



**Perspectives on Network
Routing: Lessons Learned and
Challenges for the Future – An
Operational Overview**

IEEE Infocom

Barcelona, April 25, 2006

Vijay Gill <vgill@vijaygill.com>

Agenda

- Routing in context
 - The rubber meets the road
- Design of a Backbone
 - Taking protocols and turning them into operational reality
 - Cost structure
 - Tools of the Trade
 - How do we manage this monster we have created?
- Going forward

Do not be so proud of this technological terror you have constructed. The ability to criticize Star Wars is insignificant next to power of the Fans

-Brandon David Short

Operational Rule To Live By

When you want it bad, you get it bad,
and most people want it in the worst
way.

- Heidi Heiden's First Law

Design Goals – Operational Simplicity

- When to touch the network
 - Routing policy based on simple performance metrics and cost (price/Mb)
 - No fancy tricks – Traffic engineering based on using common policy across similar peers (e.g., free or transit) rather than exception policy
- Achieved through engineering simplicity
 - Focus on reducing OPEX
 - It's not about building a network that's cool or is a challenge to engineer
 - "Make it simple, but no simpler"

Ask not what evolution can do for you, ask what you can do for evolution.

- Jimbo Kukla



ATDN Speeds & Feeds

- 162 Routers (Juniper T640s & M40Es, Cisco 124xxs & 7500s)
- 36 POPs on 4 Continents:
 - 24 US
 - 10 European
 - 1 Asia (Tokyo, Japan)
 - 1 South America (Sao Paulo, Brazil)
- Perspective on Throughput
 - 320 Gigabit/sec of edge traffic
 - AOL Datacenters (106GB/s Daily Peak)
 - Streaming (24GB/s Daily Peak)
- Interconnections:
 - 189 interconnects to 44 ISP peers with
 - **350 Gbps total capacity & 60 Gbps typical peak utilization**
 - 97 unique customer AS peerings
 - **857 Gbps capacity & 193 Gbps typical peak utilization**
 - **Total edge capacity of 1,180 Gbps with peak utilization of ~320 Gbps**
- Human Resources:
 - 6 Eng, 1 Tech Mgr, 1 PM, allocations from NOC, Arch, Sr. Staff & Finance
 - **Grand total of about 15 FTE**

▶ Planning Parameters (examples)

	Network	Host
Usage Drivers	<ul style="list-style-type: none"> • Membership & % logged in at peak • Bit Intensity of Session (Kbps /user) • Number of hosts supported • Bits streamed • New projects (XM Radio, Live 8, WAGS) 	<ul style="list-style-type: none"> • Software performance • Application profile balancing • Format choices (streaming & web) • Replication / failover requirements
Measurements	<ul style="list-style-type: none"> • 95% percentile Mbps (edge & Customer) • Traffic matrix analysis (failover) 	<ul style="list-style-type: none"> • URLs/sec • Mbps • Cache effectiveness
Utilization Limits	<ul style="list-style-type: none"> • Lit Wavelength versus SONET protected • ~45% max on ATDN BB links • Uplink utilization of <50% in LANs 	<ul style="list-style-type: none"> • Cache managed to 80% URLs/sec • CDN sized for peak events + overflow • BOC simultaneous channels
Implementation	<ul style="list-style-type: none"> • Build versus buy analysis • ATDN - \$XX per mbps • AOLWave - \$YY per provisioned mbps 	<ul style="list-style-type: none"> • Build versus buy analysis • CDN - \$0.XX /GB Transferred • BOC - \$0.YY / min. encoded • Cache - \$0.ZZ / member month
Inventory/Capital Management	<ul style="list-style-type: none"> • "Let it rust" • EOL switch and DWDM management • Redeployment of legacy technology Multiple vendor strategy 	<ul style="list-style-type: none"> • "Let it rust" • Rotate & reuse • '04 migration to low cost storage (reduced cost by an order of magnitude)

