



**Open Source Routing**  
**KTH CSD Kick-Off Workshop**

Robert Olsson Uppsala University

2008-09-02

# Why Open Source?

- Reclaim research and development to universities etc
- To be a part in the development loop
- Open for wide collaboration
  - No national boundaries
  - No organizational boundaries
- Easy experimentation to prototype new ideas
  - Next-Generation Internet take-off
  - Other ideas we can't even think of right now

# Why Open Source?

- Possibilities for superior quality  
Work can be reviewed by many people
- Very fast development can be achieved
- Process can be independent from business or politics
- Non-discriminate
- Economical possibilities
- Idea started in computer science

# Relation to Open Source

- You are getting other people's work for "free"  
Respect
- Open Source does not work without contributions  
Compare a relay race. Reuse and recycle work.
- Open Source has strong momentum  
Business models are developed etc

# Open Source Networking Now

- Interesting suitable hardware
  - Technological breakthrough
  - Multi-Core CPU , other silicon
  - Fiber Optics
  - Fast buses PCI-Express
- There are interesting applications
- Open source OS has come a long way

# Open Source Competitors



- MIT click modular router
- Berkeley, CA  
XORP
- Vyatta

# Over 10 years in production

## UU facts

Dual ISP BGP connect GIGE

Local BGP peering GIGE

Ipv4/Ipv6

OSPFv2/OSPFv3

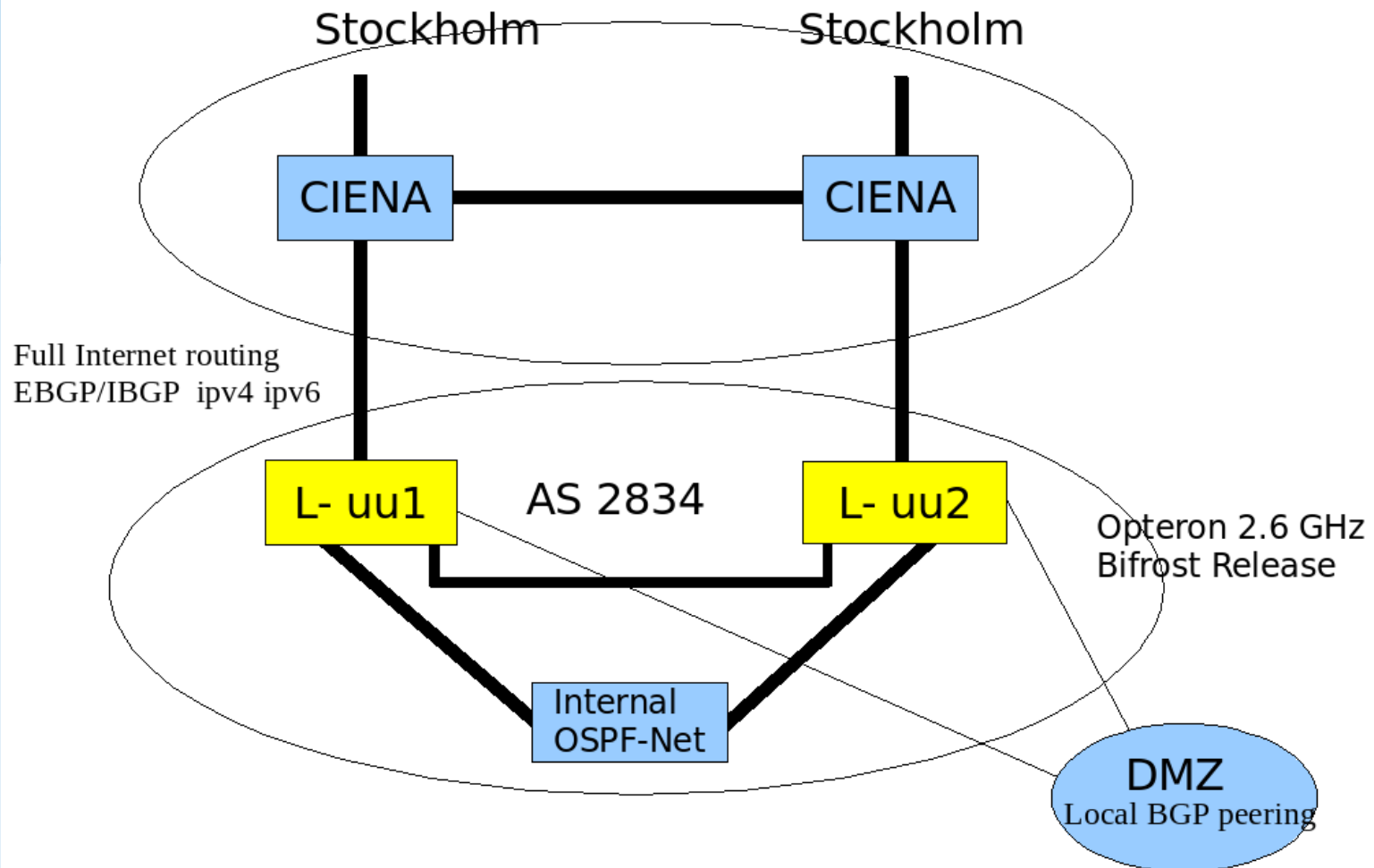
netfilter

Cisco 6500

10g planned.

