

# (A b )U s i n g R o u t e S e r v e r s

E l i s a J a s i n s k a <elisa.jasinska@ams-ix.net>

C h r i s M a l a y t e r <cmalayer@switchanddata.com>

A r n o l d N i p p e r <arnold.nipper@de-cix.net>

# A g e n d a

- W h y R o u t e S e r v e r s ?
- W h a t d o R o u t e S e r v e r s d o ?
- C u r r e n t i m p l e m e n t a t i o n s a n d R o u t e S e r v e r  
W o r k i n g G r o u p
- F u n c t i o n a l i t y a n d s c a l a b i l i t y t e s t i n g

# A g e n d a

- W h y R o u t e S e r v e r s ?
- W h a t d o R o u t e S e r v e r s d o ?
- C u r r e n t i m p l e m e n t a t i o n s a n d r o u t e S e r v e r  
W o r k i n g G r o u p
- F u n c t i o n a l i t y a n d s c a l a b i l i t y t e s t i n g

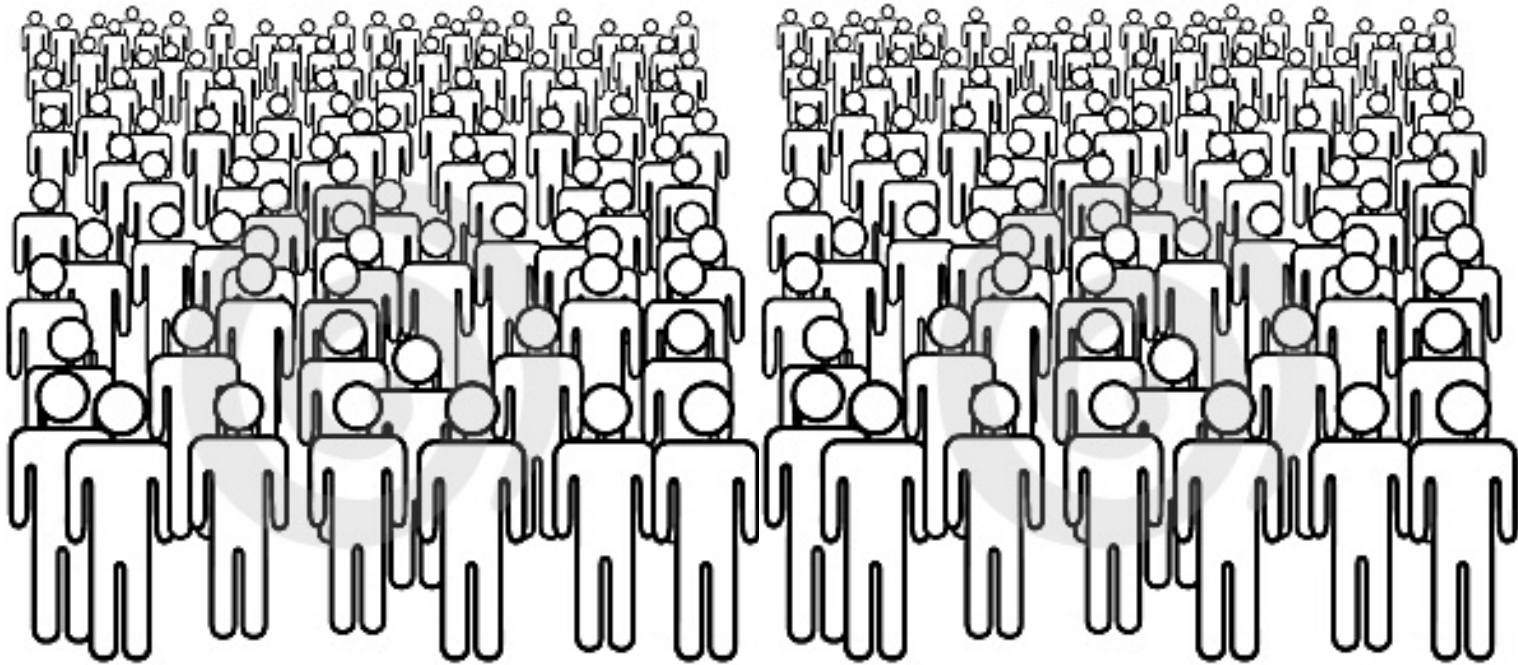
# W h y R o u t e S e r v e r s ?

- I n t e r n e t E x c h a n g e ( e . g . A M S - I X , D E - C I X , L I N X )
- P e e r i n g P l a t f o r m f o r m a n y P a r t i e s
- R o u t e S e r v e r s f o r t h e P a r t i c i p a n t s