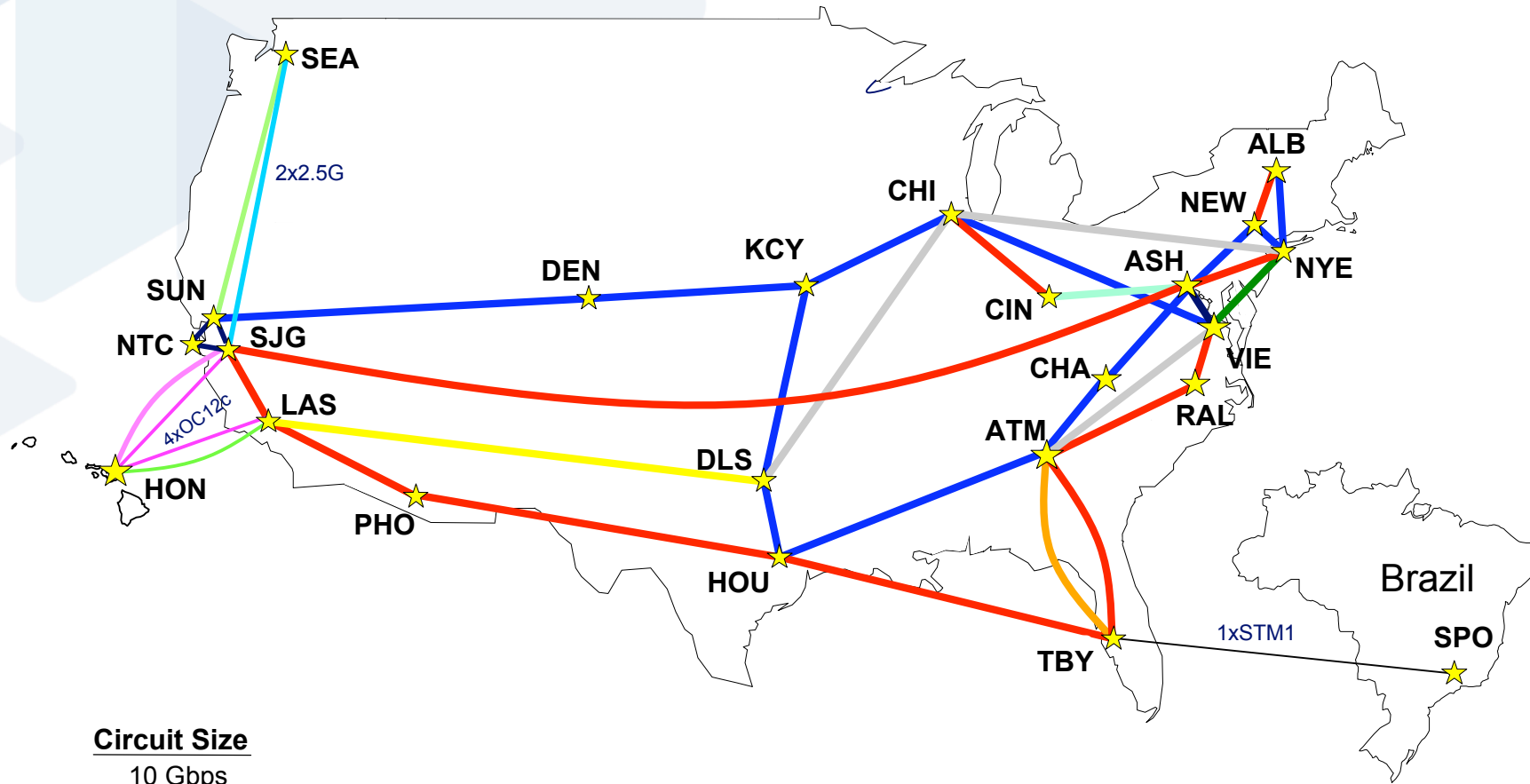
Move to the edge – Mini datacenters

- Historically, AOL was entirely centralized in Northern Virginia
- Only recently has distribution become part of our thinking
- We need to constantly evaluate the core vs. edge spectrum and strike that delicate balance between centralized & decentralized resources
- Bigger issue as
 - Our applications & content become richer (more bits)
 - Our users become more latency sensitive & accustomed to “Internet Speed”
 - Our competitors move to get topologically closer to the consumer
- Appropriate distribution
 - Improves performance, end-to-end reliability & failure resiliency
 - Avoids network costs & improves our SFP positioning
 - Simplifies maintenance

Engineering Issues

- Utilization (peak engineering)
- Vendor & route diversity; wavelength vs. SONET protected networking
- Design to avoid single points of failure
- We rely heavily on tools, stats & automation to manage these systems
 - More than 120 tools being actively managed (1/3rd in major revisions; 1/3rd minor)
- Resiliency through “all-active diverse” replication
 - At scale, we deploy clusters of machines to perform discrete functions
 - Resiliency is generally provided by deploying sufficiently diverse clusters
 - All clusters are active
 - Failure of one cluster spreads load to remaining cluster
- Model build versus buy (Unit cost/Benchmark)
 - ATDN - \$XX/Mbps versus open market of \$YY/Mbps
 - AOL Wave – \$ZZ per provisioned Mbps versus \$QQ (in Metro)

