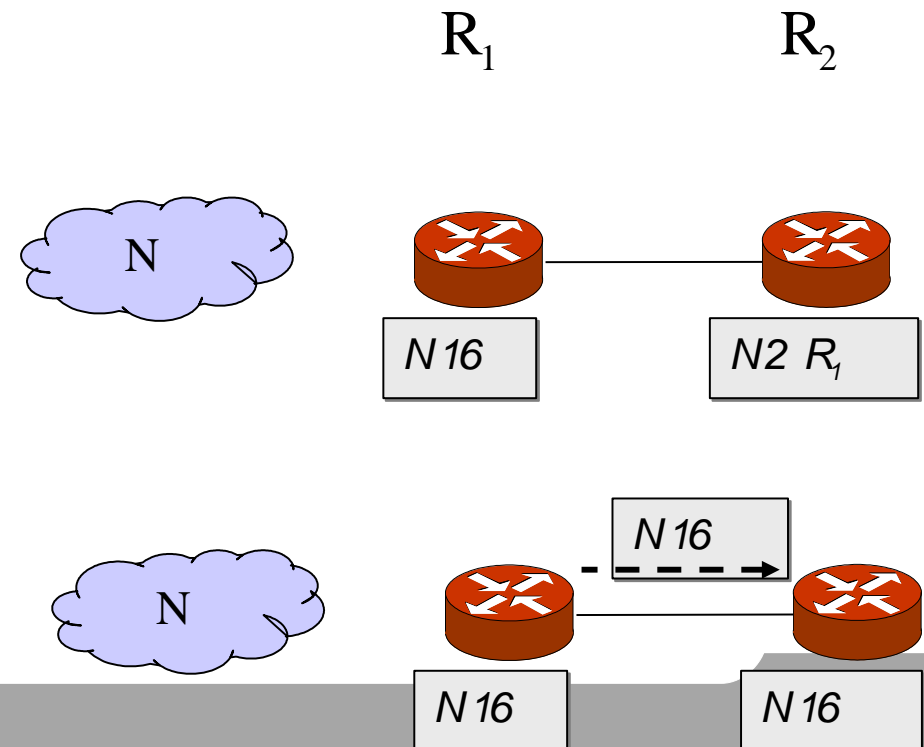
N

N 16  N 16

# Solution: Split Horizon

- *Do not send routes back over the same interface from where the route 'arrived'.*

- This helps in avoiding "mutual deception": two routers tell each other they can reach a destination via each other.

$R_1$          $R_2$

$R_2$, does not announce the route to N to $R_1$ since that is where it was learnt.

N

N 16          N2  $R_1$

Eventually, $R_1$ reports its route to $R_2$ and everything is fine.

N 16

N

N 16          N 16

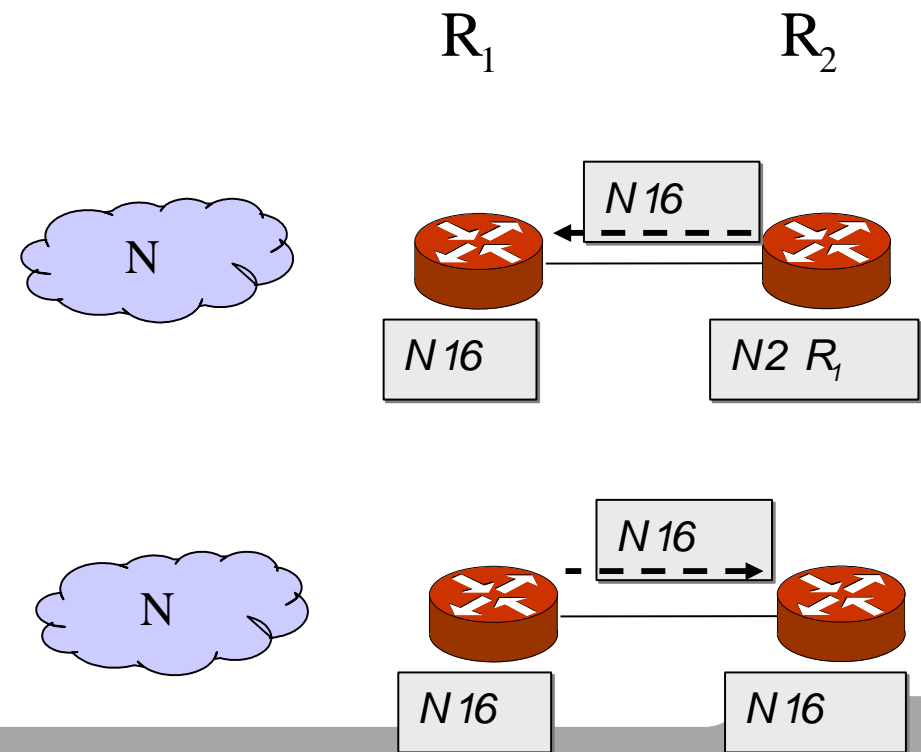# Solution: Split Horizon + Poison Reverse

- *Advertise reverse routes with a metric of 16 (i.e., unreachable).*

- Does not add information but breaks loops faster

- Adds protocol overhead

$R_1$     $R_2$

$R_2$ always announces an unreachable route to N to $R_1$.

N

N 16

N 16          N2  $R_1$

Eventually, $R_1$ reports its route to $R_2$ and everything is fine.

N

N 16

N 16          N 16

# Remaining problems

- More than two routers involved in mutual deception
    - A may believe it has a route through B, B through C, and C through A
- In this case, split horizon with poison reverse does not help



A

B

C

# Solution: Triggered Update

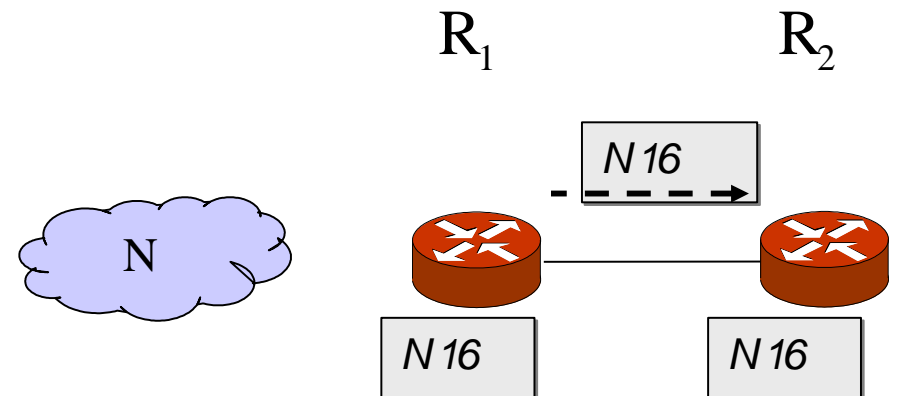- *Send out update immediately when metrics change*
- But only the changed route, not the complete table
- This may lead to a cascade of updates
  - Apply the rule above recursively!
  - Therefore, triggered updates are not allowed more often than, for example 1-5 seconds.
- A router may use triggered update only when deleting routes (16).

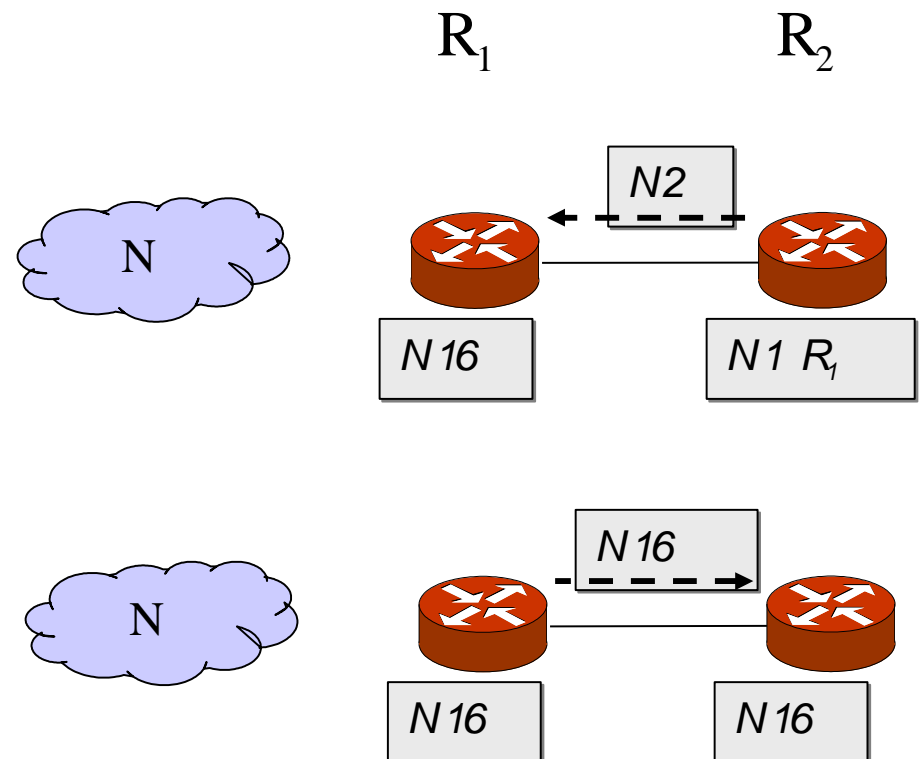$R_1$ immediately announces the broken link when it happens.