# Why Route Servers?

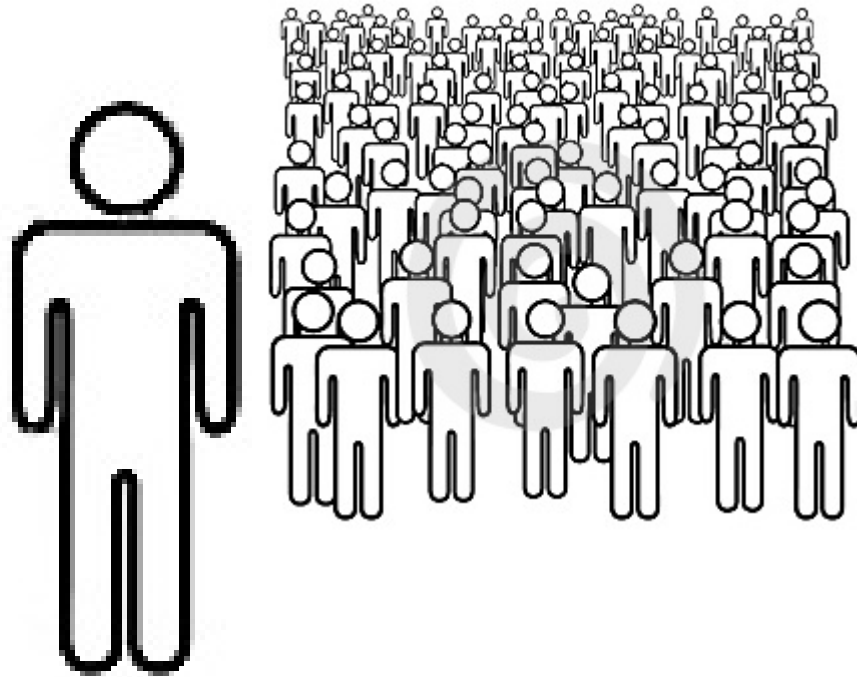
- Peer with as many parties as possible

➔ Maintaining lots of BGP Sessions



# Why Route Servers?

- Reach a lot of Parties with just one BGP session



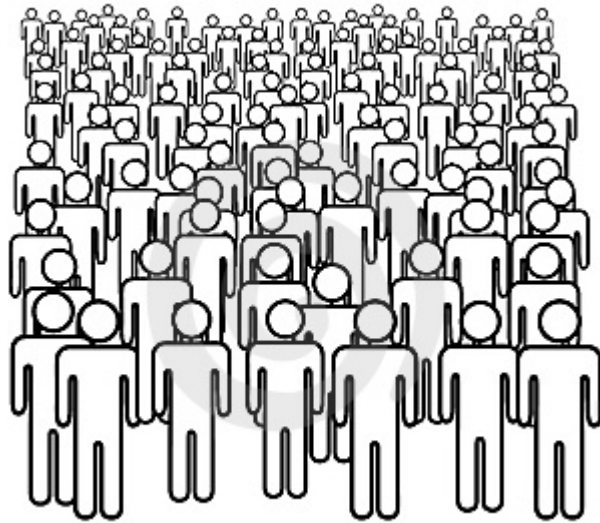
# W h y R o u t e S e r v e r s ?

- R e d u n d a n c y ... i n c a s e y o u r s e s s i o n s d i e ...



# Why Route Servers?

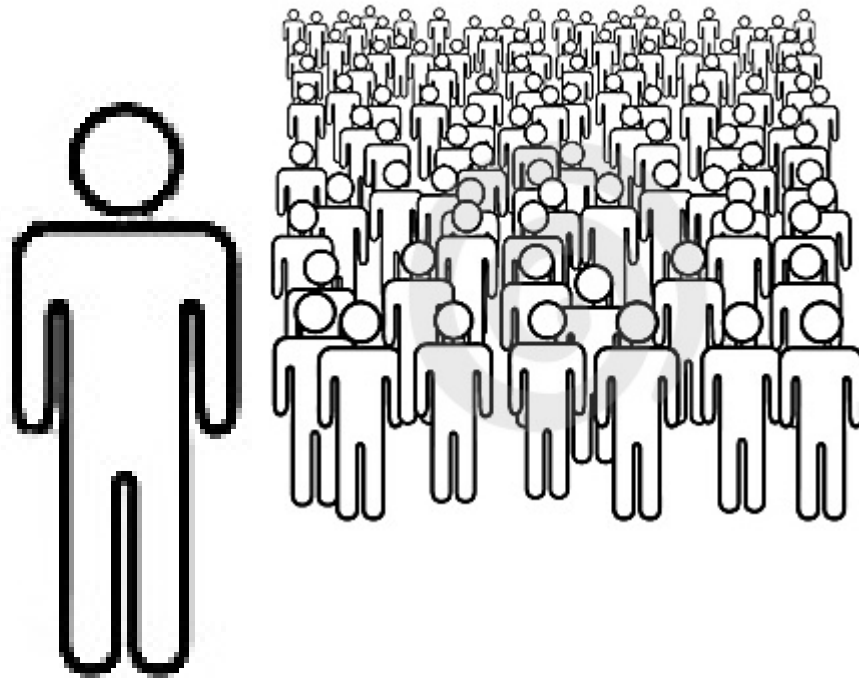
- Redundancy ... in case the Route Server dies ...





# W h y R o u t e S e r v e r s ?

- E a s y e n t r y p o i n t f o r n e w M e m b e r s t o t h e E x c h a n g e - i m m e d i a t e t r a f f i c



# A g e n d a

- W h y R o u t e S e r v e r s ?
- W h a t d o R o u t e S e r v e r s d o ?
- C u r r e n t i m p l e m e n t a t i o n s a n d r o u t e S e r v e r  
W o r k i n g G r o u p
- F u n c t i o n a l i t y a n d s c a l a b i l i t y t e s t i n g

# What do route servers do?

- Receive UPDATES from every participant

```
19:58:33.721679 IP (tos 0x0, ttl 64, id 44576, offset 0, flags [DF], proto TCP (6), length 117)
10.23.0.5.58880 > 10.23.0.1.179: P, cksum 0xb892 (correct), 48:113(65) ack 61 win 1460
<nop,nop,timestamp 1976783474 3177206804>: BGP, length: 65
  Update Message (2), length: 65
```

```
...
  AS Path (2), length: 10, Flags [T]: 65499 11 12 13
```

```
...
  Next Hop (3), length: 4, Flags [T]: 10.23.0.5
```

```
...
  Updated routes:
    2.0.5.0/24
```

```
19:58:33.723897 IP (tos 0x0, ttl 64, id 42762, offset 0, flags [DF], proto TCP (6), length 117)
10.23.0.4.33349 > 10.23.0.1.179: P, cksum 0xb033 (correct), 48:113(65) ack 61 win 1460
<nop,nop,timestamp 1976783474 1916183085>: BGP, length: 65
  Update Message (2), length: 65
```

```
...
  AS Path (2), length: 10, Flags [T]: 65500 11 12 13
```

```
...
  Next Hop (3), length: 4, Flags [T]: 10.23.0.4
```

```
...
  Updated routes:
    2.0.4.0/24
```

