# Operational routing experience in NTT/OCN

**Routing-SIG @ APNIC19**

**NTT Communications / OCN**
**Tomoya Yoshida <yoshida@ocn.ad.jp>**

# Our History from "OCN Economy"

- We started "OCN Economy" Service in 1996
  - This is the epoch-making service
    - » The Price was very cheap at that time : ¥38,000 128Kbps
  - We distribute|assign /28 or /29 to users
- /28 or /29 is redistributed to OSPF by external route
  - Static route information on the edge router is redistributed to OSPF
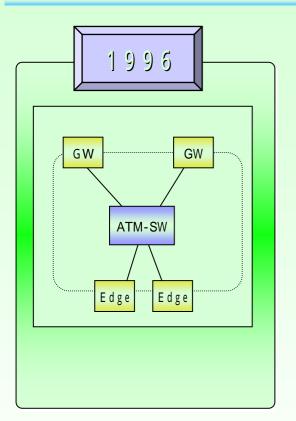- Many OSPF external routes is growing
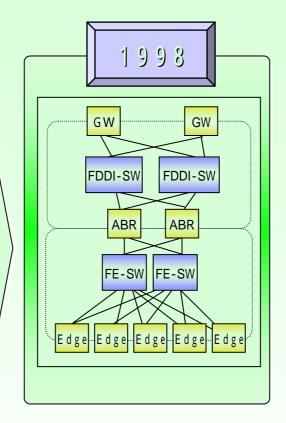
# Our History Cont.

- When the OSPF external route reached around 20,000, OSPF convergence time needed more and more
  - We tried to separate OSPF domain
    - » Operation would be complicated
    - » Extension would be difficult
  - We changed from OSPF to BGP around 2000
  - iBGP route is growing and growing very fast
  - Then we use route reflector hierarchy
- Address problem
  - We could not get enough address to assign at once
  - As the result it was difficult to aggregate the route

# The changes of OCN Backbone Topology



full-mesh topology

using ATM-SW etc

divided OSPF area
FDDI  FE-SW
Reduction of routing

Clustering topology
according to the service or
routing

# Backbone Topology in 1999

**Global Internet**
Domestic **600M**, International:**1.2G**

Osaka

Tokyo

Spporo

Sendai

Tachikawa

**Kyoto**

Karagasaki

Yokohama

Fukuoka

Kobe

Nagoya

Hiroshima

Doujima

**NTT** Communications

OCN

# Current OCN Backbone Topology

# Square Backbone



Square Topology
(2 redundant square)

dispersion of the External Link
(Hot-Potato Routing)

Regional POP dispersion

International

Domestic

Osaka

Tokyo

Osaka

Tokyo

Regional POP

Regional POP

Osaka

Tokyo

| | |
|---|---|
| Tokyo-Osaka : | 80G |
| International : | 48G |
| Domestic : | over 100G |

# Routing (OSPF/BGP)

- **OSPF:IGP**
  - Backbone area and many other areas : normal design
  - Cost design is basically equal cost load balancing
  - Distribution the function of DR/BDR in the same router for more than two segment
  - Restriction of the number of router in the same area
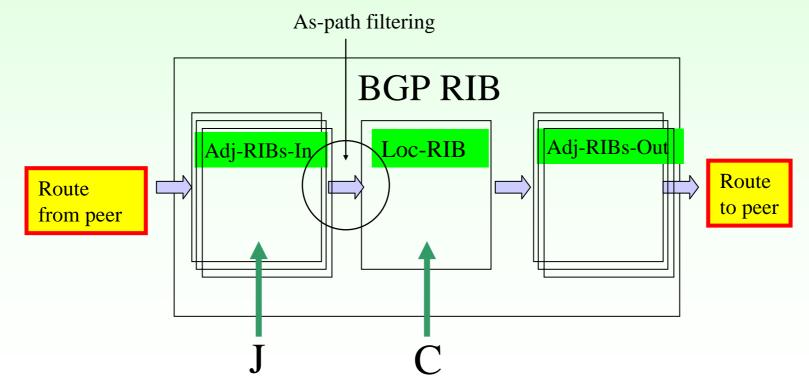
- **BGP:EGP**
  - Route reflector hierarchy topology
  - Distribute for needed cluster

