# Carrier-Grade NAT

Josef Ungerman

Consulting Systems Engineer, Cisco, CCIE#6167

Email, Jabber, TP: josef.ungerman@cisco.com

# Agenda

CGN (Carrier-Grade NAT)

Definition and purpose

NAT vs. Firewall

Mapping/filtering, app traversal

Other CGN Behaviors

Pooling, port limits, etc.

Session Logging

Netflow/Syslog, formats, variations

CGN Design

Performance, placement

# Network Address & Port Translation
## Most of Broadband users are behind NAT today!

*When say "NAT", they typically mean "NAPT"*

- NAT

  First described in 1991 (draft-tsuchiya-addrtrans), RFC1631

  1:1 translation: Does not conserve IPv4 addresses

  Per-flow stateless

  Today's primary use is inside of enterprise networks

  - Connect overlapping RFC1918 address space

  Note: NAT66 is stateful or stateless, but it is not NAPT

- NAPT

  Described in 2001 (RFC3022)

  1:N translation

  Conserves IPv4 addresses

  Allows multiple hosts to share one IPv4 address

  Only TCP, UDP, and ICMP

  Connection has to be initiated from 'inside'

  Per-flow stateful

  Commonly used in home gateways and enterprise NAT

*"NAT44" is used to differentiate IPv4-IPv4 NAPT from Address Family Translation, typically referred to as NAT64 and NAT46"*
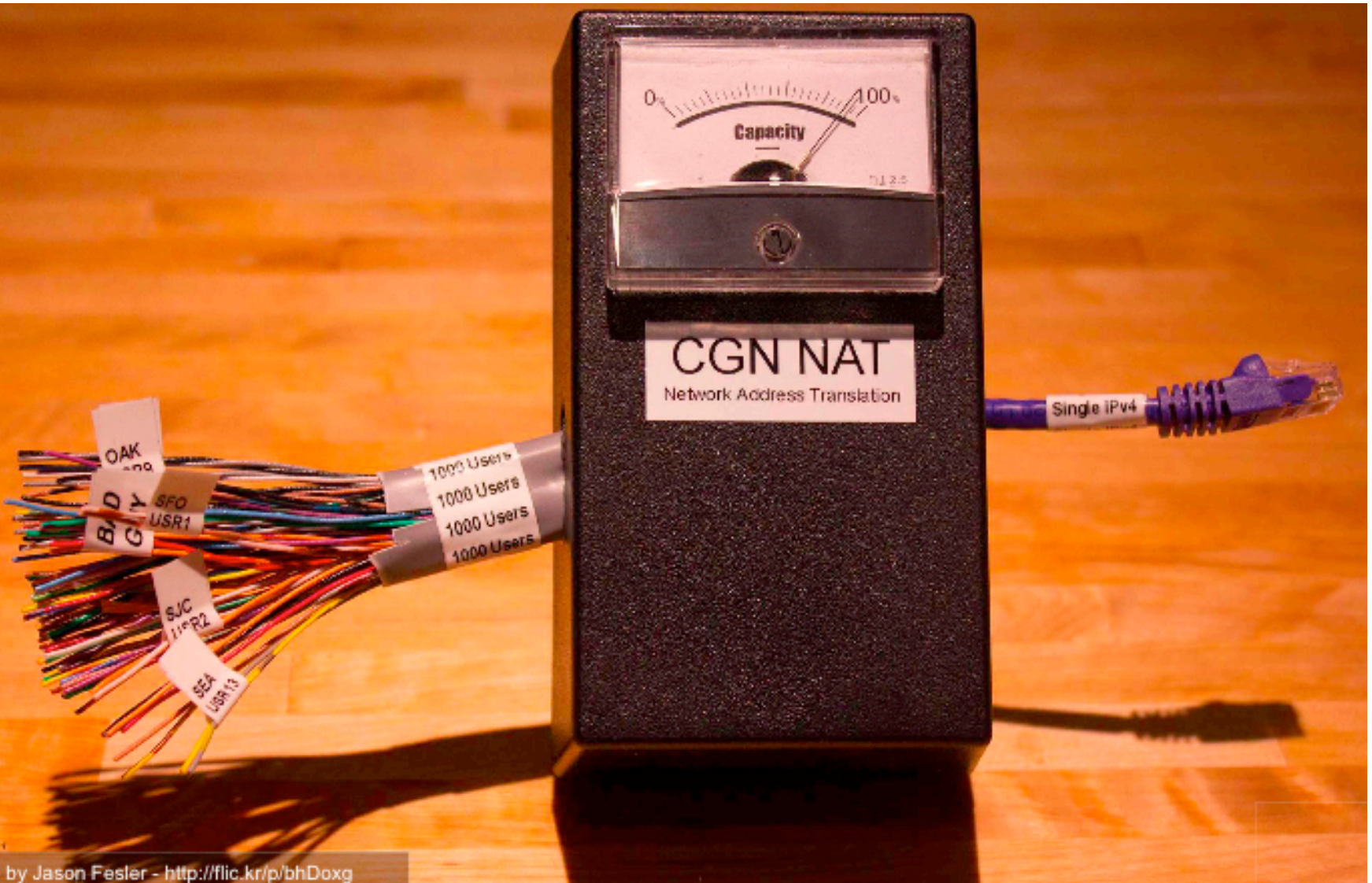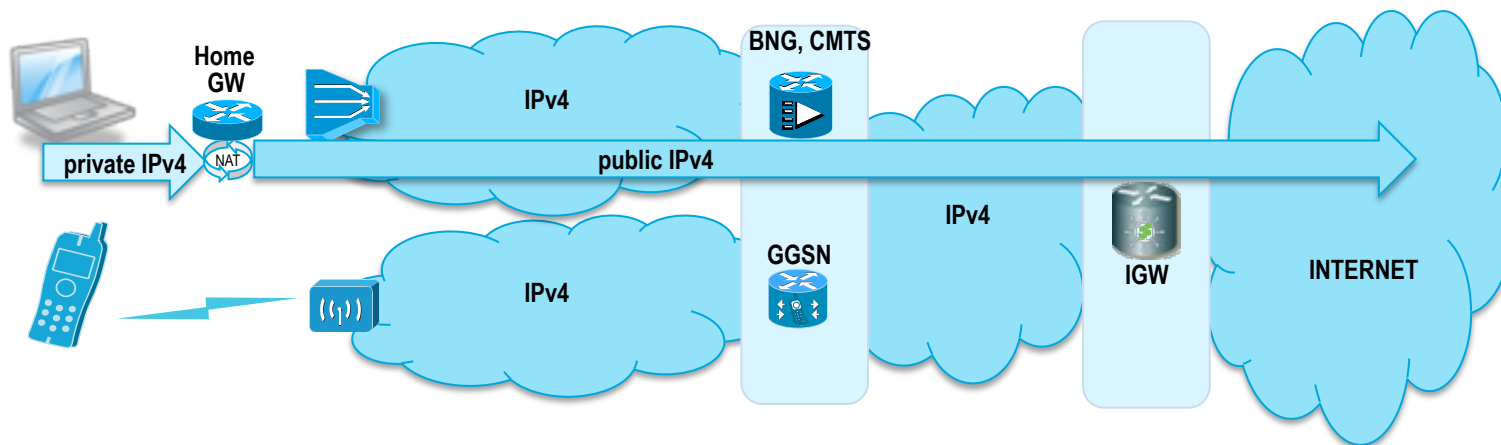
# What is CGN?