# W h a t d o r o u t e s e r v e r s d o ?

- A p p l y f i l t e r s f o r t h e r e c e i v i n g p e e r s

```
from AS65500 accept ANY  
to AS65500 announce AS65499
```

```
from AS65499 accept ANY  
to AS65499 announce AS65500
```

# W h a t d o r o u t e s e r v e r s d o ?

- Perform “best path” selection for each peer
- A S N , M E D and next-hop transparent
- Store Routing Information Base (R I B ) for every peer

```
flags destination          gateway          lpref    med aspath origin
*>  2.0.4.0/24             10.23.0.4       100      200 65500 11 12 13 i
*>  2.0.5.0/24             10.23.0.5       100      200 65499 11 12 13 i
```

# What do route servers do?

- Forward the RIB contents to the desired peer

```
19:58:33.901718 IP (tos 0xc0, ttl 1, id 15745, offset 0, flags [DF], proto TCP (6), length 103)
10.23.0.1.179 > 10.23.0.4.33349: P, cksum 0x2b21 (correct), 61:112(51) ack 114 win 17376
<nop,nop,timestamp 1916183105 1976783474>: BGP, length: 51
  Update Message (2), length: 51
```

...

```
AS Path (2), length: 10, Flags [T]: 65499 11 12 13
```

...

```
Next Hop (3), length: 4, Flags [T]: 10.23.0.5
```

...

```
Updated routes:
```

```
  2.0.5.0/24
```

```
19:58:33.903463 IP (tos 0xc0, ttl 1, id 12268, offset 0, flags [DF], proto TCP (6), length 103)
10.23.0.1.179 > 10.23.0.5.58880: P, cksum 0x377e (correct), 61:112(51) ack 114 win 17376
<nop,nop,timestamp 3177206824 1976783474>: BGP, length: 51
  Update Message (2), length: 51
```

...

```
AS Path (2), length: 10, Flags [T]: 65500 11 12 13
```

...

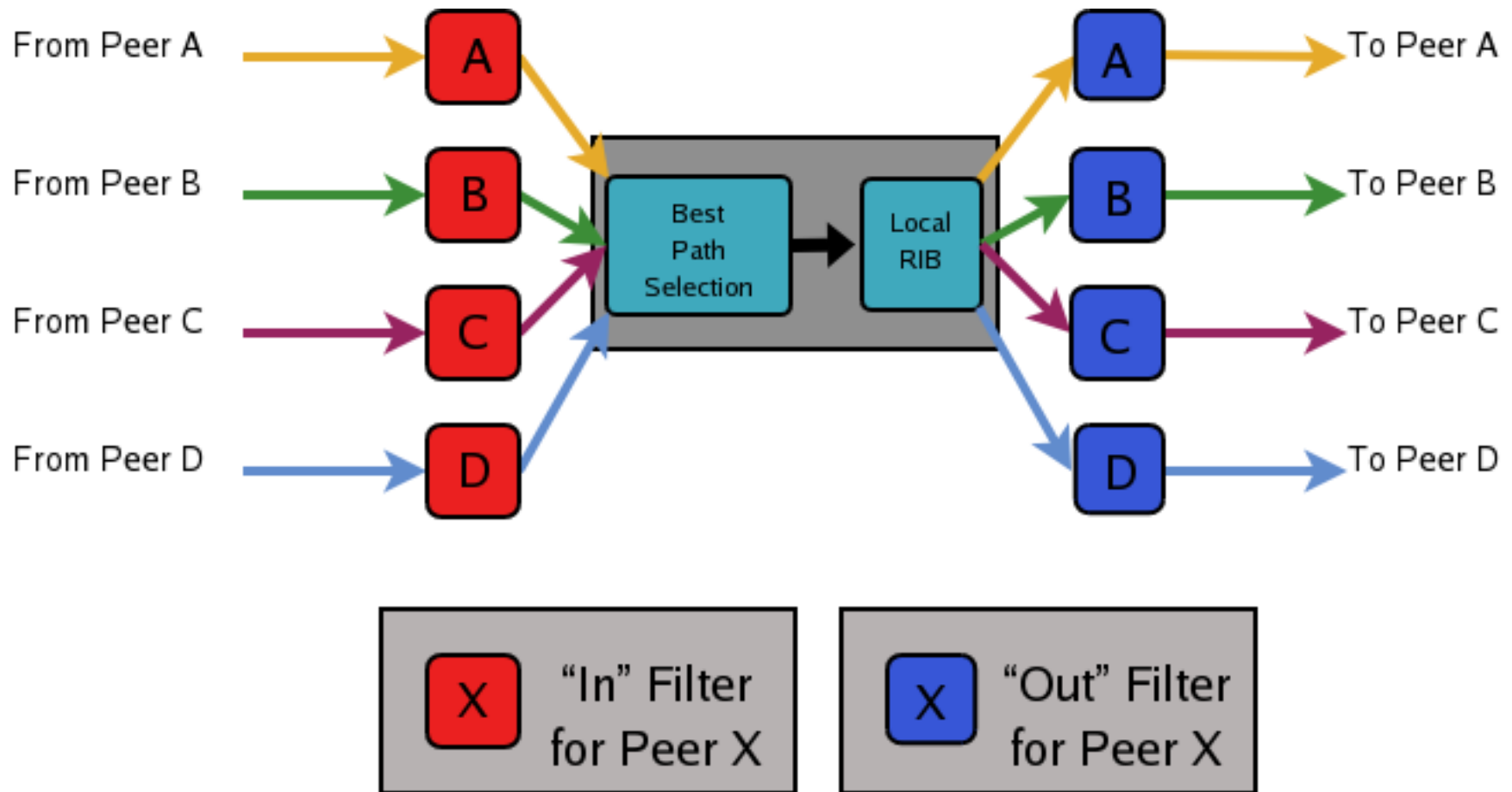
```
Next Hop (3), length: 4, Flags [T]: 10.23.0.4
```