# BGP prefix limitation experience

- Both Cisco and Juniper have a limitation mechanism of the BGP route from peer
- But those implementation are different

As-path filtering

BGP RIB

Adj-RIBs-In    Loc-RIB    Adj-RIBs-Out

Route from peer    Route to peer

J    C

# Next-hop self / redistribution

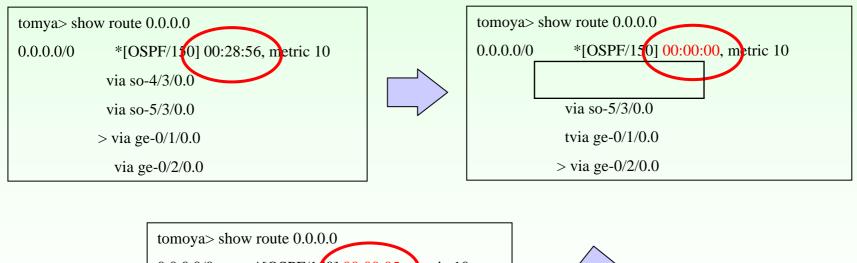- If you forget next-hop-self at the eXchange border route and not redistributed to your backbone the IX segment around /24

- In Japan, 3 major IXs is announcing around /20 the part of the IX's segment IP like /24, so when some ISP forget the next-hop-self and not redistribute those segment to IGP, traffic will go to the IX's AS (dix-ie, JPIX, JPNAP's AS)

**NTT** Communications

**OCN**
OPEN COMPUTER NETWORK

# LSA refresh experience

- ## Some LSA was flapping
  - ### Default refresh timer is different
    - » Cisco is 30 minutes, Juniper is 50 minutes

```
tomya> show route 0.0.0.0

0.0.0.0/0        *[OSPF/150] 00:28:56, metric 10

                  via so-4/3/0.0

                  via so-5/3/0.0

                > via ge-0/1/0.0

                  via ge-0/2/0.0
```

```
tomoya> show route 0.0.0.0

0.0.0.0/0        *[OSPF/150] 00:00:00, metric 10



                  via so-5/3/0.0

                  tvia ge-0/1/0.0

                > via ge-0/2/0.0
```

```
tomoya> show route 0.0.0.0

0.0.0.0/0        *[OSPF/160] 00:00:05, metric 10

                > via so-4/3/0.0

                  via so-5/3/0.0

                  via ge-0/1/0.0

                  via ge-0/2/0.0
```

# Route cache is very useful

- Currently almost vendor is implemented "route refresh capability"

- But soft-reconfiguration inbound ( for crs-1 need always keyword ) is very useful

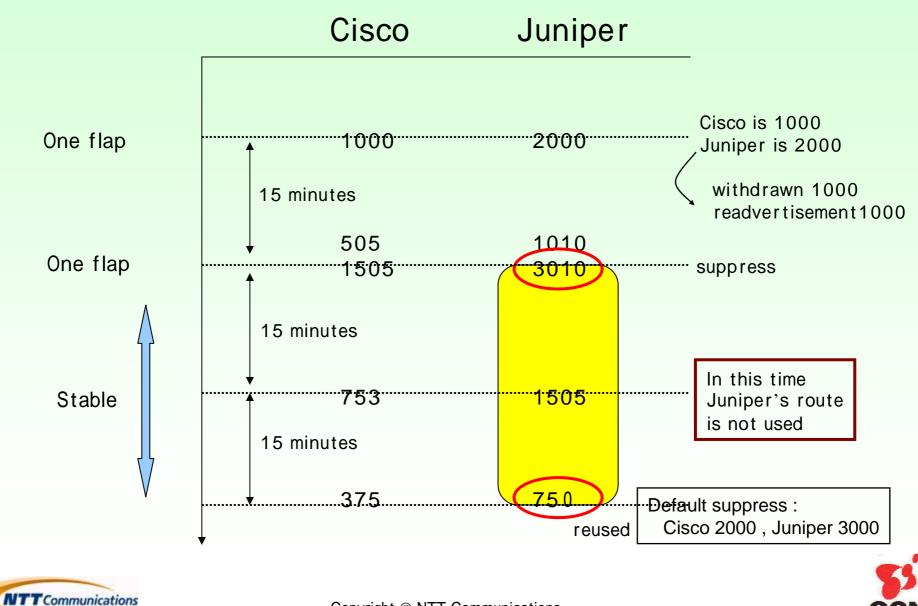- When you set a new peer, you set low priority to this new peer, but more specific is strong!
  - Firstly check the route not receiving any route, only monitor the route from peer by using cache then receive