AOL Transit Data Network (ATDN) Domestic and South America



Circuit Size

10 Gbps

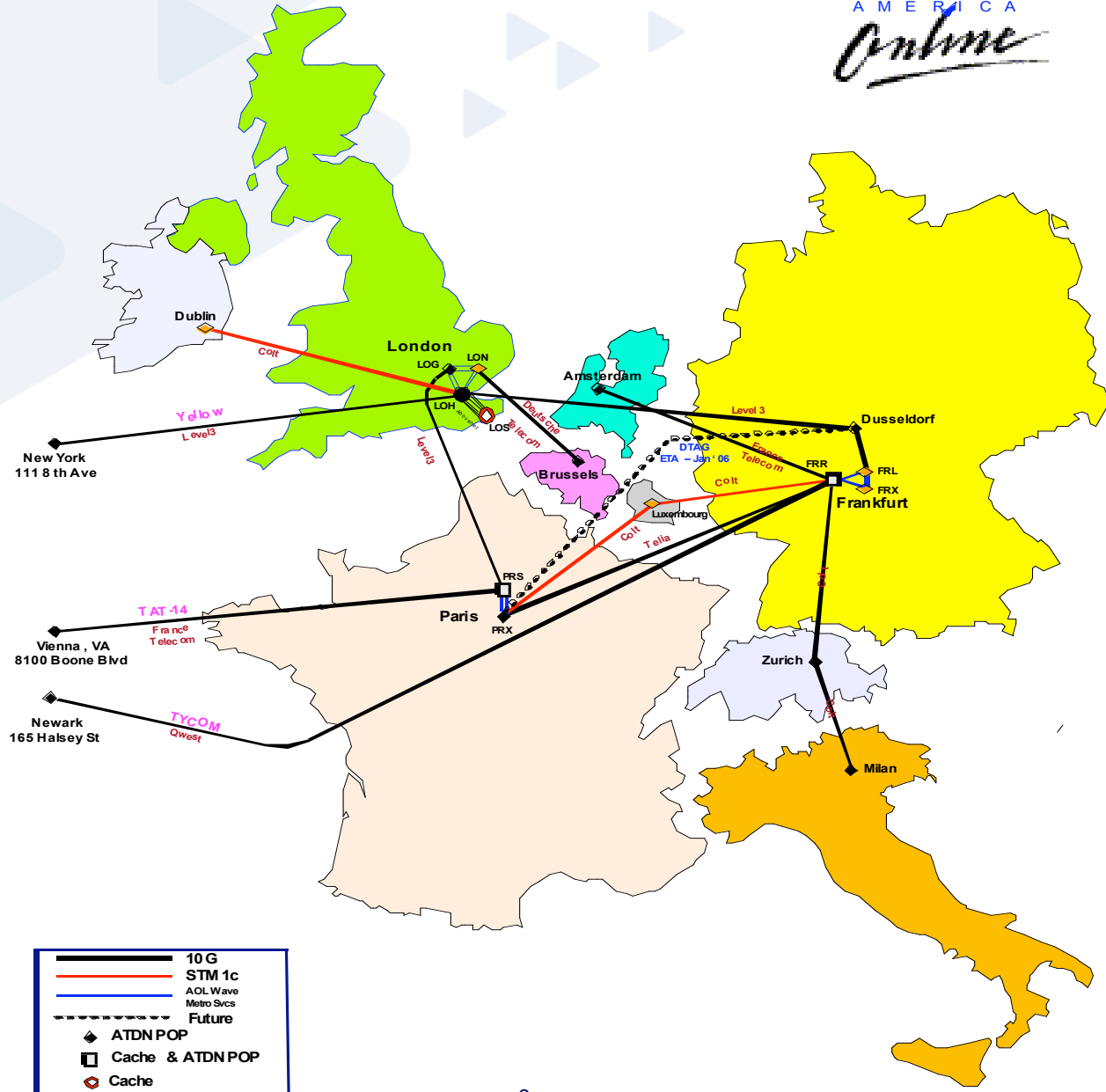
2.5 Gbps

622 Mbps

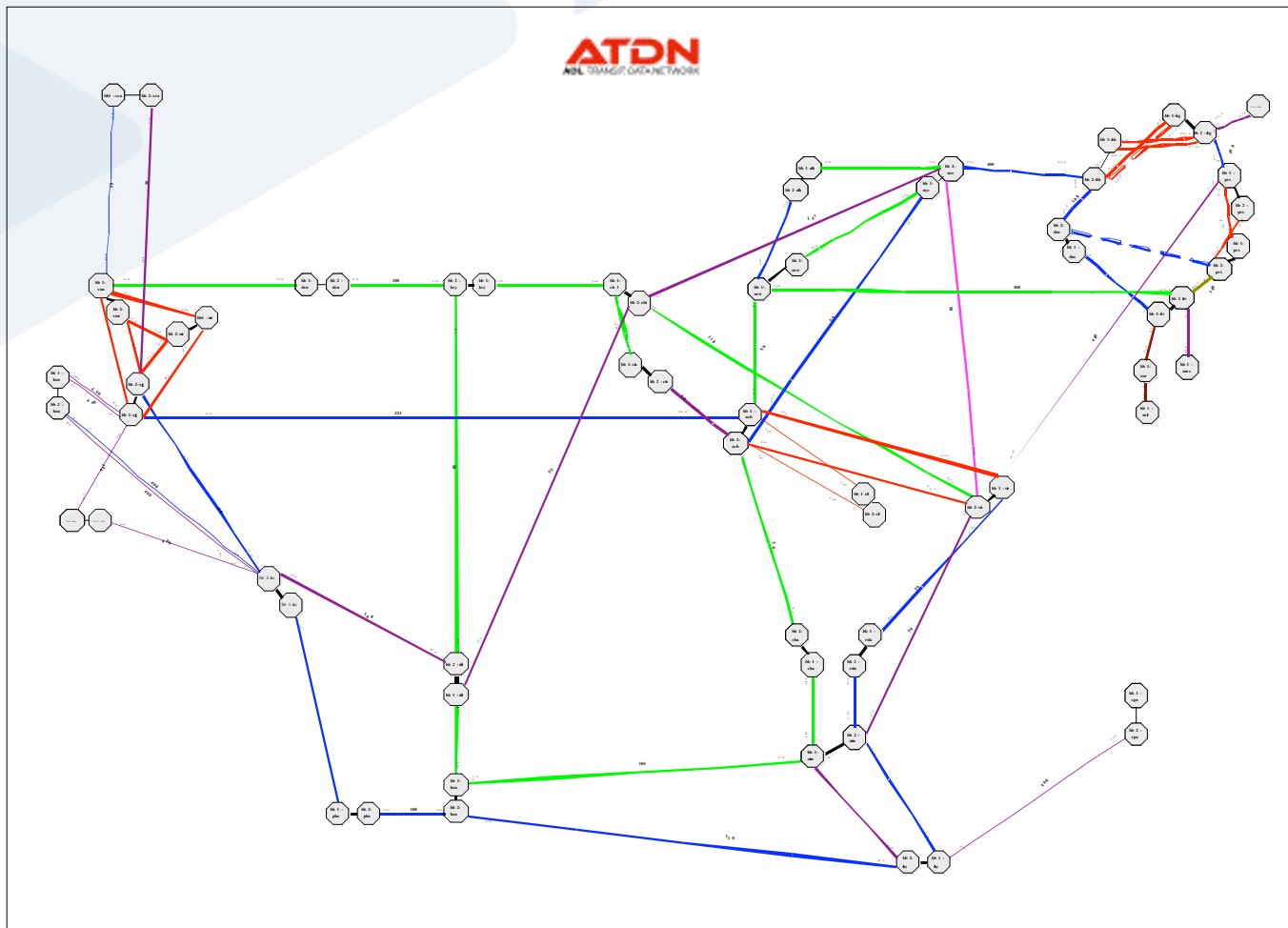
155 Mbps



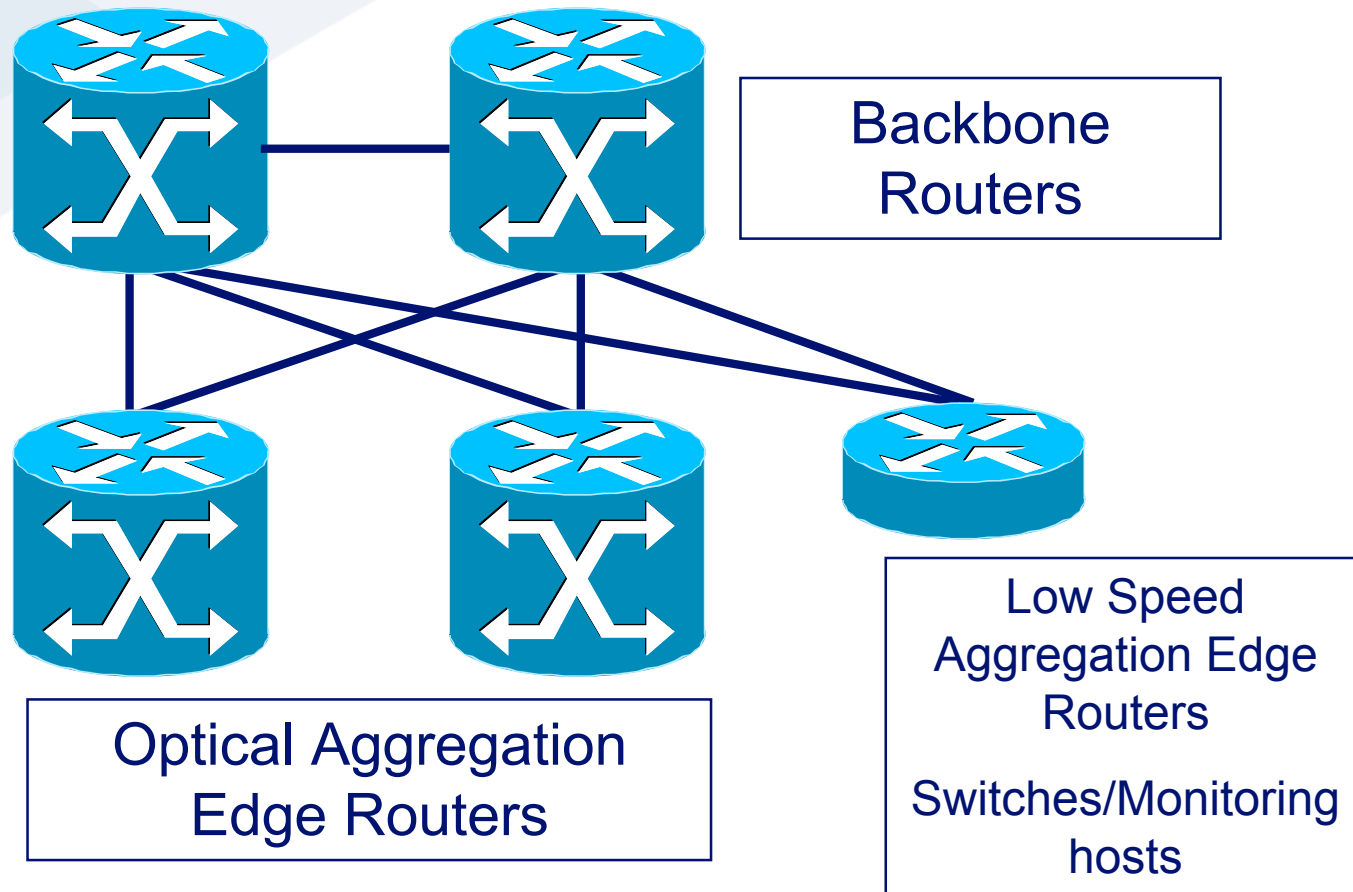
ATDN Europe – 12/1/2005