...

```
Updated routes:
```

```
  2.0.4.0/24
```

# What do route servers do?



# A g e n d a

- W h y R o u t e S e r v e r s ?
- W h a t d o R o u t e S e r v e r s d o ?
- C u r r e n t i m p l e m e n t a t i o n s a n d R o u t e S e r v e r  
W o r k i n g G r o u p
- F u n c t i o n a l i t y a n d s c a l a b i l i t y t e s t i n g



# Route Server Virtual Working Group

- Formed at 14<sup>th</sup> Euro-IX in April 2009 in Prague as a subgroup with support from AMS-IX, BCIX, CATNIX, CIX, DE-CIX, INEX, LINX, LONAP, MSK-IX, NaMeX, Netnod, NIX-CZ, PacketExchange, TOP-IX and VIX
- Goal is to have at least two exchange-ready routeservers
- Management team comprising of Arnold Nipper, Nick Hilliard and John Souter

# C u r r e n t I m p l e m e n t a t i o n s

- Q u a g g a
  - T h i s i s w h e r e m o s t o f t h e f u n d s h a v e g o n e s o f a r
- O p e n B G P D
  - M o s t l y d r i v e n b y A M S - I X
- B I R D
  - S e l f - f u n d e d , N I X - C Z

# Route Server Testing Working Group

- Andy Davidson - LONAP
- Chris Malayter - Switch & Data
- Elisa Jasinska - AMS-IX
- Mo Shivji - LINX
- Robert Wozny - PL-IX
- Sebastian Spies - DE-CIX
- Wolfgang Hennerbichler - VIX



# A g e n d a

- W h y R o u t e S e r v e r s ?
- W h a t d o R o u t e S e r v e r s d o ?
- C u r r e n t i m p l e m e n t a t i o n s a n d r o u t e S e r v e r  
W o r k i n g G r o u p
- F u n c t i o n a l i t y a n d s c a l a b i l i t y t e s t i n g

# Functional Testing

# A S 4 / 3 2 B i t A S N

- All three implementations support A S 4
- All three versions tested as of 2009-12-04 to properly implement A S 4

# IP v 6

- All three implementations support IP v 6
- We highly recommend running a current version of any of three implementations
- M A N Y bugs fixed between 2009-10-01 and 2010-01-01
- Running a port of a route server is ill-advised and can leave a bad taste in your mouth !

# Scalability



# Testing

- 100 sessions, set up from IXIA
- 500 or 1000 prefixes per session
- Additional random flapping

# Q u a g g a

- S i n g l e t h r e a d e d i m p l e m e n t a t i o n
- I s s u e s w i t h p e r f o r m i n g i t s t a s k s o n t i m e
- C P U t h r a s h i n g d u r i n g p e r i o d s o f i n s t a b i l i t y
- B u g c a u s i n g c r a s h d u r i n g f l a p p i n g

# O p e n B G P D

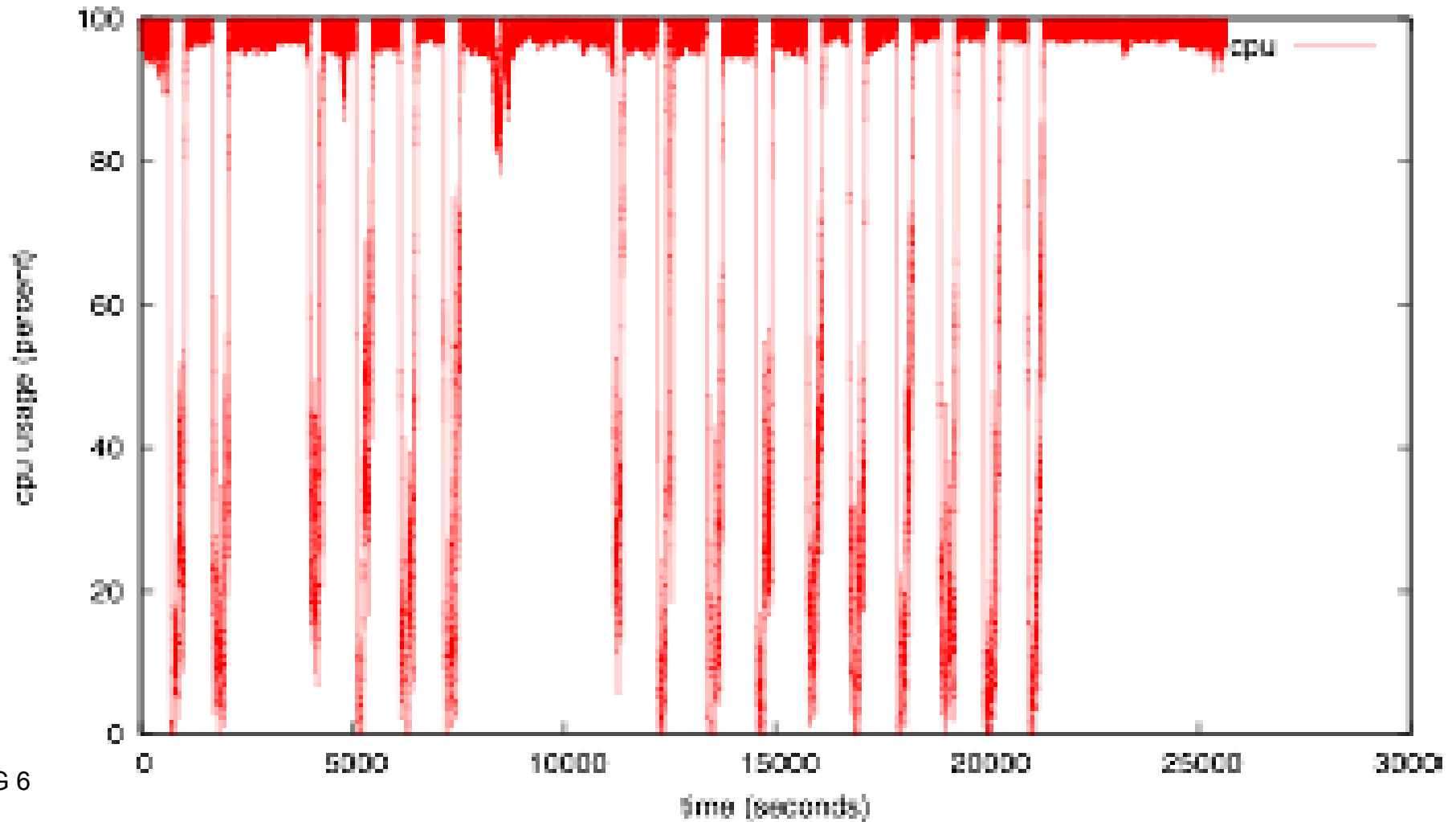
- M u l t i - t h r e a d e d i m p l e m e n t a t i o n
- S e s s i o n t h r e a d k e e p s s e s s i o n s a c t i v e w h i l e i n s t a b i l i t y i s o c c u r r i n g
- 1 G B m e m o r y l i m i t a t i o n p e r p r o c e s s o n i 3 8 6 a n d a 4 G B m e m o r y l i m i t a t i o n o n a m d 6 4
- O n l y a b l e t o u s e > 4 G B m e m o r y w i t h O p e n B S D 4 . 7
- V e r y w e a k o n f i l t e r i n g