# Over 10 years in production

## BGP topology at Uppsala University



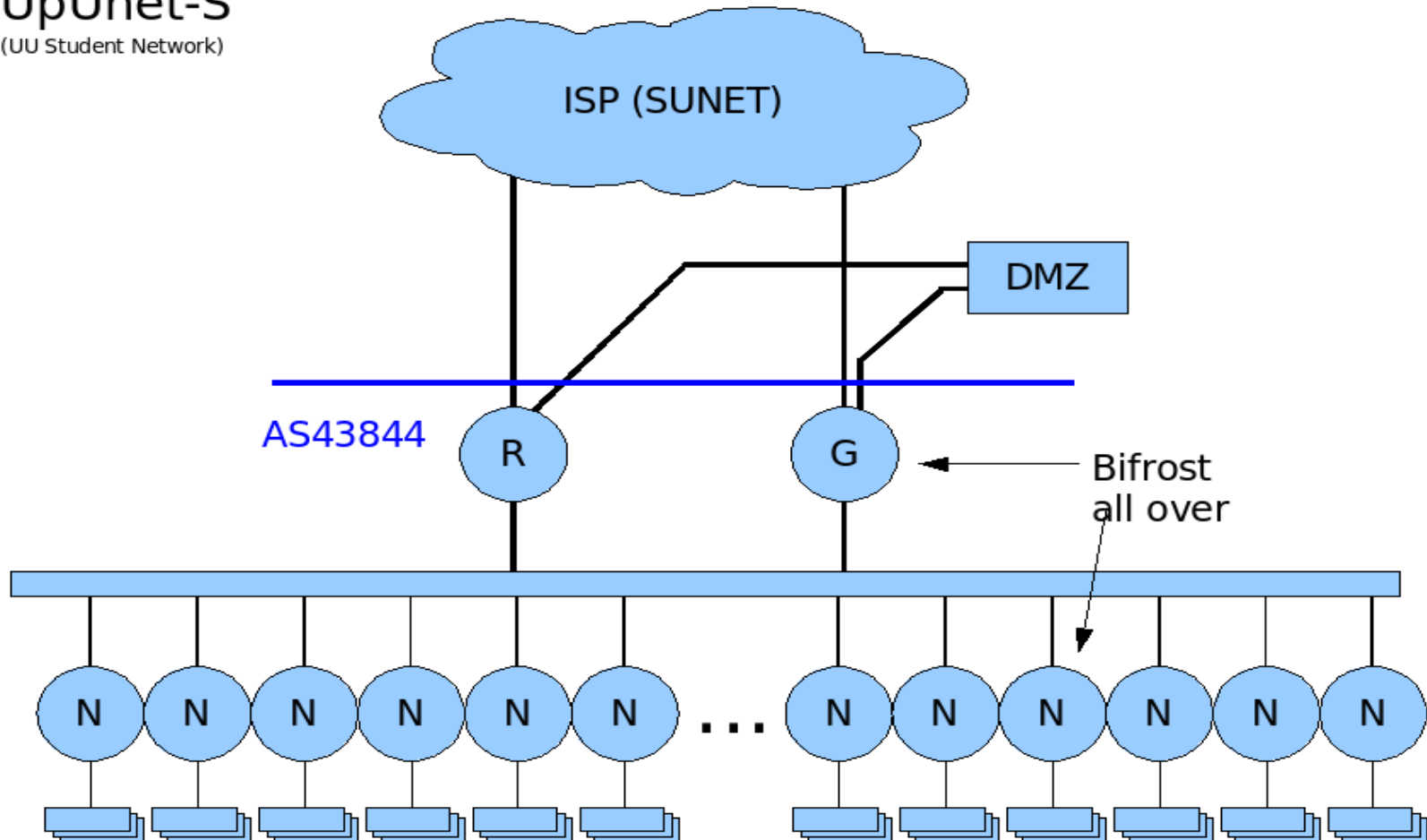


# Over 10 years in production

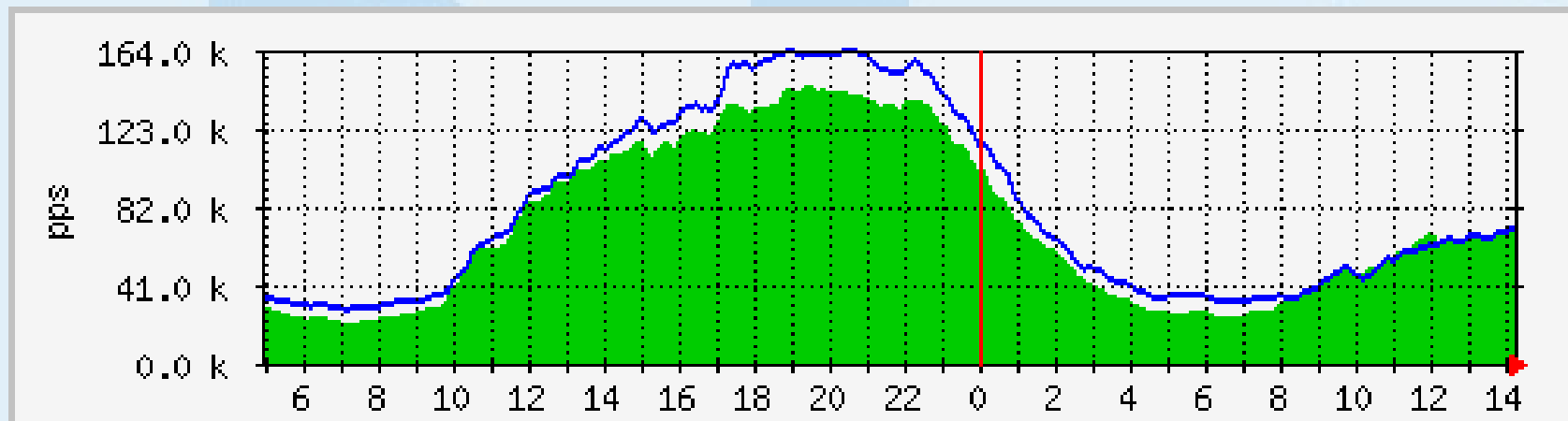
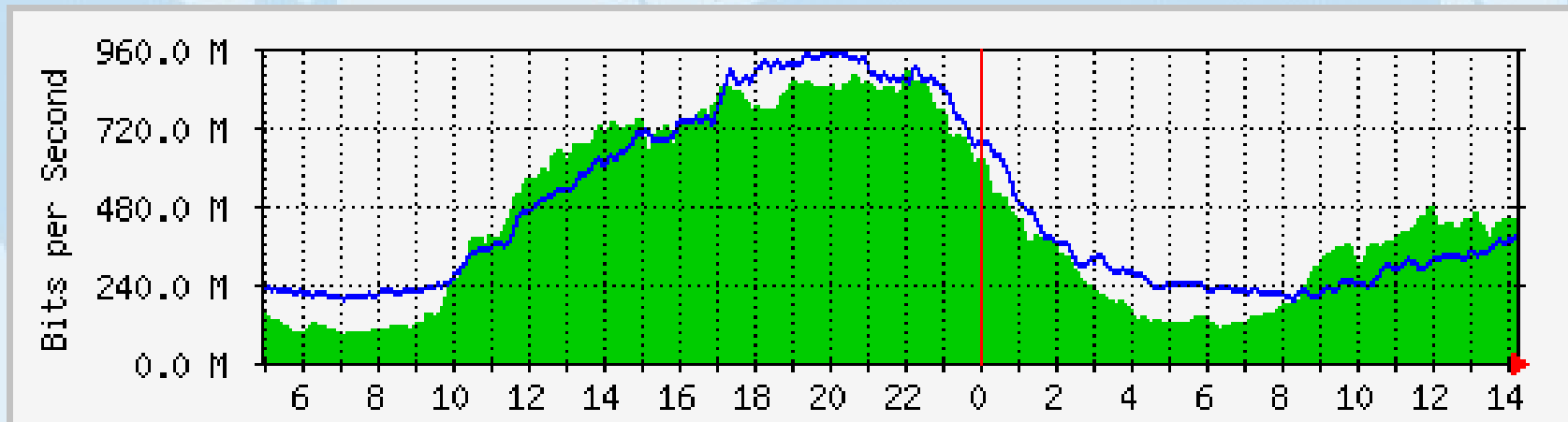
- Three major installations
- UU core routers towards SUNET
- UU StudentNetwork 30.000 students
- `ftp.sunet.se`

# Over 10 years in production

UpUnet-S  
(UU Student Network)



# Over 10 years in production Student Network Core Router



# Over 10 years in production

## Student Network facts

Dual ISP BGP connect GIGE

Local BGP peering GIGE

Ipv4

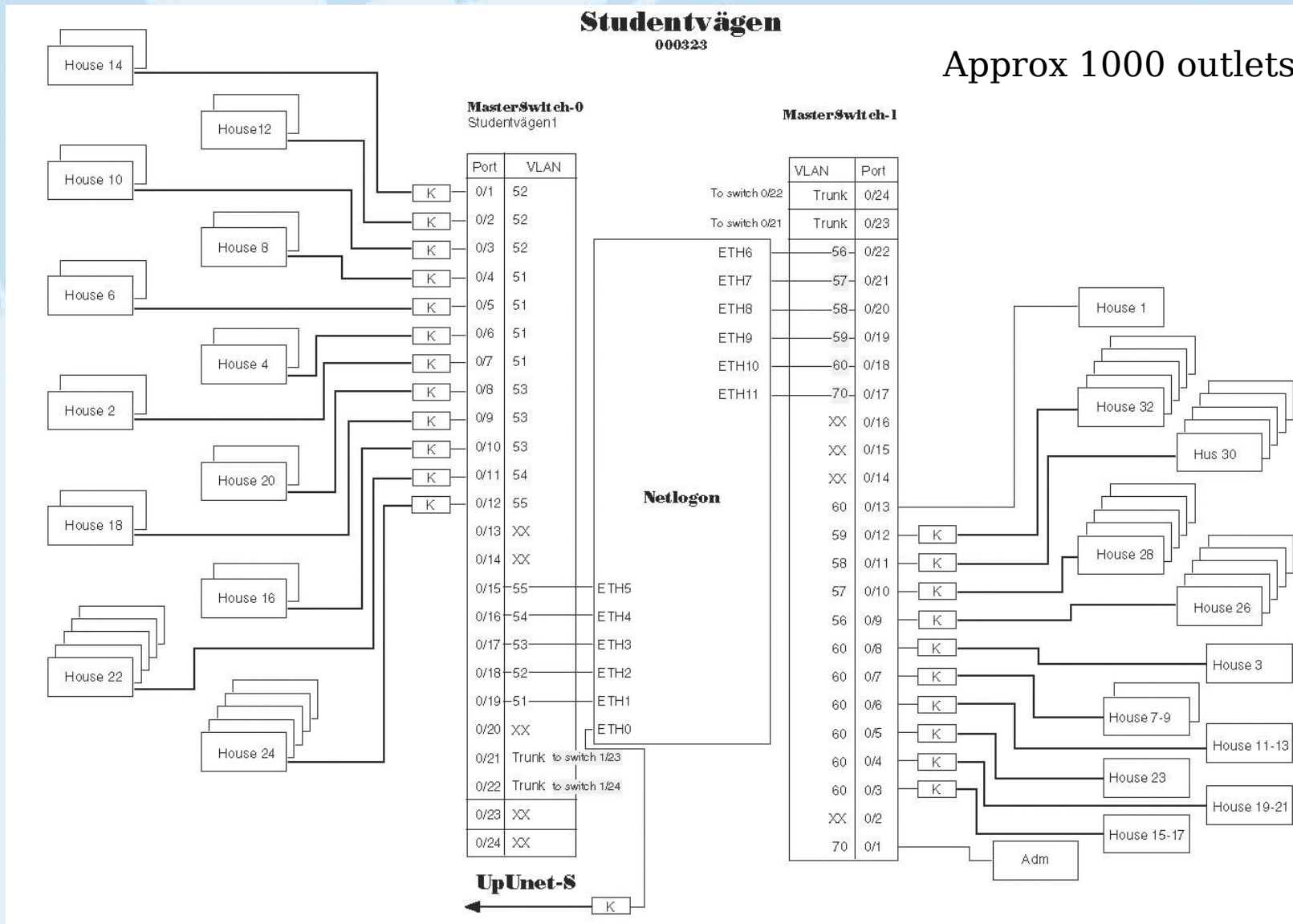
OSPFv2

xxx netfilter rules

netlogin-service at premises

10g planned.

# IP-login installation at Uppsala University





# Testing, Verification Development & Research

- Started out as simple testing.
- Curiosity, Open Source, Collaboration
- Relatively freedom, the idea to use in own infrastructure. No need for external funding.
- OS was intended for desktops.