# Solution: Hold Down

- *When a route is removed, no update of this route is accepted* for some period of time (hold-down time)- to give everyone a chance to remove the route.

$R_1$ ignores updates to N from $R_2$ for some period of time.

Eventually, $R_1$ sends the update to $R_2$.

# Dijkstra's shortest path first

From the link-state database, compute a shortest path delivery tree using a permanent set S and a tentative set Q:

1. Define the root of the tree: the router
2. Assign a cost of 0 to this node and make it the first permanent node.
3. Examine each neighbor node of the last permanent node.
4. Assign a cumulative cost to each node and make it tentative.
5. Among the list of tentative nodes:
   - ·Find the node with the smallest cumulative cost and make it permanent.
   - ·If a node can be reached from more than one direction, select the direction with the smallest cumulative cost.
6. Repeat steps 3 to 5 until every node is permanent.

# Example network

# Example graph

# Exercise: Dijkstra from A

| Permanent set | Tentative set |
|---|---|
| A 0 – | 10.0.3.0/24 1 –<br>10.0.1.0/24 1 – |

# Exercise: Dijkstra

| Permanent set | Tentative set |
|---|---|
| A 0 – | ~~10.0.3.0/24 1 –~~ |
| 10.0.3.0/24 1 – | 10.0.1.0/24 1 – |
| | B 1 – |

# Exercise: Dijkstra

| Permanent set | Tentative set |
|---|---|
| A 0 – | |
| 10.0.3.0/24 1 – | 10.0.1.0/24 1 – |
| B 1 – | ~~B 1 –~~ |
| | 10.0.2.0/24 2 B |
| | 10.0.6.0/24 2 B |

# Exercise: Dijkstra

| Permanent set | Tentative set |
|---|---|
| A 0 – | |
| 10.0.3.0/24 1 – | ~~10.0.1.0/24 1 –~~ |
| B 1 – | |
| 10.0.1.0/24 1 – | 10.0.2.0/24 2 B |
| | 10.0.6.0/24 2 B |
| | C 1 – |

# Exercise: Dijkstra

| Permanent set | Tentative set |
|---|---|
| A 0 – | |
| 10.0.3.0/24 1 – | |
| B 1 – | |
| 10.0.1.0/24 1 – | 10.0.2.0/24 2 B |
| C 1 – | 10.0.6.0/24 2 B |
| | ~~C 1 –~~ |
| | 10.0.2.0/24 2 C |
| | 10.0.4.0/24 2 C |

# Exercise: Dijkstra

|  Permanent set | Tentative set |
|---|---|
| A 0 – | |
| 10.0.3.0/24 1 – | |
| B 1 – | |
| 10.0.1.0/24 1 – | ~~10.0.2.0/24 2 B~~ |
| C 1 – | 10.0.6.0/24 2 B |
| 10.0.2.0/24 2 B | |
| 10.0.2.0/24 2 C | ~~10.0.2.0/24 2 C~~ |
| | 10.0.4.0/24 2 C |

Note: ECMP

ECMP: Equal Cost MultiPath. More than

# Exercise: Dijkstra

| Permanent set | Tentative set |
|---|---|
| A 0 – | |
| 10.0.3.0/24 1 – | |
| B 1 – | |
| 10.0.1.0/24 1 – | |
| C 1 – | 10.0.6.0/24 2 B |
| 10.0.2.0/24 2 B | |
| 10.0.2.0/24 2 C | |
| 10.0.4.0/24 2 C | ~~10.0.4.0/24 2 C~~ |
| | D 2 C |
| | E 2 C |

# Exercise: Dijkstra

| Permanent set | Tentative set |
|---|---|
| A 0 –<br>10.0.3.0/24 1 –<br>B 1 –<br>10.0.1.0/24 1 –<br>C 1 –<br>10.0.2.0/24 2 B<br>10.0.2.0/24 2 C<br>10.0.4.0/24 2 C<br>E 2 C<br>D 2 C | 10.0.6.0/24 2 B<br><br>~~D 2 C~~<br>~~E 2 C~~<br>10.0.5.0/24 3 C |

# Exercise: Dijkstra

| Permanent set | Tentative set |
|---|---|
| A 0 – | |
| 10.0.3.0/24 1 – | |
| B 1 – | |
| 10.0.1.0/24 1 – | |
| C 1 – | ~~10.0.6.0/24 2 B~~ |
| 10.0.2.0/24 2 B | |
| 10.0.2.0/24 2 C | |
| 10.0.4.0/24 2 C | |
| E 2 C | |
| D 2 C | |
| 10.0.6.0/24 2 B | 10.0.5.0/24 3 C |
| | F 2 B |

# Exercise: Dijkstra

| Permanent set | Tentative set |
|---|---|
| A 0 – | |
| 10.0.3.0/24 1 – | |
| B 1 – | |
| 10.0.1.0/24 1 – | |
| C 1 – | |
| 10.0.2.0/24 2 B | |
| 10.0.2.0/24 2 C | |
| 10.0.4.0/24 2 C | |
| E 2 C | |
| D 2 C | |
| 10.0.6.0/24 2 B | 10.0.5.0/24 3 C |
| <span style="color:red">F 2 B</span> | ~~F 2 B~~ |
| | 10.0.5.0/24 3 B |