Courtesy of Jason Fesler, Yahoo (V6 World Congress 2012)
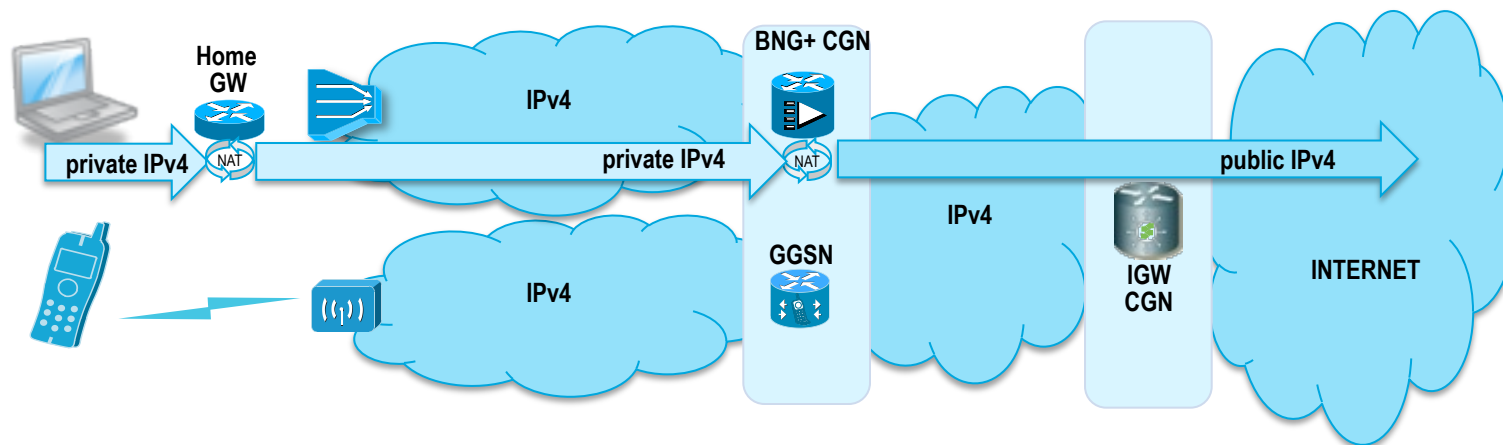
# NAT in Internet Access
*typical deployment today (wireline)*

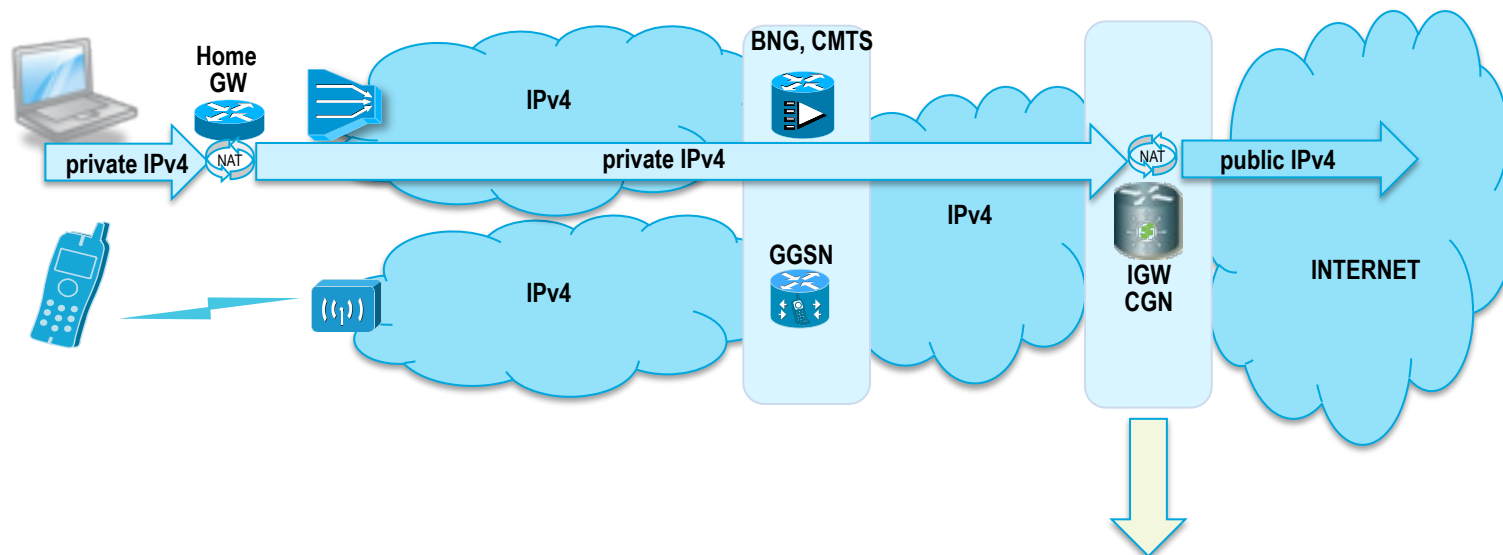# NAT in Internet Access
## *CGN – NAT444 (wireline)*

RFC1918 or RFC6598 (100.64.0.0/10)

# NAT in Internet Access
## *CGN – NAT444 (wireline)*

RFC1918 or RFC6598 (100.64.0.0/10)



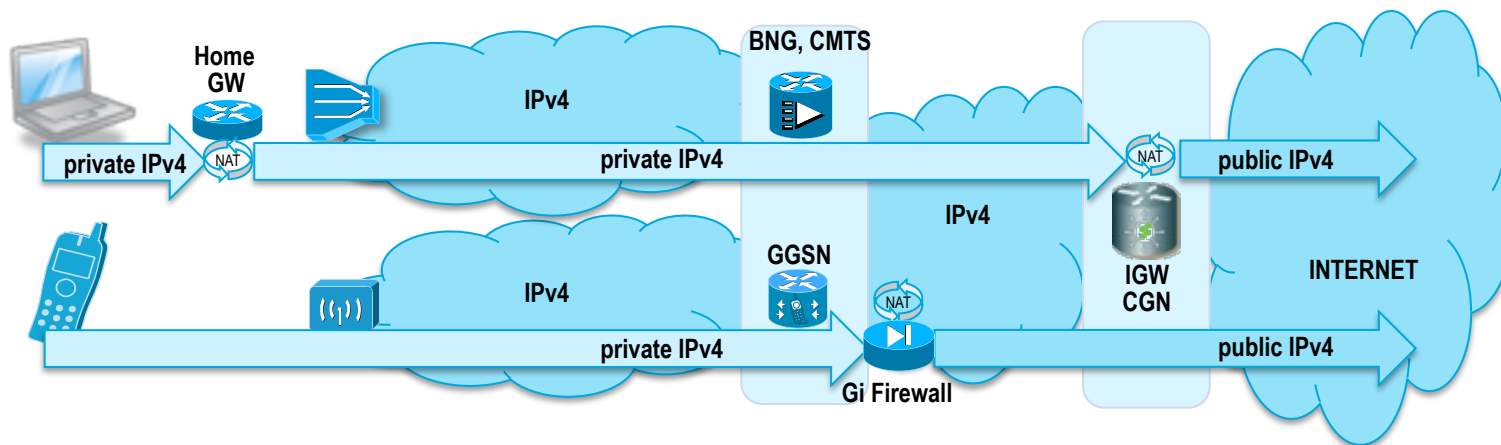**Large Scale**
- **100G+ throughput**
- **100M+ concurrent bidirectional sessions**
- **1M+ sessions per second setup rate**

# NAT in Internet Access
## *Mobile*

# NAT in Internet Access
## *Mobile*



Home GW

BNG, CMTS

private IPv4

private IPv4

IPv4

public IPv4

NAT

IPv4

GGSN

IGW CGN

INTERNET

private IPv4

public IPv4

NAT

Gi Firewall

**Virtualization Support (VRF)**

# NAT in Internet Access
## *IPv6 ultimately bypasses CGN*



Home GW

BNG, CMTS

private IPv4

private IPv4

public IPv4

IPv4

IPv6 bypass

GGSN

IGW CGN

INTERNET

IPv4

private IPv4

public IPv4

Gi Firewall

NAT

**Dual-Stack**
**6in4 tunneling – 6rd BR**
**4in6 tunneling – DS-Lite AFTR, MAP**
**NAT64 for v6-only hosts**

# CGN = IP Address Sharing
## IPv4 Exhaust workaround: move part of Internet from L3 to L4

- Adds capital and operations cost [$/Gbps] → decreases ARPU

  CGN is here not to stay, but to be replaced (by IPv6)

- Inherent issues (general NAPT issues, mostly not NAT444 related)

  draft-ford-shared-addressing-issues
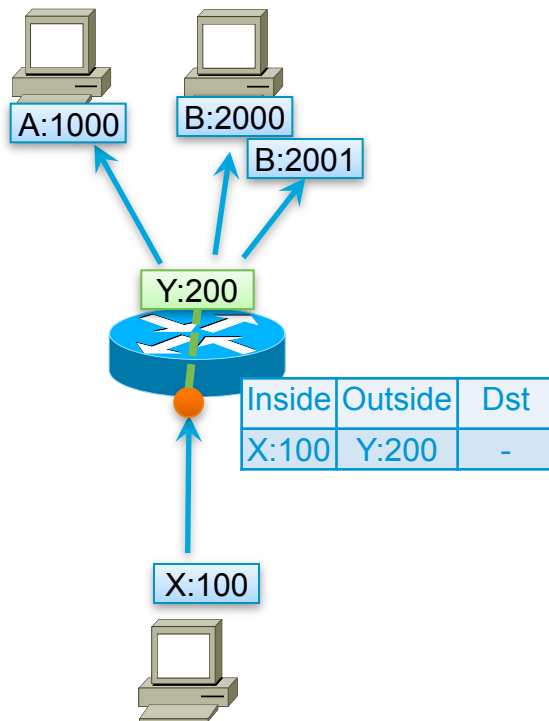
- Servers must log also source port numbers

  Shared IP address = shared suffering (blacklisting, spam,…)

  Tracking and Law Enforcement, draft-ietf-intarea-server-logging-recommendations

  Otherwise CGN must do log also Destination IP:port (privacy issue!)