# Testing, Verification Development & Research

No need for test network. We could test in own infrastructure. (Or SLU)

Problem oriented vs Project oriented

We could work on complicated issues

- NAPI 3 years
- pktgen 2 years
- fib\_trie 1 year
- TRASH 1 year

# Building Blocks

## Hardware:

PC

Motherboard/CPU/Memory

Network Interfaces

GIGE/10gWiFi/etc

## Software

Operating System

Linux/BSD/Microsoft

Applications

Routing Daemons

Quagga/XORP

IP-login/netlogon

## Network

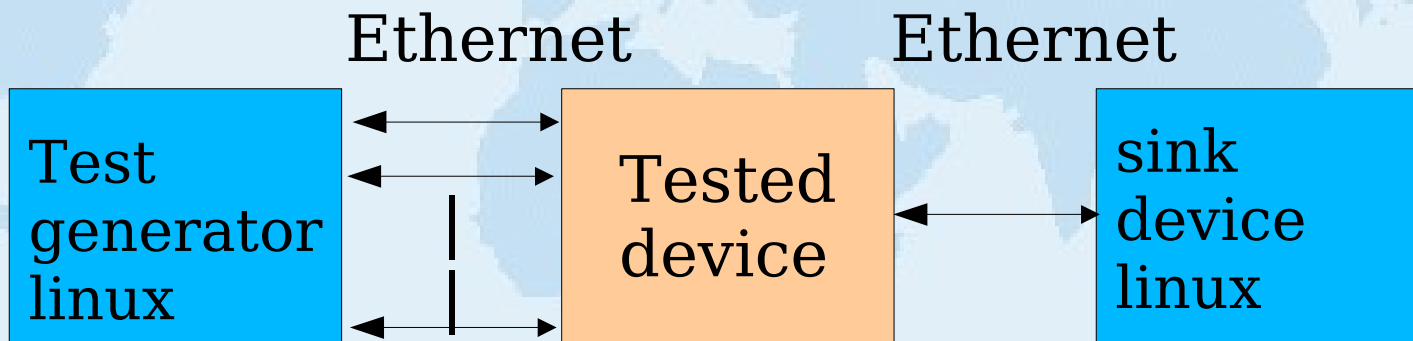
Cable, Fiber, Copper

Equipment, Switches



# Flexible netlab at Uppsala University

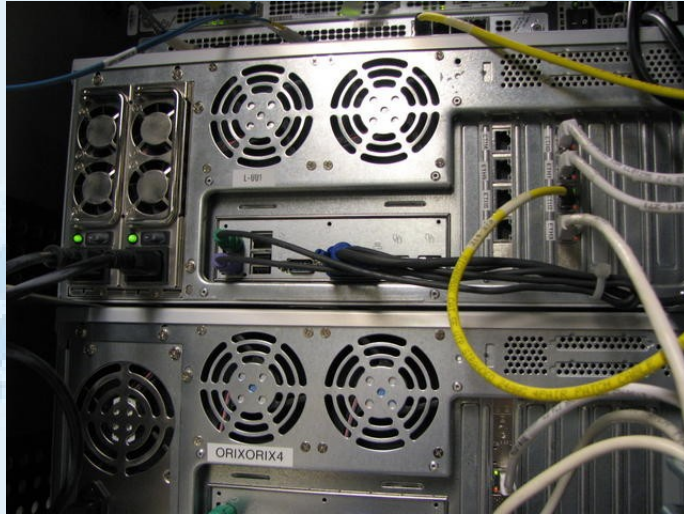
El cheapo-- High customizable -- We write code :-)



- \* Raw packet performance
- \* TCP
- \* Timing
- \* Variants

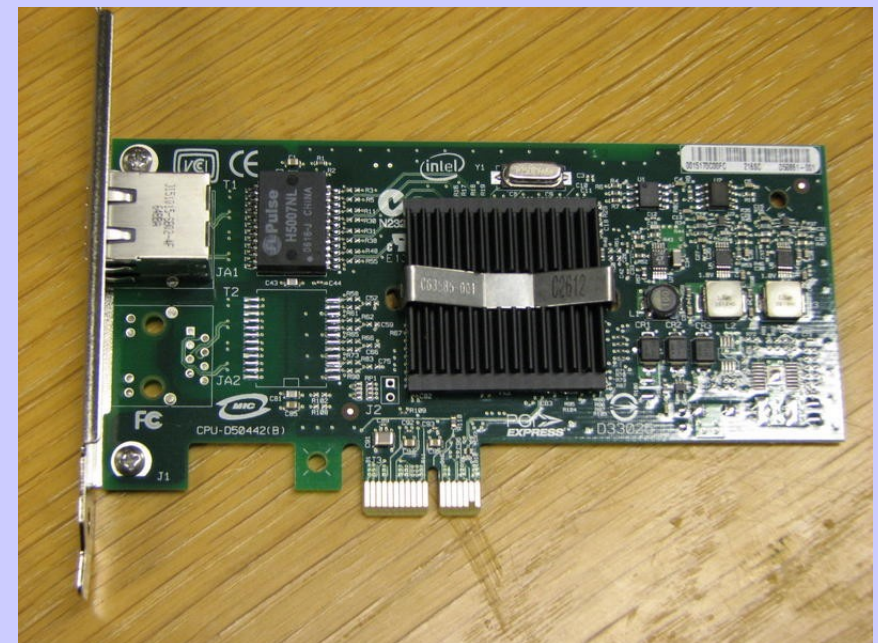
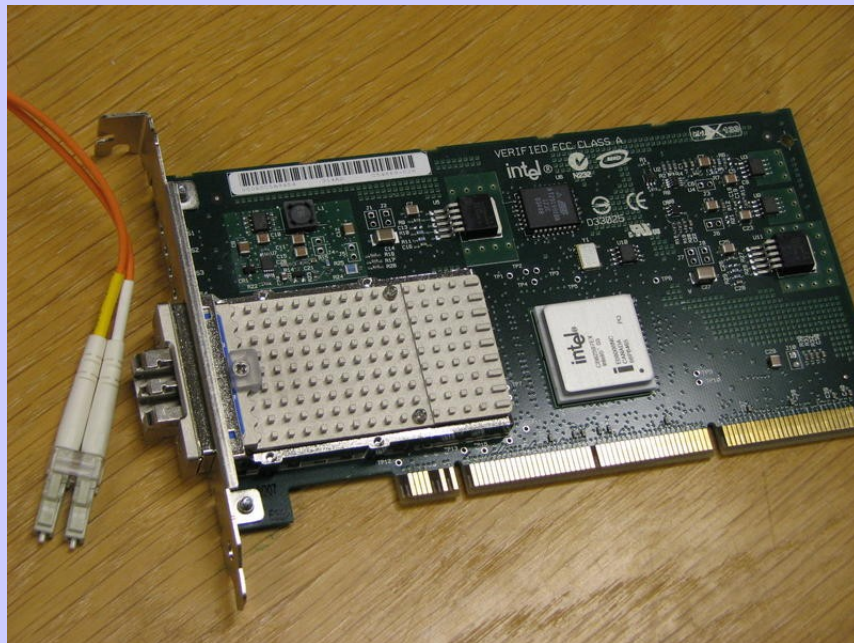
# *Lab*







# Intel NIC's





*Not or were blessed ...*



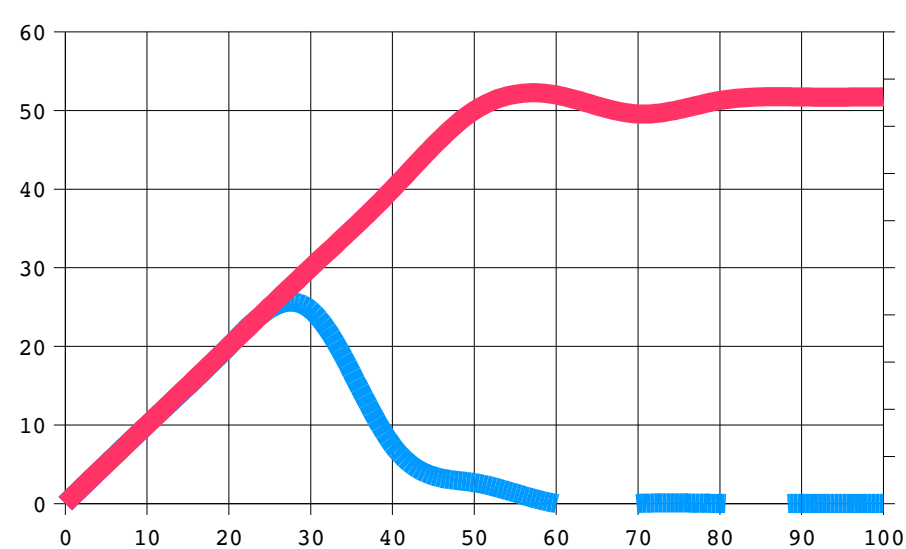
# Bifrost concept

- Linux kernel collaboration
- Performance testing, development of tools and testing techniques
- Hardware validation, support from big vendors
- Detect and cure problems in lab not in the network infrastructure.
- Test deploy (Often in own network)