# Exercise: Dijkstra (complete)

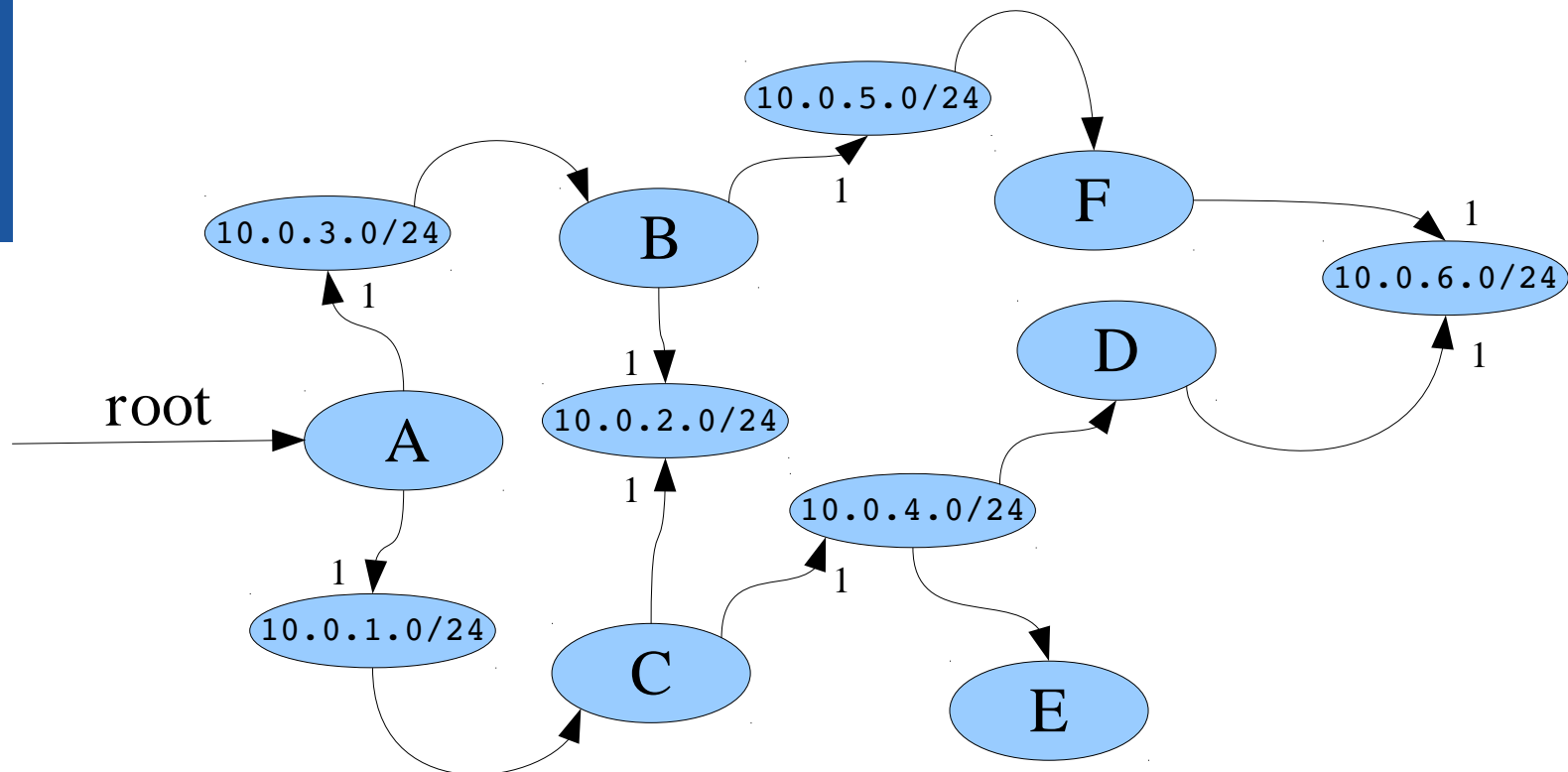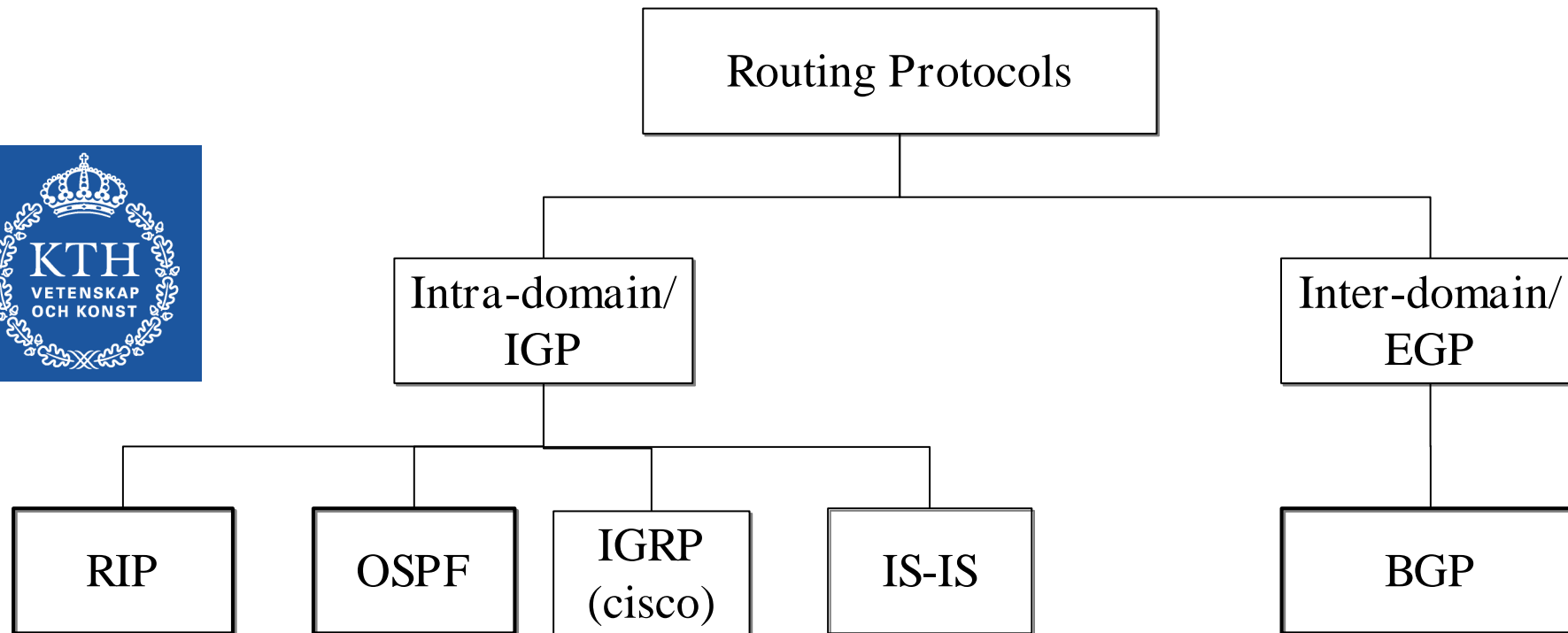| Permanent set | Tentative set |
|---|---|
| A 0 –<br>10.0.3.0/24 1 –<br>B 1 –<br>10.0.1.0/24 1 –<br>C 1 –<br>10.0.2.0/24 2 B<br>10.0.2.0/24 2 C<br>10.0.4.0/24 2 C<br>E 2 C<br>D 2 C<br>10.0.6.0/24 2 B<br>F 2 B<br>10.0.5.0/24 3 B<br>10.0.5.0/24 3 C | <br><br><br><br><br><br><br><br><br><br>~~10.0.5.0/24 3 C~~<br><br>~~10.0.5.0/24 3 B~~ |

Note: ECMP

# Exercise: Dijkstra tree graph view

- Compare with table view in the previous slide
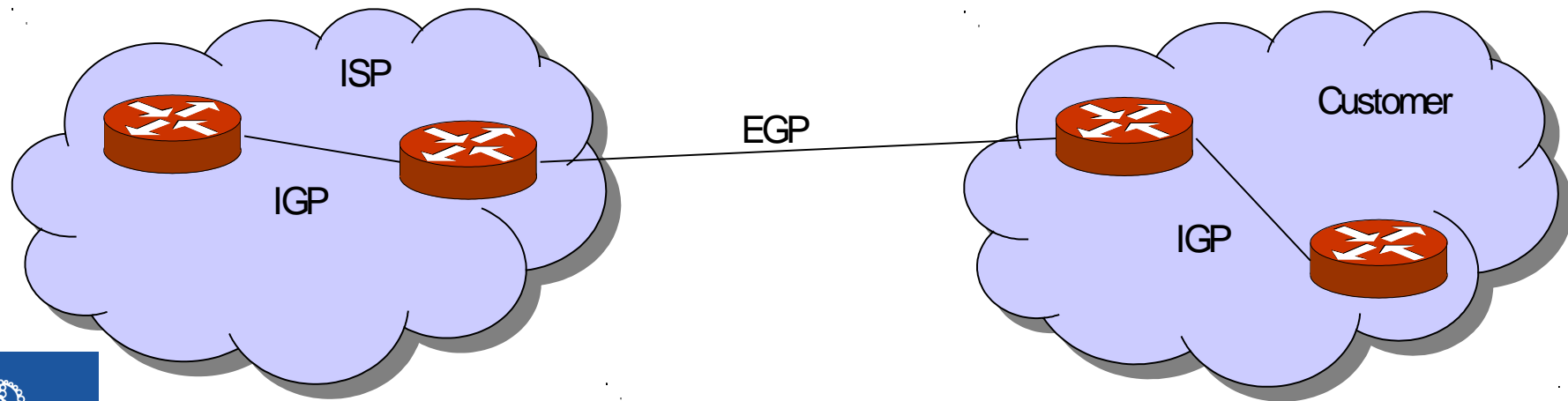- Note the ECMP routes to 10.0.2.0/24 and 10.0.6.0/24

# Intra-domain routing protocols

# Popular Unicast Routing Protocols

# IGP/EGP



## IGP

– Interior Gateway Protocol.

– Runs within a network/domain (intra-domain)

– Handles *internal* routes within a domain

– Examples: RIP, OSPF, IS-IS.

## EGP

– Exterior Gateway Protocol.

– Primarily exchanges routes between networks/domains (inter-domain)

– Handles *external* routes

– Examples: BGP, static routing

– Note: an EGP can handle *external* routes *within* a domain (IBGP)

# Routing Information Protocol - RIP

- RIP-1 (RFC 1058), RIP-2 (RFC 2453)
- Metric is hop counts
  - 1: directly connected
  - 16: infinity
  - RIP cannot support networks with diameter > 14.
- RIP uses distance vector
- RIP messages are carried via UDP datagrams.
  - IP Multicast (RIP-2): 224.0.0.9
  - Broadcast (RIP-1)

# Disadvantages with RIP

- Slow convergence
    - Changes propagate slowly
    - Each neighbor only speaks ~every 30 seconds; information propagation time over several hops is long
- Instability
    - After a router or link failure RIP takes minutes to stabilize.
- Hops count may not be the best indication for which is the best route.
- The maximum useful metric value is 15
    - Network diameter must be less than or equal to 15.
- RIP uses lots of bandwidth
    - It sends the whole routing table in updates.

# Why would anyone use RIP?

–It is easy to implement

–It is generally available

–Implementations have been rigorously tested

–It is simple to configure.