- Requesting specific ports – "Not everyone can get port 80"

- Geo-Location issues ("get me the nearest ATM")

- Complicates inbound access to media

- Keepalives → power consumption, mobile battery drain

# NAPT internals: Mapping

**Endpoint Independent**   Address Dependent   Address and port Dependent

A:1000   B:2000   B:2001

| Inside | Outside | Dst |
|--------|---------|-----|
| X:100  | Y:200   | -   |

Y:200

X:100

A:1000   B:2000   B:2001

Y:200   Y:300

| Inside | Outside | Dst    |
|--------|---------|--------|
| X:100  | Y:200   | A:any  |
| X:100  | Y:300   | B:any  |

X:100

A:1000   B:2000   B:2001

Y:200   Y:300   Y:400

| Inside | Outside | Dst    |
|--------|---------|--------|
| X:100  | Y:200   | A:1000 |
| X:100  | Y:300   | B:2000 |
| X:100  | Y:400   | B:2001 |

X:100

IP Addres: Port Number

# NAPT internals: Filtering

**Endpoint Independent**

| Inside | Outside | from |
|--------|---------|------|
| X:100  | Y:200   | -    |

Address Dependent

| Inside | Outside | from |
|--------|---------|------|
| X:100  | Y:200   | A    |

Address and Port Dependent

| Inside | Outside | from   |
|--------|---------|--------|
| X:100  | Y:200   | A:1000 |

A:1000  A:1001  B:2000  Y:200  X:100

IP Addres: Port Number

# NAT mapping/filtering behavior

| | behavior | Filtering | | |
|---|---|---|---|---|
| | | Independent | Address Dependent | Address:Port Dependent |
| **Mapping** | Independent | Full Cone NAT | Address Restricted NAT | Port Restricted NAT |
| | Address Dependent | Symmetric NAT | | |
| | Address:Port Dependent | | | |

Restricted

Internet NAT (CGN)

Home Gateway NAT

Firewall NAT

STUN NAT Types
- Classic STUN : simple traversal of UDP through NAT(RFC3489)
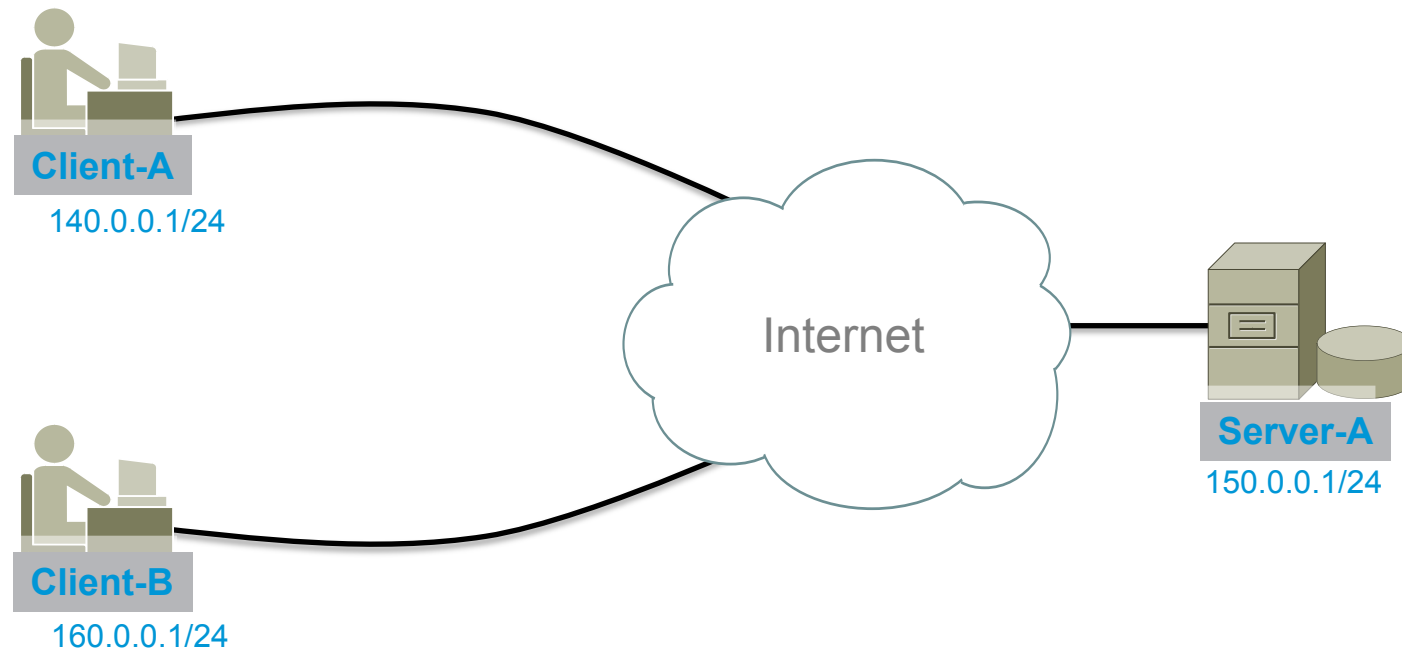- Now : Session Traversal Utilities for NAT(RFC5389)

# NAT != Firewall
## Application transparency behaviors

**Client-A**

140.0.0.1/24

**Client-B**

160.0.0.1/24

Internet

**Server-A**

150.0.0.1/24

# What is Symmetric NAT
## Address:port restricted NAT = Firewalling behavior

NAT/PAT Entry:

| Inside local | Inside global | Outside local | Outside global |
|---|---|---|---|
| 192.168.1.1 :5000 | 140.0.0.1 :6000 | 150.0.0.1 :6000 | 150.0.0.1 :6000 |

**Client-A**

192.168.1.1/24

**Firewall+NAT**

NAT Pool
140.0.0.0/24

Internet

**Server-A**

150.0.0.1/24

**Client-B**

192.168.1.1/24

**Firewall+NAT**

NAT Pool
160.0.0.0/24

## Endpoint Dependent Mapping, Endpoint Dependent Filtering

# What is Full cone NAT
## Pure NAT with no Firewalling behavior

NAT/PAT Entry:

| Inside local | Inside global | Outside local | Outside global |
|---|---|---|---|
| 192.168.1.1 :5000 | 140.0.0.1 :6000 | any | any |

**Client-A**

192.168.1.1/24

**NAT**

NAT Pool
140.0.0.0/24

Internet

**Server-A**

150.0.0.1/24

**Client-B**

192.168.1.1/24

**NAT**

NAT Pool
160.0.0.0/24

**Endpoint Independent Mapping, Endpoint Independent Filtering – EIM/EIF**

# NAT Traversal with Full cone NAT
## Peer-to-peer Applications Transparency

Client Identities Table

| Client-A | 192.168.1.1 :5000 | 140.0.0.1 :6000 |
|----------|-------------------|-----------------|
| Client-B | 192.168.1.1 :5000 | 160.0.0.1 :6000 |

