# Going with the (s)Flow at 200G

Richard Yule<richard@linx.net>

Nigel Titley<nigel@titley.com>
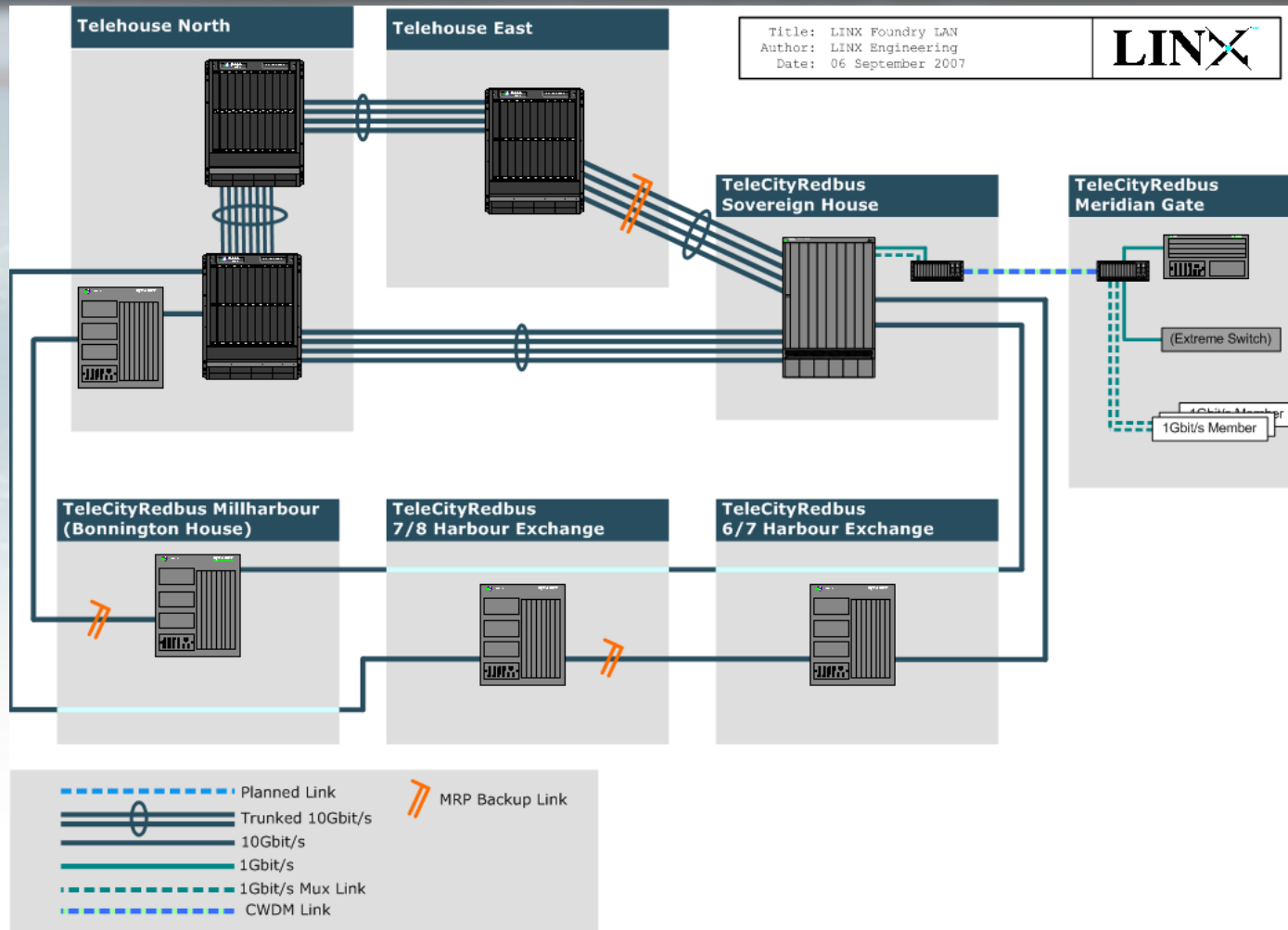
Mike Hughes <mike@linx.net>