# Overall Effect

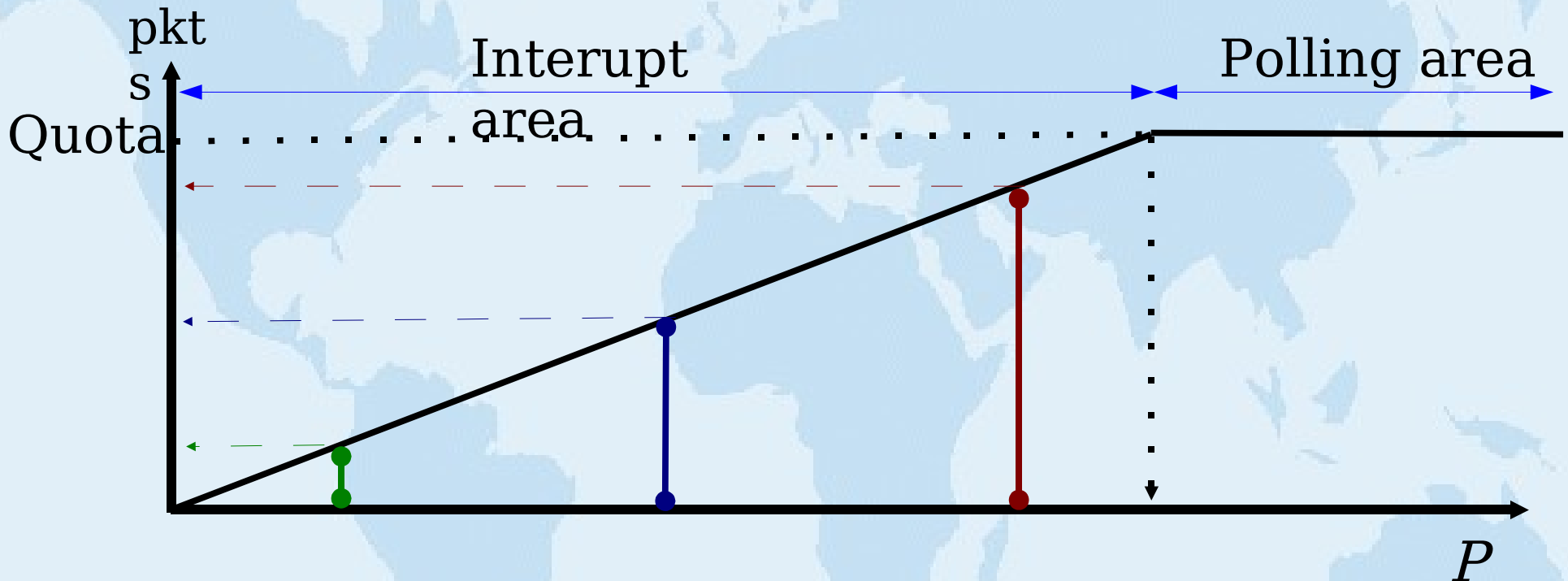
- Inelegant handling of heavy net loads
- System collapse
- Scalability affected
- System and number of NICs
  - A single hogger netdev can bring the system to its knees and deny service to others

Summary 2.4 vs feedback



March 15 report on lkml  
Thread: "How to optimize routing performance"  
reported by  
[Marten.Wikstron@framsfab.se](mailto:Marten.Wikstron@framsfab.se)  
- Linux 2.4 peaks at 27Kpps  
- Pentium Pro 200, 64MB RAM

# A high level view of new system



- $P$  packets to deliver to the stack (on the RX ring)
- Horizontal line shows different netdevs with different in
- Area under curve shows how many packets before next
- *Quota* enforces fair share



# Kernel support

NAPI kernel part was included in:  
2.5.7 and back ported to 2.4.20

Current driver support:

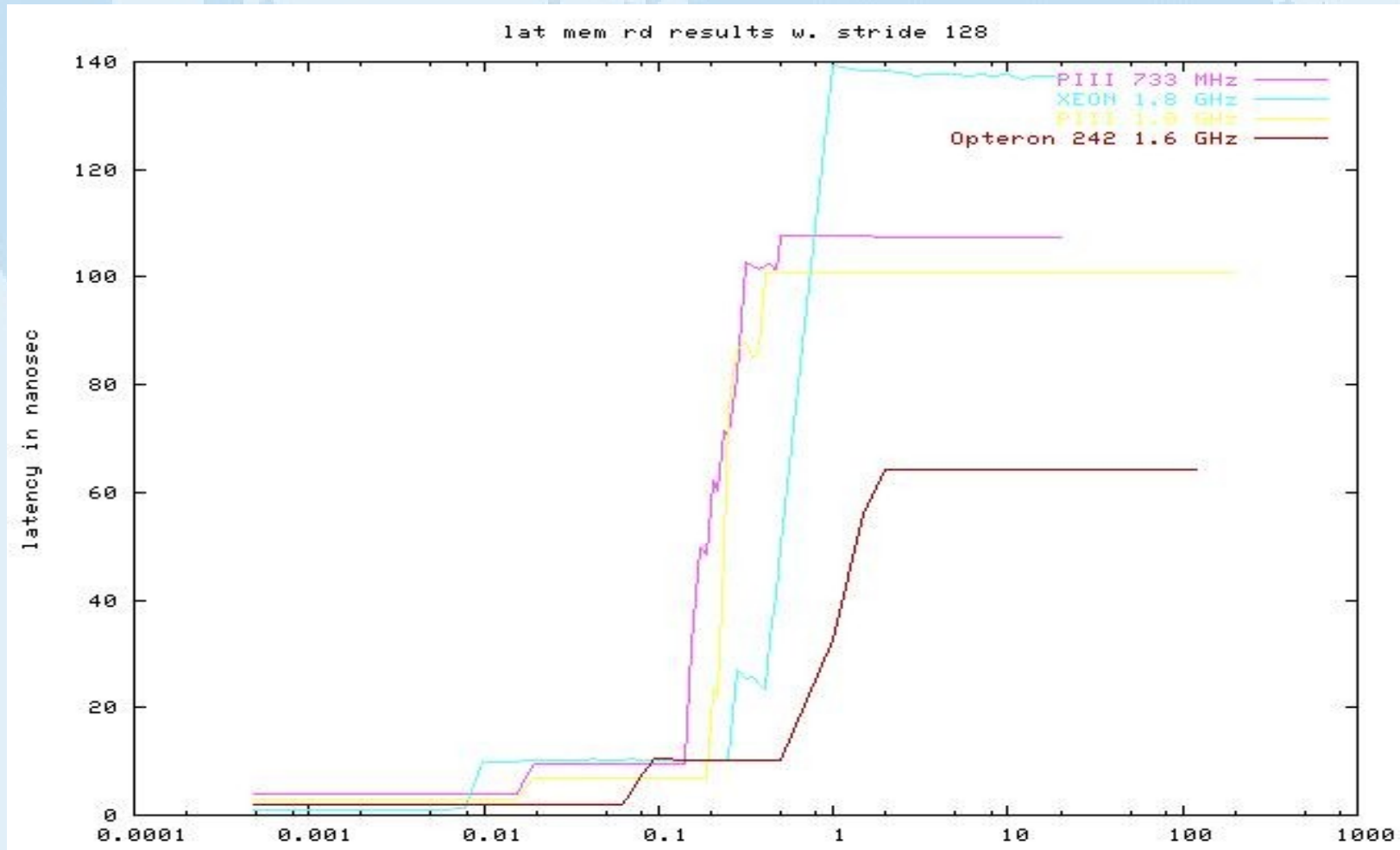
e1000 Intel GIGE NIC's

tg3 BroadCom GIGE NIC's

dl2k D-Link GIGE NIC's

tulip (pending) 100 Mbs

# Cache effect/Performance



# Cache effect/Performance

Cache line 32 – 128 bytes

Optimize struct for cache and multiprocessors  
usage

PIO even worse than cache miss

PIO READ stalls CPU

PIO WRITE can be posted

DMA copies of data into RAM

Does prefetch solve problems?

A new network symbol has been seen...

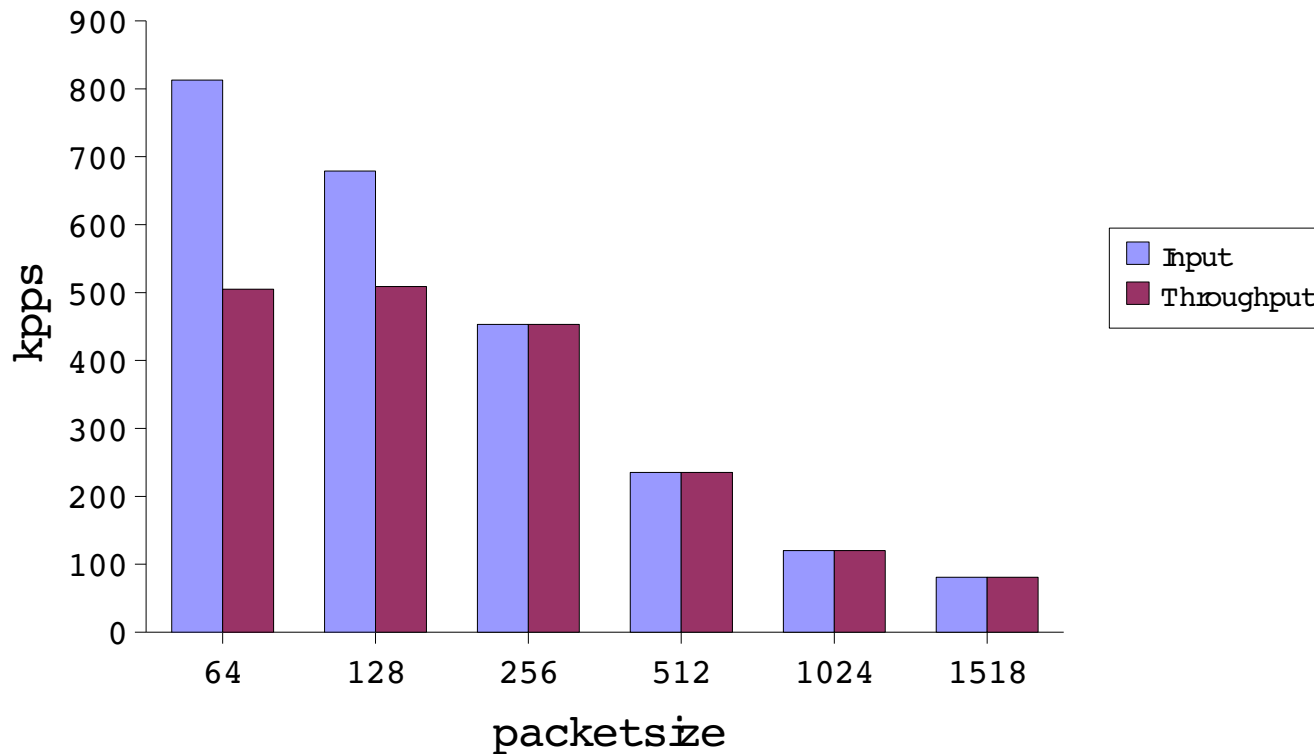
*The Penguin Has Landed*



# Forwarding performance

Linux forwarding rate at different pkt sizes

Linux 2.5.58 UP /skb recycling 1.8 GHz XEON



Fills a GIGE pipe -- starting from 256 byte pkts

# Other activities

## informal linux agenda

Ericsson is willing to open patent for Linux  
Jamal have the contacs via Ericsson Montreal