ATDN Logical Design



ATDN Hub Design



The Basics

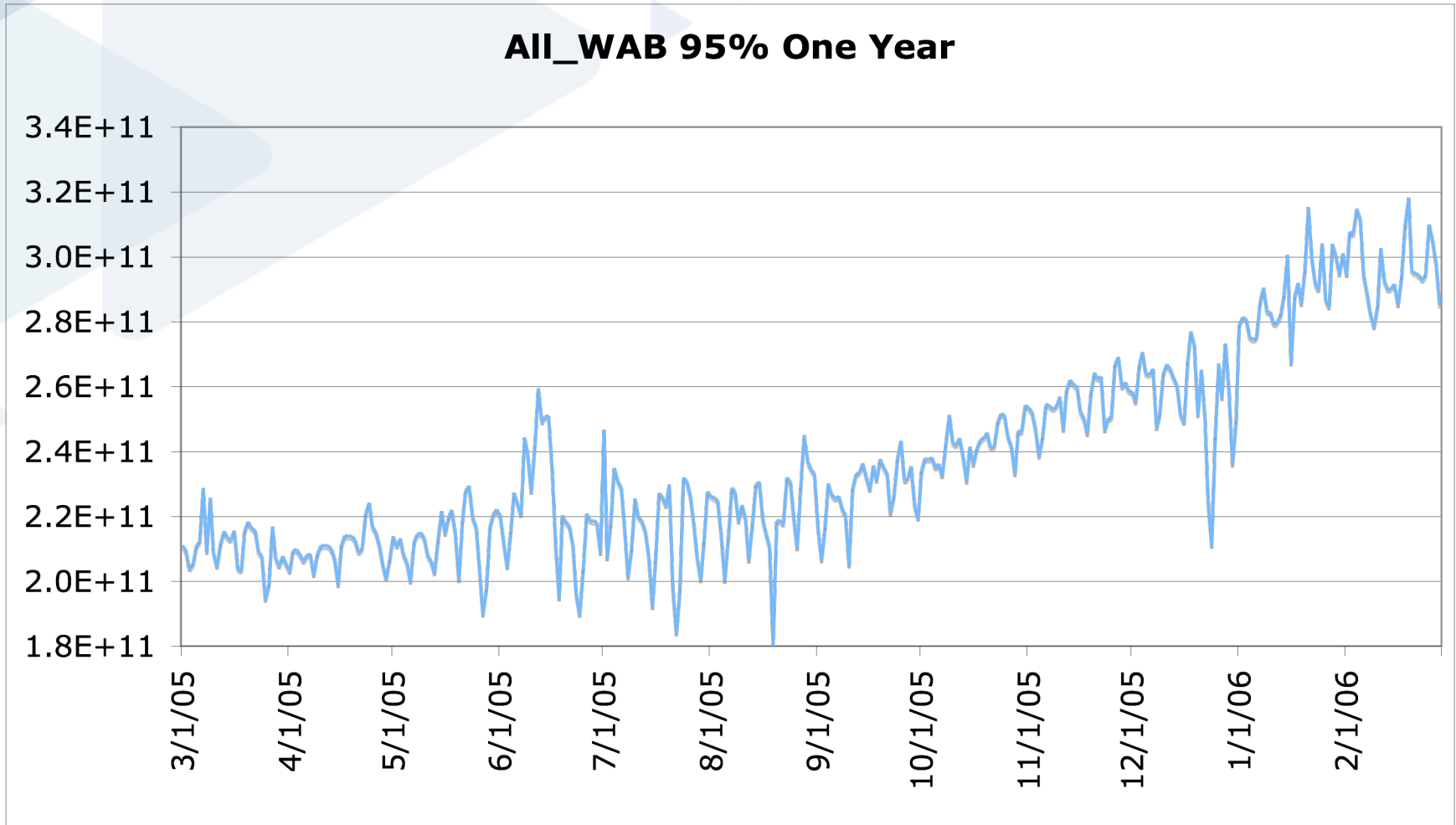
- Environmentals
 - Front-to-back cooling
 - Fire is bad for equipment
- Physical Plant
 - 7 foot rack
 - 19 inch rack mountable
- Chassis
 - No dependency on spinning media
 - Seen more flash card failures than hard disk however



Routing Architecture

- BGP/IS-IS
- Run IGP lean, carry everything in iBGP
- IGP cost structure:
 - Local
 - Regional
 - Long-haul
 - International

Traffic Volume



▶ What Are The Bits?

