# ISP Edge design
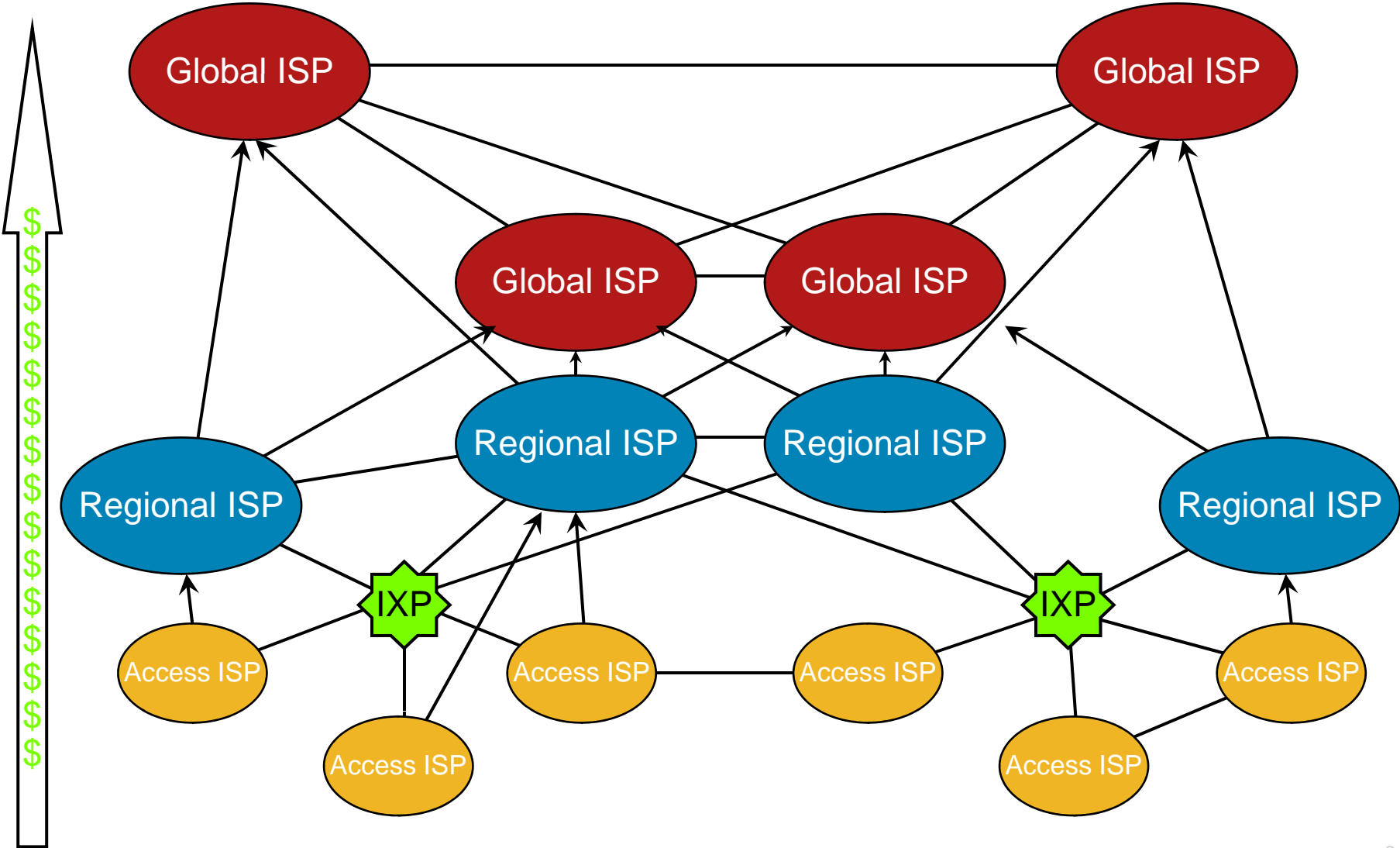
**Josef Ungerman**

CCIE #6167

# Agenda

- **The Internet**

- **IXP Intro**

- **Euro-IX**

- **Technical Details**

- **Live Examples**

- **OTT, Video and IXP**

- **Summary & Resources**

# Categorising ISPs

# Peering and Transit

- **Transit**

    Carrying traffic across a network

    **Usually for a fee**

    Example: Access provider connects to a regional provider

- **Peering**

    Exchanging routing information and traffic

    **Usually for no fee**

    Sometimes called **settlement free peering**

    Example: Regional provider connects to another regional provider

# Private Interconnect

- Two ISPs connect their networks over a **private link**
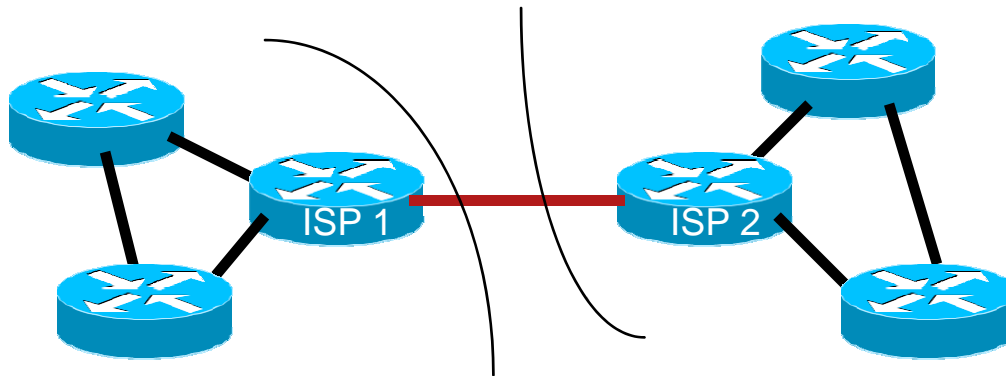    - Can be peering arrangement
        - No charge for traffic
        - Share cost of the link
    - Can be transit arrangement
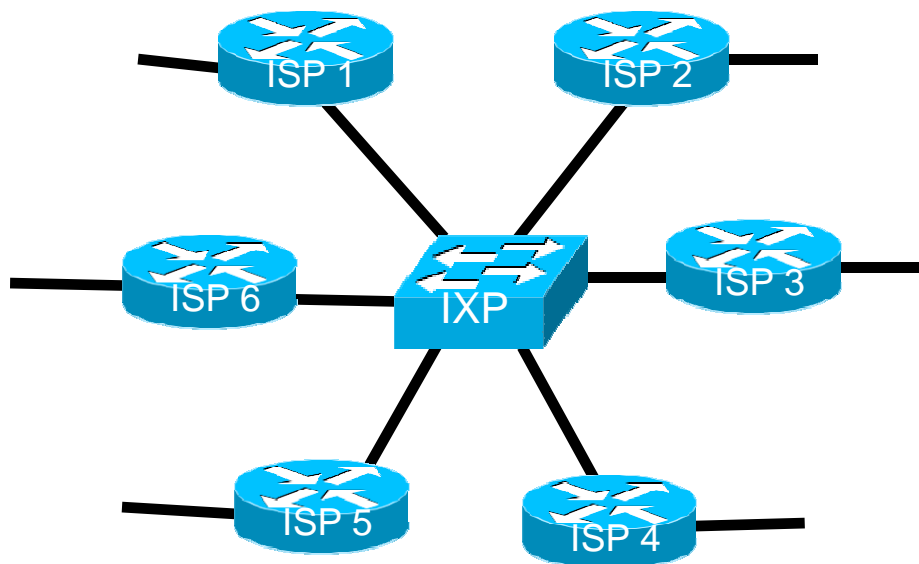        - One ISP charges the other for traffic
        - One ISP (the customer) pays for the link

ISP 1

ISP 2

# Public Interconnect

- Several ISPs meeting in a common neutral location and interconnect their networks