–It has little overhead (for small networks)

# Link-state routing

- Each router spreads information about its links to its neighbours.
- This information is flooded to every router in the routing domain so that every router has knowledge of the entire network topology.
- Using Dijkstra's algorithm, the shortest path to each prefix in the network is calculated
- OSPF and IS-IS are two well-known link-state routing protocols
- OSPF is popular among organizations (KTH uses OSPF)
- IS-IS is popular among operators (SUNET uses IS-IS)

# Comparison with distance-vector

- Link-state uses a distributed database model
- Distance-vector uses a distributed processing model
- Link-state pros:
  - More functionality due to distribution of original data, no dependency on intermediate routers
    - Easier to troubleshoot
  - Fast convergence: when the network changes, new routes are computed quickly
  - Less bandwidth consuming
- Distance-vector pros:
  - Less complex – easier to implement and administrate
  - Needs less memory

# The OSPF protocol

1) The *hello* protocol
   - Is there anybody out there?
   - Detection of neighboring routers
   - Election of designated routers

2) The *exchange* protocol
   - Exchange database between neighbours

3) Reliable *flooding*
   - When links change/age send: update to neighbours and flood *recursively*.

4) *Shortest path* calculation
   - Dijkstra's algorithm
   - Compute shortest path tree to all destinations
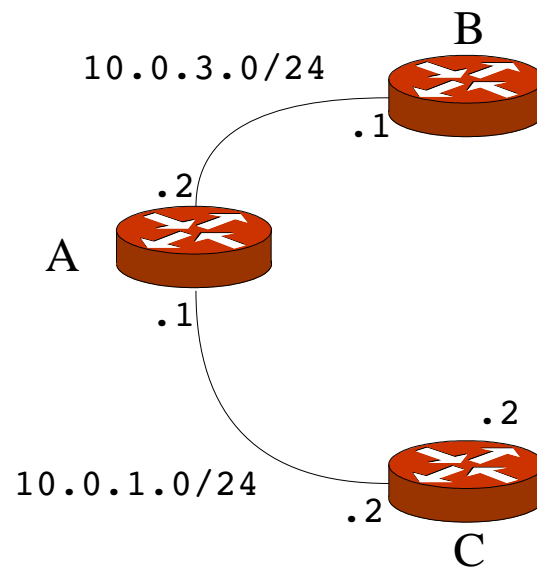
# The link-state

- Each router describes its environment in the form of networks (links) that it is connected to
- These link-states are the elements of the distributed database
- Fundamental task in OSPF is to distribute the link-states to all nodes in a reliable way
- Then, each node can compute Dijkstra on the *same* database
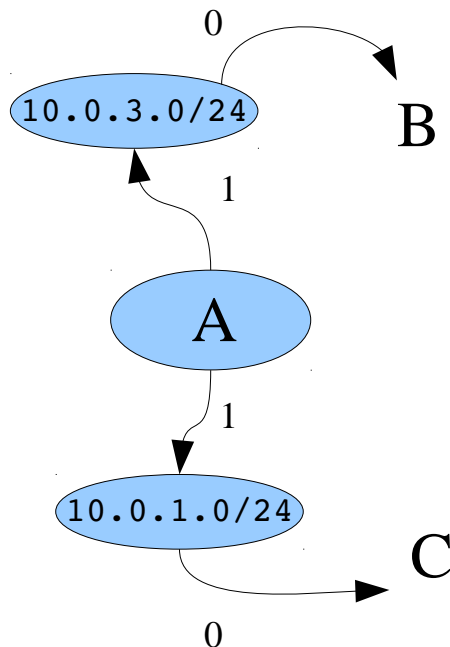  - The result (shortest path) is consistent everywhere

# Example: OSPF link state

- Translate the network below to link states (from As point of view)



B

10.0.3.0/24

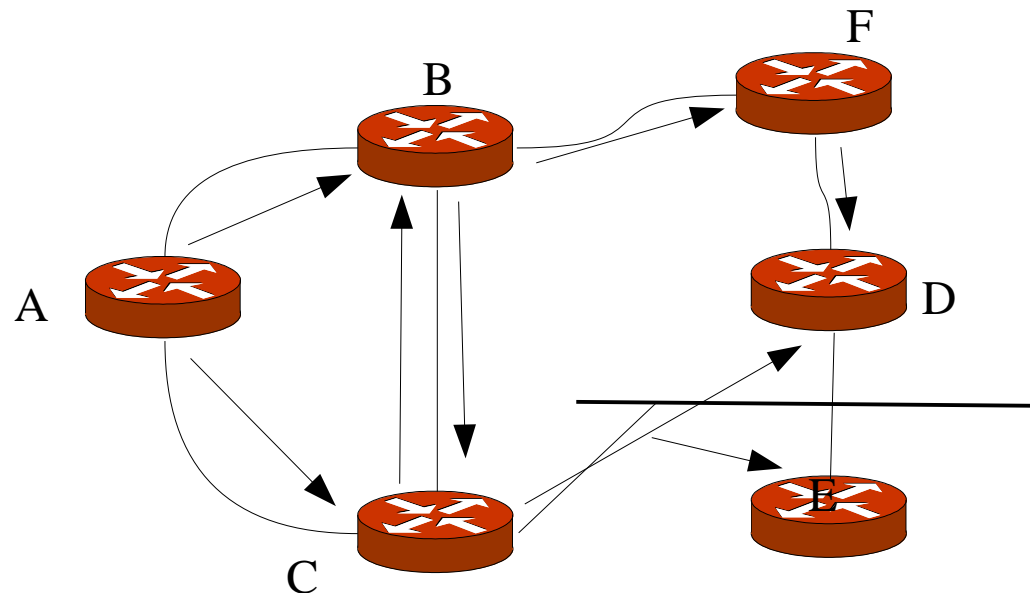.1

.2

A

.1

.2

10.0.1.0/24

.2

C

# Example: OSPF link state

- Every node creates the link-state of its connected links
- Example: A is connected to two 'transit' links, which are in-turn connected to B and C respectively
- The transit links in this case 'belongs' to A, since A is *designated router* of these sub-networks.
- A therefore distributes the three link-states in the figure:
  - One for A itself (it is connected to two transit networks)
  - Two transit links which are connected to A and B, and A and C respectively.

# Flooding link-state

- Every node distributes its link-state to all others
- Initially and after link changes (also every ~30 mins)
- Example: 'A' floods its link-state by sending it to its neighbors, who in turn distributes it to their neighbors, etc
- Flooding is made reliably and to all routers
  - No need for periodic retransmit – no waste of bandwidth
- The flooding protocol is the most complex part of OSPF (not Dijkstra!)

# Inter-domain routing

# Inter-domain routing

- The objective of inter-domain routing is to bind together the hundreds of thousands of independent IP networks that constitute the Internet

  Perspective from one network:
- Spread routing information to the outside world
  - Originate and aggregate address prefixes
  - Announce prefixes to other domains
  - Tag prefixes with routing information
- Receive information from the outside world
  - Receive and choose (filter) between prefixes from other domains
- Transfer information through your routing domain
  - Received information from one domain may be transferred (and possibly modified) to other domains