- 320 Gigabit/sec of edge traffic
- 26+ million AOL subscribers (including EU)
- 3+ million CompuServe subscribers
- 3+ million TW Cable subscribers
- 125+ million AIM users worldwide
- 110+ million ICQ registered users
- Peak simultaneous usage:
 - 3.1M AOL users online
 - 7.3M active AIM users
 - 2.4M active ICQ users

Infrastructure:

- 25,000+ host servers
- ~ 500,000 sq ft raised floor space
- ~ 800 optical backbone routers
 - ATDN – 66+ 10-gig capable routers in the backbone
 - ATDN – 100+ edge routers
- ~ 3000 L2/L4-7 Switches
- 66,000 interfaces polled every 5 mins
- Over 1 million network variables captured

Speedboat

- **Purpose:**

- Speedboat is a collection of components that implement the Argus protocol specification, and together provide a framework to quickly and cheaply store and retrieve time-series data across NW&DC Ops.

- **Benefits:**

- Stored metrics are used by System Admins to debug problems.
- Metric points are used for alarming, trending, KPI.

- **Implementation Strategy: (By component)**

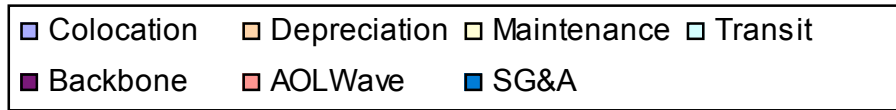
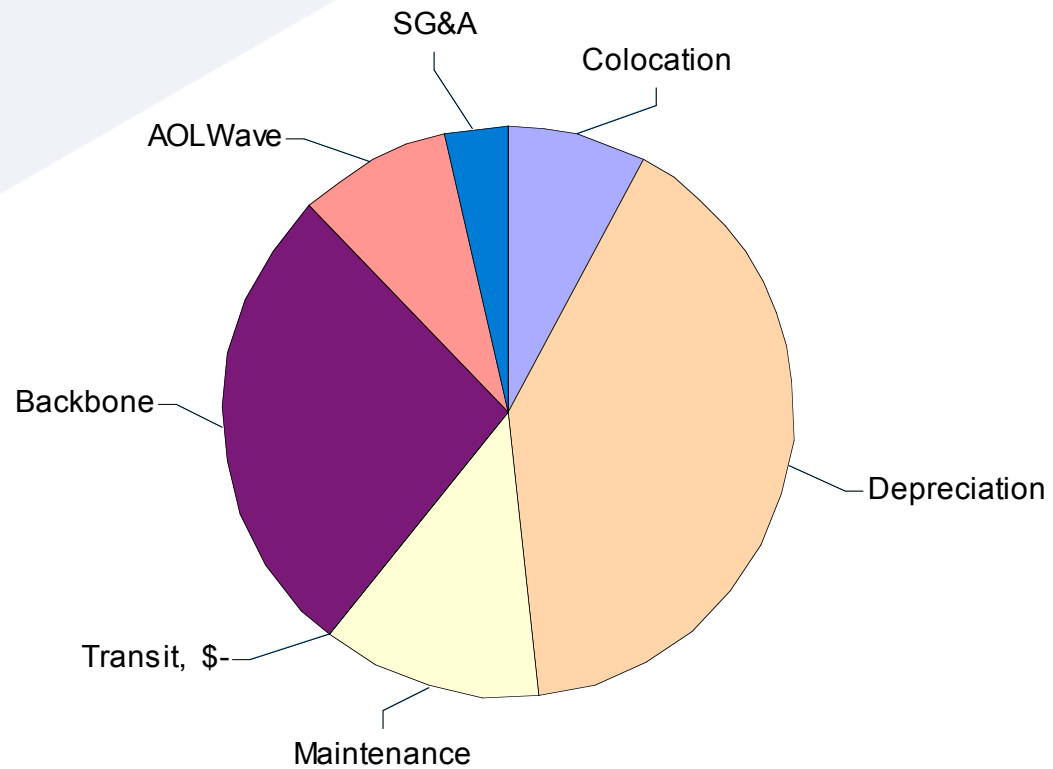
- Simple Temporal Database (STDB): Berkeley DB based data storage, designed for time-series data.
- Meta-data Server (MDS): Abstracts textual names from numeric storage Ids.
- Publisher: Provides an interface to clients that write data points.
- Unification Layer (UL): Provides the interface to clients that read data. Proxies for all other servers.
- Network Stats Grapher (NSG): Provides a “fat” client for data analysis.
- Heimdall: Provides a web front end.

- **Status:**

- In production storing 6.7B metrics/day @ \$0.0X/million.
- Minor version development for some components.
- Working with Sysops tools team to share more components beyond STDB.
- Improving Grouping and Aggregation.