  Usually is a peering arrangement between their networks
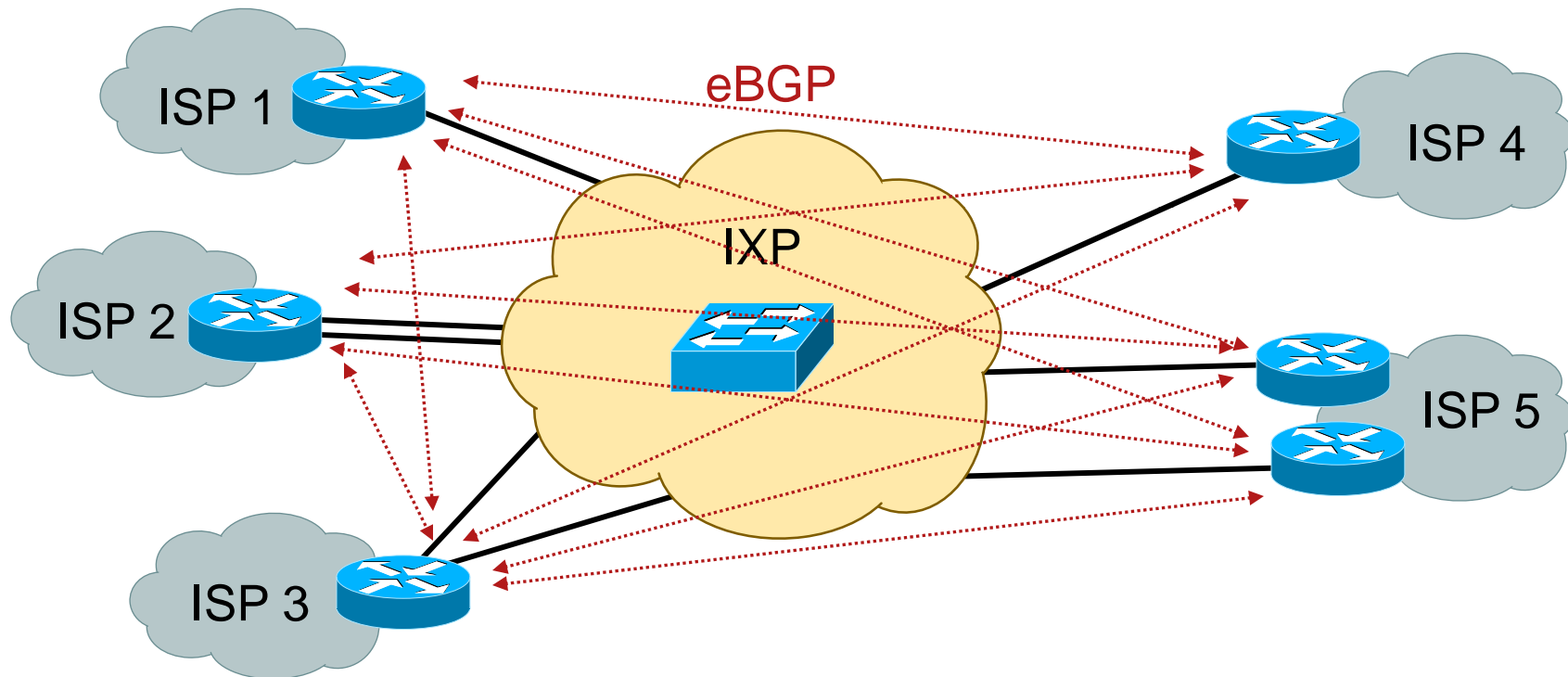
# IXP
# (Internet Exchange Points)

# IXP (Internet eXchange Point)

A physical network infrastructure operated by a single entity with the purpose to facilitate the exchange of Internet traffic between Autonomous Systems. The number of Autonomous Systems connected should at least be three and there must be a clear and open policy for others to join.
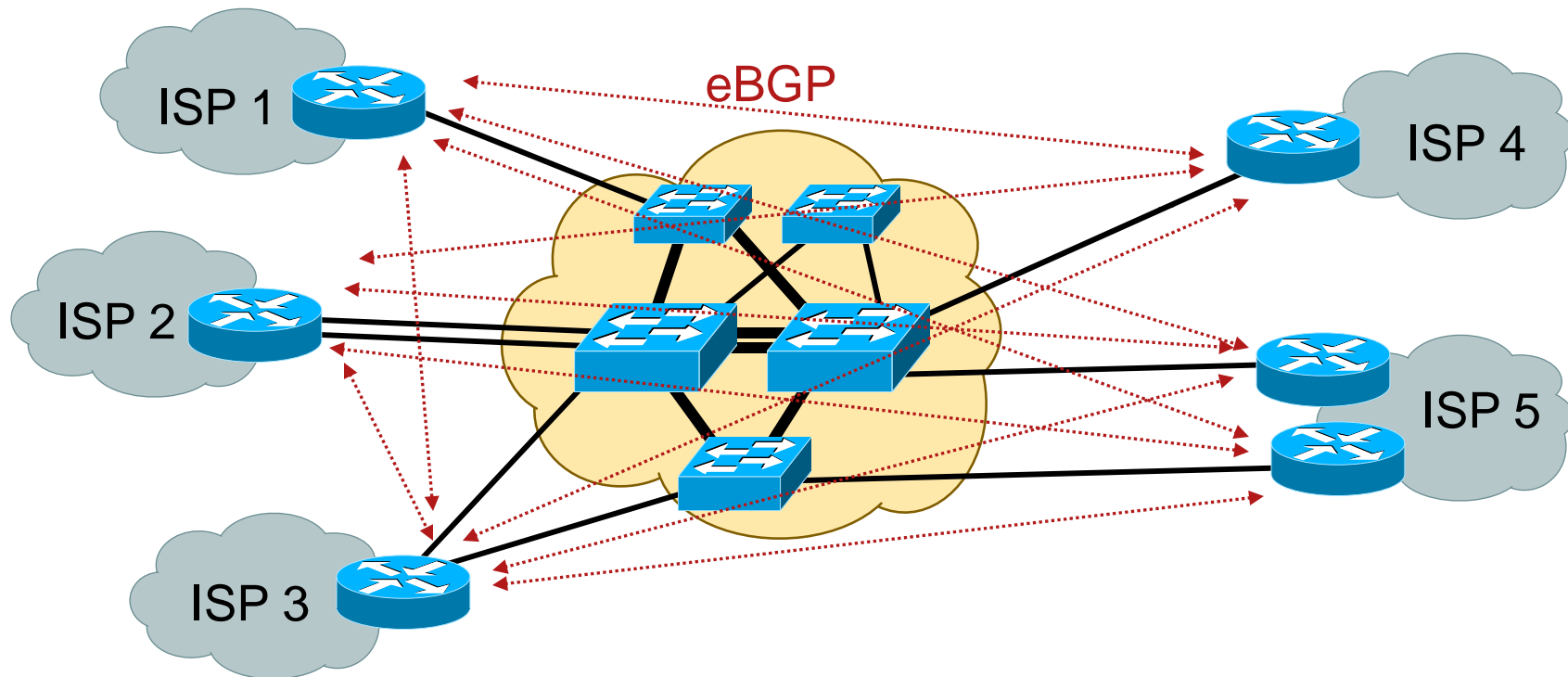
- High-speed/Low-cost Internet Traffic Exchange

- A.k.a. Public Peering or Settlement-Free Peering

- Non-Profit Associations or Commercial Datacenters
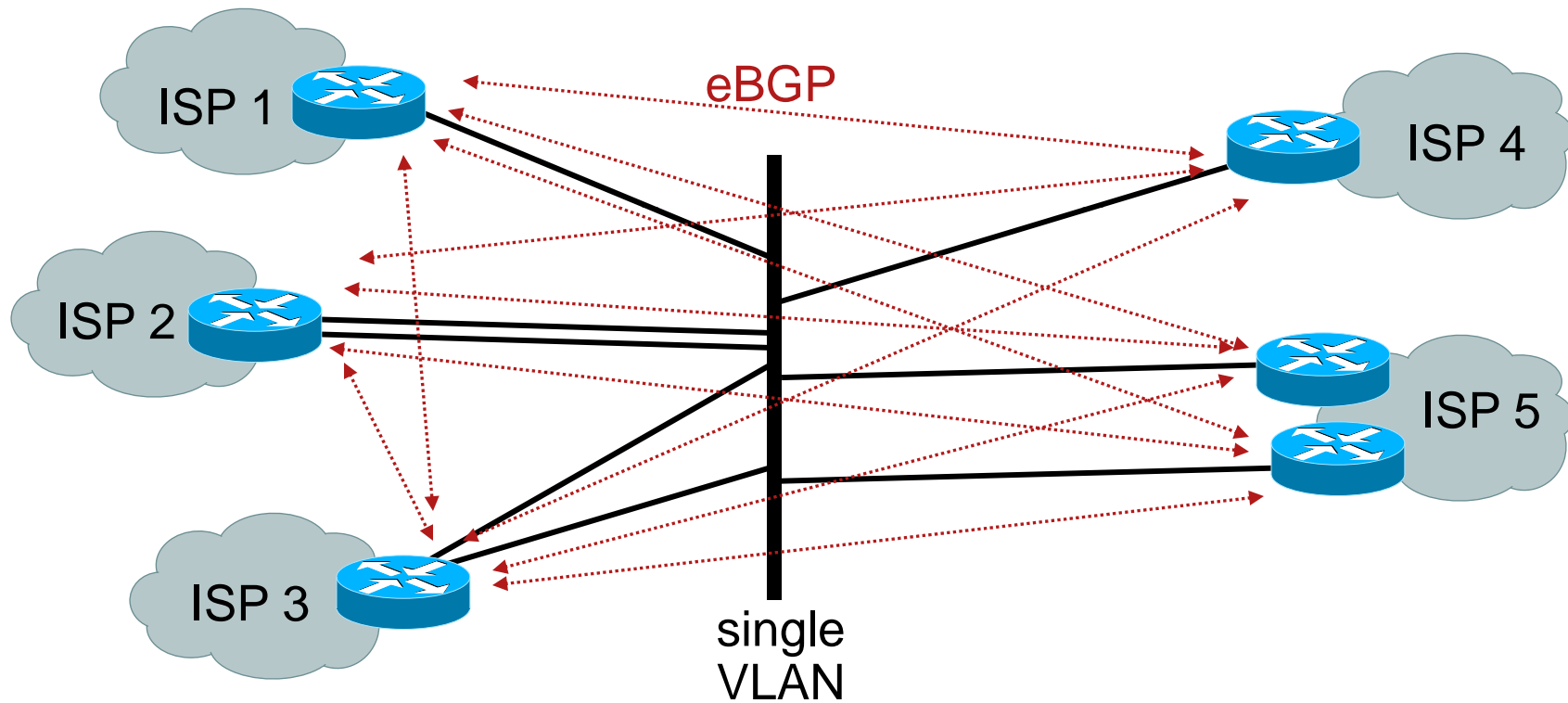
- Around 300 big IXPs in the world