# BIRD

- Single threaded implementation
- Amazing scheduling system
- The most stable route server we tested
- Discovered odd memory freeing issues in Linux glibc
- Currently 3 out of 4 of the largest IXP are running BIRD

# Quagga CPU

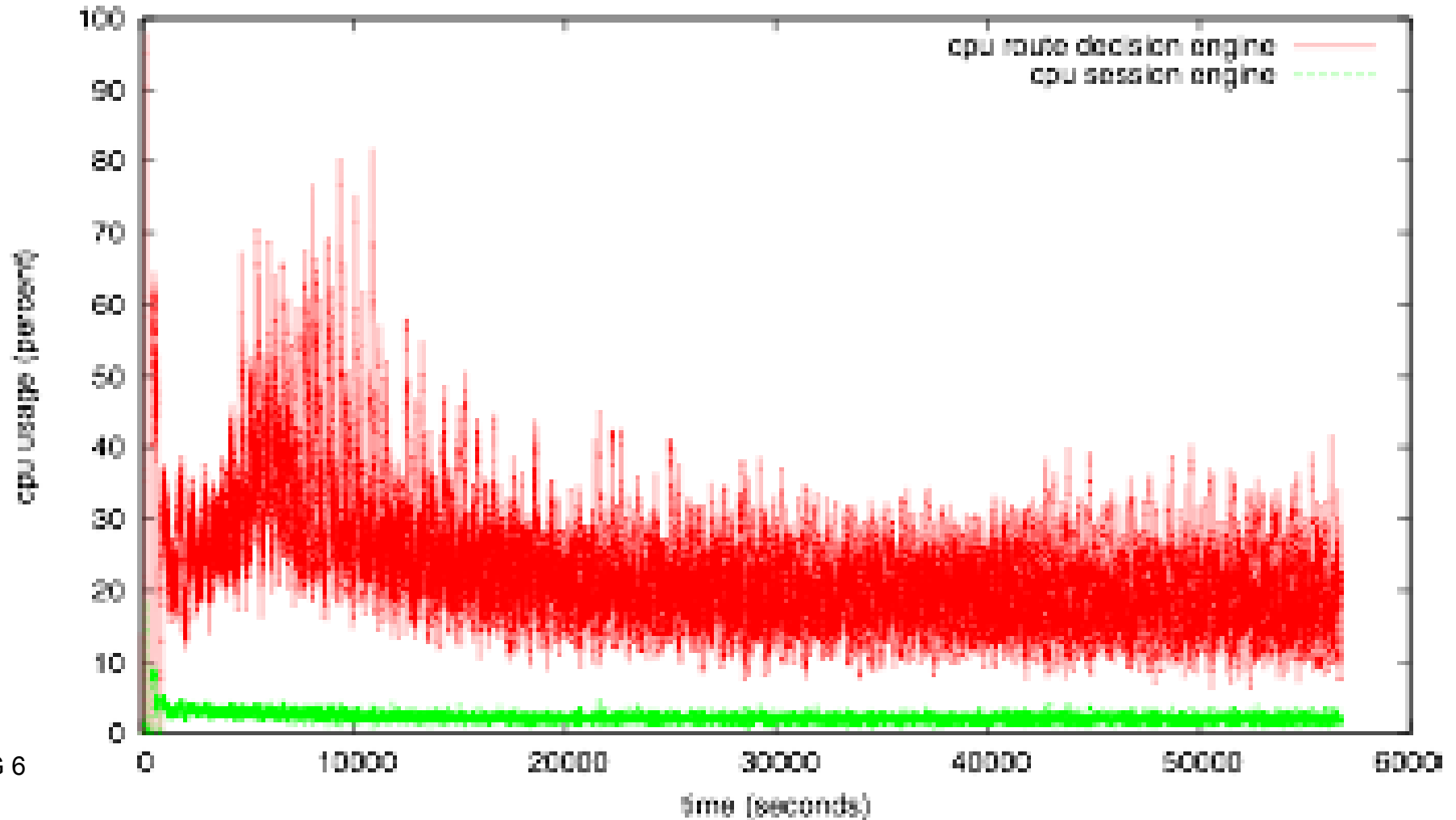
Quagga cpu usage  
multiple-rib; 100 sessions  
IXIA vs. Quagga  
500 prefixes per session with random flapping  
lab4  
2 x Intel(R) Xeon(R) CPU 3050 @ 2.13GHz