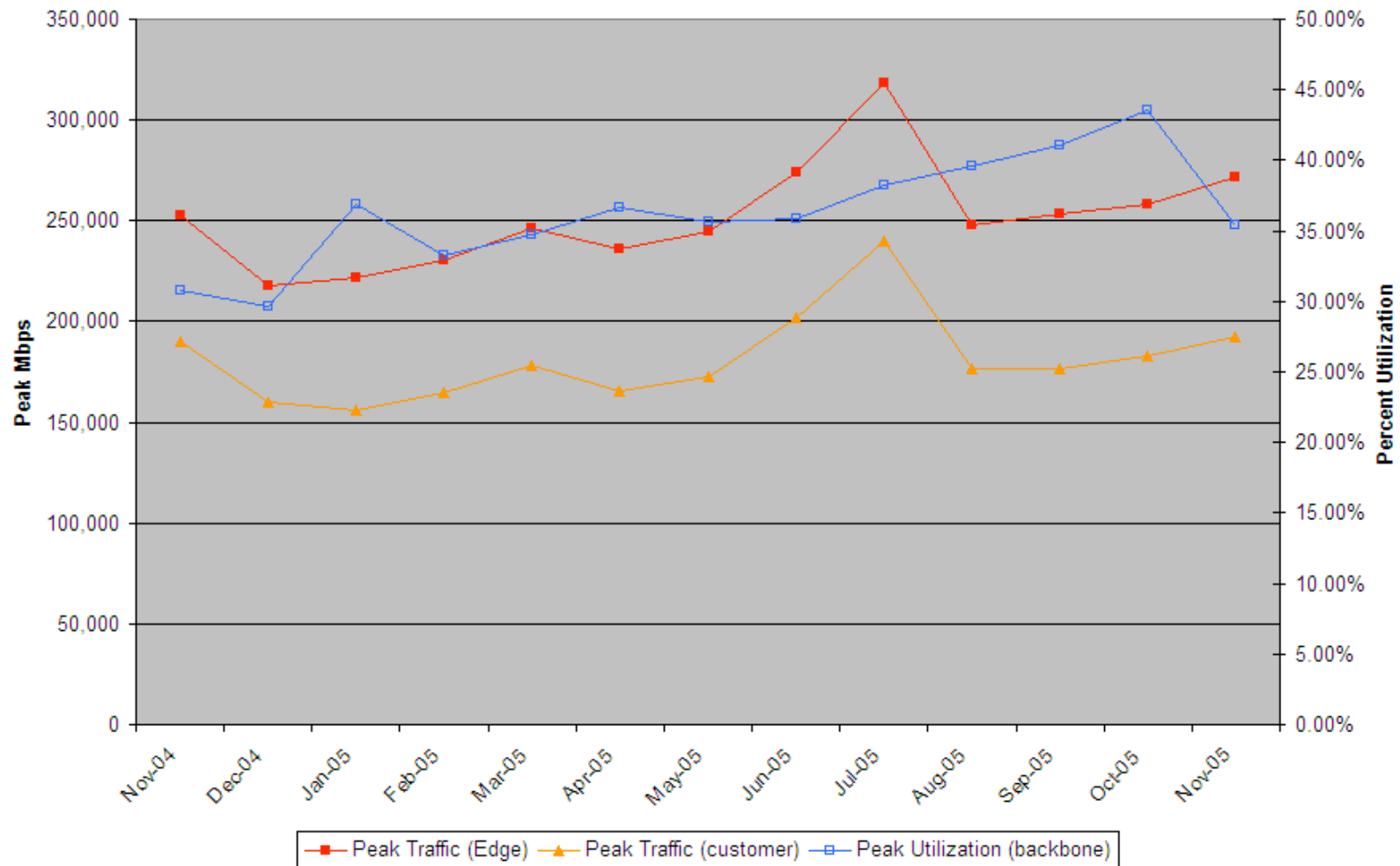
FY05 ATDN Expenses by Type

2005 ATDN Expenses



ATDN Utilization Snapshot

ATDN Peak Traffic



Key ATDN Issues

- Route Resiliency & vendor diversity
 - Wiltel purchase by Level 3
- Settlement Free Peering
 - Maintain **peering ratios** through active management and P2P
 - Improve **geographic traffic balance** through smart distribution of servers (mainly streaming today) in Mini Data Centers (MDCs)
 - New York, Chicago, Los Angeles, Dallas, Atlanta
 - Will soon include pieces of AOL.com and other Audience properties
 - 40 Gbps deployment may become Tier-1 requirement
- Deployment & adoption of advanced services (e.g. MPLS) for cost savings & to keep pace with features available on commercial backbones
- Expand capacity for XX usage from YY Gbps to 2xYY Gbps
- Renewal of Southern transcontinental path in late 2006

ATDN Security Initiatives – unicast Reverse Packet Forwarding (uRPF)

- uRPF loose mode has been tested on Cisco & Junipers and under deployment across network.

Excessive uRPF Drops Limit: 6.94444 drops/min 08/08/2005 16:51:38

Router	Total uRPF Drops	Drops/Min	Elapsed Secs
pop1-dls	1205151625	20413.55	1023
pop1-hon	4055627	45.18	1020
pop1-hou	73359500	285.28	1023
pop1-kcy	32125766	289.03	1023

```
grao@pop1-ash> show interfaces so-1/3/0 detail | match "RPF Failures"
```

```
RPF Failures: Packets: 2675, Bytes: 157188
```

```
grao@pop1-ash> show firewall log
```