# IXP (Internet eXchange Point)

# IXP (Internet eXchange Point)

eBGP

ISP 1

ISP 2

ISP 3

ISP 4

ISP 5

# IXP (Internet eXchange Point)



eBGP

ISP 1

ISP 2

ISP 3

ISP 4

ISP 5

single
VLAN

# Euro-IX

**Euro-IX (European Internet Exchange Association)** was formed in May 2001 with the intention to further develop, strengthen and improve the Internet Exchange Point (IXP) community

- **105 IXPs in 102 cities in 31 countries**

- 9 non-european members

- www.euro-ix.net

## European IXP growth

# Euro-IX Report 2008

## 2002- 2008 Traffic History (Euro-IX IXPs only)

# Euro-IX Report 2008



Aggregated Peak Traffic per country

IXPs and their peak traffic

# Euro-IX Report 2008

## Total number of IXP particpants per country



LV, PL, UA –
• highly fragmented ISP market
• maybe a lot of Hosting DC's

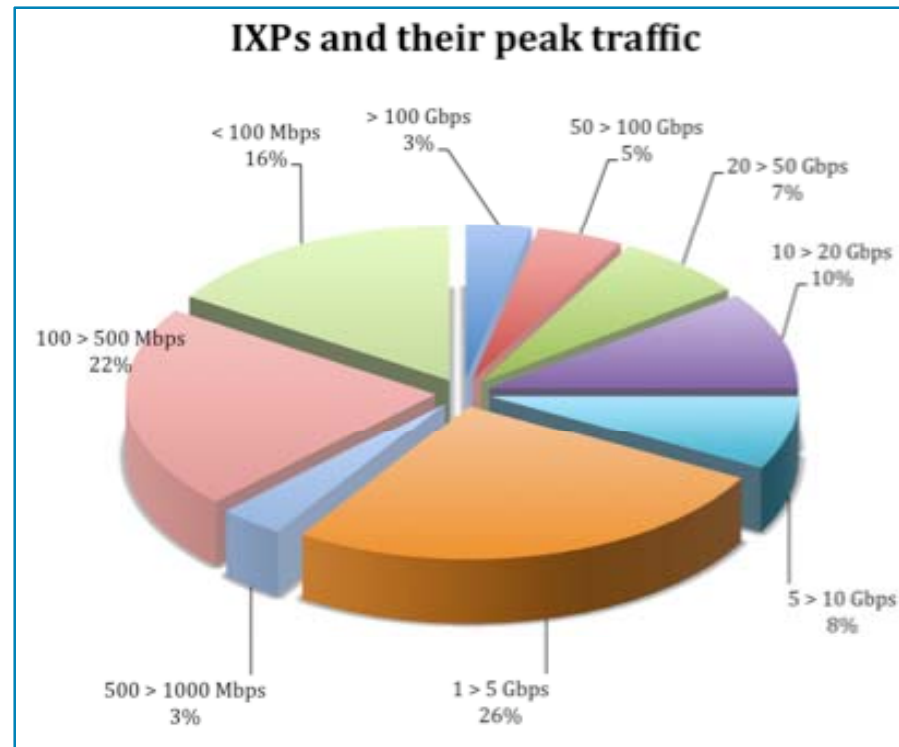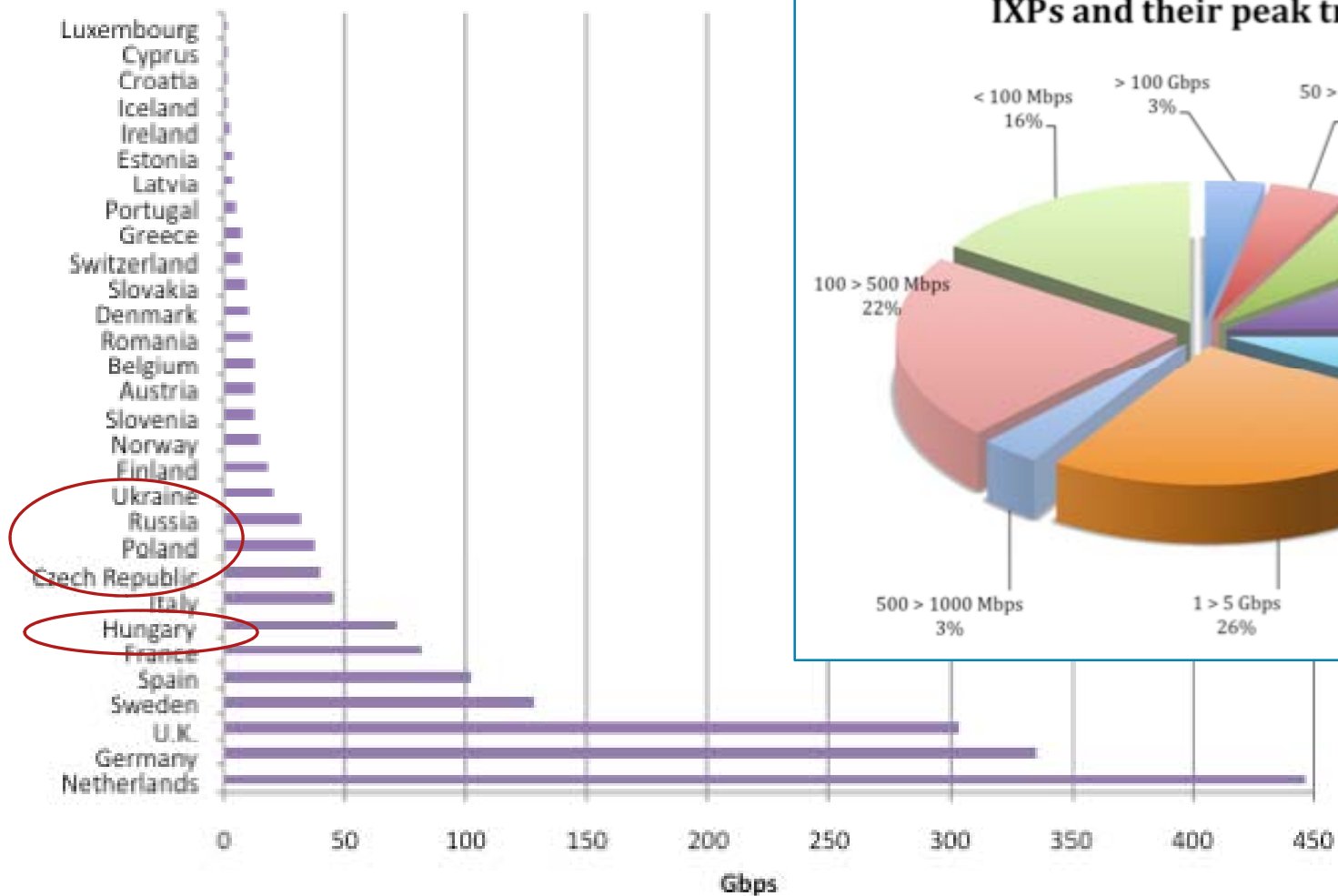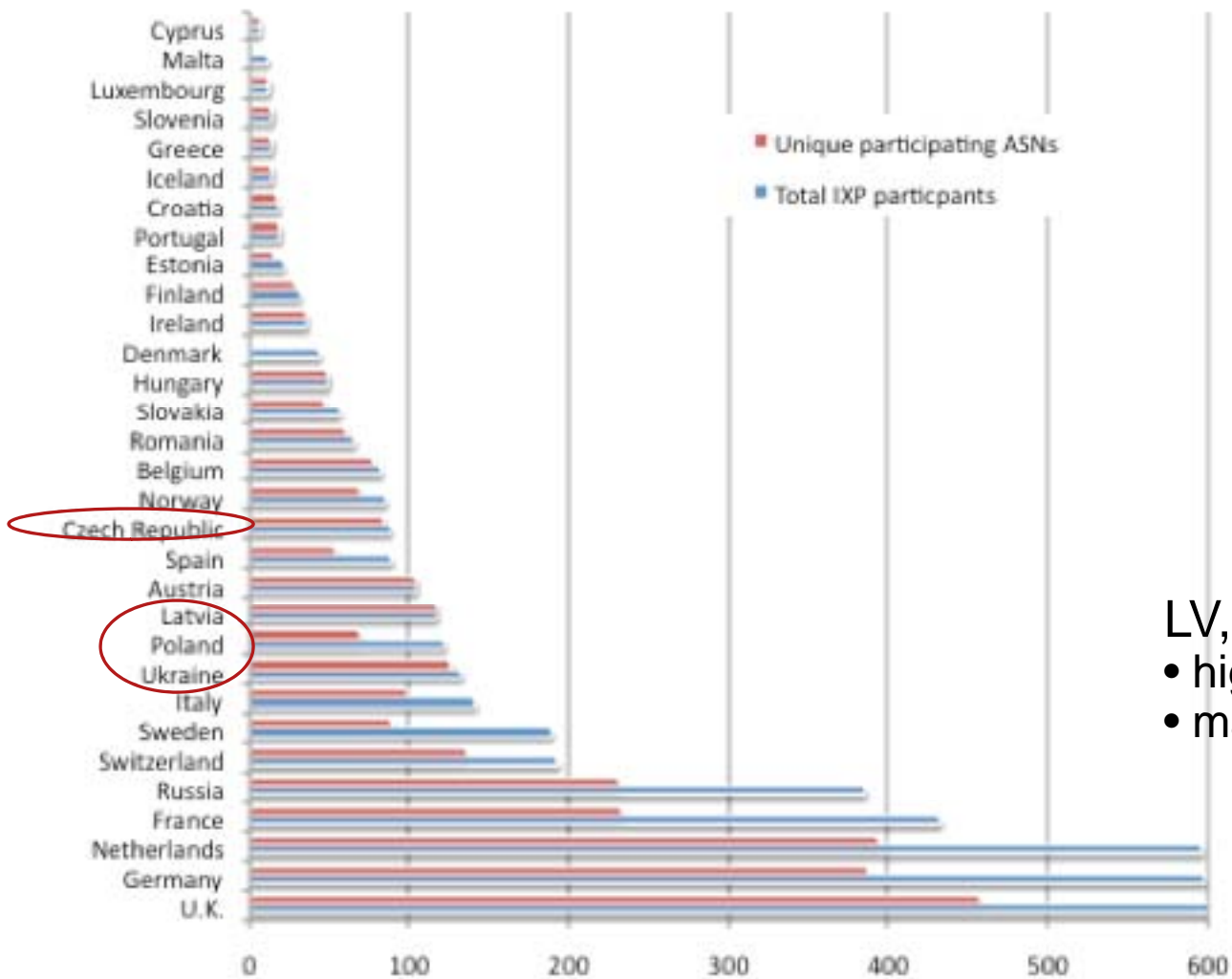# Example: GoogleNet…
## A PortalNet… Dedicated CDN… Parallel Internet BackBone