**LINX**

# LINX Overview

- Founded in 1994 by 5 Members
- Non-profit, Mutual Ownership
- Now
  - 280 Members from 43 countries
  - 60% of the routing table peered
  - ~240Gbps peak (5 minute average)
  - Over 100 managed private interconnects
    - For large traffic flows

LINX

# LINX Architecture

- Dual LAN Architecture
  - One LAN using Foundry switches
  - One LAN Extreme switches
- 7 sites in London Docklands
  - Connected by multiple diverse fibre rings
  - 8x10GE trunk ISL between top sites
  - 3 new sites to be connected in 2008
- Bi-lateral or multilateral peering
- 100M, 1GE and 10GE member ports

LINX

# LINX Foundry Network



**Engineering**

# Enter sFlow

- What is sFlow?
  - Defined in RFC 3176
  - A means of taking sampled traffic data from within a network
  - Works in Layer 2 networks (e.g. IXP)!
- sFlow agent on switch/router sends sFlow datagrams (UDP) to a sFlow collector
- sFlow collector runs analysis software

LINX

# Purposes of LINX sFlow project

- Provide member to member statistics
  - For use by LINX members
  - For use by LINX engineering staff
- Provide engineering staff with tools such as
  - Traffic matrix (i.e. between nodes)
  - Peering matrix
  - Spot traffic anomalies
- More intelligence about our network

LINX

# Challenges

- Most existing sFlow tools didn't do what we wanted
  - Commercial: expensive, inflexible
  - Too Lightweight: couldn't handle the data
  - Too Exhaustive: tried to extrapolate the data down to a specific L3 flow, scaling issues
  - Incomplete: only did part of the job

- Most importantly, little or no concept of "a member"

- We had to write our own tools

# Phase 0.9: Proof of Concept

- First attempt was done using sflowtool feeding to pmacct

- Write the output into RRD files

- Problems
  - Constrained by disk I/O
  - Produced large unwieldy page of graphs

- Not very flexible
  - Borne out by minimal use

**LINX**

# Switch Config

- All switches send sFlow packets over a VLAN interface in a specific sFlow VLAN – removes mgmt i/f concerns

- sFlow collector has interface to VLAN

- 1 in 2048 sample rate

```
sflow enable
sflow destination 172.22.0.90 6301
!
interface ethernet 4/4
 sflow forwarding
```

**LINX**

# Phase 1: Cleansheet

- We knew what we wanted to achieve
- We brought in a programmer with good DB and web programming skills
- Allowed him to come at it from his own direction
    - No dictation about type of technology to use: Other than sFlow in, traffic data out.
- Took a fairly "minimalist" approach
    - Throwing away data we don't need

**LINX**

# Overview

# Database Layout



Database

| 1 Minute samples | 5 Minute averages | 15 Minute averages | 60 Minute averages |
|---|---|---|---|
| Raw Data from sfacctd, 48 hours of data. | 7 days of data | 28 days of data | Forever and ever... |

**LINX**

# Advantages and Limitations

- Upsides
  - We get more than just graphs from the same data sets
  - pmacctd gives us huge scope for functionality (pre-processing before insert)
- Downsides
  - Deleting data
  - Temporary tables
  - Joining tables

**LINX**

# Hardware

- 2x dual-core CPU's
- 16 GB of RAM
  - Can hold table indexes in memory
- 820GB RAID6 Array (8 x 146GB disks)
- Possible to scale hardware by running a distributed system
  - Mirror of system to allow for maintenance
  - Archive databases on different boxes

LIN✕

# A few numbers

- 36 Million rows for 48 hours of  1 minute samples
- 32 Million rows for 7 days of 5 minute averages
- 43 Million rows for 4 weeks of 15 minute averages
- 40 Million rows for 3 months of 1 hour averages
- 147GB of data collected over the last 3 months