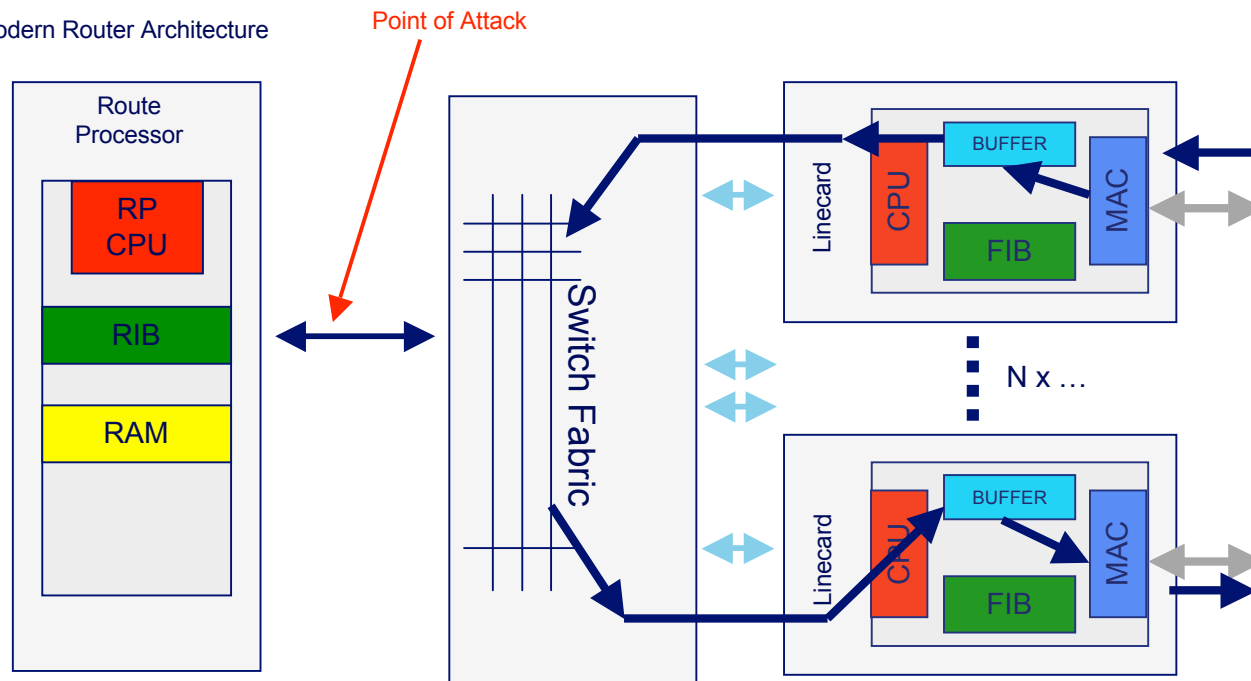
```
Log :
```

Time	Filter	Action	Interface	Protocol	Src Addr	Dest Addr
16:53:25	pfe	D	so-1/3/0.0	ICMP	10.34.1.1	68.202.210.155
16:52:51	pfe	D	so-1/3/0.0	ICMP	10.34.1.1	68.202.210.155
16:52:38	pfe	D	so-1/3/0.0	ICMP	10.34.23.1	64.236.56.22

Security

- Routers are optimized for traffic **through** the hardware
 - Not traffic **for** the hardware
- Designing a cost efficient router implies:
 - Cross-sectional bandwidth capacity dominates budget
 - No cost-effective way to engineer a router that can absorb and usefully process data at the rate it can arrive

Modern Router Architecture



Hardware – Queuing of Control Plane Traffic

- This one should be easy to get but surprisingly few can do it
- Simple, unambiguous parsing
 - Filter on stuff that is for the router
 - What I deem interesting goes onto the high priority queue
 - Everything else goes onto the low priority queue
- Simple discriminator function/ACL etc.
- Rate-limit on low priority queues
- Apply discriminator on linecard/forwarding engines BEFORE it hits the brain
- Why?

Outside Context Problem

- Attackers are seizing this weak link as a point of attack
 - DoS attacks targeted at infrastructure are increasing
 - Hackers will evolve – Have seen port 179 attacks already (and MSDP can't be far behind)
- Problem
 - Need some way to disambiguate between invalid and valid control traffic (e.g. BGP updates)
 - Rate-limiting on control traffic is not sufficient
 - Enough false data will swamp legitimate data
 - Connection flaps/resets
 - Need to focus on BGP (MSDP)– other traffic is not control, thus will not cause control plane issues

Security

- IGP traffic can be safely blocked
- MD5 on neighbors will not prevent the Router CPU from being inundated with packets that must be processed
- Solution
 - Short term - Dynamic Filtering on the line cards
 - Long term – outboard processing of SHA1/HMAC-MD5
 - This is very long term indeed – not necessarily solving a known problem today (replay or wire sniffing)
 - Vendors have to implement priority queuing for control traffic from line cards to control plane

Network Tuning

IGP tuning for faster convergence

- How do we measure this?
- Routing policy analysis
- Diffserv-based QoS
 - Protect control plane traffic
 - Provide classes of service for various products

Take Aways

- Central thrust is manageability
 - Operators are managing to things that are sometimes OPAQUE to the end-user
 - E.G. Traffic distribution, SFI, projected maintenance
 - Cost
- New protocols must
 - Be incrementally deployable – Flag days are over
 - Must show visible benefit in terms of ROI
 - Networks are now controlled by the finance folk
 - Must not make life harder (see above)



Q&A

There is a difference between making something foolproof and reducing the number of fools"

-Bill Barns