**GoogleNet** (Faster, Cheaper, More Reliable)

DataCenters can be colocated at Peering Points

10G    N*10G    40G    100G    N*100G

Full Mesh Peering    Tier 1 ISPs    Transit Free

**Some 300 Exchanges Worldwide**

IX    IX    IX    IX    IX

Partial Mesh Peering    Tier 2 ISPs    Must Buy Transit

IPTV Local Loop

Upgrades

**Google-WIFI**    **Mobile**

Content / Enterprise Companies / Users

Generally No Peering    Must Buy Transit

- Google has been buying Fiber on a Worldwide basis

- Google builds it's own worldwide IP Backbone.

- Google peers locally, often on a Settlement Free Basis, with Eyeball Carriers.

- Google can send any amount of traffic into the Internet without paying anyone, they are Nobody's Customer.

- Google distributes it's DataCenters to be virtually ONnet to Eyeball networks. Google is now only a few Hops away from Any User on the Internet.

- Tier2 ISP's invest in massive Local Loop upgrades to support IPTV.

- Google drives Net Neutrality so that whatever Traffic they send, can't be impaired.

- Google can now addresses Service Substitution (Google TV, Voice…)

# Internet Edge

# ISP design – peering layer



P

MPLS Core

# ISP design – peering layer

INTERNET

Upstream ISP's

eBGP

International
IGW

P

MPLS Core

# ISP design – peering layer



INTERNET

Upstream ISP's

eBGP

International
IGW

iBGP

IPv4 Route
Reflectors

P

MPLS Core

# ISP design – peering layer

INTERNET

Upstream ISP's

IXP

eBGP

**National IGW**

**International IGW**

iBGP

**IPv4 Route Reflectors**

P

MPLS Core

# ISP design – peering layer



INTERNET

IXP

Upstream ISP's

National IGW

International IGW

IPv4 Route Reflectors

P

MPLS Core

ISP Transit Routers

ISP Customers

eBGP

iBGP

eBGP

22

# ISP design – peering layer



INTERNET

IXP

Other ISP's

IGW

P

IPv4 Route Reflectors

MPLS Core

ISP Transit Routers

ISP Customers

eBGP

iBGP

eBGP

# ISP design – peering layer



INTERNET

IXP

Other ISP's

Internet GW
+ ISP Transit

N-PE

MPLS

EoMPLS
pseudowire

U-PE

ISP Customers

eBGP

eBGP

# Internet Gateway

Cisco Internal

# Cisco Internet Gateway Routers

| | ASR 1000 | CRS-1/4 | CRS-1/8 | CRS-1/16 | CRS-1 MC |
|---|---|---|---|---|---|
| Throughput | 20 Gbps | 320 Gbps | 640 Gbps | 1.28 Tbps | 10 Tbps |
| Scalability | 40 Gbps | 960 Gbps | 1.92 Tbps | 3.84 Tbps | 100+ Tbps |
| FIB entries | 2 Million | 2 Million | 2 Million | 2 Million | 2 Million |
| Netflow entries | 2 Million | 4 Million | 8 Million | 16 Million | 100+ Million |

**Existing deployments (~60% marketshare)**
- The most used ISP GW is Cisco 12000 (GSR)
- Many deployments are based on Cisco 7600
- Many small IGW's are still Cisco 7200

# IGW – Essential Feature set

**Broad LAN and WAN interfaces support**
- international links – POS STM-1/4/16/64
- national links – GE, 10GE, future full-rate 100GE

**IPv4 and IPv6 Routing and Forwarding**
- 2M hardware entries (IPv4 + IPv6) – no compression tricks!
- BGP, OSPF/ISIS, BFD – fast, prefix-independent convergence

**IPv4 and IPv6 filters (access-lists)**
- thousands of L3/L4 entries (IPv4 + IPv6) – no impact on forwarding rate!
- loose uRPF (Unicast Reverse Path Forwarding)

**IPv4 and IPv6 netflow monotoring**
- at least 1:1000 sampling rate, V9 export

**DDoS attack protection and Control Plane protection**
- in-hardware protection of router's brain
- anti-hacking tools – management plane protection

# IGW – some optional features

**MPLS support**
- rarely used on IGW, but sometimes yes
- MPLS Netflow is required too

**Traffic Shaping with RED – per-interface or per-VLAN**
- if the circuit runs over MAN or ISP subrate service
- shaping prevents unnecessary drops and improves TCP goodput

**Accounting**
- BGP Policy Accounting – per-AS accounting for large networks
- BGP Policy Propagation – packet marking based on BGP Communities
- MAC accounting – for peering/transit via IXP

**Secure Virtualization of the router**
- Logical Routers with secure resources allocation

**Carrier Grade NAT**
- IPv4 exhaustion is close!
- large scale IPv4 NAT and IPv6 AFT with V6 Tunneling is desirable

**LI (Lawful Intercept)**
- if used as a ISP Transit, LI may be mandatory

# ISP Security

# Anti-spoofing
## *RFC2827/BCP38 Ingress Packet Filtering*

Anti-spoofing filter (ingress filter on source IP)

allow only source addresses from the customer's 96.0.X.X/24
RFC2827 and RFC3704 (BCP 38 and 84)

Bogon filter (ingress filter on destination IP)

Drops packets with "insane" destination IP address
RFC1918, own block, internal IP core, NMS

96.0.20.0/24

96.0.21.0/24

Internet        ISP

ISP's Customer
Allocation Block:
96.0.0.0/19

96.0.19.0/24

96.0.18.0/24

Anti-spoofing Filter Applied
ingress on Downstream
Aggregation or NAS Routers

# uRPF (Unicast Reverse Path Forwarding)
## "Strict Mode" (v1) and "Loose Mode" (v2)

router(config-if)# ip verify unicast source reachable-via rx

i/f 2

i/f 1          i/f 3

S D data

FIB:
. . .
. . .
S -> i/f 1
D -> i/f 3
. . .

Same i/f:
FORWARD

i/f 2

i/f 1          i/f 3

S D data

FIB:
. . .
. . .
S -> i/f 2
D -> i/f 3
. . .