**Client-A**

192.168.1.1/24

**NAT**

NAT Pool
140.0.0.0/24

Internet

**Rendezvous Server
(aka. Super-Node, Hub,…)**

150.0.0.1/24

**Client-B**

192.168.1.1/24

**NAT**

NAT Pool
160.0.0.0/24

# NAT Traversal with Full cone NAT
## Connectivity Probes and the Shortest Path



One of the PROBE replies comes back first
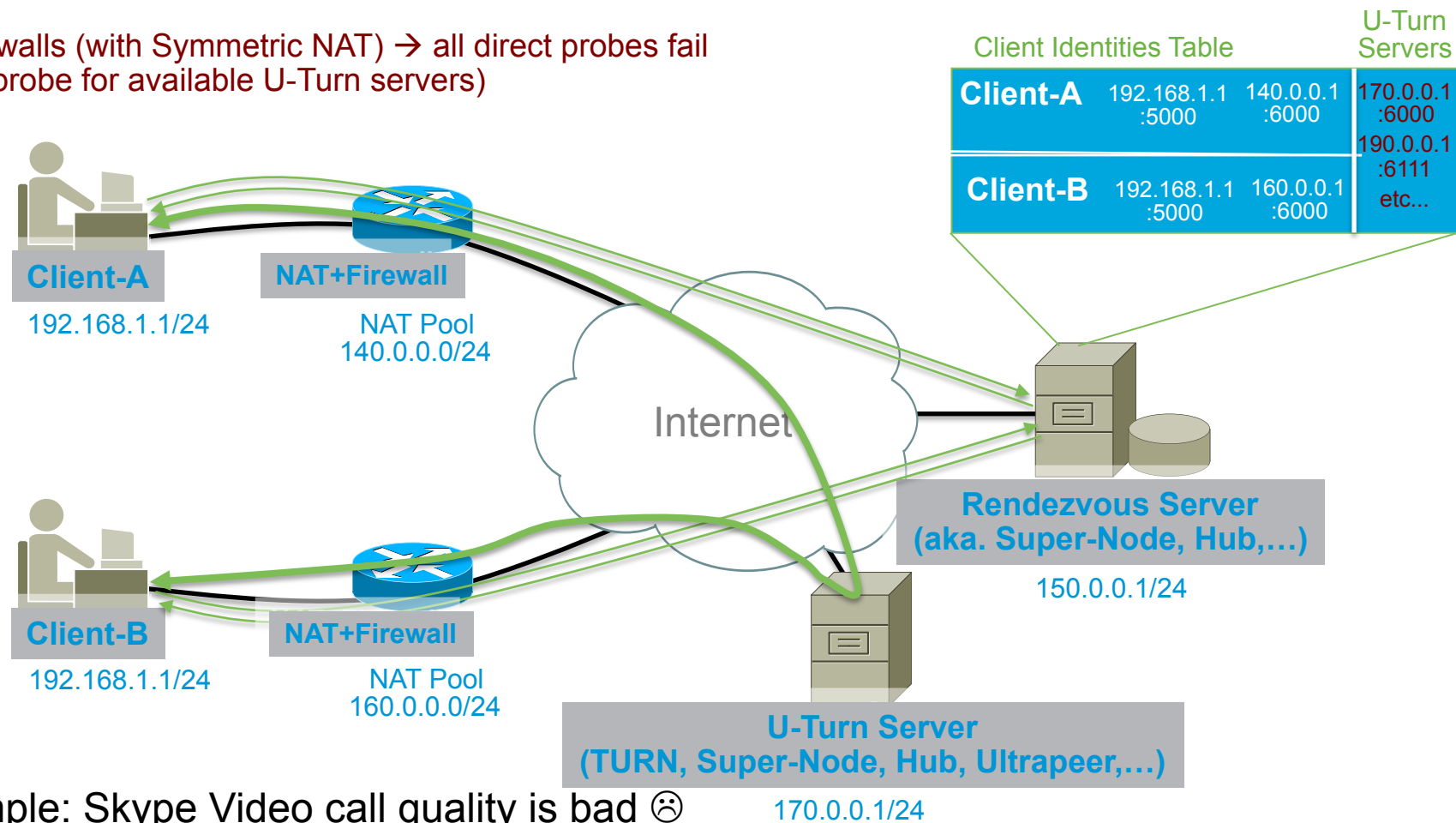(probably on-net one → fastest path)

Client-A
192.168.1.1/24

PROBE-1
192.168.2.1

NAT
NAT Pool
140.0.0.0/24

PROBE-2
160.0.0.1

Local Network

Internet

Client-B
192.168.2.1/24

NAT
NAT Pool
160.0.0.0/24

Client Identities Table

| Client-A | 192.168.1.1 :5000 | 140.0.0.1 :6000 |
|----------|-------------------|-----------------|
| Client-B | 192.168.2.1 :5000 | 160.0.0.1 :6000 |

Rendezvous Server
(aka. Super-Node, Hub,…)
150.0.0.1/24

Example: Skype Video call quality is good ☺

# NAT Traversal with Firewalls
## U-Turn server (Super-Node)

Firewalls (with Symmetric NAT) → all direct probes fail
(→ probe for available U-Turn servers)

**Client Identities Table**

**U-Turn Servers**

| | | | |
|---|---|---|---|
| **Client-A** | 192.168.1.1 :5000 | 140.0.0.1 :6000 | 170.0.0.1 :6000 |
| | | | 190.0.0.1 :6111 |
| **Client-B** | 192.168.1.1 :5000 | 160.0.0.1 :6000 | etc... |

**Client-A**

192.168.1.1/24

**NAT+Firewall**

NAT Pool
140.0.0.0/24

Internet

**Rendezvous Server
(aka. Super-Node, Hub,…)**

150.0.0.1/24

**Client-B**

192.168.1.1/24

**NAT+Firewall**

NAT Pool
160.0.0.0/24

**U-Turn Server
(TURN, Super-Node, Hub, Ultrapeer,…)**

Example: Skype Video call quality is bad ☹         170.0.0.1/24

# Avoiding U-Turn Servers with Firewalls
## Opening port holes in Firewall/NAT

- Application talks to the Firewall

  UPnP-IGD (app to home router)

  Common in Home Environment

  Not applicable in Internet (scale, no suitable protocols defined)

- ALG's (Application Level Gateway) aka. Fixup

  App proxy in the firewall – we jump from L4 to L6+

  Common in Enterprise (closed environment – we know the list of supported apps)

  Hardly applicable in Internet (open environment – we don't know all the apps)

  How about encrypted and integrity-protected protocols?

  How about URL with literals https://1.2.3.4

# Summary: CGN behavior best practices
## draft-behave-lsn-requirements, RFC4787, RFC6145

- Internet way: keep it simple

  EIM/EIF → Full Cone NAT → no ALG's (with well controlled exceptions)

  CGN role is IPv4 exhaust solution, not security, LI, traffic monitoring, etc.

  Respect OSI model – stay at L3, at L4 if you have to, not higher layers


- Firewall = ALG's = going above L4

  Breaks Net Neutrality (OTT regulatory pushback)

  ALG for Vendor-A breaks app from Vendor-B (same port, different traversal)

  Undefined performance impact of ALG's → numerous DoS attack vectors

  Bugs, ISP liable for 3rd party apps?

# Reality Check: Internet Apps work with NAT



iTunes

Google Maps

Playstation Network

iPhone App Store

Windows Live Messenger

Google Talk

# Reality Check: Apps and NAT Traversal

- STUN, ICE, TURN

  NAT EIM/EIF – Intelligence in endpoint

  Useful for offer/answer protocols
  (SIP, XMPP, probably more)

  Standardized in MMUSIC and BEHAVE

**ICE apps – exaples:**

- **Google chat (XMPP)**

- **Microsoft MSN (SIP inside of XML)**

- **Yahoo (SIP)**

- **Counterpath softphone (SIP)**

STUN: "Session Traversal Utilities for NAT" – RFC 5389
ICE: "Interactive Connectivity Establishment" – RFC 5245
TURN: "Traversal Using Relays around NAT" – RFC 5766

- Other examples

  IPSec over TCP/UDP

  FTP PASV – data connection always to server

  RTSPv1 → RSTPv2 (effectively replaced HTTP Video, ABR,…)

  Skype – encrypted, does its own NAT traversal

  Port 80/443 apps

**Known Problems:**
Active FTP (old browsers), RTSPv1 (old m.youtube.com), MS PPTP (old PC VPN)
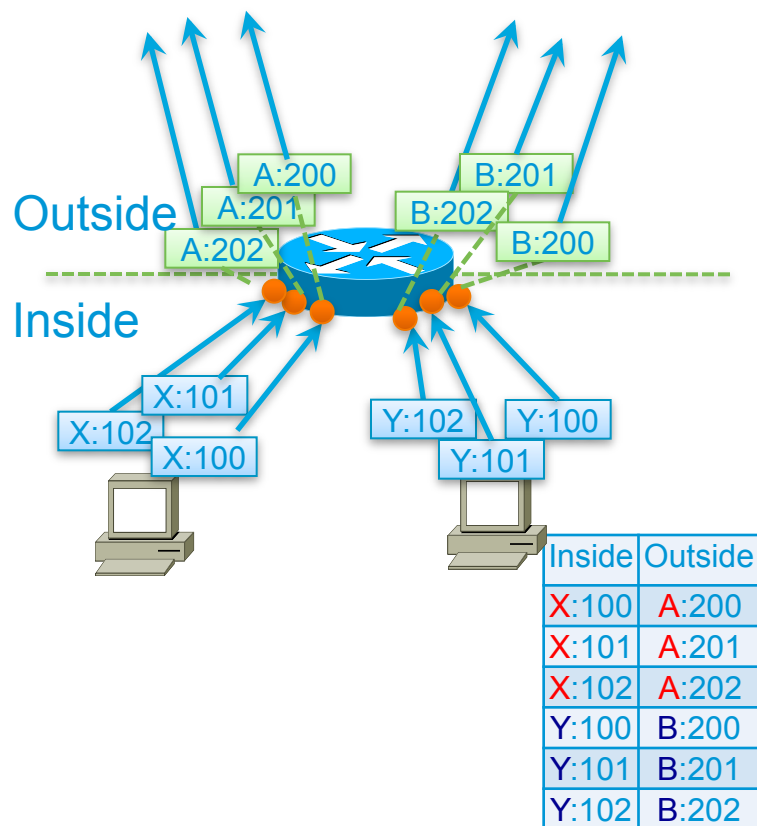
# CGN Behaviors
## *draft-ietf-behave-lsn-requirements*