# OpenBGPD CPU

OpenBGPD cpu usage  
multiple-rib; 100 sessions  
IXIA vs. OpenBGPD  
1000 prefixes per session  
lab2.pax.net

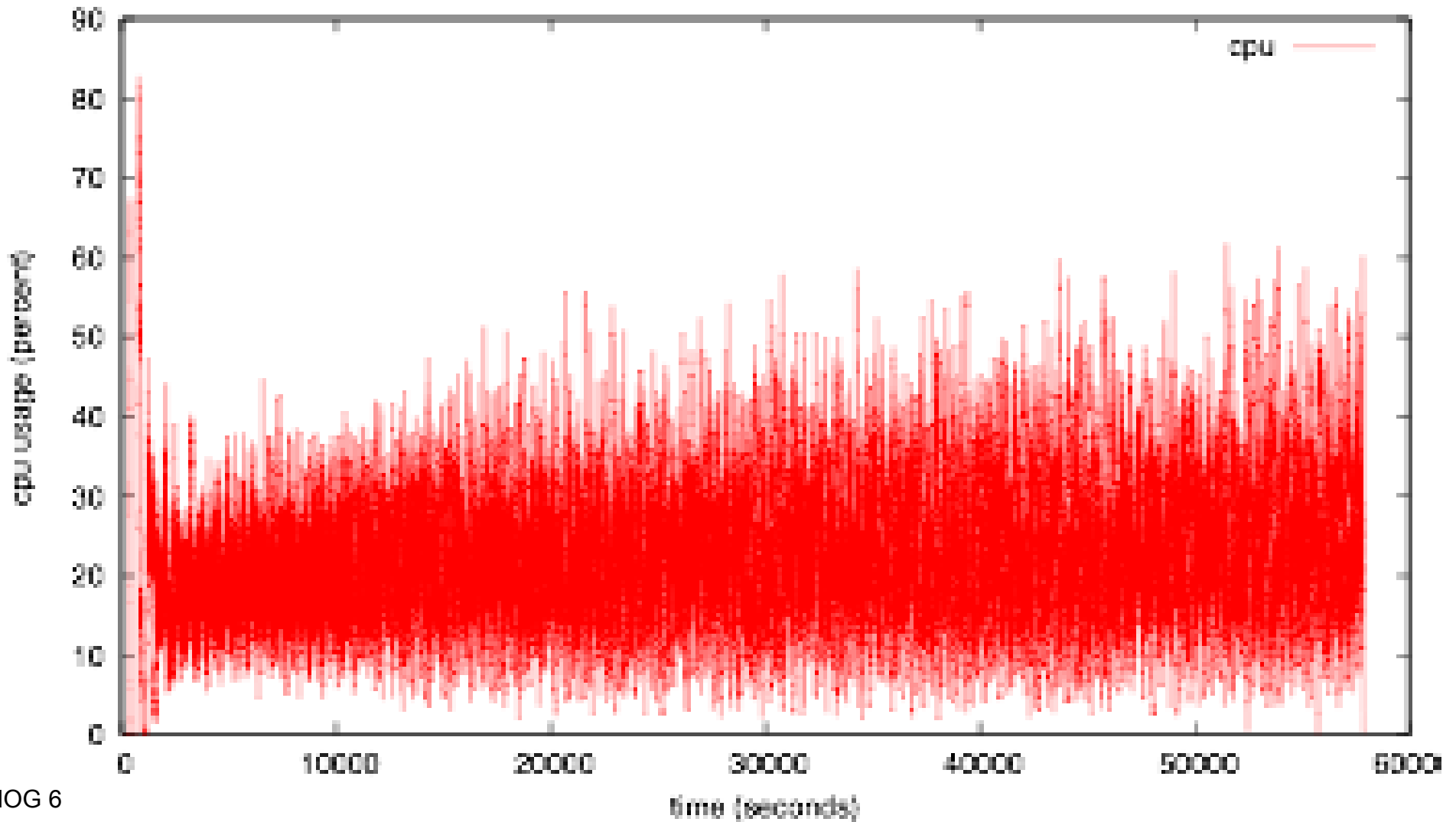
4 x Intel(R) Xeon(TM) CPU 3.60GHz (GenuineIntel 686-class) 3.61 GHz