Other i/f:
DROP

**"Strict Mode"**
(aka "v1")

router(config-if)# ip verify unicast source reachable-via any

i/f 2

i/f 1          i/f 3

S D data

FIB:
. . .
. . .
S -> i/f x
D -> i/f 3
. . .

Any i/f:
FORWARD

i/f 2

i/f 1          i/f 3

S D data

FIB:
. . .   ?
. . .
D -> i/f 3
. . .

Src not in FIB
or route = null0:
DROP

**"Loose Mode"**
(aka "v2")

# Bogons

- A Bogon prefix is a route that should never appear in the Internet routing table

- Different from DSUA.

    Bogons are defined as Martians (private and reserved addresses defined by RFC 1918 and RFC 3330) and netblocks that have not been allocated to a (RIR) by IANA

- CYMRU maintains list of Bogons, works with IANA and RIR etc.

- http://www.cymru.com/Bogons/index.html

- BOGON List Keeps on Changing as IANA allocates routes.

    BE AWARE!

    The bogon prefixes are announced unaggregated by the bogon route-servers is **65333:888**; as of 14 JUL 2008 this includes **45** prefixes

- BOGON Router Server.

    Peer with CYMRU Route Server keep BOGON list upto date.

# Hardware protection against DOS attacks CRS-1 Control Plane Protection



RP

CPU
Input processes

CoPP

CSAR queue

Ingress LC

CPU

raw queues

To RP queue

ASI

4: Multiple queues to LC and RP CPU

3: LPTS in iFIB police traffic

2b: Skip LC CPU!

2a: LPTS iFIB lookup (Match, BTSH/GTSM)

1: Ingress iACL, uRPF

# IOS XR – Dynamic Control Plane Protection

```
Router bgp
  neighbor 202.4.48.99
  ··ttl_security
!


  mpls ldp
  …
!
```

## LC 1 PreIFIB TCAM HW Entries

| Local | port | Remote | port | Rate | Priority |
|---|---|---|---|---|---|
| Any | ICMP | ANY | ANY | 1000 | low |
| any | 179 | any | any | 100 | medium |
| any | 179 | 202.4.48.99 | any | 1000 | medium |
| 202.4.48.1 | 179 | 202.4.48.99 | 2223 | 10000 | medium |
| 200.200.0.2 | 13232 | 200.200.0.1 | 646 | 100 | medium |

ttl 255

## LC 2 PreIFIB TCAM HW Entries …

LPTS

Socket

bgp

ldp

**TCP Handshake**

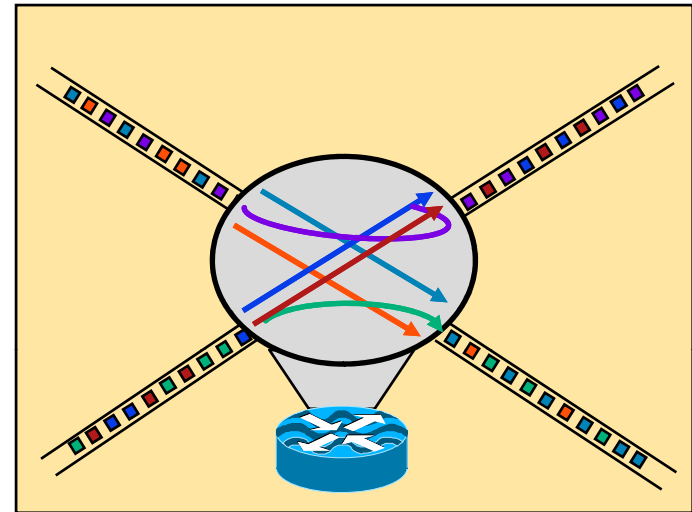# Detecting an attack:

# Netflow

# Netflow is a Security tool #1 today!

## 7 Keys define a flow

Source Address, Destination Address, Source Port, Destination Port, Layer 3 Protocol Type, TOS byte (DSCP), Input Logical Interface (ifIndex)

## A flow is unidirectional



Turning it on (generic):
```
interface GigabitEthernet 1/1/1
   ip route-cache flow [sampled]
```
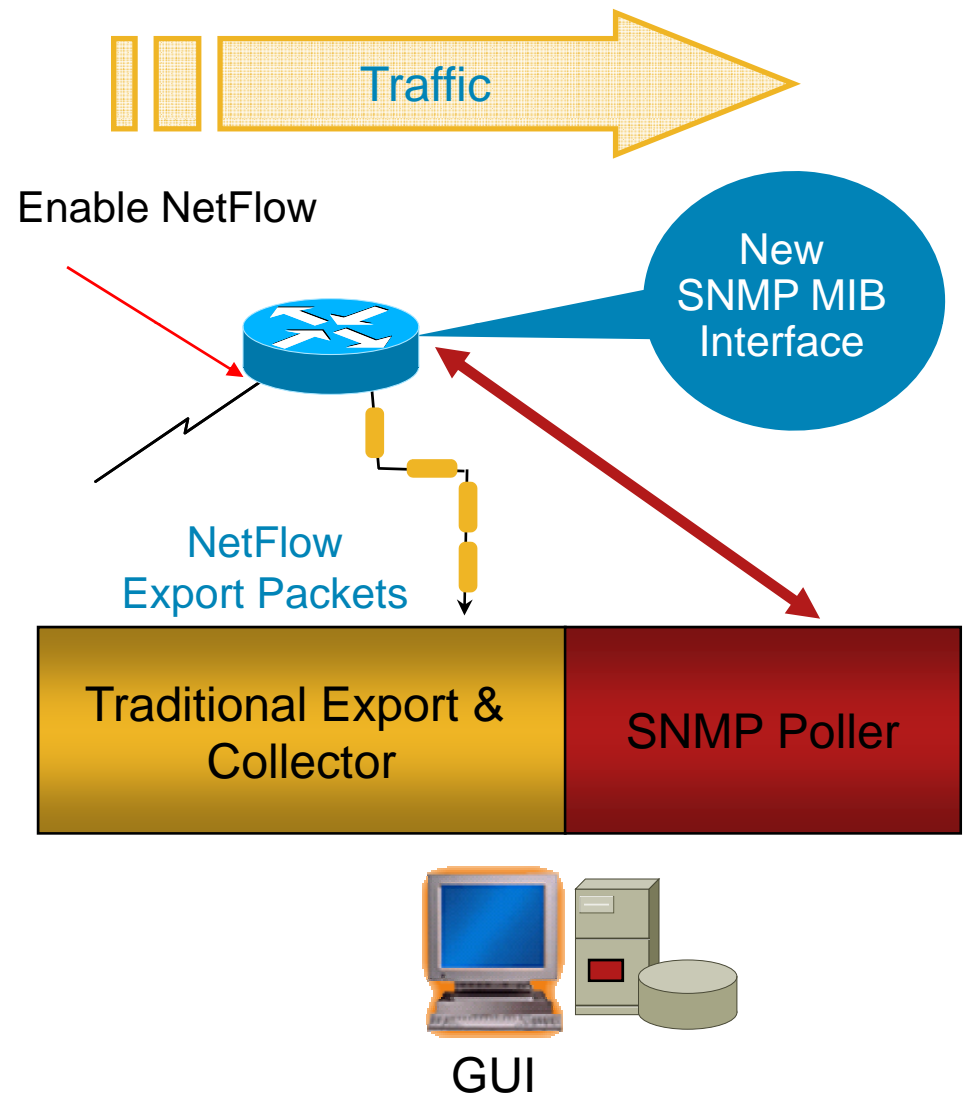
Export (optional):
```
ip flow-export destination 172.17.246.225 9995
```

Sampled Netflow (mostly used for Security):
```
ip flow-sampling-mode packet-interval x
```

# Flow Is Defined By Seven Unique Keys

- Source IP address
- Destination IP address
- Source port
- Destination port
- Layer 3 protocol type
- TOS byte (DSCP)
- Input logical interface (ifIndex)

Traffic

Enable NetFlow

New SNMP MIB Interface

NetFlow Export Packets

Traditional Export & Collector

SNMP Poller

GUI

# NetFlow Cache Example

1. Create and update flows in NetFlow cache

| SrcIf | SrcIPadd | DstIf | DstIPadd | Protocol | TOS | Flgs | Pkts | Src Port | Src Msk | Src AS | Dst Port | Dst Msk | Dst AS | NextHop | Bytes/Pkt | Active | Idle |
|-------|----------|-------|----------|----------|-----|------|------|----------|---------|--------|----------|---------|--------|---------|-----------|--------|------|
| Fa1/0 | 173.100.21.2 | Fa0/0 | 10.0.227.12 | 11 | 80 | 10 | 11000 | 00A2 | /24 | 5 | 00A2 | /24 | 15 | 10.0.23.2 | 1528 | 1745 | 4 |
| Fa1/0 | 173.100.3.2 | Fa0/0 | 10.0.227.12 | 6 | 40 | 0 | 2491 | 15 | /26 | 196 | 15 | /24 | 15 | 10.0.23.2 | 740 | 41.5 | 1 |
| Fa1/0 | 173.100.20.2 | Fa0/0 | 10.0.227.12 | 11 | 80 | 10 | 10000 | 00A1 | /24 | 180 | 00A1 | /24 | 15 | 10.0.23.2 | 1428 | 1145.5 | 3 |
| Fa1/0 | 173.100.6.2 | Fa0/0 | 10.0.227.12 | 6 | 40 | 0 | 2210 | 19 | /30 | 180 | 19 | /24 | 15 | 10.0.23.2 | 1040 | 24.5 | 14 |