LIN✕

# sFlow Portal Entry Screen

# Select Your Switch Port

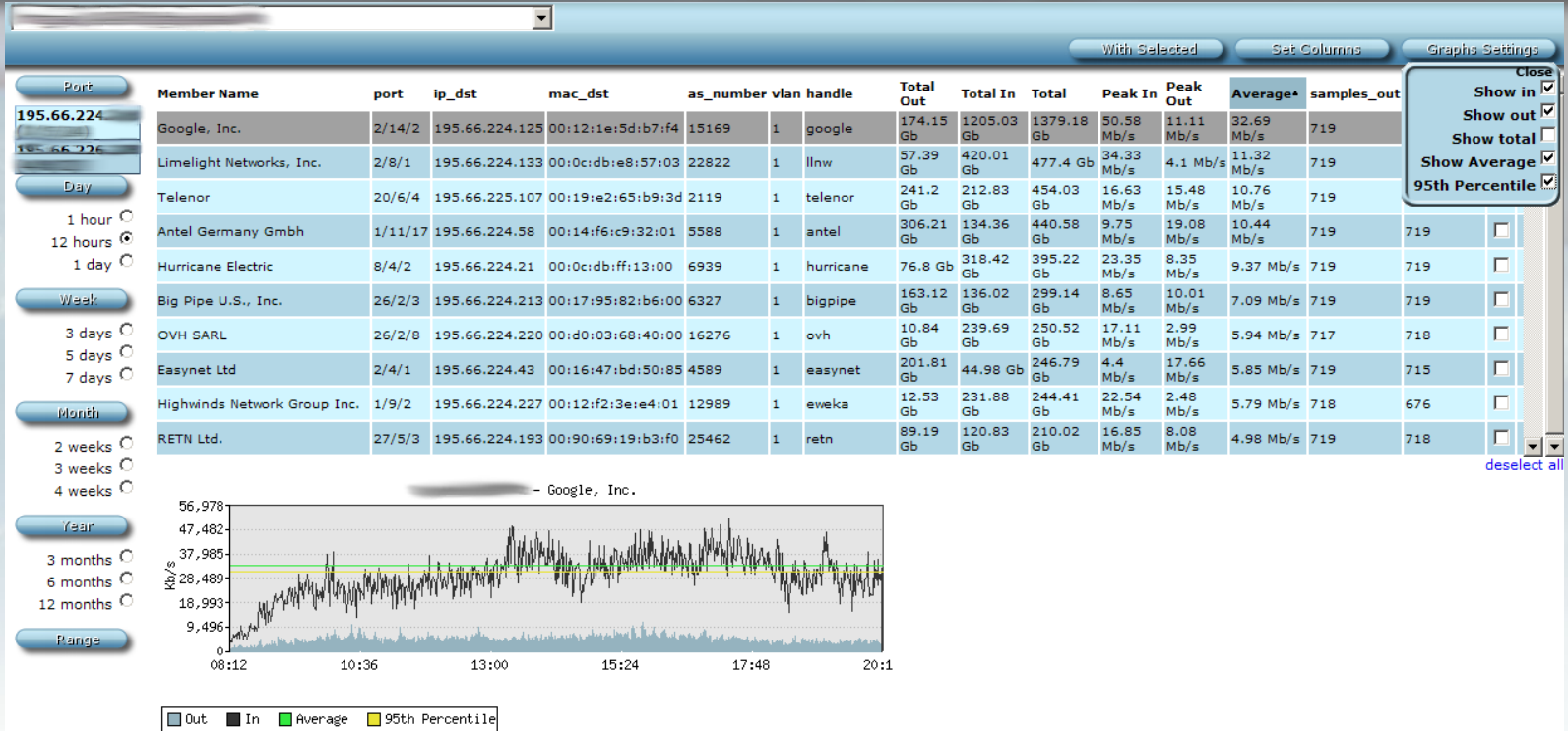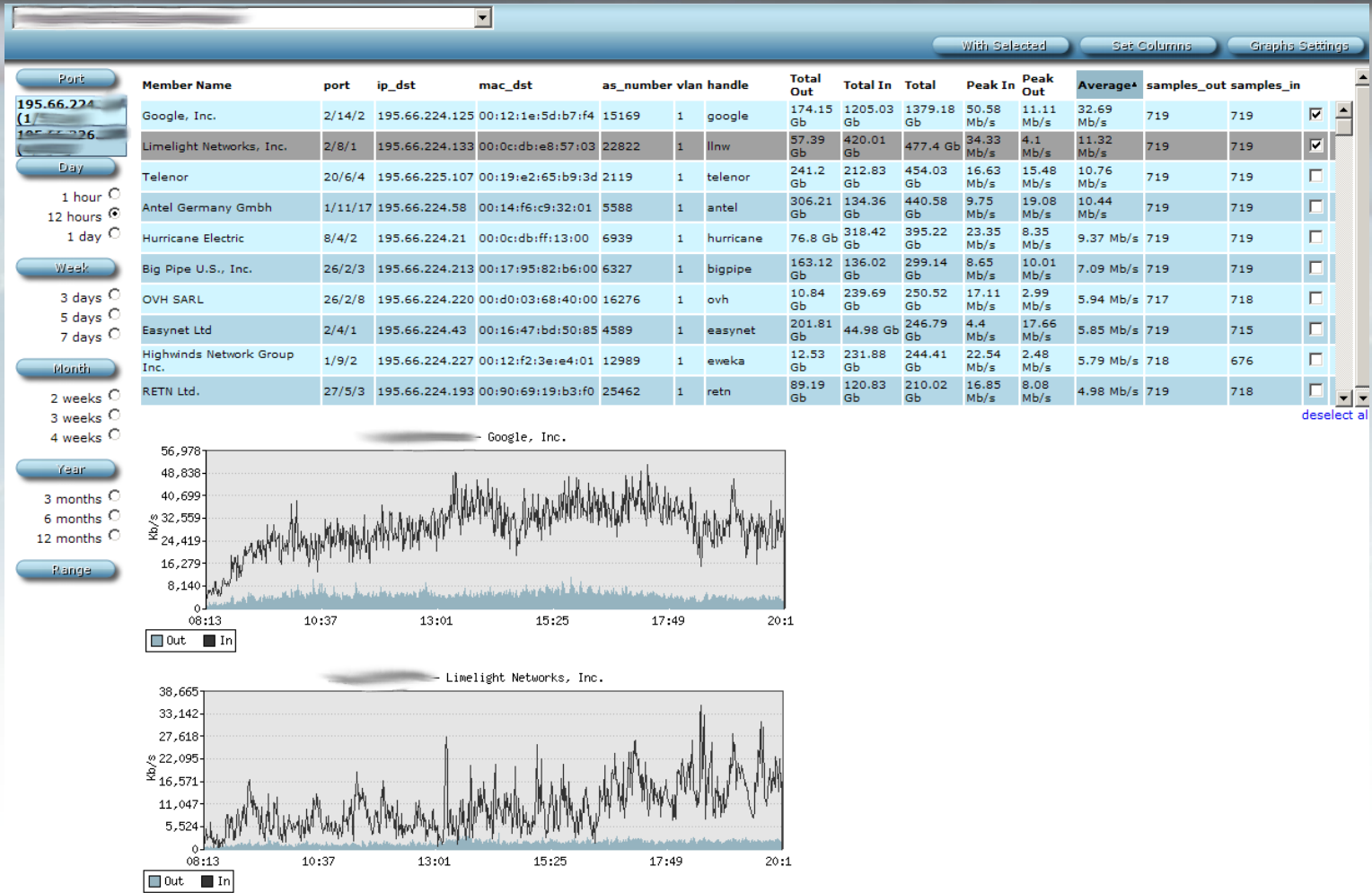# Select/Deselect Columns

LINX

# Select Time Window, Show Graph

# Show Average & 95ᵗʰ Percentile

# Select multiple peers & compare



**Engineering**

# Future Wishlist

- Add-in processing of Extreme data
  - Held up due to implementation issues
- XML schema and authenticated direct XML interface for members
  - Integrate directly into their own systems
  - Perform direct queries of the db
- Peering Matrix
- Improve engineering tools
  - "Toptalkers", interswitch traffic matrix

**LINX**

# Other wacky ideas?

- Pro-active notification agent
  - Be able to configure various thresholds, receive alerts
- Weekly "overview report"

LINX

# End Results

- Allow members to manage their peerings more intelligently
- Allow LINX to better understand flows inside the peering networks
- Identify traffic flows for optimisation
    - Switch platform relief
    - Through PNI or regrooming of member connections onto same switch

**LINX**

# Questions?

- 

**LINX**