# What is BGP?

- Border Gateway Protocol version 4
- An inter-domain routing protocol
- Uses the *destination-based* forwarding paradigm
  - Other relations are not expressed: sources, tos, link load
- Uses path-vector routing
- Views the Internet as a collection of autonomous systems
- Exchanges information between peers using TCP as underlying protocol
- Maintains a database (RIBs) of network layer reachability information
- Tags destinations with *path attributes* which describes different properties of the destination
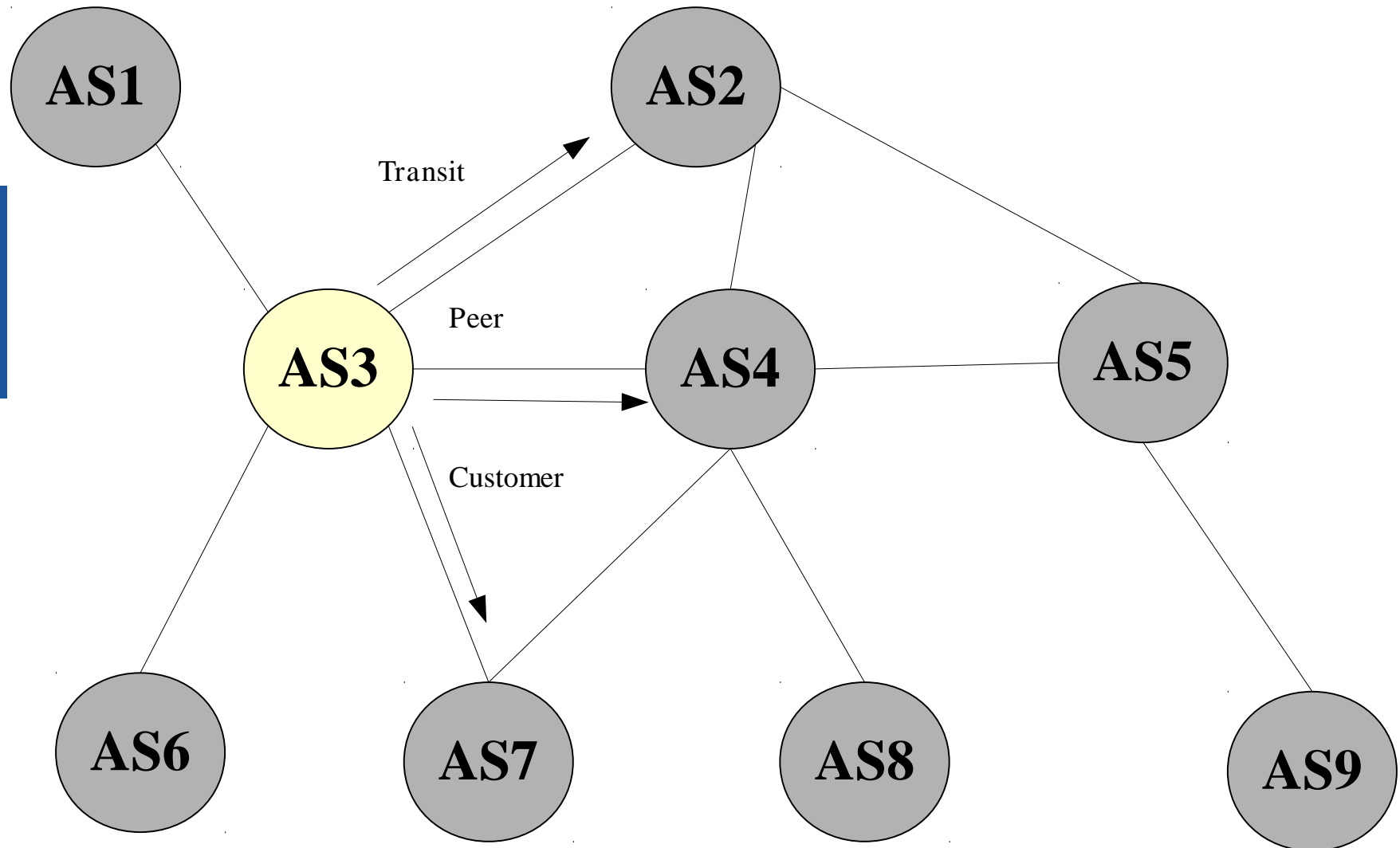- Supports a toolkit of mechanisms to express and enforce policy decisions at AS level

# Autonomous Systems (AS)

- A set of routers that has a single routing policy, that run under a single technical administration
  - A single network or group of networks
  - University, business, organization, operator
- This is viewed by the outside world as an Autonomous System
  - All interior policies, protocols, etc are hidden within the AS
- Represented in the Internet by an Autonomous System Number (ASN). 0-65535
  - Example: ASN 1653 for SUNET
- Currently, operators are switching to four-byte ASNs
  - RFC 4893: BGP Support for Four-octet AS Number Space

# AS peering relations

# Whois example

```
gelimer.kthnoc.net> whois -h whois.ripe.net AS1653
aut-num:        AS1653
as-name:        SUNET
descr:          SUNET Swedish University Network
import:         from AS42 accept AS42
export:         to AS42 announce AS-SUNET
import:         from AS702 accept AS702:RS-EURO AS702:RS-CUSTOMER
export:         to AS702 announce AS-SUNET
import:         from AS2603 accept any
export:         to AS2603 announce AS-SUNET
import:         from AS2831 accept AS2831 AS2832
export:         to AS2831 announce any
import:         from AS2833 accept AS2833
export:         to AS2833 announce any
import:         from AS2834 accept AS2834
export:         to AS2834 announce any

gelimer.kthnoc.net> whois -h whois.ripe.net AS-SUNET
as-set:         AS-SUNET
descr:          SUNET AS Macro
descr:          ASes served by SUNET
members:        AS1653, AS2831, AS2832, AS2833, AS2834, AS2835, AS2837
members:        AS2838, AS2839, AS2840, AS2841, AS2842, AS2843, AS2844
members:        AS2845, AS2846, AS3224, AS5601, AS8748, AS8973, AS9088
members:        AS12384, AS15980, AS16251, AS20513, AS25072, AS28726
members:        AS-NETNOD
```
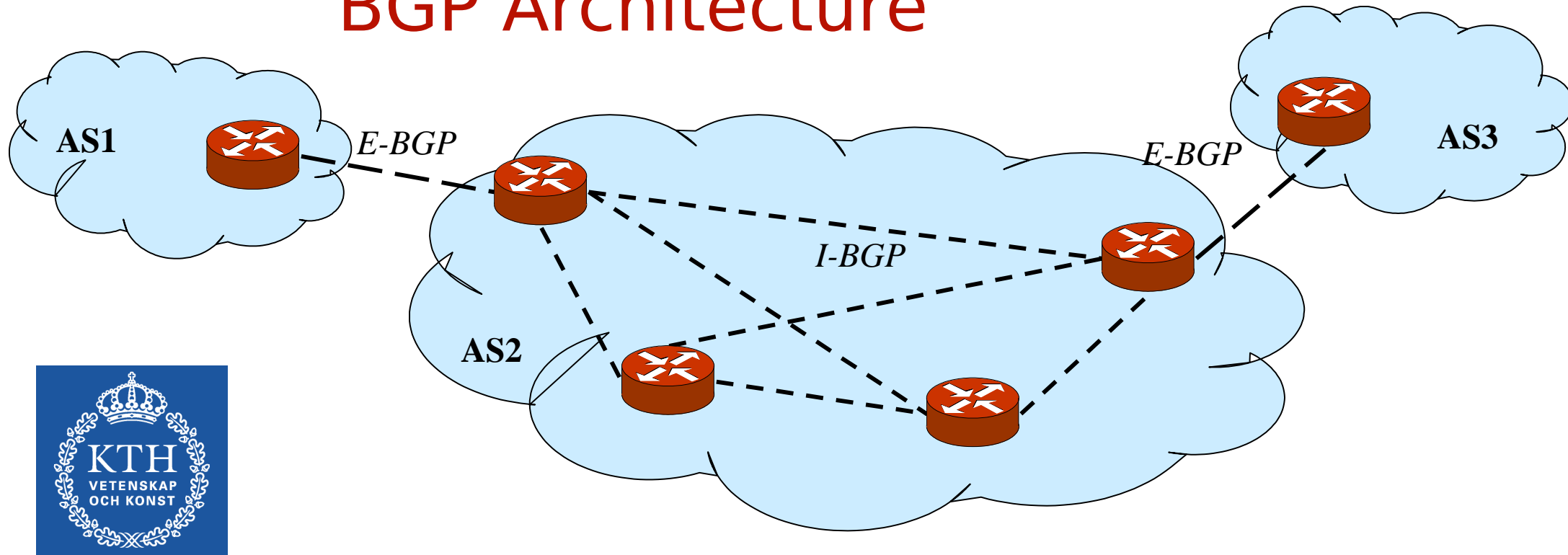
(Edited example)