DaveM has discussions with Washington university  
about who is willing to grant another patent for use  
with Linux

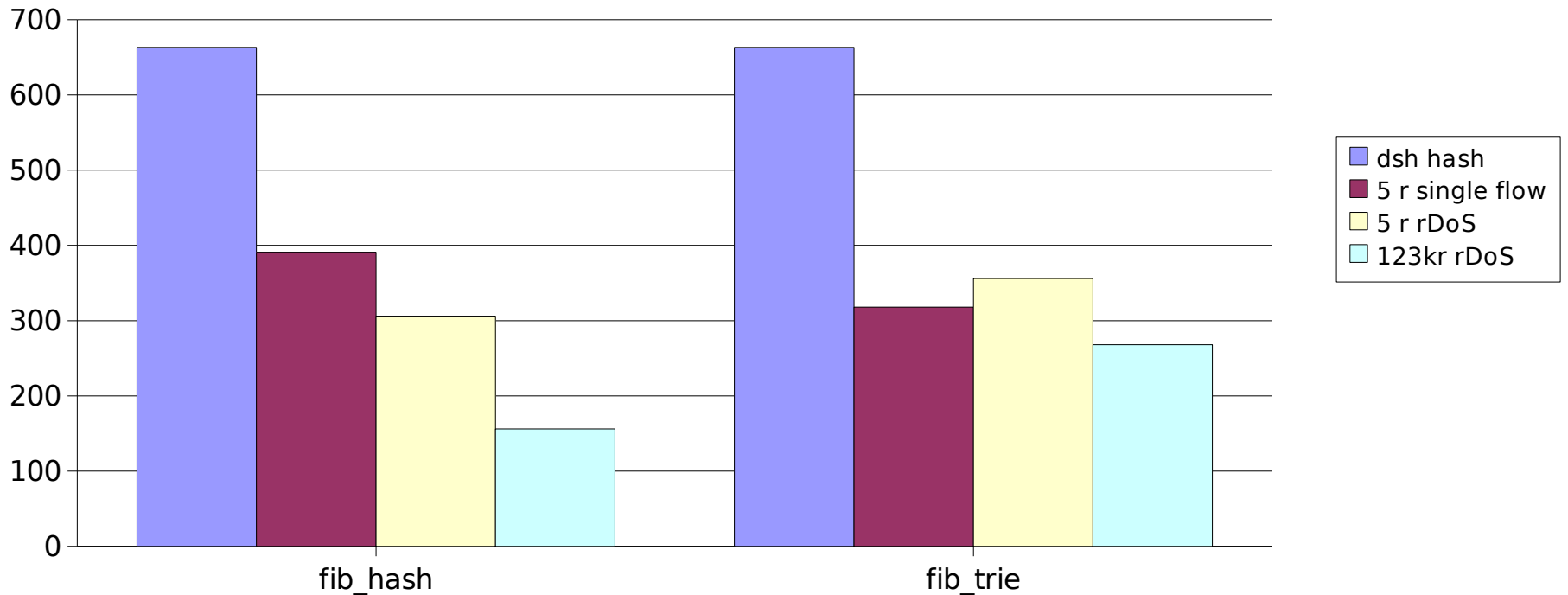
Discussed LC-trie with Alexey Kuznetsov.

LC-trie investigations. Got GPL from authors.

# fib\_trie performance comparison

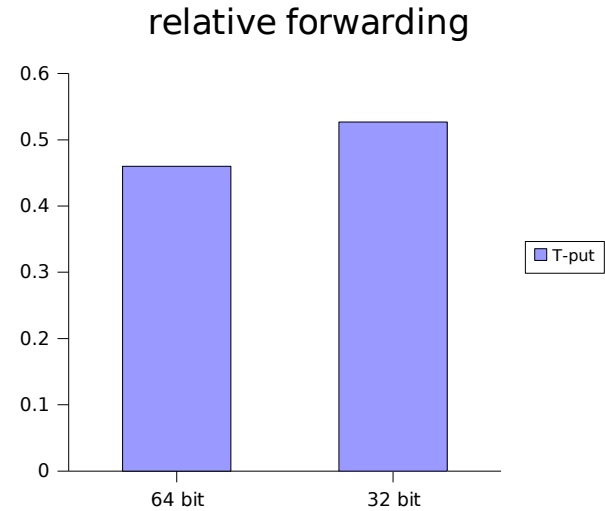
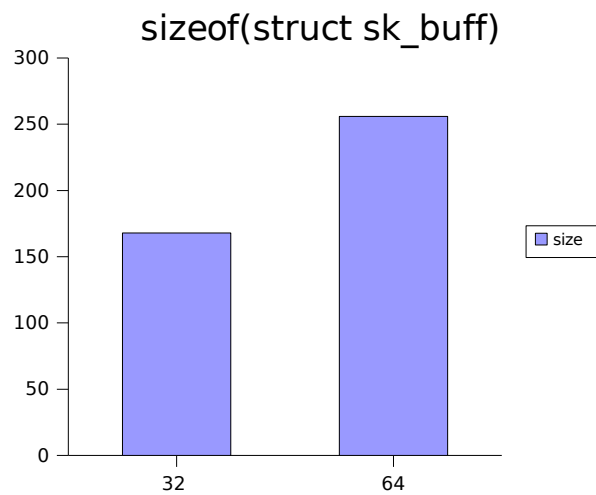
forwarding kpps

Linux 2.6.16 1 CPU used(SMP) Opteron 1.6 GHz e1000



Preroute pathes to disable route hash

# 32/64 bit || sizeof(sk\_buff)



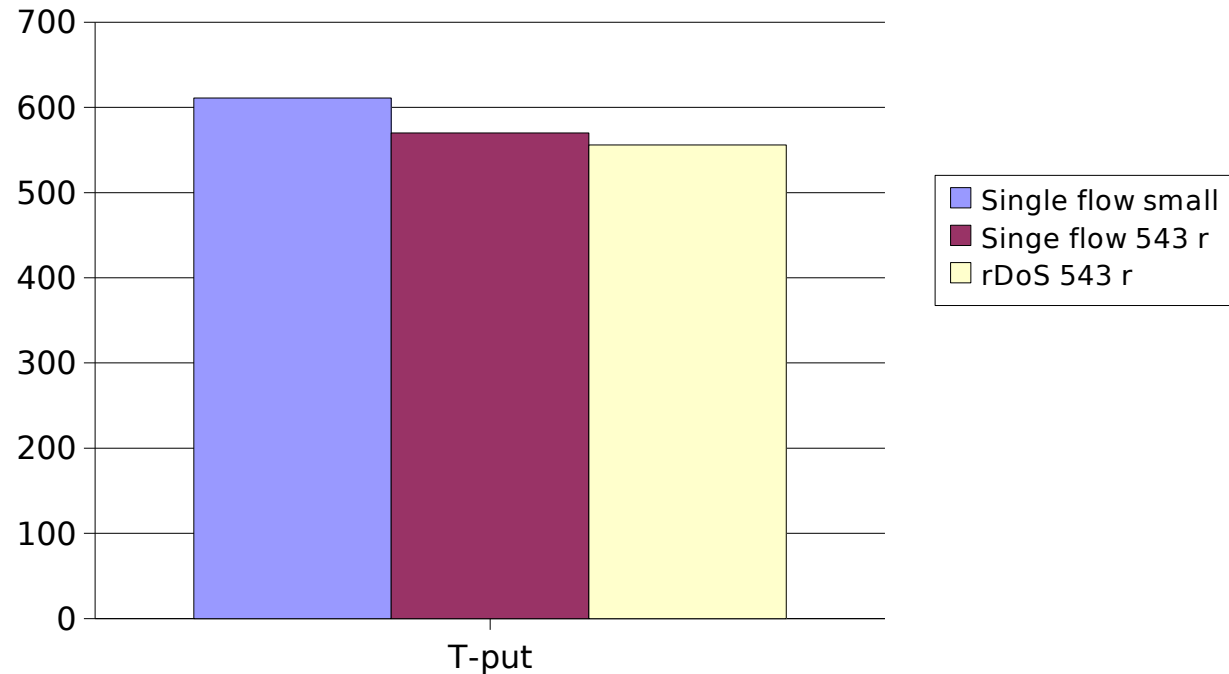
Gcc 3.4 x86\_64 vs i686 on same HW



# ipv6 performance

Forwarding kpps 76 byte pkt.