# BIRD CPU

BIRD cpu usage  
multiple-rib; 100 sessions  
IXIA vs. BIRD

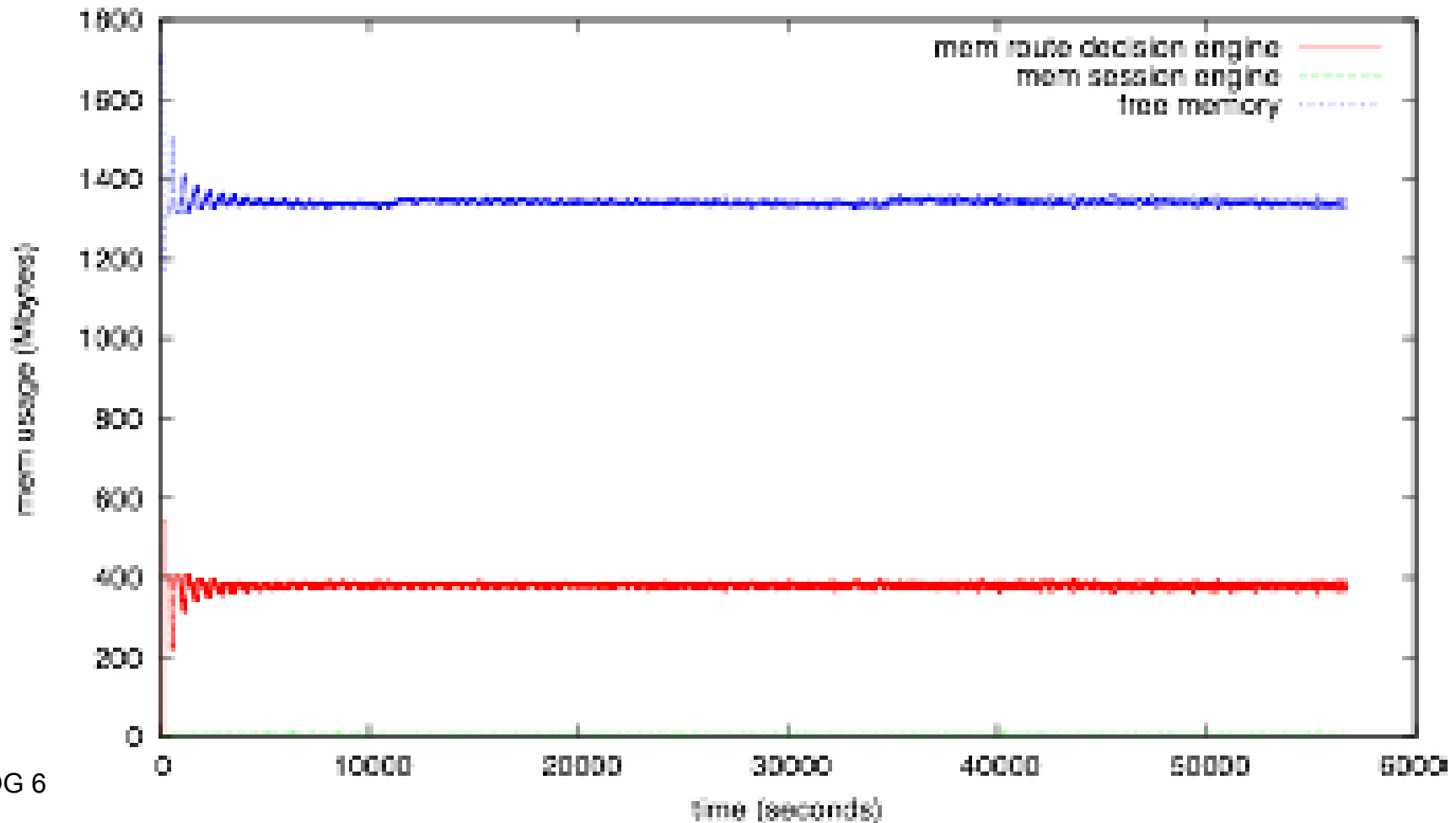
500 prefixes per session with random flapping  
lab5.paix.net  
4 x Intel(R) Xeon(TM) CPU 3.80GHz



# OpenBGPD Mem

OpenBGPD mem usage  
multiple-rib; 100 sessions  
IXIA vs. OpenBGPD  
1000 prefixes per session  
lab2.pax.net

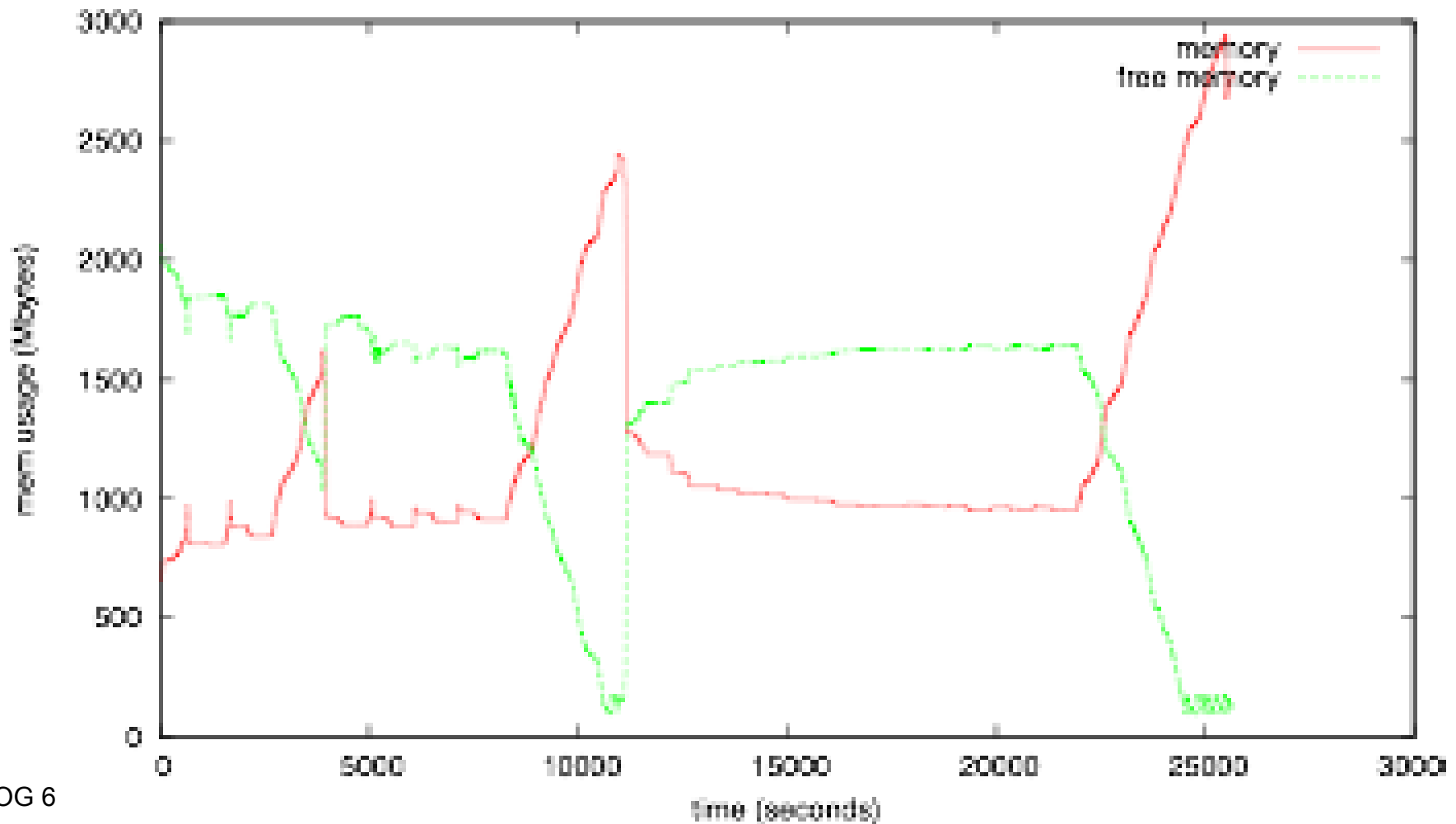
4 x Intel(R) Xeon(TM) CPU 3.50GHz (GenuineIntel 686-class) 3.51 GHz