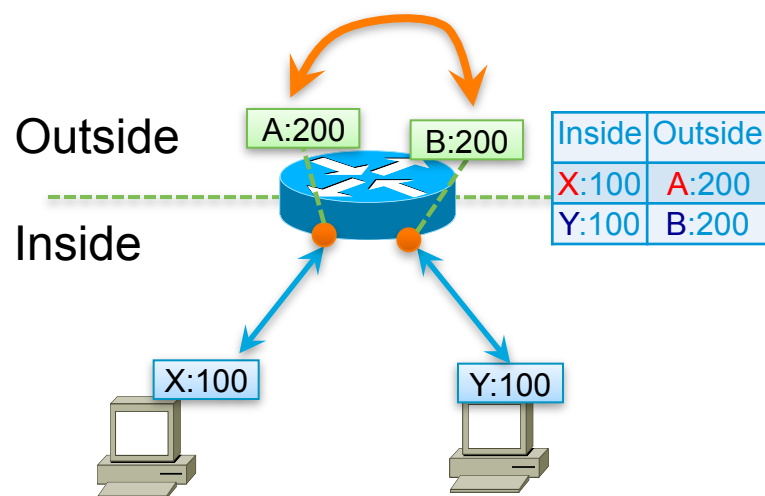A CGN is defined by constrained behavior:

- NAT Behavior Compliance (RFC4787, RFC5382, RFC5508)

    Endpoint Independent Mapping and Filtering (Full Cone NAT)

    ALG's (fixups) should not be used (exceptions like A-FTP)

    Paired IP address pooling behavior

    Port Parity preservation for UDP

    Hairpinning behavior

    Static Port Forwarding (PCP)

- Management

    Port Limit per subscriber

    Mapping Refresh

    Very scalable NAT logging (binary Netflow)

- Redundancy – Intra-box Active/Standby, Inter-box Active/Active

- Scale – 10M's concurrent sessions, 100K's sessions per second, Virtualization (VRF-aware)

- IPv6 Transition Tool-set – dual-stack, NAT64, 6RD, DS-Lite, MAP-T…

# Paired IP Address Pooling Behavior



| Inside | Outside |
|--------|---------|
| X:100  | A:200   |
| X:101  | A:201   |
| X:102  | A:202   |
| Y:100  | B:200   |
| Y:101  | B:201   |
| Y:102  | B:202   |

- Paired (recommended) : use the same external IP address mapping for all sessions associated with the same internal IP address

- Some peer to peer applications don't negotiate the IP address for multiple sessions (eg. apps that are not able to negotiate the IP address for RTP and RTCP separately)

# Hairpinning Behavior



Outside

| Inside | Outside |
|--------|---------|
| X:100 | A:200 |
| Y:100 | B:200 |

Inside

- Use Case: Allow communications between two endpoints behind the same NAT when they are trying each other's external IP addresses
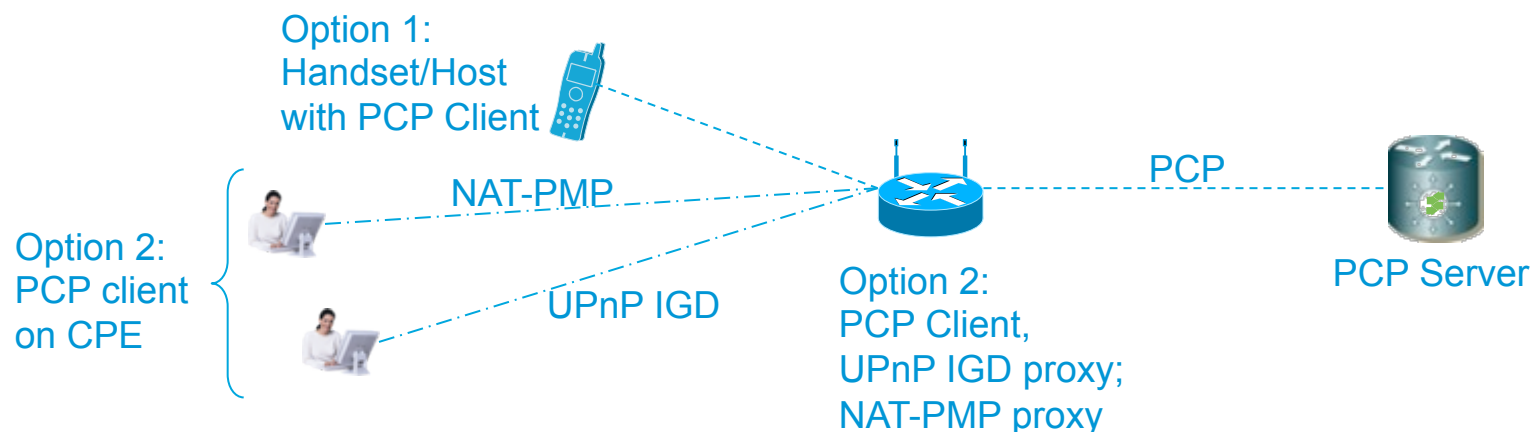
Notation | X:100 | IPv4 address:Port *

\* TCP/UDP port or Query ID for ICMP

# Static Port Forwarding

- Requirement: Ability to configure, a fixed private (internal) IP address:port associated with a particular subscriber while CGN allocates a free public IP address:port

- New protocol: PCP (Port Control Protocol)

Option 1:
Handset/Host
with PCP Client

NAT-PMP

Option 2:
PCP client
on CPE

UPnP IGD

PCP

Option 2:
PCP Client,
UPnP IGD proxy;
NAT-PMP proxy

PCP Server

- Delegate port numbers to requesting applications/hosts to avoid requirement for ALGs

- draft-ietf-pcp-base

# Other Port Behaviors

## No Port Overloading

• A NAT must not have a "Port assignment" behavior of "Port overloading"( i.e. use port preservation even in the case of collision). Most applications will fail if this is used.

## Port Parity Preservation

• An even port will be mapped to an even port, and an odd port will be mapped to an odd port. This behavior respects the [RFC3550] rule that RTP use even ports, and RTCP use odd ports.
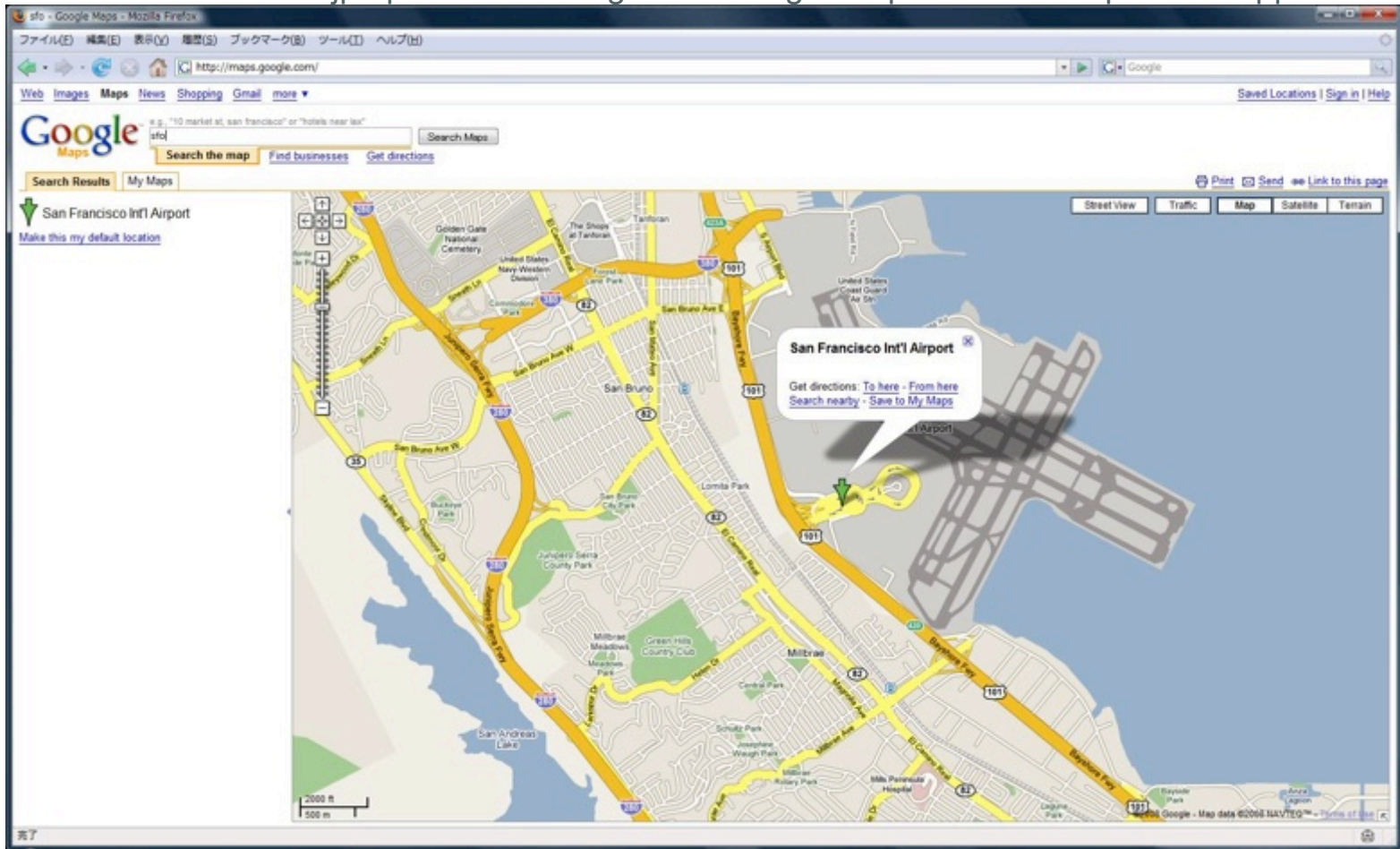
## Port Limit Per Subscriber