OCN

# Route flapping experience

Cisco            Juniper

One flap          1000            2000

Cisco is 1000
Juniper is 2000

withdrawn 1000
readvertisement1000

15 minutes

                  505             1010
One flap          1505            3010          suppress

15 minutes

Stable            753             1505

In this time
Juniper's route
is not used

15 minutes

                  37.5            75

reused

Default suppress :
Cisco 2000 , Juniper 3000

# Routing Hijack

- We have around /10 IP blocks

- Sometimes our prefix hijacked

- When we hijacked our route, we announce more specific prefix to the internet

  – But When someone hijack /24, it is very difficult

    » We announce two /25s but almost ISPs cannot receive

    » Also we announce /24 in addition to /16 our PA

- We need BGP origin validation security mechanism
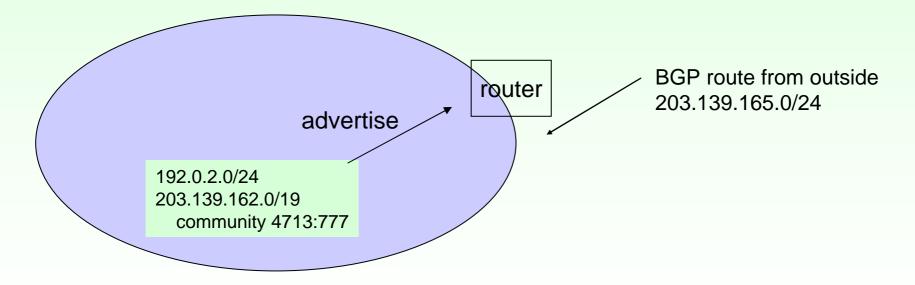
  – sBGP/So-BGP or IRR etc.

# We need…

- **TTL hack security mechanism for many vendor**
- **Prefix limitation by using LOC-RIB for Juniper**
- **Mac accounting for 10G**
- **Feasible path reverse path forwarding for uRPF**
  - Strict mode is dangerous
  - Loose mode is just loose…
- **BGP Inactive reason for Cisco is coming**
  - Cisco implemented for CRS-1
  - Operational additional information is very important
- **Dynmic filtering by using bgp community, just my idea**