# Peering relations

- An abstract way of defining peering relations is for example:
- Prefix sets:
    - Define a customer set, a peering set and a transit set
- Example rules:
    - Customer prefixes should be announced to transit and peers
    - Peer and transit prefixes should be announced to customers
    - Prefer prefixes from peers over prefixes from transit
    - Do not accept illegal prefixes (RFC 1918 for example), or unknown prefixes from customers
    - Load balance over several transit providers
    - Filter traffic (eg src addresses) according to the prefixes announced
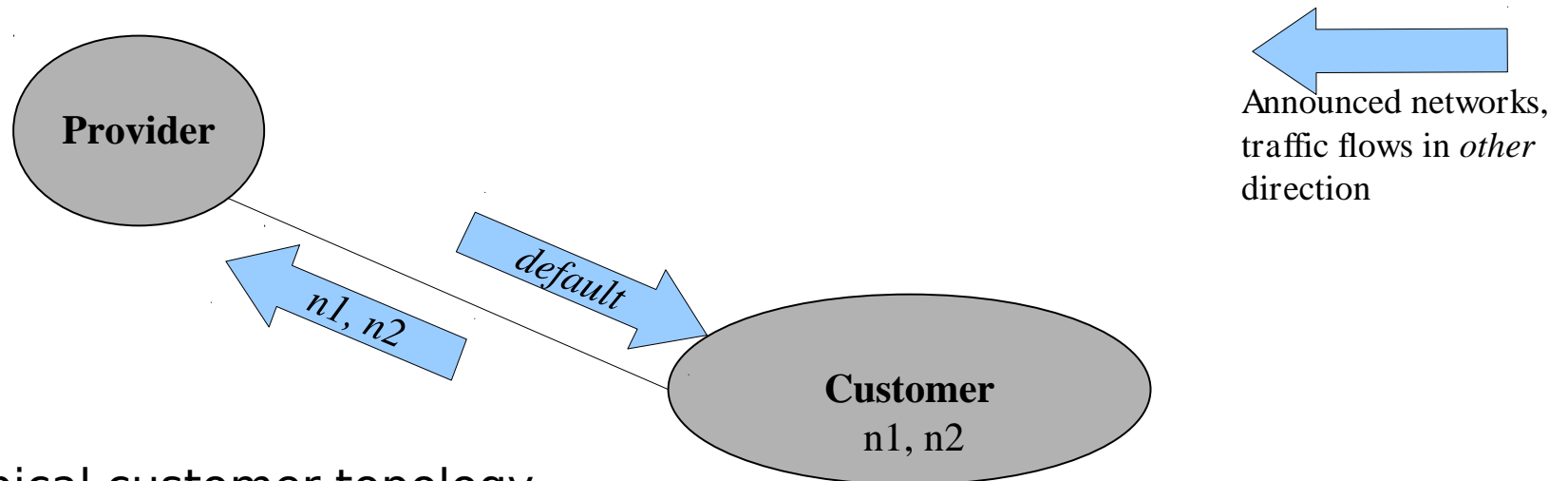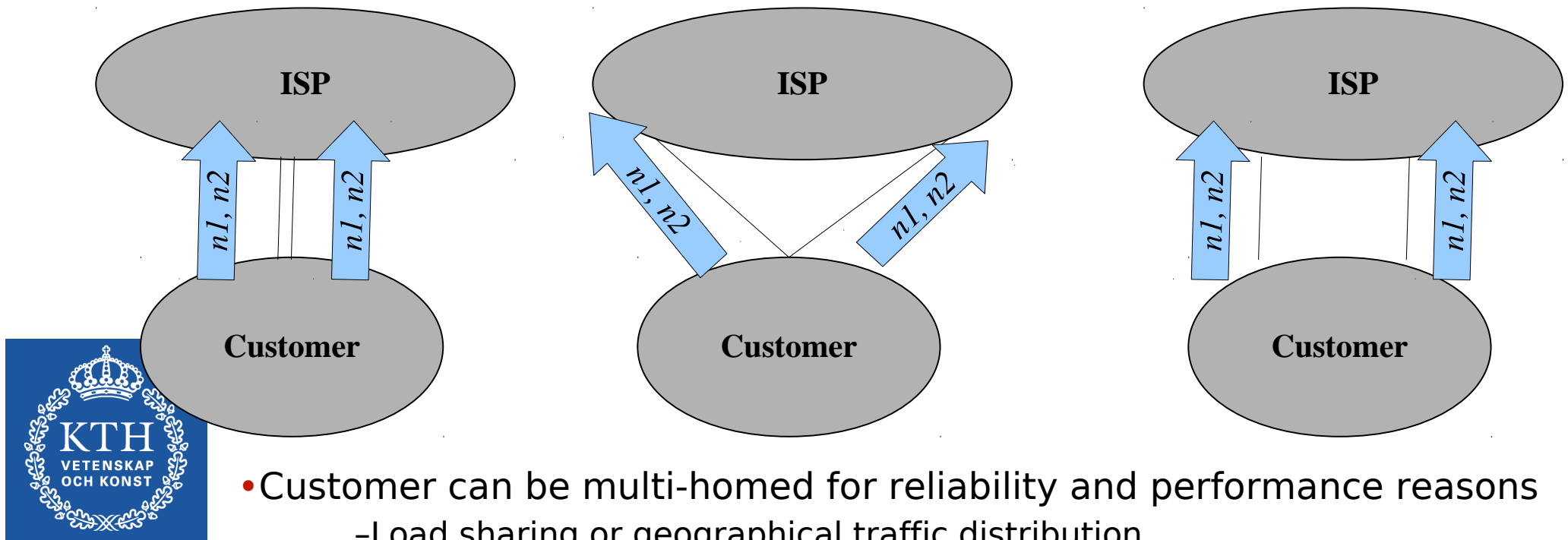
# BGP Architecture



- BGP interacts with the Internal routing (OSPF/IS-IS/RIP/...)
    - *Redistributes* internal / external routes between the two protocols
- BGP really consists of two variants:
    - E-BGP: exchanges external routes between border routers *between* AS:s
    - I-BGP : synchronises *external* routes *within* an AS (IGP takes care of internal routes)

# Customer / ISP Relations: Stub AS


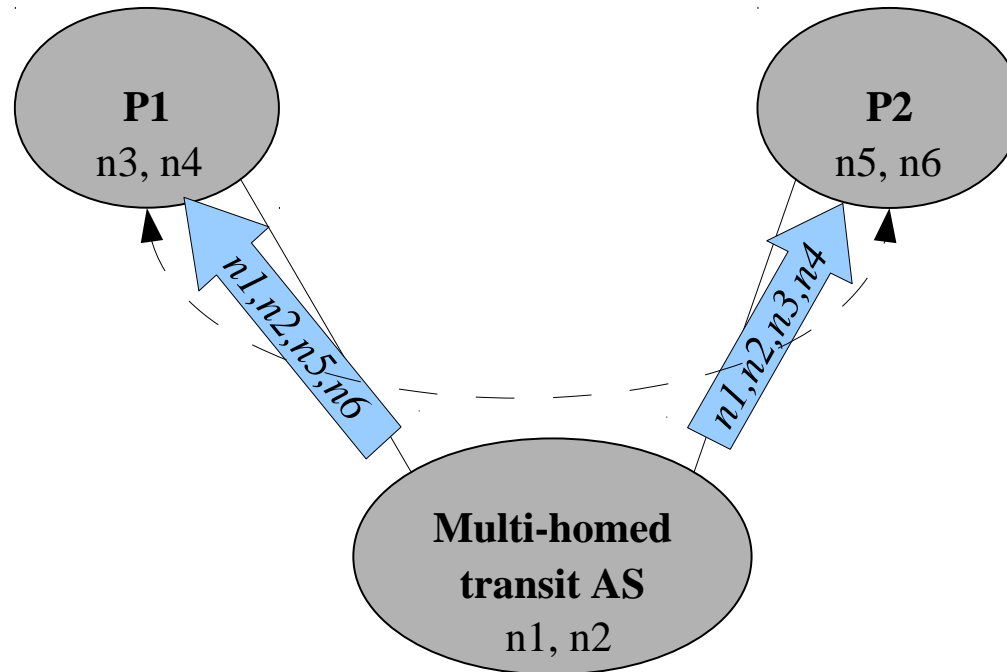
Announced networks, traffic flows in *other* direction

- Typical customer topology
- Can use default route to reach the Provider and Internet
- Customer can use address block of provider
- Customer does not need to be a separate AS
- Typically use static routing but can also use BGP

# Multi-homed customer



- Customer can be multi-homed for reliability and performance reasons
  - Load sharing or geographical traffic distribution
- Multi-homed non-transit AS
  - Non Transit AS does now allow external traffic to pass through
- What to think about:
  - How to announce the prefixes
  - Default routes
  - Symmetrical routing
  - Packet filtering,
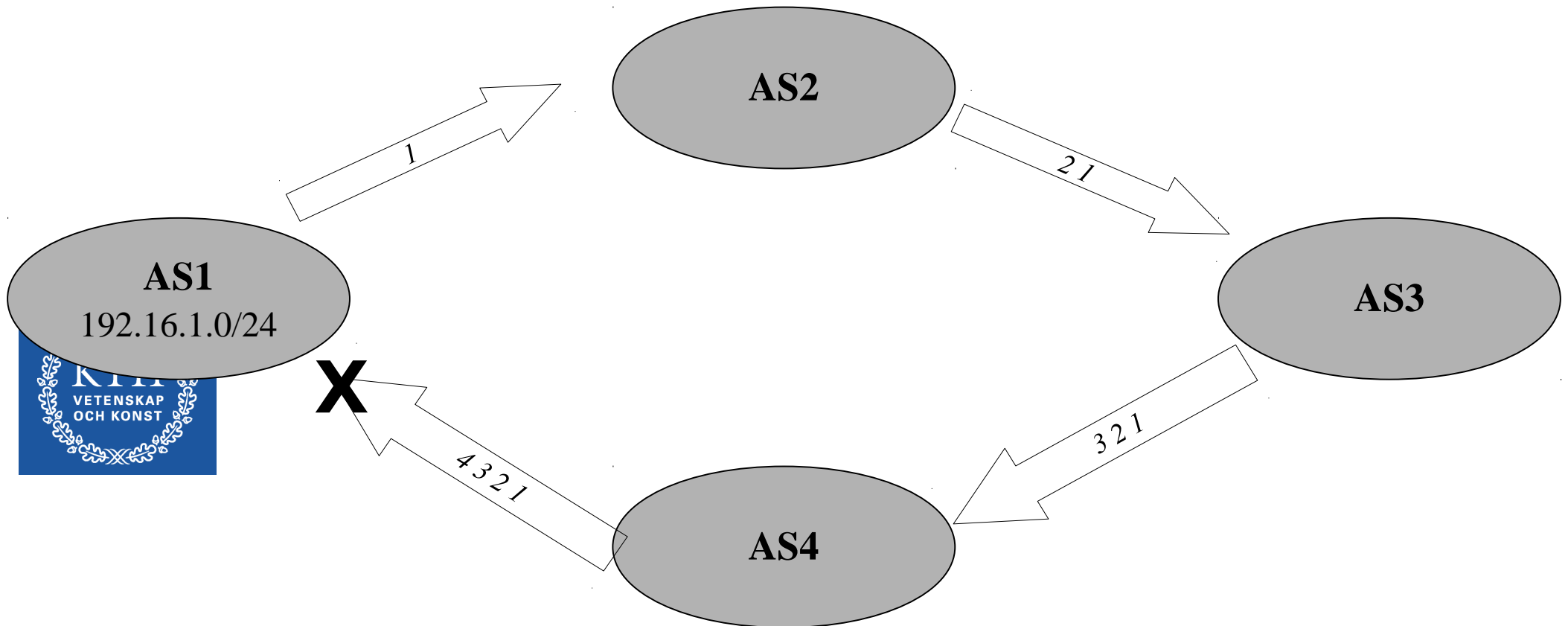  - Address aggregation, etc

# Provider: Multi-homed Transit AS



- Transits traffic within own network
- This the most general configuration and is how an internet provider works.