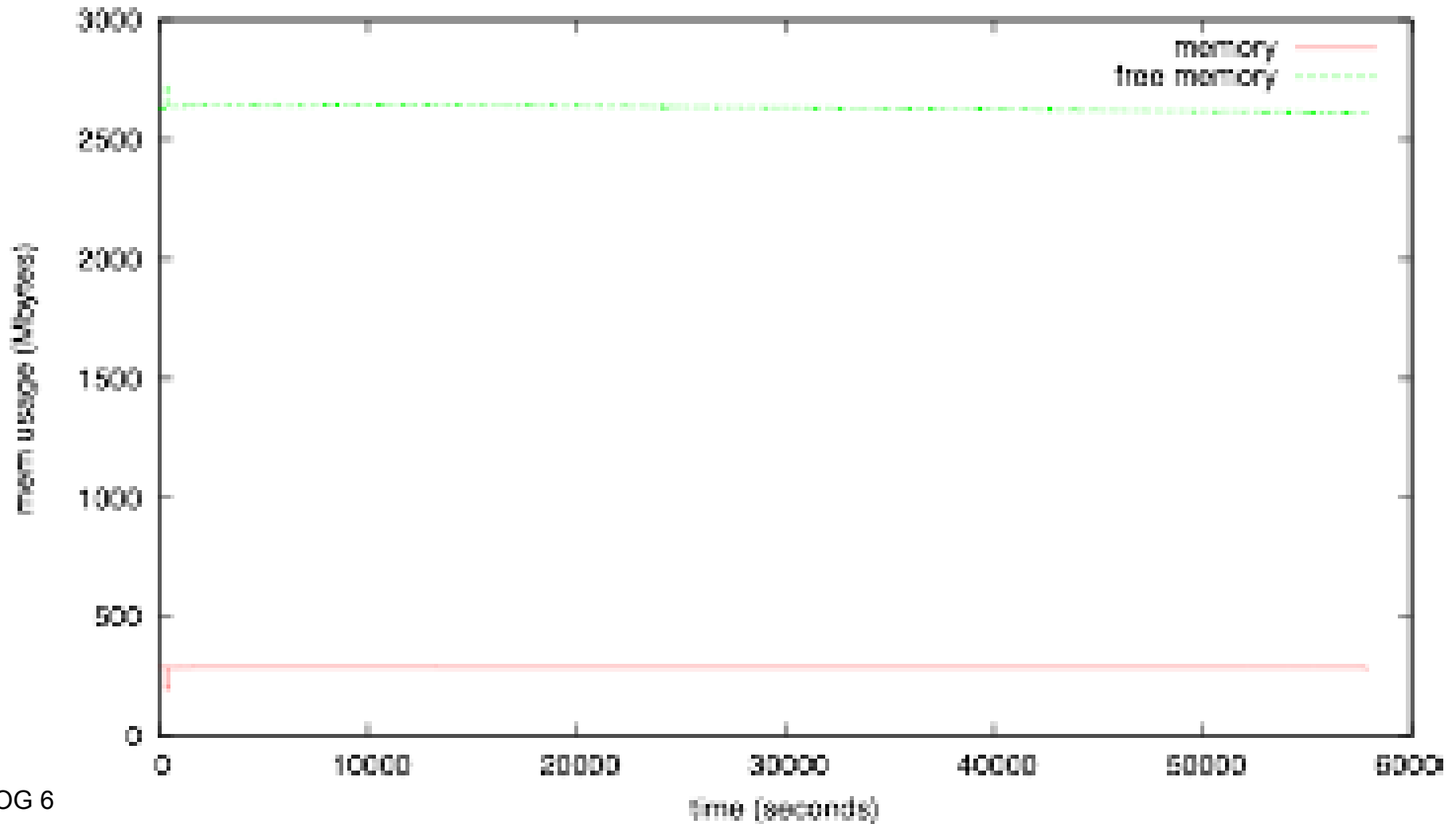
# Quagga Mem

Quagga mem usage  
multiple-rib; 100 sessions  
IXIA vs. Quagga  
500 prefixes per session with random flapping  
lab4  
2 x Intel(R) Xeon(R) CPU 3050 @ 2.13GHz



# BIRD Mem

BIRD mem usage  
multiple-rib; 100 sessions  
IXIA vs. BIRD  
500 prefixes per session with random flapping  
lab6.pax.net  
4 x Intel(R) Xeon(TM) CPU 3.80GHz



Thank you!  
Questions?

<elisa.jasinska@ams-ix.net>

<cmalayer@switchanddata.com>

<arnold.nipper@de-cix.net>