Linux 2.5.12 1 CPU(SMP) Opteron 1.6 GHz e1000



How rDoS work on sparse routing table?

# *Trash data structure*

Interesting novel approach. Trie-Hash → Trash

When extending the LC-trie

Paper with Stefan Nilsson/KTH

Exploits that key len does not affect tree depth

When lengthen the so key it can be better compressed.

Implemented in Linux forwarding patch as a replacement to the route hash.

# *Trash datastructure*

Can do fullkey bokup . src /dst /sport /dport /proto /if etc and later socket.

For even ip6 with little performance degradation

Could be a candidate for the grand unified bokup

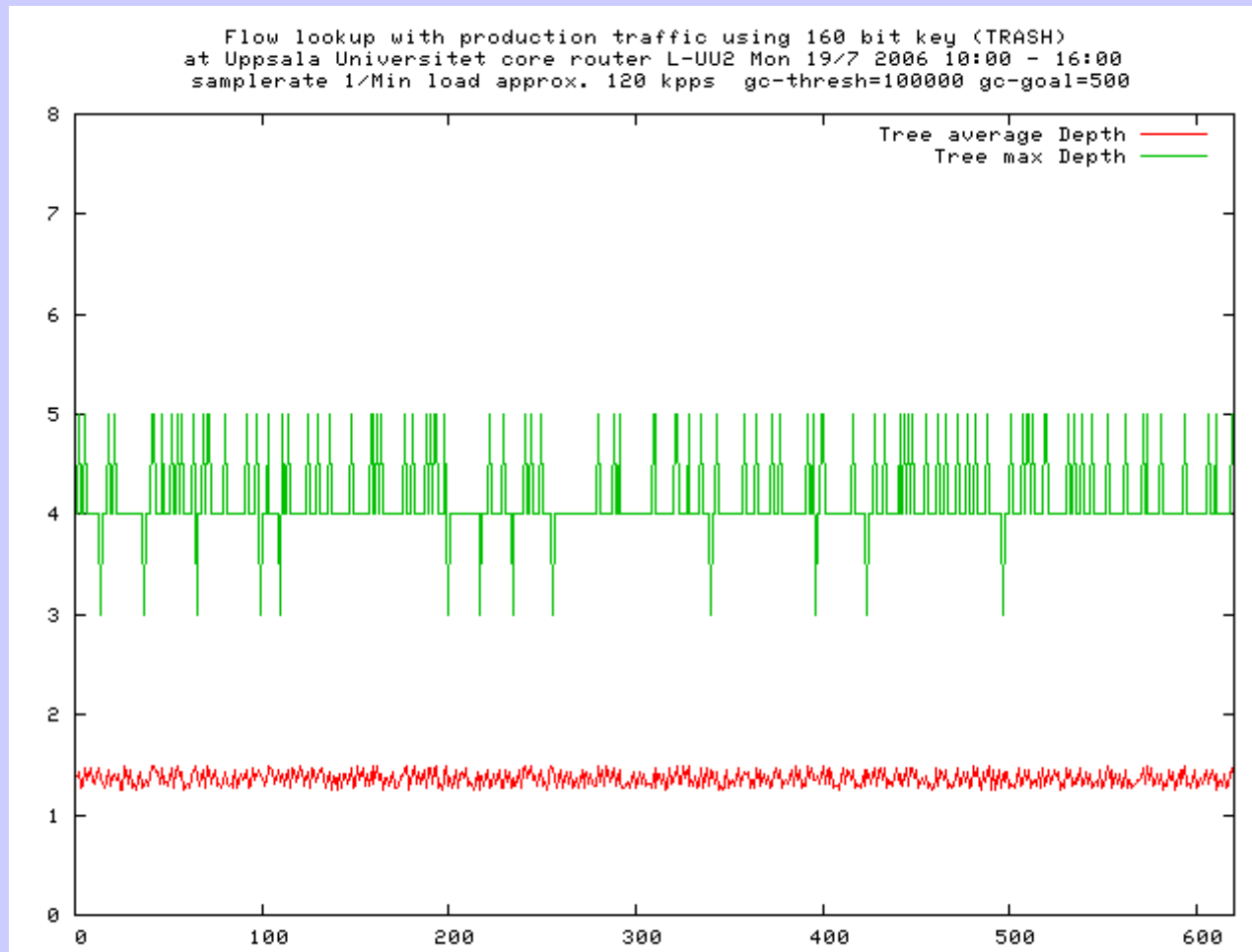
Fullfw bokup can understand connections.

Free fw bgging etc

New garbage collection (GC) possible. Active GC stated AGC in the paper. Listen to TCP SYN, FIN and RST Show to be performance winner.