# Path vector protocol

- In a distance-vector protocol, vectors with destination information are distributed between routers:
- Example:
    - <dst: 10.1.10/24, metric: 5, nexthop: 10.2.3.4>
- Distance-vector has problems with converging
    - Example: count-to-infinity
- Path-vector extends the information with a *path* to the destination
    - This enables immediate loop detection
    - Several other attributes associated with path
- Also, in BGP, the path vector uses AS:s, not IP addresses
    - This hides internal structure in the domains
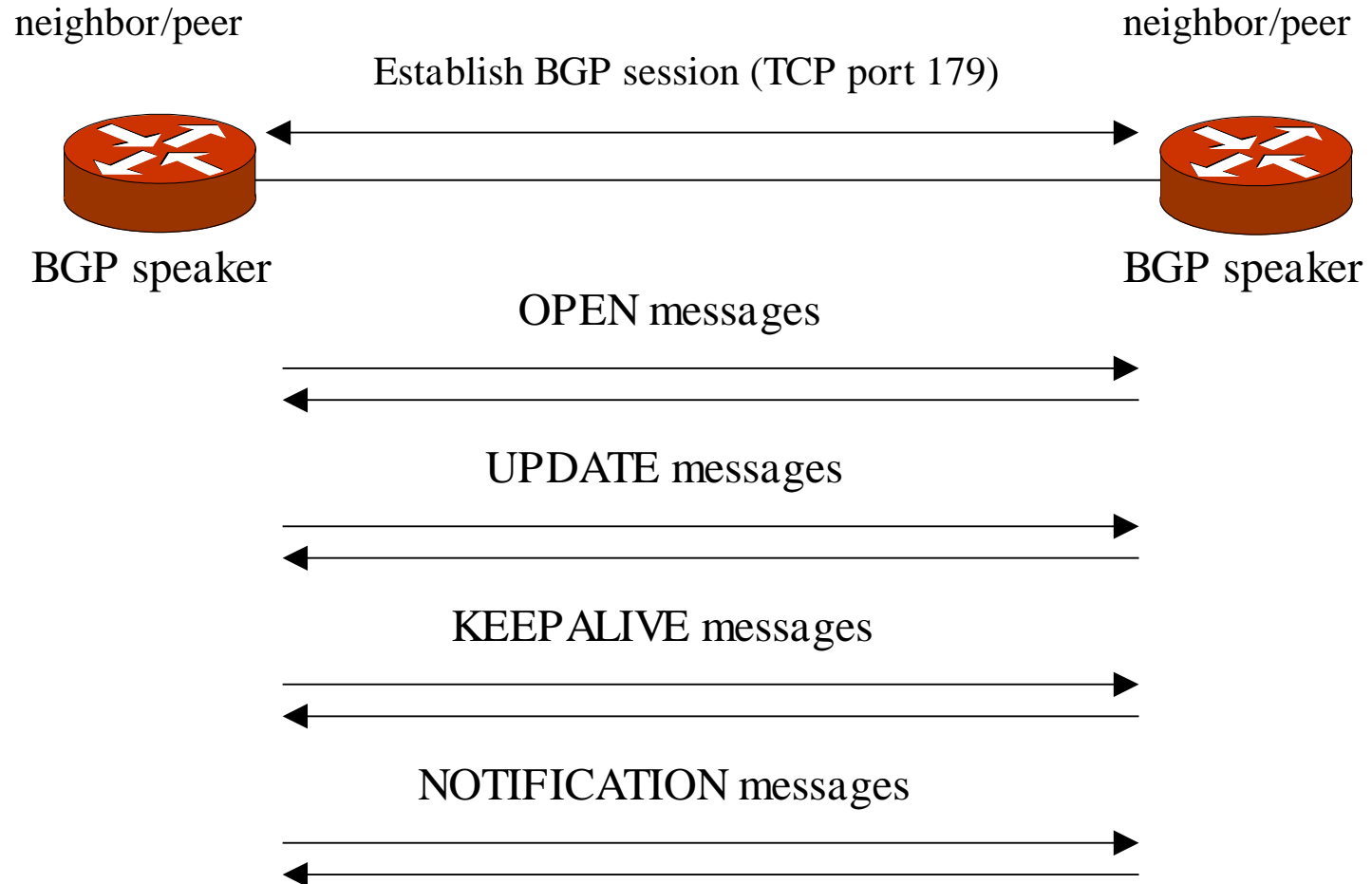    - Loop detection only on AS-numbers!
- Example:
    - <dst: 10.1.10/24, path: AS1:AS3:AS5, nexthop: 10.2.3.4>

# AS_PATH



- AS-PATH is used to break loops (between AS:s)
- AS1 announces 192.16.1.0/24 to AS2 and detects its own ASN when received from AS4
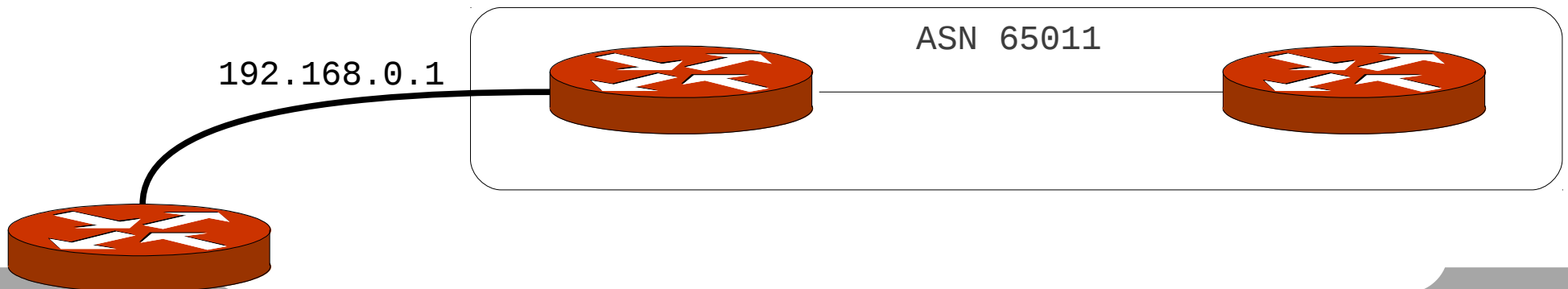- AS-PATH is the most well-known path-attribute, there are several others

# BGP Operation

neighbor/peer                                                     neighbor/peer

Establish BGP session (TCP port 179)

BGP speaker                                                       BGP speaker

OPEN messages

UPDATE messages

KEEPALIVE messages

NOTIFICATION messages

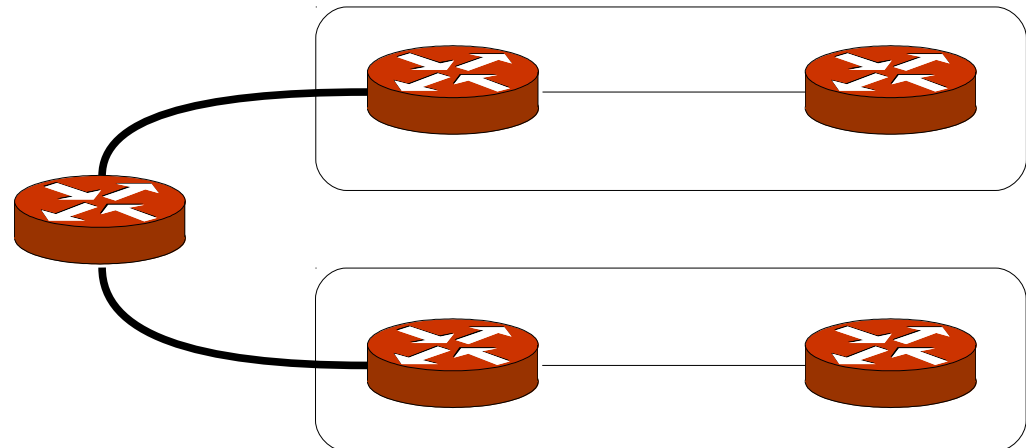# Example: JunOS BGP configuration

```
routing-options {
  autonomous-system 65011;
}
protocols {
  bgp {
      group EXTERN {
        type external;
        peer-as 65001;
        export MYNETWORK;
        neighbour 192.168.0.1;
        }
      }
  }
}
```