# Dynamic Filtering : just idea

■ If you receive the BGP route with this community (4713:777 attribute), the route which in scope of this community will be rejected automatically
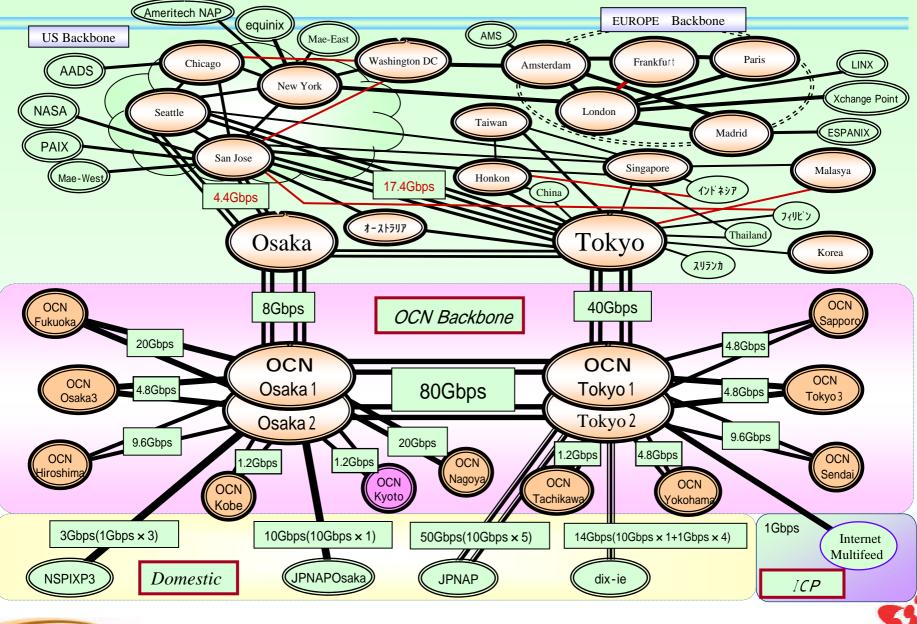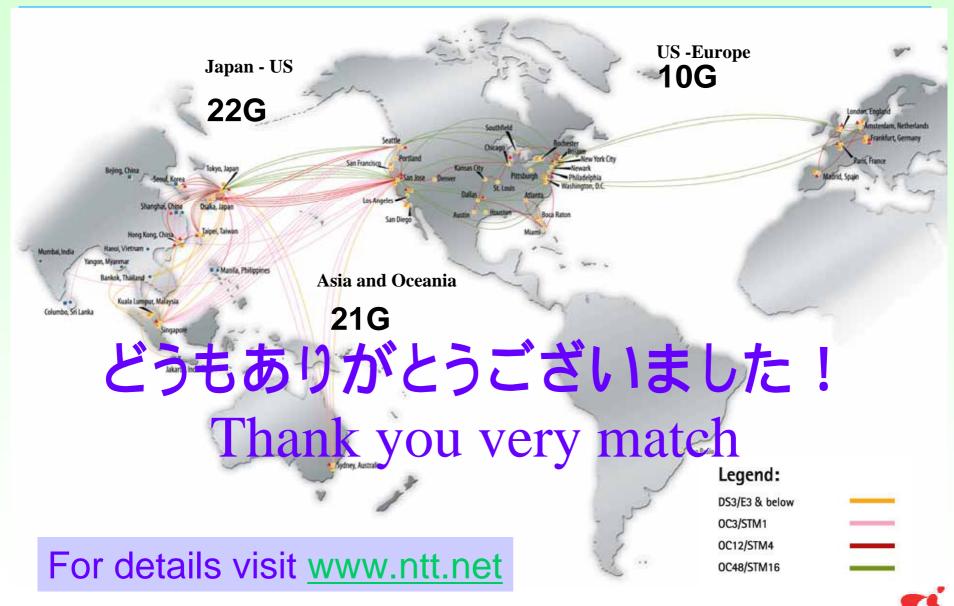  – Useful for filtering for your PA

router

advertise

192.0.2.0/24
203.139.162.0/19
    community 4713:777

BGP route from outside
203.139.165.0/24

# Our Backbone



US Backbone

Ameritech NAP
equinix
Mae-East
Washington DC
AADS
Chicago
New York
NASA
Seattle
PAIX
San Jose
Mae-West
4.4Gbps
17.4Gbps

EUROPE Backbone

AMS
Amsterdam
Frankfu
Paris
LINX
London
Madrid
Xchange Point
ESPANIX
Taiwan
Honkon
China
Singapore
Malasya
Thailand
Korea

Osaka
Tokyo

**OCN Backbone**

OCN
Fukuoka
8Gbps
40Gbps
OCN
Sapporo
20Gbps
4.8Gbps
OCN
Osaka
OCN
Tokyo
OCN
Osaka3
4.8Gbps
**80Gbps**
4.8Gbps
OCN
Tokyo
Osaka
Tokyo
9.6Gbps
9.6Gbps
OCN
Hiroshima
20Gbps
OCN
Sendai
1.2Gbps
1.2Gbps
OCN
Nagoya
1.2Gbps
4.8Gbps
OCN
Kobe
OCN
Kyoto
OCN
Tachikawa
OCN
Yokohama

3Gbps(1Gbps× 3)
10Gbps(10Gbps× 1)
50Gbps(10Gbps× 5)
14Gbps(10Gbps× 1+1Gbps× 4)
1Gbps
Internet
Multifeed

NSPIXP3
*Domestic*
JPNAPOsaka
JPNAP
dix-ie
*P*

**NTT** Communications

Copyright © NTT Communications

OCN

# NTT Communications Global IP Network



**Japan - US**

**22G**

**US -Europe**
**10G**

**Asia and Oceania**

**21G**

Thank you very match

Legend:

DS3/E3 & below
OC3/STM1
OC12/STM4
OC48/STM16

For details visit www.ntt.net