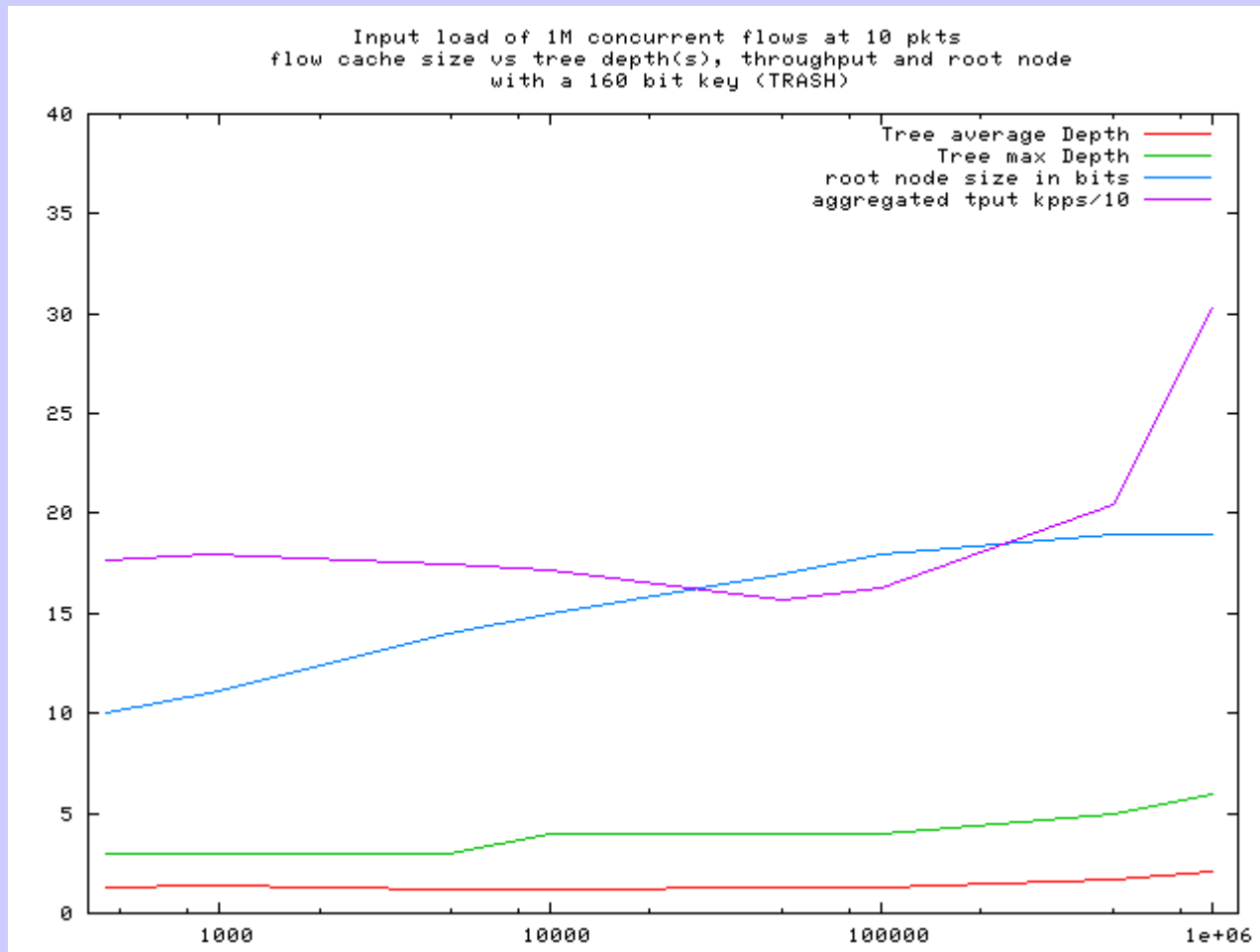
# Trash datastructure

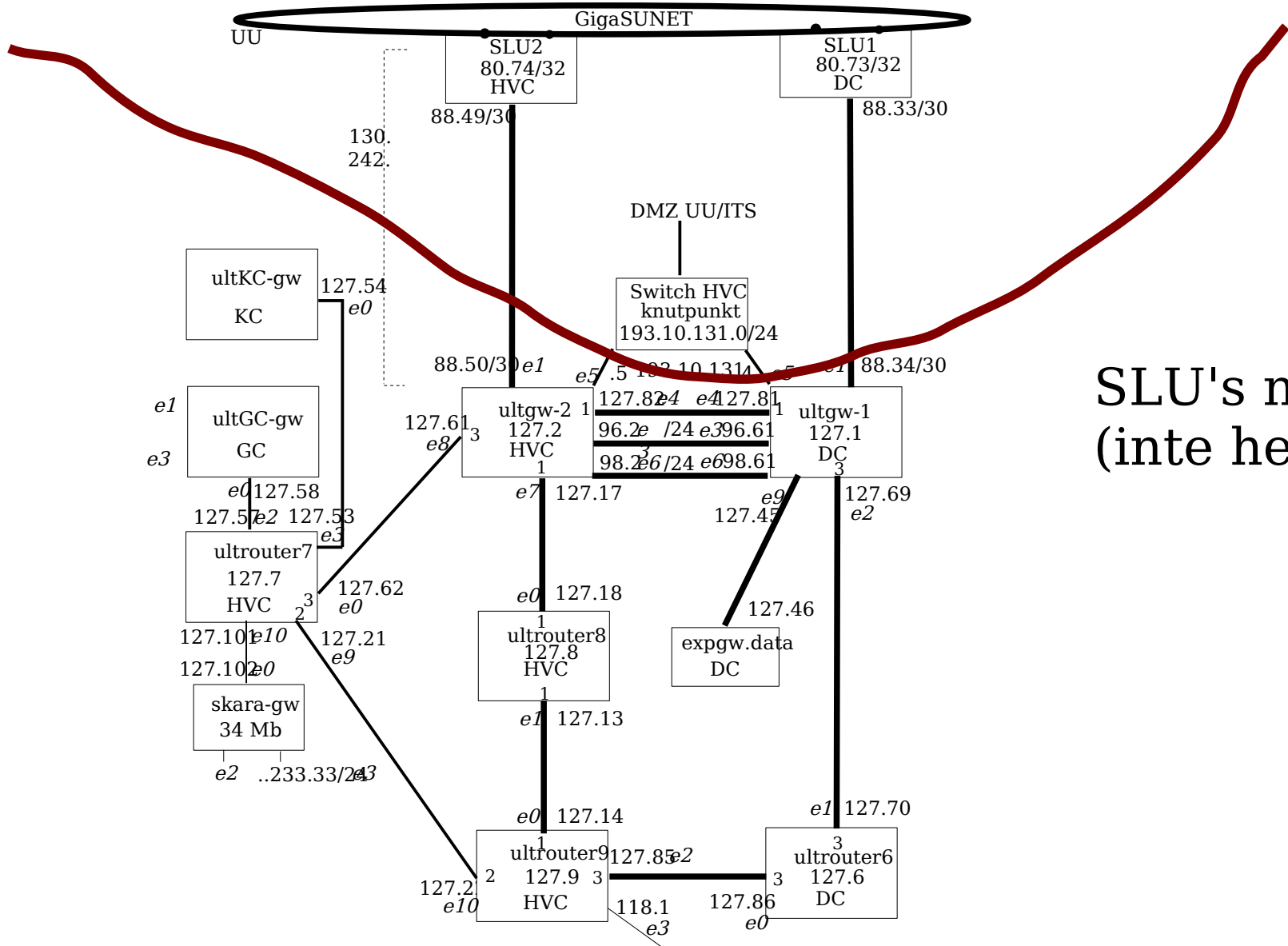
## Uppsala Universitet core router



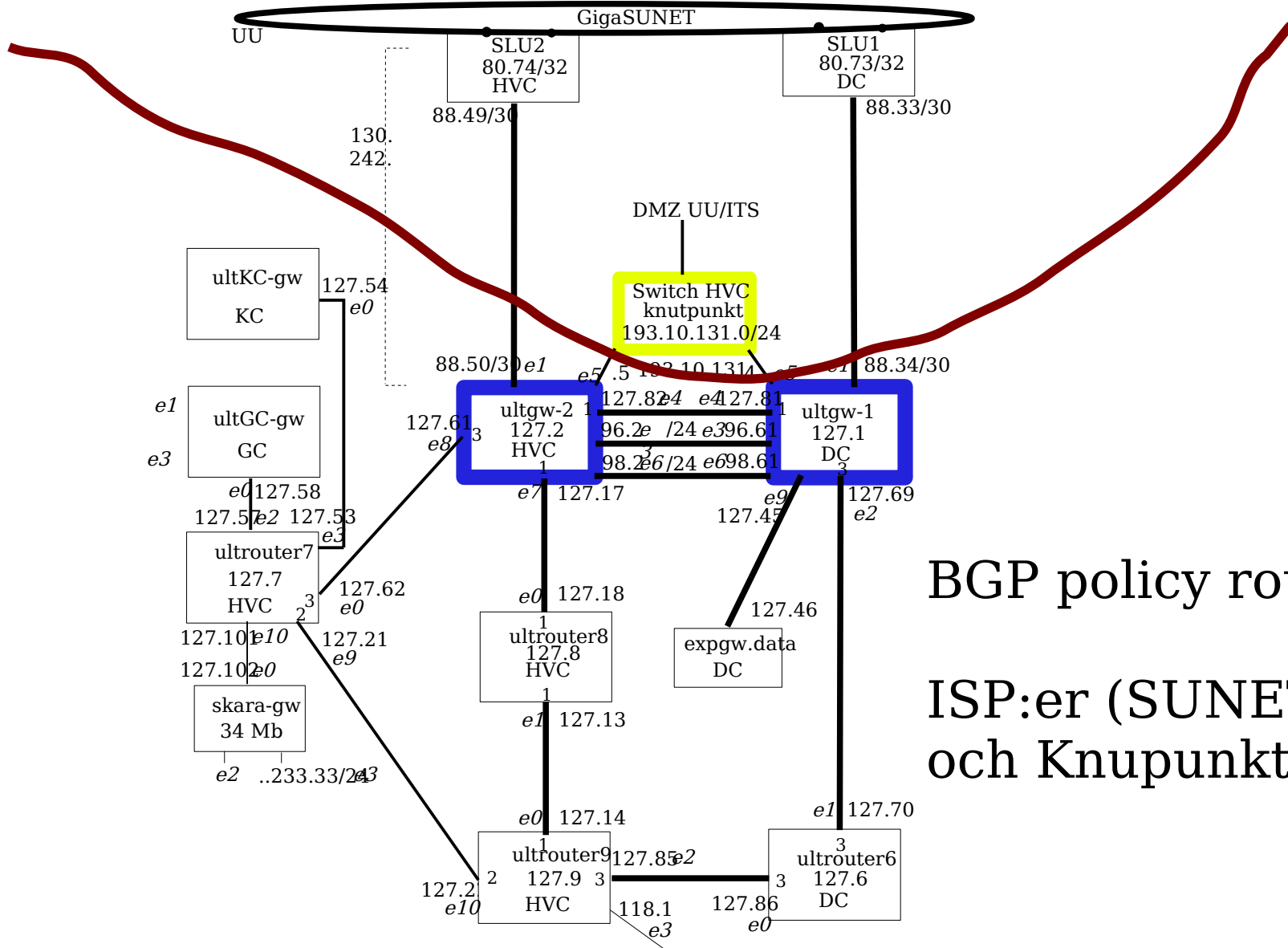
# Trash datastructure

## Very flat(fast) trees



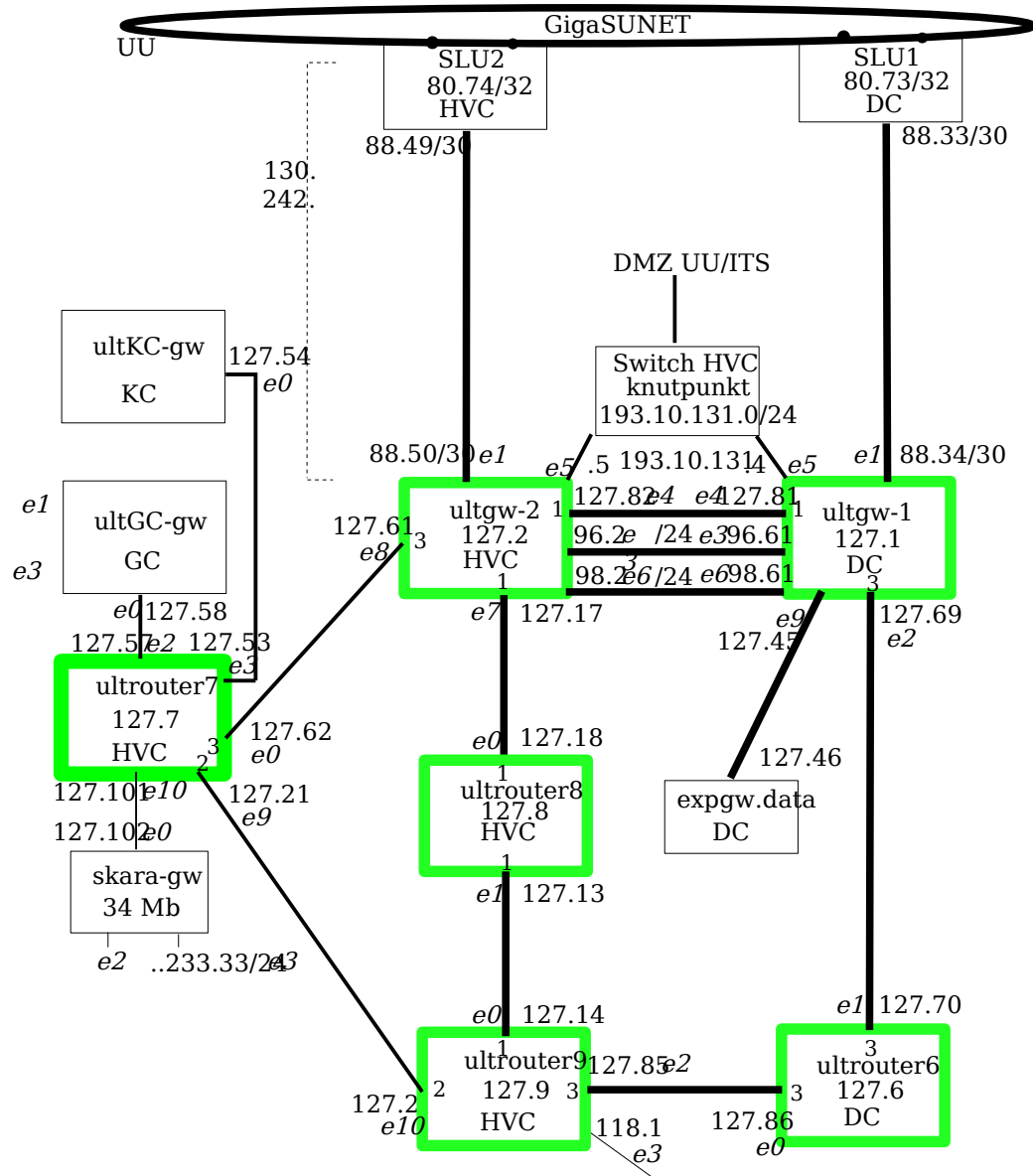


SLU's nät  