• Configurable port limit per subscriber for the system (includes TCP, UDP and ICMP). NAT Security – DoS attack/virus exhaust prevention.
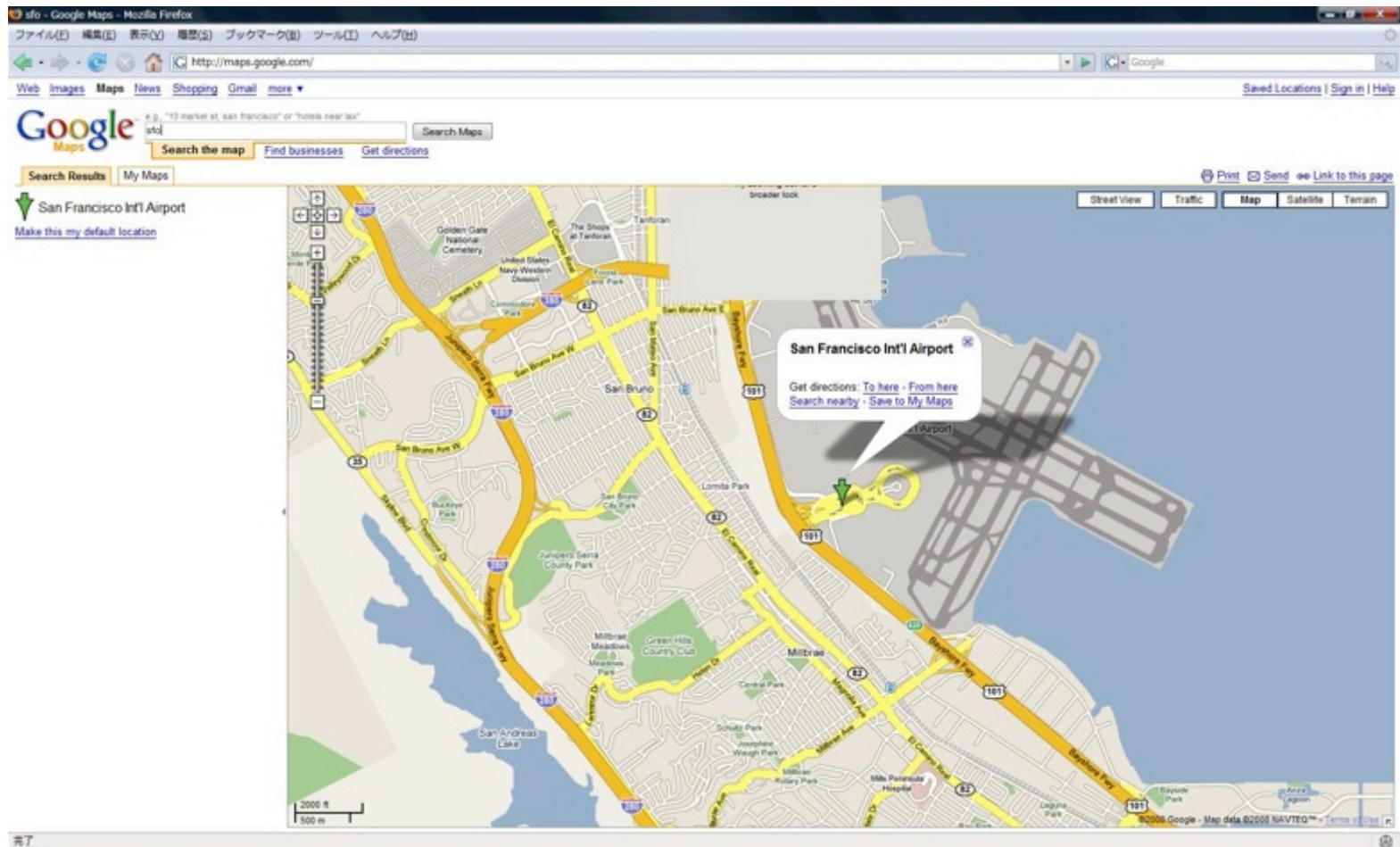
# Impact of NAT Port Limits
## *Example: GoogleMaps with Max 30 Connections*
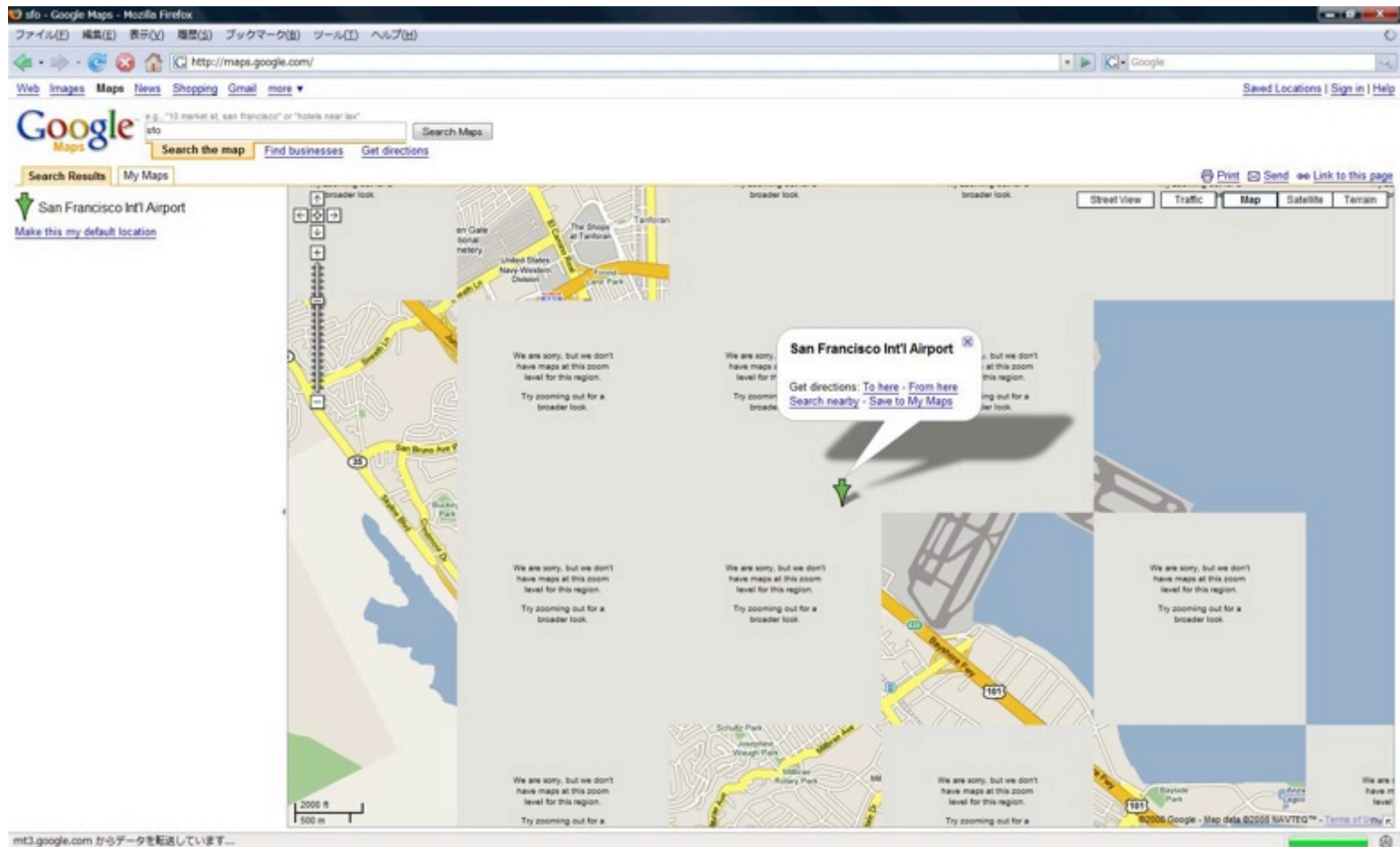
Example/Slides Courtesy of NTT, See Also:
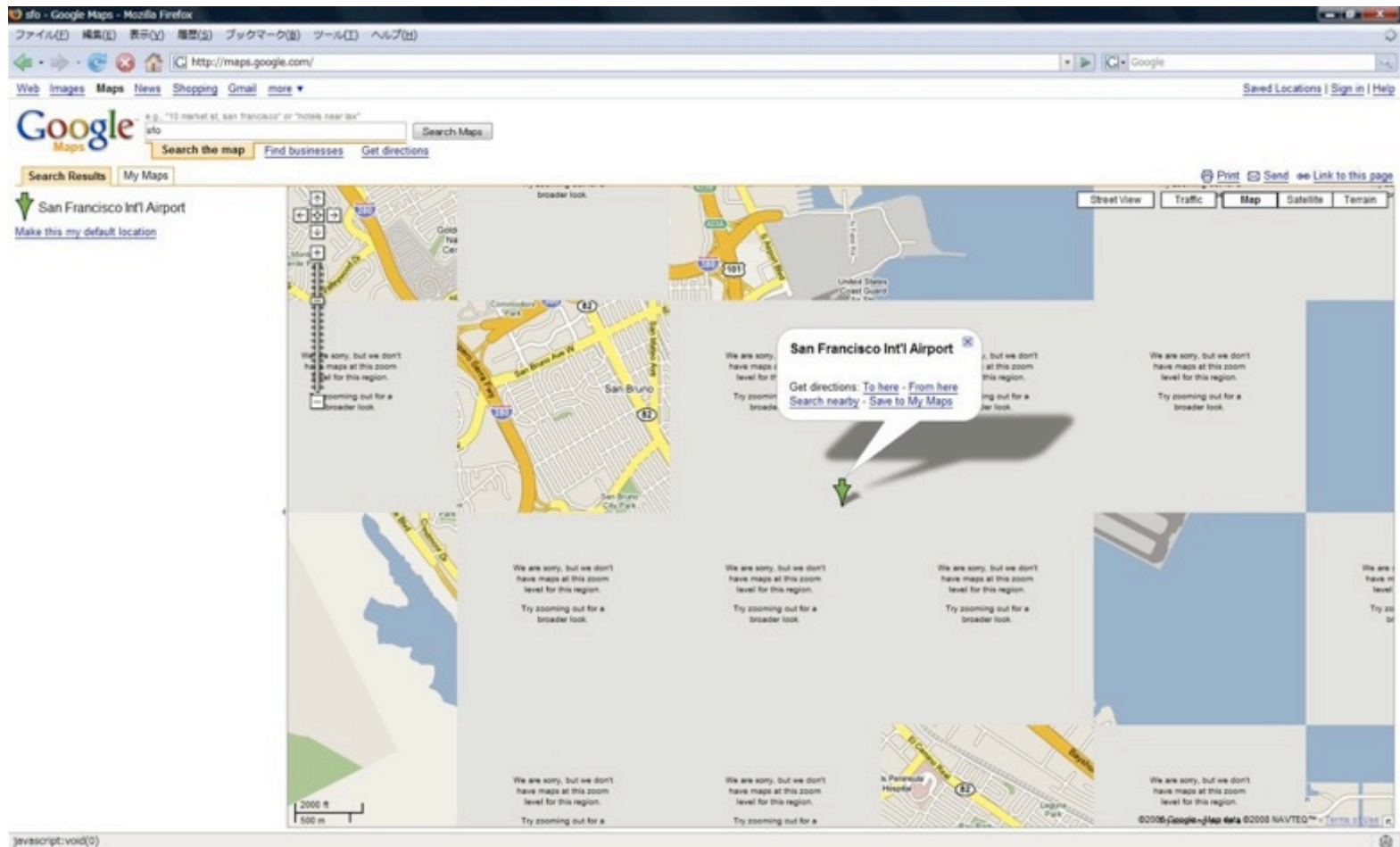Hiroshi Esaki: www2.jp.apan.net/meetings/kaohsiung2009/presentations/ipv6/esaki.ppt