2. Expiration

- Inactive timer expired (15 sec is default)
- Active timer expired (30 min (1800 sec) is default)
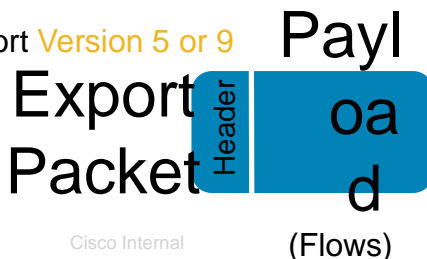- NetFlow cache is full (oldest flows are expired)
- RST or FIN TCP Flag

| SrcIf | SrcIPadd | DstIf | DstIPadd | Protocol | TOS | Flgs | Pkts | Src Port | Src Msk | Src AS | Dst Port | Dst Msk | Dst AS | NextHop | Bytes/Pkt | Active | Idle |
|-------|----------|-------|----------|----------|-----|------|------|----------|---------|--------|----------|---------|--------|---------|-----------|--------|------|
| Fa1/0 | 173.100.21.2 | Fa0/0 | 10.0.227.12 | 11 | 80 | 10 | 11000 | 00A2 | /24 | 5 | 00A2 | /24 | 15 | 10.0.23.2 | 1528 | 1800 | 4 |

No

Yes

3. Aggregation

4. Export version

Non-Aggregated Flows—Export Version 5 or 9
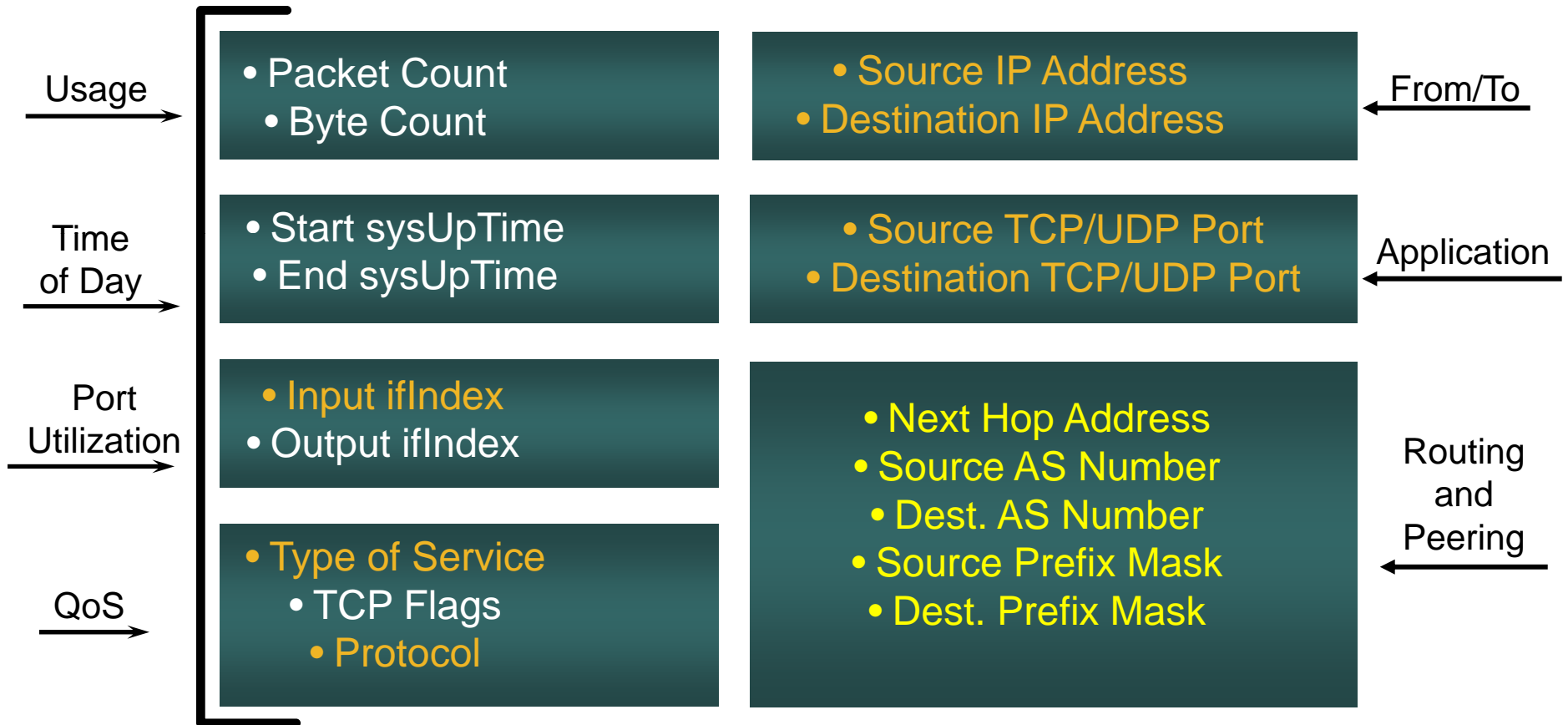
5. Transport protocol

Export Packet
Header | Payload
(Flows)

e.g. Protocol-Port Aggregation Scheme Becomes

| Protocol | Pkts | SrcPort | DstPort | Bytes/Pkt |
|----------|------|---------|---------|-----------|
| 11 | 11000 | 00A2 | 00A2 | 1528 |

Aggregated Flows—Export Version 8 or 9

# Netlow Export – V5 fixed format

Usage →

Time of Day →

Port Utilization →

QoS →

- Packet Count
  - Byte Count

- Start sysUpTime
- End sysUpTime

- Input ifIndex
- Output ifIndex

- Type of Service
  - TCP Flags
    - Protocol

- Source IP Address
- Destination IP Address

← From/To

- Source TCP/UDP Port
- Destination TCP/UDP Port

← Application

- Next Hop Address
- Source AS Number
- Dest. AS Number
- Source Prefix Mask
- Dest. Prefix Mask

← Routing and Peering

## Version 5 used extensively today

# NetFlow Export – V9 flexible format

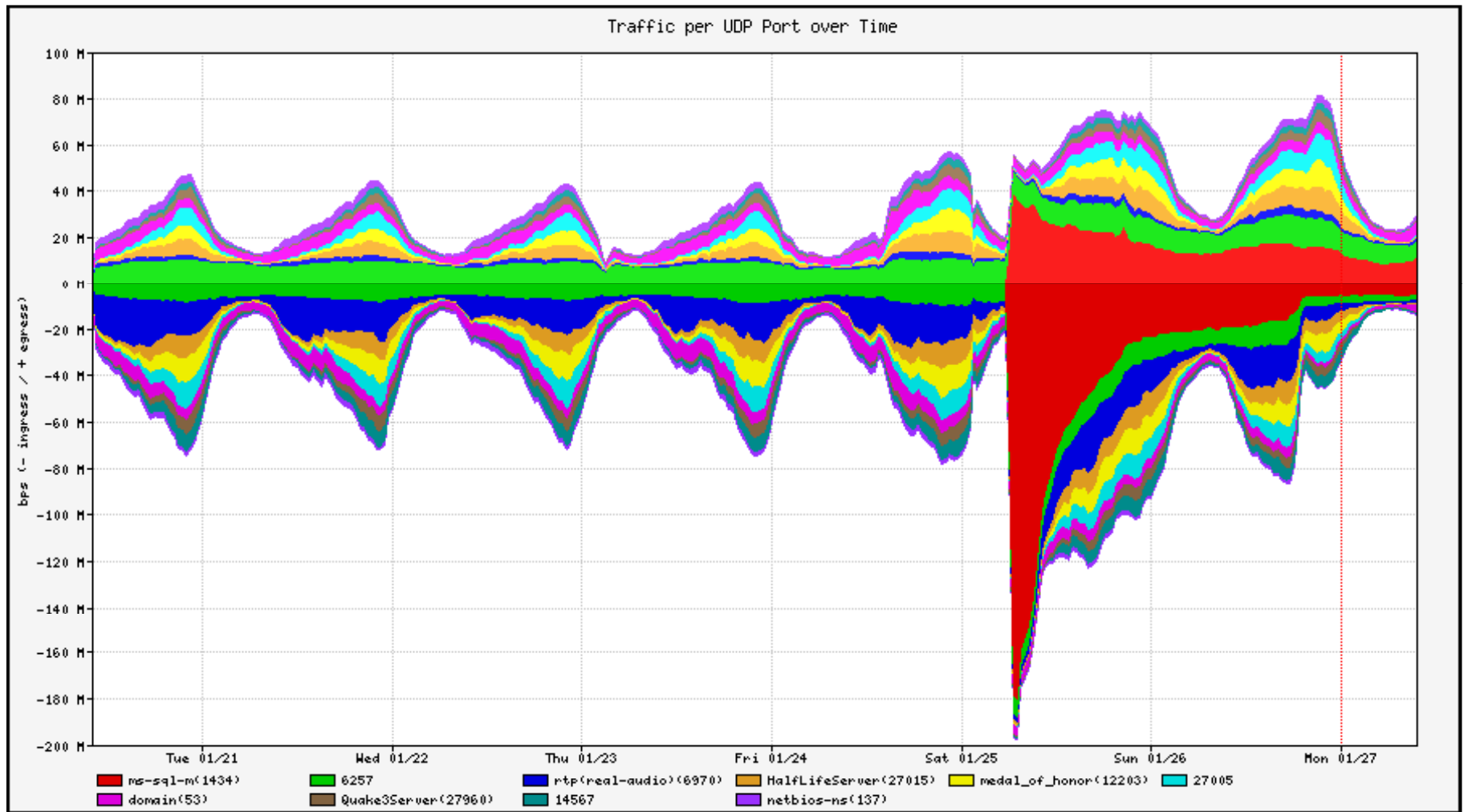Example of Export Packet right after router boot or NetFlow configuration