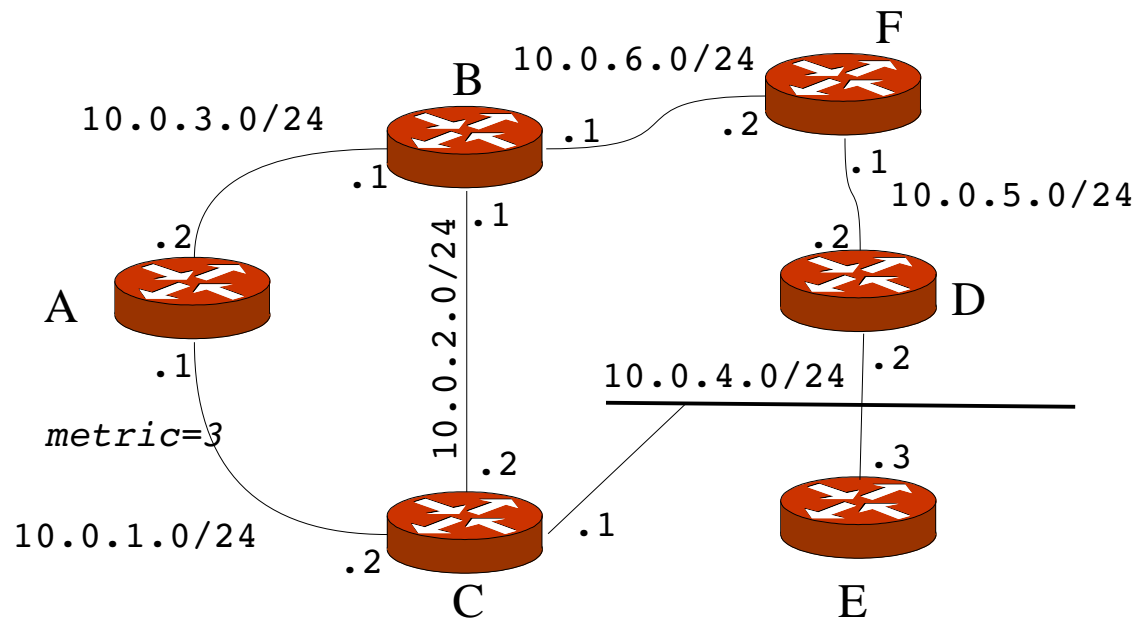ASN 65011

192.168.0.1

# (EP2120 only) Lab 4: Dynamic routing

- Four parts:
  - JunOS CLI
  - Static configuration
  - OSPF
  - BGP
- Make preparation questions in advance!

# Homework

- Second part of homework 3 is a distance-vector exercise.
- Fill in routing tables and distance-vectors from A's perspective on the following network:
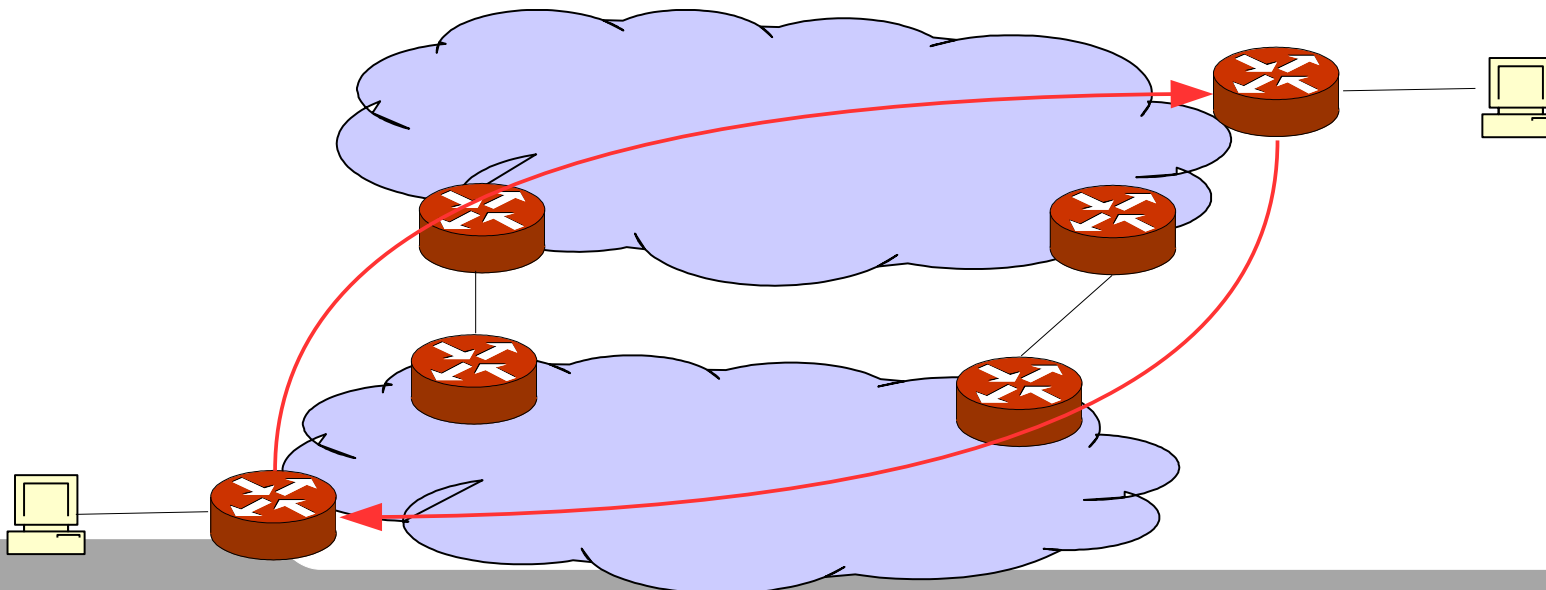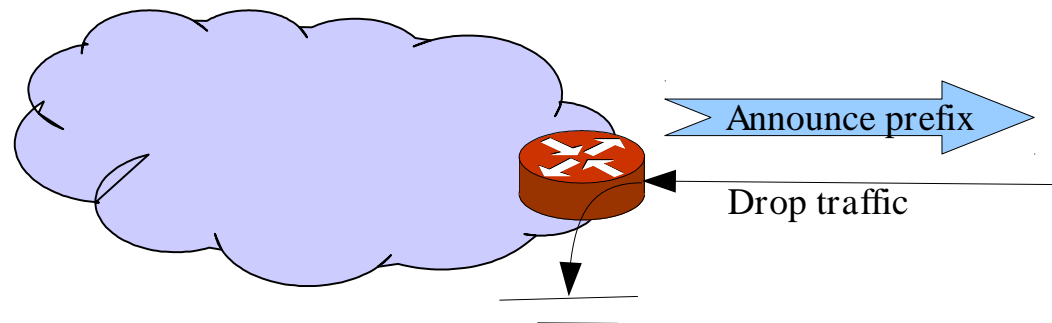
# Extra material

# Asymmetric Routing

- A rule rather than an exception:
  - To- traffic and from- traffic take different paths
- Hot-potato routing
  - Send traffic out of your AS as soon as possible
- Cold-potato
  - Try to keep your traffic as long as possible.

# Black-holing

- Black-holing: announce prefix, but traffic to the prefix is dropped (not delivered)
- Loops: circular announcements causing packet loops
    - TTL is decremented until packet drops -> same symptom as black-holing
- Reasons:
    - Transient errors due to long convergence (see count-to-infinity in distance-vector)
    - Misconfigurations
    - Attacks (DOS, man-in-the-middle)
    - Response to attacks: create a black-hole for attacked prefixes which removes DOS traffic



Announce prefix

Drop traffic

# Redistribution of routing information

- If several protocols are running on the same router
  - E.g., an OSPF as interior and BGP as exterior
- The router can distribute routes from one protocol to another
  - Interior routes need to be advertized to the Internet
    - Typically these routes are aggregated
  - Exterior routes (or a default) may need to be injected into the interior network
    - But only a subset – the backbone tables are very large
    - Necessary for domain carrying *transit* traffic
    - Not necessary for a domain using only a default route
- Typically, redistributed routes are filtered in different ways due to routing policies