# Max 20 Connections

# Max 15 Connections

# Max 10 Connections

# Max 5 Connections

# Number of Sessions for Some Applications

## Some examples of major Web site

| Application | # of TCP sessions |
|---|---|
| No operation | 5～10 |
| Yahoo top page | 10～20 |
| Google image search | 30～60 |
| ニコニコ動画 | 50～80 |
| OCN photo friend | 170～200+ |
| iTunes | 230～270 |
| iGoogle | 80～100 |
| 楽天(Rakuten) | 50～60 |
| Amazon | 90 |
| HMV | 100 |
| YouTube | 90 |

* Courtesy of NTT, Hiroshi Esaki

**Port Consumption can be big**
- Eg. AJAX-based applications with tens-hundreds of TCP sessions
- Eg. Relaunching Firefox with Tabs opens hundreds of sessions



Port Consumption Comparision Between Mobile Broswer and PC Broswer

Source:
Application behaviors in in terms of port/session consumptions on NAT
http://opensourceaplusp.weebly.com/experiments-results.html
See also "An Experimental Study of Home Gateway Characteristi
https://fit.nokia.com/lars/papers/2010-imc-hgw-study.pdf
http://www.ietf.org/proceedings/78/slides/behave-8.pdf

# Default Session Timers – example

| Type | Default Value |
|---|---|
| ICMP | 60 sec |
| UDP init | 30 sec |
| UDP active | 120 sec |
| TCP Init | 120 sec |
| TCP active | 30 min |

**\*) Default Refresh Direction is Bidirectional (configurable to OutBound only)**

# CGN Session Logging

- Data Retention Law compliance, user trackability

  Who posted a content to a server on Tue at 8:09:10pm?

  - Global IP:port → CGN Log → Private IP:port → MSISDN
  - Directive 2006/24/EC - Data Retention

- Logging Format

  Must be fast and efficient (think 1Msps)

  - ASCII format (Syslog) – very chatty (113B add-event), inefficient, no sequence #
  - Binary Format (Netflow) – efficient (21B add event), sequencing

- Netflow V9 Logging

  21B add-event, 11B delete-event

  Up to 68 add-events per 1500B export packet

  - Dynamic, template-based format (1 Msps = cca 176 Mbps, 14.7 Kpps)

  Future evolution → IPFIX (SCTP – reliable streaming, multi-core CPU's)

# Netflow v9 logging

Tip: IsarFlow – tested CGN NFv9 Collector
www.isarnet.de

Add Event
Template 256
(21B)

| Field ID | Attribute | Value |
|----------|-----------|-------|
| 234 | Incoming VRF ID | 32 bit ID |
| 235 | Outgoing VRF ID | 32 bit ID |
| 8 | Source IP Address | IPv4 Address |
| 225 | Translated Source IP Address | IPv4 Address |
| 7 | Source Port | 16 bit port |
| 227 | Translated Source Port | 16 bit port |
| 4 | Protocol | 8bit value |

Delete Event
Template 257
(11B)

| Field ID | Attribute | Value |
|----------|-----------|-------|
| 234 | Incoming VRF ID | 32 bit ID |
| 8 | Source IP Address | IPv4 Address |
| 7 | Source Port | 16 bit port |
| 4 | Protocol | 8bit value |