(inte hela)



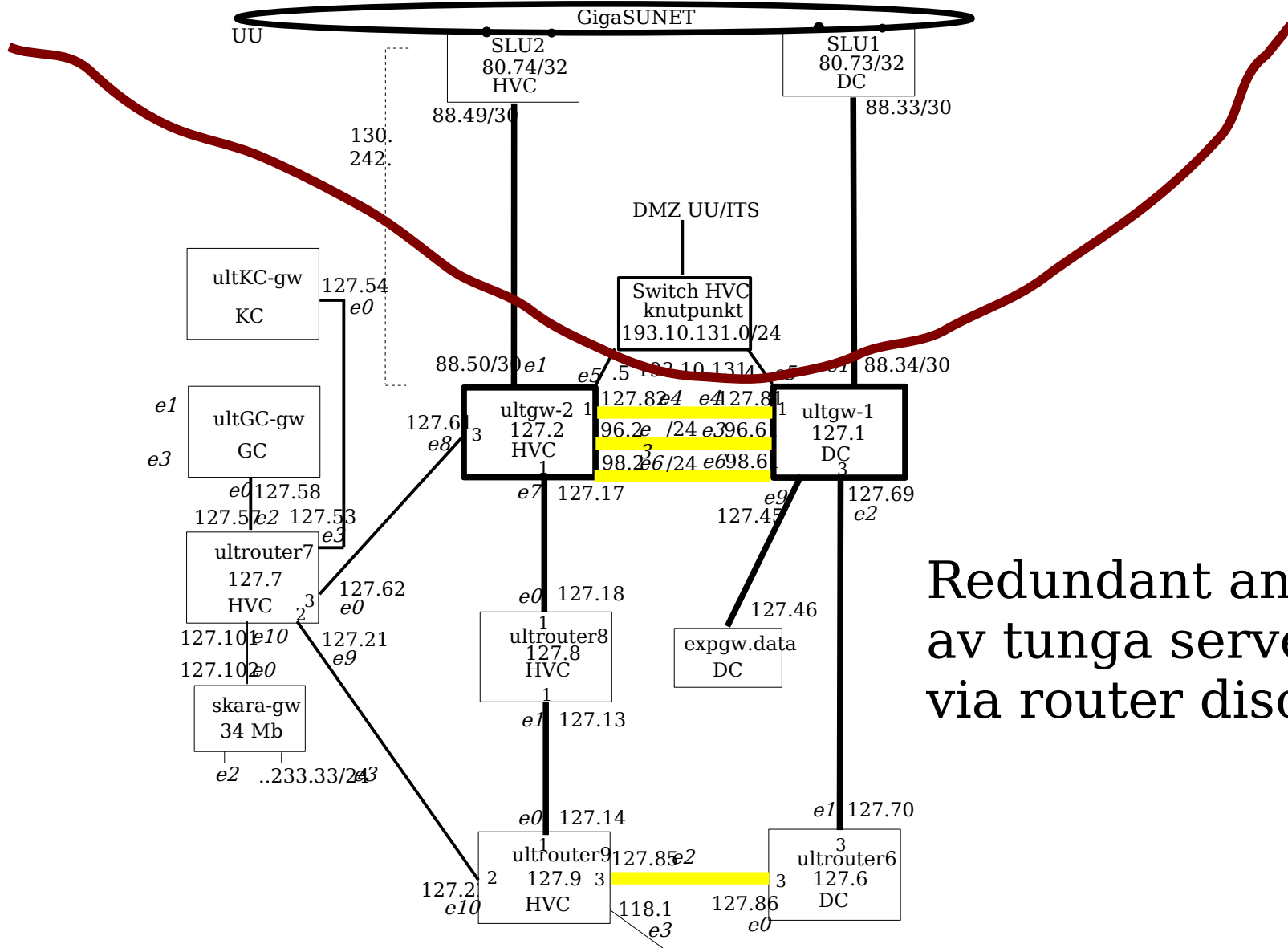
BGP policy routing

ISP:er (SUNET)  
och Knutpunkt.

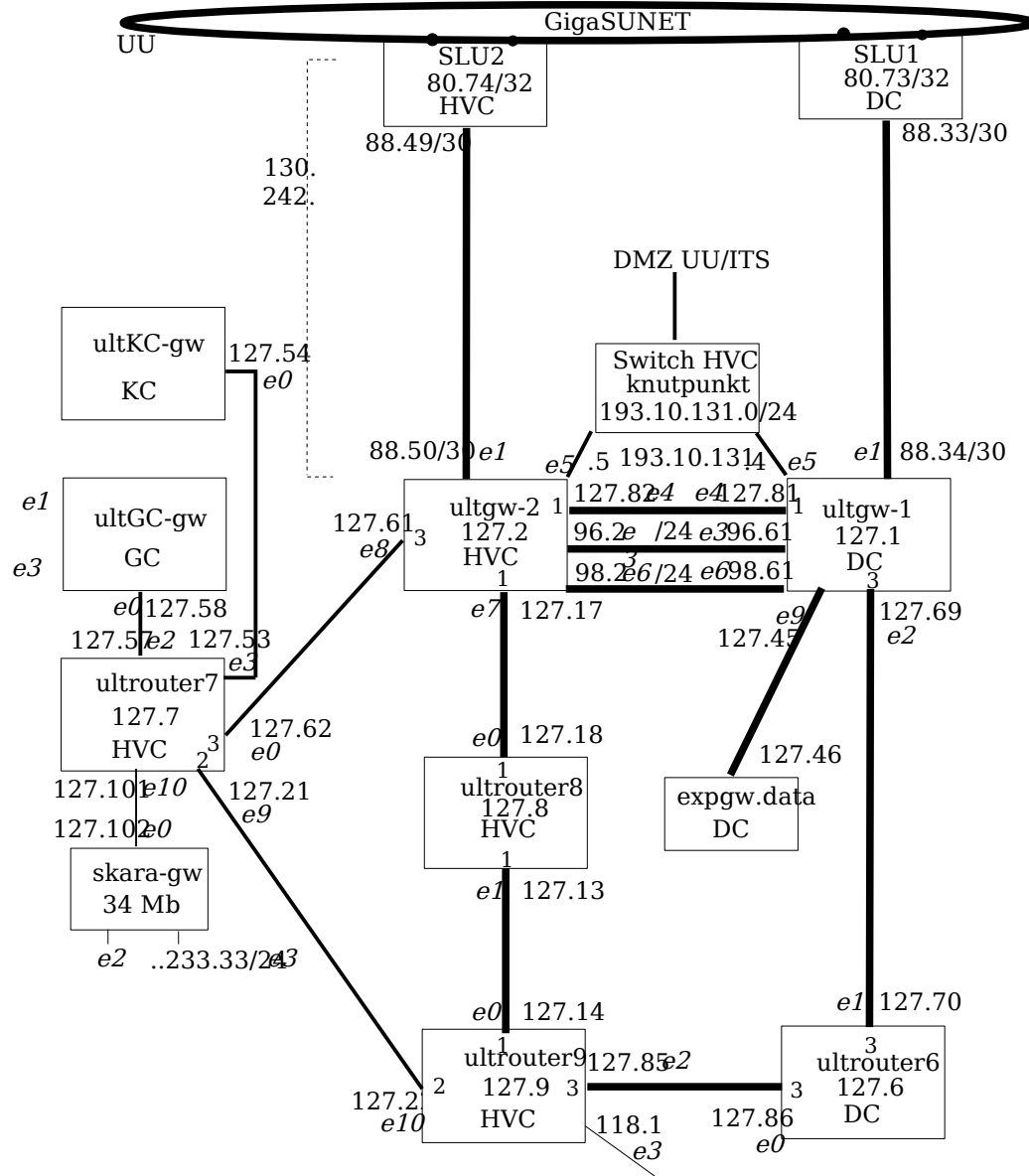


Redundant  
inre kärna





Redundant anslutning  
av tunga servernät  
via router discovery



Almost all...

OpenSource  
bifrost  
zebra/quagga