| (version, # packets, sequence #, Source ID) | Template FlowSet | | | | Option Template FlowSet | Option Data FlowSet |
|---|---|---|---|---|---|---|

Template FlowSet

Template Record **Template ID** (specific Field types and lengths)

Template Record **Template ID** (specific Field types and lengths)

Template Record **Template ID** (specific Field types and lengths)

Template Record **Template ID** (specific Field types and lengths)

Option Template FlowSet **Template ID** (specific Field types and lengths)

Option Data FlowSet **FlowSet ID**

Option Data Record (Field values)

Option Data Record (Field values)

Example of Export Packets containing mostly flow information

Header (version, # packets, sequence #, Source ID)

Data FlowSet **FlowSet ID**

Data Record (Field values)

Data Record (Field values)

Data Record (Field values)

Data Record (Field values)

Data Record (Field values)

Data Record (Field values)

Data FlowSet **FlowSet ID**

Data Record (Field values)

# Example—What is an Anomaly?

# NetFlow—nfdump and nfsen



Source: http://nfsen.sourceforge.net, ev. http://software.uninett.no/stager/

# Arbor Peakflow SP — Application Distribution

# Example—Arbor Peakflow SP DoS Module

# BGP Next Hop TOS Aggregation

## Typical Example
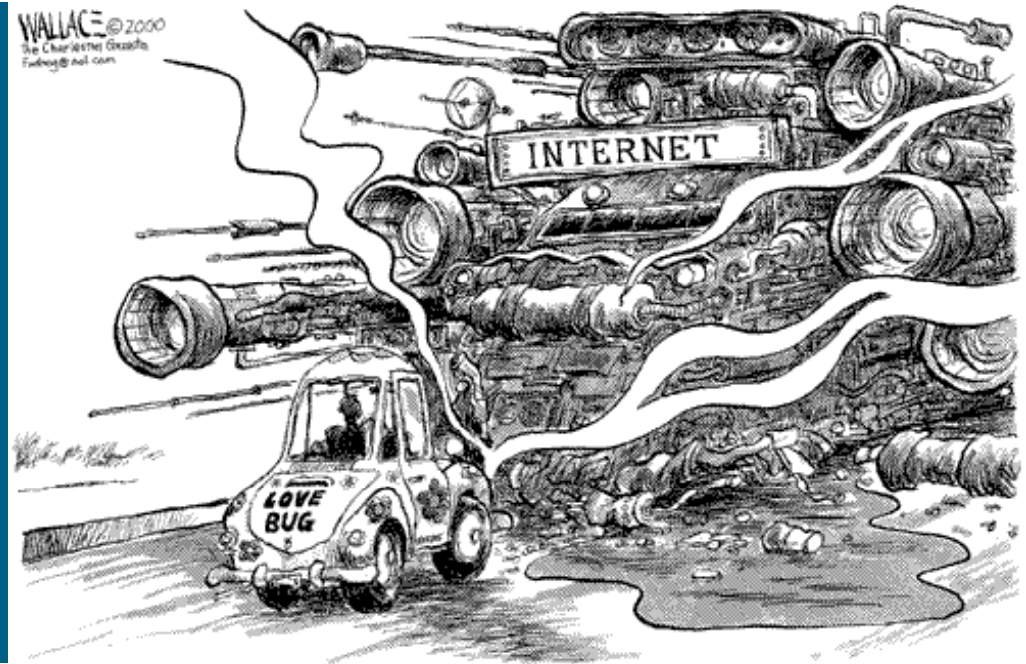


Internal Traffic: "PoP to PoP"
External Traffic Matrix PoP to BGP AS

# Dropping a DDoS attack:

# BGP Blackholing

# Customer is DOSed
*Before*



Peer A

Peer B

IXP-W

IXP-E

A

B

C

D

E

Upstream A

Upstream B

Upstream A

Upstream B

Target

F   POP

NOC

G

Target is *taken out*

# Customer is DOSed
## *Before – Co-Lateral Damage*



IXP-W

Peer A

Peer B

IXP-E

A

Upstream A

D

Upstream A

B

C

Upstream B

Upstream B

E

Target

Customers

F    POP

Attack causes
Co-Lateral
Damage

G    NOC

48

# Customer is DOSed
## *After – Packet Drops Pushed to the Edge*



Peer A

Peer B

IXP-W

IXP-E

A

B

C

D

E

Upstream A

Upstream A

Upstream B

Upstream B

Target

F   POP

G   NOC

iBGP Advertises List of Black Holed Prefixes

# BGP Blackholing: Reacting to an Attack

BGP Sent – 171.68.1.0/24 Next-Hop = 192.0.2.1

Static Route in Edge Router – 192.0.2.1 = Null0

171.68.1.0/24 = 192.0.2.1 = Null0

Next hop of 171.68.1.0/24 is now        equal to Null0

- Remote Triggered Black Hole filtering is the foundation for a whole series of techniques to traceback and react to DDOS attacks on an ISP's network.

- Easy preparation, does not effect ISP operations or performance.

- It does adds the option to an ISP's *security toolkit.*

# BGP Blackholing: IOS configuration

- place a host-route to Null on <u>every BGP router</u>

```
ip route 192.0.2.1 255.255.255.255 Null0
```

- prepare a injection into BGP with the blackhole next-hop

```
router bgp 10
  redistribute static route-map set-blackhole

route-map set-blackhole permit 10
 match tag 666
 set ip next-hop 192.0.2.1
 set community 10:666 no-export
 set local-preference 50
```

- simply filter it out everywhere by one command:

```
BH(config)# ip route 1.2.2.2 255.255.255.255 Null0 tag 666
```

# BGP Blackholing: Filtering on source IP address

- loose uRPF (unicast reverse path forwarding)

```
ip route 192.0.2.2 255.255.255.255 Null0
int PoS 1/0/0
   ip verify unicast source reachable-via any
```

  !!! packet with source IP prefix pointing to Null0 will be dropped !!!

- prepare a injection into BGP with the blackhole next-hop

```
route-map set-blackhole permit 20
 match tag 667
 set ip next-hop 192.0.2.2
 set community 10:667 no-export
 set local-preference 50
```