# Netflow logging Data Volume example

**Collector Performance** – 100K users, average and peak

| Average NF Collector Load | |
|---|---|
| NF records per sec | 188.800 |
| Avg Load - Logging BW Rate (Mbps) | 33 |
| Avg Load - Logging Packets per Second (pps) | 2.776 |
| | |
| **Peak NF Collector Load during CGSE Switchover** | |
| Peak Load - NF records per sec | 1.000.000 |
| Peak Load - Logging BW Rate (Mbps) | 176 |
| Peak Load - Logging Packets per Second (pps) | 14.706 |
| Peak Load - Duration (s) | 9 |

**Storage Capacity** – includes per-day user behavior

| Storage Data Volume | |
|---|---|
| Average Subscriber active Duration (hrs/day) | 8 |
| Raw NF Data Volume per Day (Gbytes) | 29 |
| Total Data Volume per year (TBytes) | 10.5 |
| Database Overhead Factor | 2.0 |
| Database Compression Rate | 6.0 |
| Total Data Volume per year (TBytes) - compressed | 3.5 |

Reality check: 100K CGN users would consume 3.5TB storage per year
(compressed, fully SQL searchable data)
E-Shop: 4TB disk, 300 Euro…

**Usually no need to bother with logging reduction…**

# Logging reduction

- Bulk port range allocation

  - Pre-allocates a port-set per user (eg. 512 ports)

  - Logs only once (when port-set is allocated/deallocated)

- Deterministic NAT

  - Port-sets are determined algorithmically from user IP (patented)

  - No logging until port-set overflows

ISSUES

  - Breaks TCP port randomization (user security consequences)

  - Inefficient usage of global IP pools
    (active users eat hundreds of ports, inactive users eat only few ports)

  - Troubles to get acceptance in IETF (BEHAVE WG)

  - Mutually exclusive with DBL (Destination Based Logging)

→ It is not worth it and is NOT RECOMMENDED. It's better to use Netflow.

# DBL (Destination Based Logging)

Add Event Template 271 (27B)

| Field ID | Attribute | Value |
|----------|-----------|-------|
| 234 | Incoming VRF ID | 32 bit ID |
| 235 | Outgoing VRF ID | 32 bit ID |
| 8 | Source IP Address | IPv4 Address |
| 225 | Translated Source IP Address | IPv4 Address |
| 7 | Source Port | 16 bit port |
| 227 | Translated Source Port | 16 bit port |
| 12 | Destination Address | IPv4 Address |
| 11 | Destination Port | 16 bit port |
| 4 | Protocol | 8 bit value |

DBL logs also destination IP:port
- data retention vs. user privacy
- keeps EIM/EIF behavior

# CGN Scale and Performance
## *Session = full-duplex, bidirectional L4 flow*

- <u>Session Setup Rate [sps]</u> – sessions per second

  Average # of New Sessions per User, during peak hours

  - Huge load during a failover scenarios or after a power blackout

  - Failing to cope with SPS = huge TCP delays, timeouts/retransmissions


- Maximum Number of Concurrent Sessions [cs] per CGN

  Average # of Concurrent Sessions per User, during peak hours

  - UDP must not expire in less than 2 minutes (RFC4787)

  - UDP/TCP timers for Initializing and Established sessions should be configurable


- Throughput per CGN [bps]

  Aggregate (downstream + upstream) bandwidth

# CGN Scale and Performance – design

**L (Low-scale) Scenario** – 3G mobile users, smart-phones

**M (Medium-scale) Scenario** – ADSL subscribers, PC users with 3G/4G dongles, Tablets, WiFi and top smart-phone users

**H (High-scale) Scenario** – heavy Broadband users, Internet sharing

| | H | M | L |
|---|---|---|---|
| Average BW – ∑ of combined DL + UL BW per subscriber during peak hrs | 300 kbps | 150 kbps | 30 kbps |
| Average #of Concurrent Sessions – ∑ of all ports/protocol and direction | 200 | 100 | 10 |
| Maximum #of Concurrent Sessions – ∑ of all ports/protocol and direction | 1000 | 500 | 100 |
| Average Session Transaction Rate – #of inbound and outbound sessions that are established or deleted per second | 4,0 | 1,0 | 0,1 |

**100K BB users = up to 100Ksps and 10Mcs during peak hour!**

# CGN Redundancy

- High Availability scenarios

  Intra-chassis, Inter-chassis

  Active/Standby, Active/Active

- Stateful or stateless
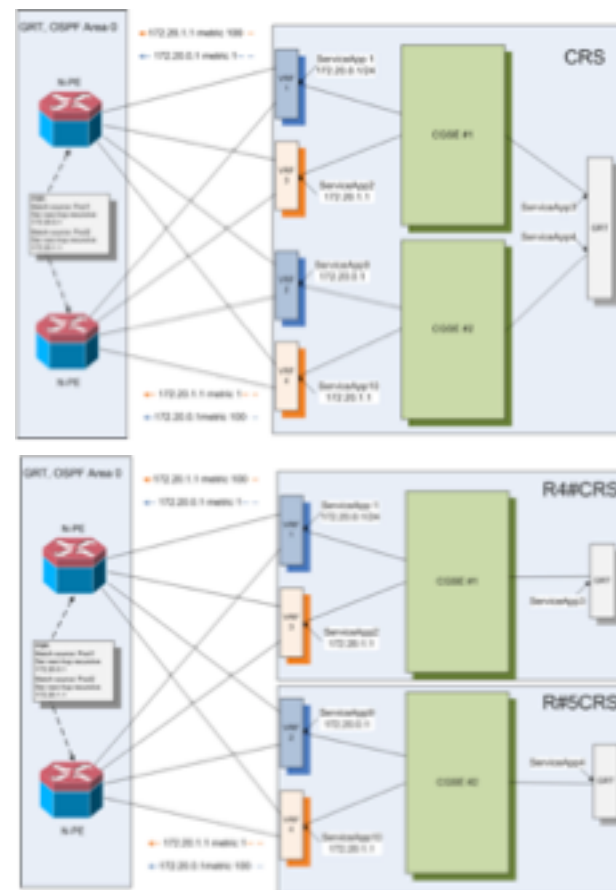
  Millions of short-lived Layer-4 session

  Stateful sync makes no sense for such
  ephemeral state (memory & CPU) – eg.



- **Stateless redundancy**

  1Msps = 100K active users (10Mcs) are up in 10s → minimal loss

  Load-sharing = simple ECMP routing

  Best Practice: Simple Non-Revertive 1:1 Warm Standby

# CGN design: Basic Scenarios



**CGN module**

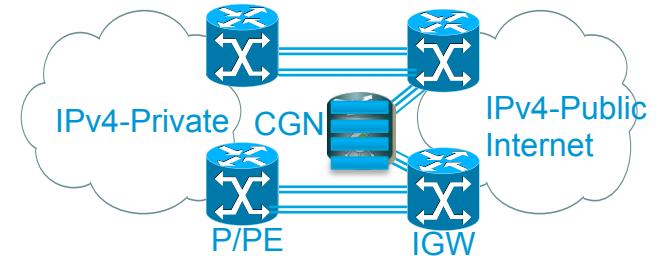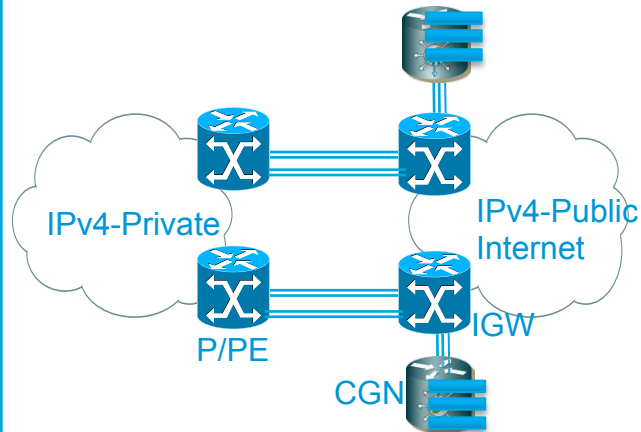|  | **1+1 Redundancy** - 3+3 in this example | **N+1 Redundancy** - 3+1 in this example |
|---|---|---|

**Bump-in-a-wire Design**
- new 4xTGE per box in this example

**Router-on-a-stick Design**
- new 3xTGE per box in this example

**Integrated Design**
- No new ports

**most efficient**

IPv4-Private

IPv4-Public Internet

CGN

P/PE

IGW

IGW with CGN module

# Summary

- CGN is here to overcome IPv4 exhaust before IPv6 migration

- CGN is ISP element, focus on transparency

- CGN is not firewall

- CGN behavioral requirements

- CGN logging

- CGN performance – SPS, # of sessions, deployment options

Thank you.