- simply filter it out everywhere by one command:

```
BH(config)# ip route 1.2.2.3 255.255.255.255 Null0 tag 667
```
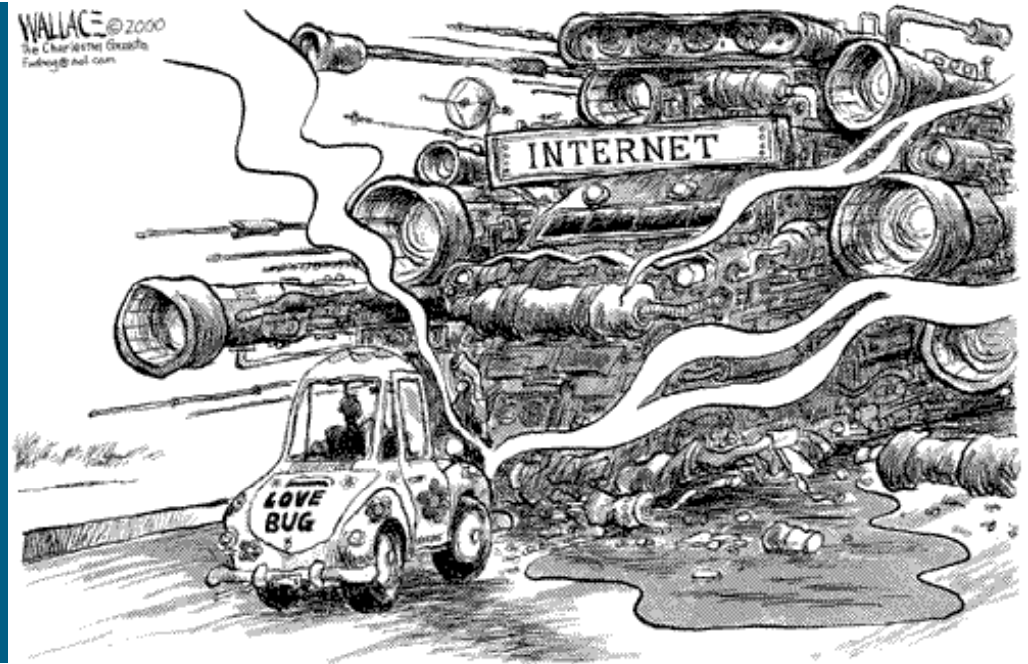
# BGP Triggered Rate Limiting
## QPPB (QoS Policy Propagation via BGP)

```
router bgp 10
  table-map DOS-Activate
  neighbor 200.200.14.4 remote-as 10
  neighbor 200.200.14.4 update-source Loopback 0

  neighbor 200.200.14.4 send-community
!
ip bgp-community new-format
!
ip community-list 1 permit 10:666
!
route-map DOS-Activate permit 10
  match community 1
  set ip qos-group 66
!
route-map DOS-Activate permit 20
!
interface PoS 0/0/0

  bgp-policy source ip-qos-map

  rate-limit input qos-group 66 256000 8000 8000
      conform-action transmit
      exceed-action drop
```

- **QPPB marking is done before rate-limit or policing**

- **hardware support in Cisco 10000, 12000, CRS-1**

# Dark IP space:

# Sinkholes

# Default Route & the Internet

```
BHole(config-router)# default-information originate always
```

- Advertising Default from the Sink Hole will pull down all sort of *junk* traffic.

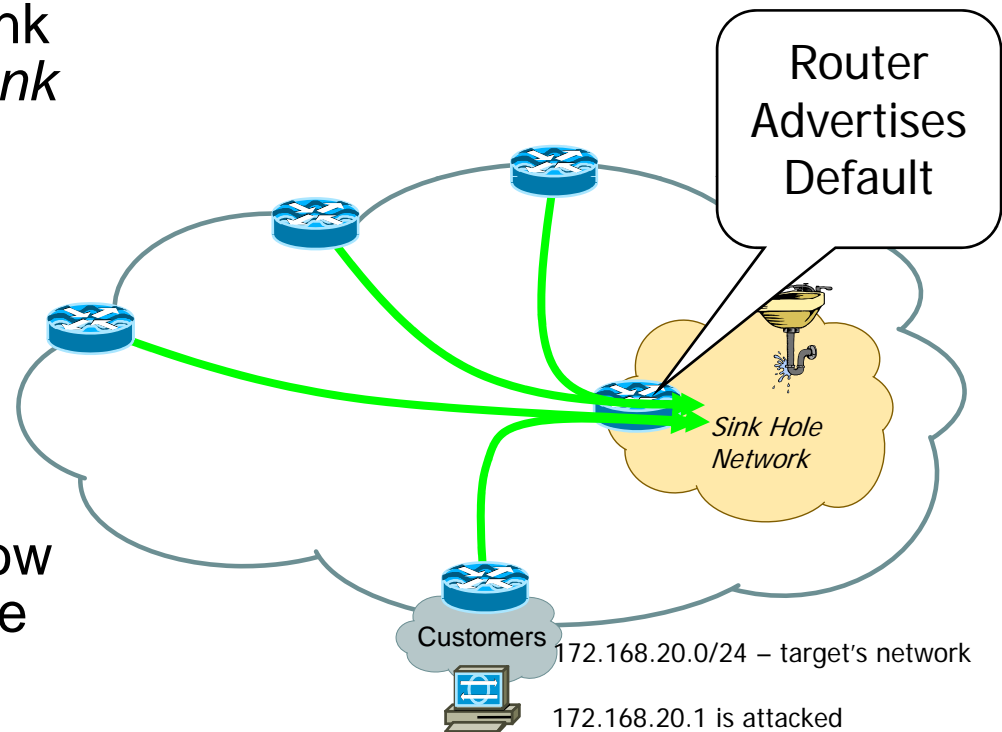  - Customer Traffic when circuits flap.
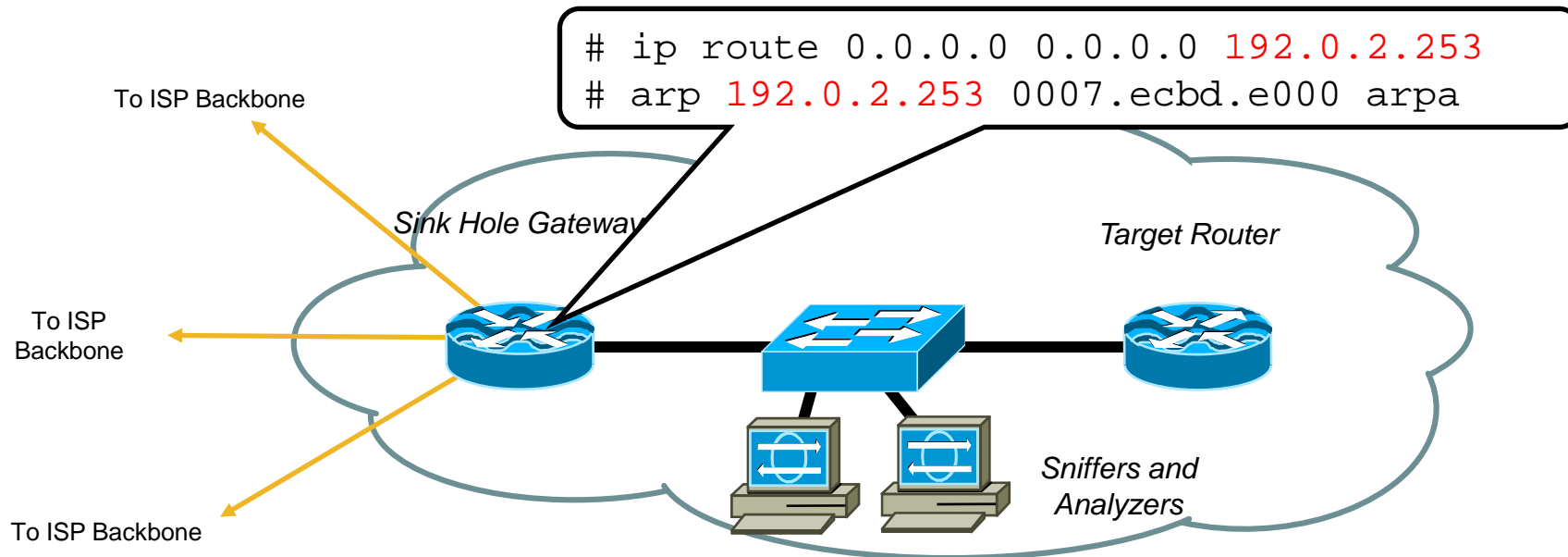
  - Network Scans

  - Failed Attacks

  - Code Red/NIMDA

  - Backscatter

- Can place tracking tools (Netflow cache) and IDS in the Sink Hole network to monitor the noise.

- BCP: Default should be always a blackhole (Null0 or Static ARP) !!

Router Advertises Default

*Sink Hole Network*

Customers

172.168.20.0/24 – target's network

172.168.20.1 is attacked

# Target Routers are Expendable

```
# ip route 0.0.0.0 0.0.0.0 192.0.2.253
# arp 192.0.2.253 0007.ecbd.e000 arpa
```

To ISP Backbone

To ISP Backbone

To ISP Backbone

*Sink Hole Gateway*

*Target Router*

*Sniffers and Analyzers*

- Sink Hole Gateway Generates the more specific iBGP Announcement.

- Pull the DOS/DDOS attack to the sink hole and forwards the attack to the target router.

- Static ARP to the target router keeps the Sink Hole Operational – Target Router can crash from the attack and the static ARP will keep the gateway forwarding traffic to the ethernet switch.

# What to Monitor in a Sinkhole?

- Scans on dark IP (allocated and announced but unassigned address space)

    Who is scoping out the network—pre-attack planning, worms…

- Scans on bogons (unallocated)

    Worms, infected machines, and Bot creation

- Backscatter from spoofed attacks

    Who is getting attacked

    > don't use "no ip icmp unreachables"
    >
    > use "ip icmp rate-limit unreachables"

- Backscatter from garbage traffic (RFC-1918 leaks)

    Which customers have mis-configuration or "leaking" networks

# Summary & Resources

# Summary

- **Transit vs. Peering**

- **The importance of IXP**


- **Anatomy of the ISP Edge**

- **Cisco peering platforms and features**


- **The importance of Netflow**

- **Basic ISP cecurity techniques**

# Cisco Networkers
# 25-28. januar 2010.
# Barselona
# Registrujte se

61