



March 20<sup>th</sup> 2012, OSG/ATLAS/CMS

Jason Zurawski, Internet2 Research Liaison

[zurawski@internet2.edu](mailto:zurawski@internet2.edu)

## **Addressing the “things that go bump in the net” – perfSONAR/DYNES/LHCONE**

# Agenda

- Current Networking
  - perfSONAR Status (ATLAS, CMS, LHCOPN, LHCONE)
  - Reaching for the Brass Ring (why we monitor)
- Future Networking
  - DYNES
  - LHCONE

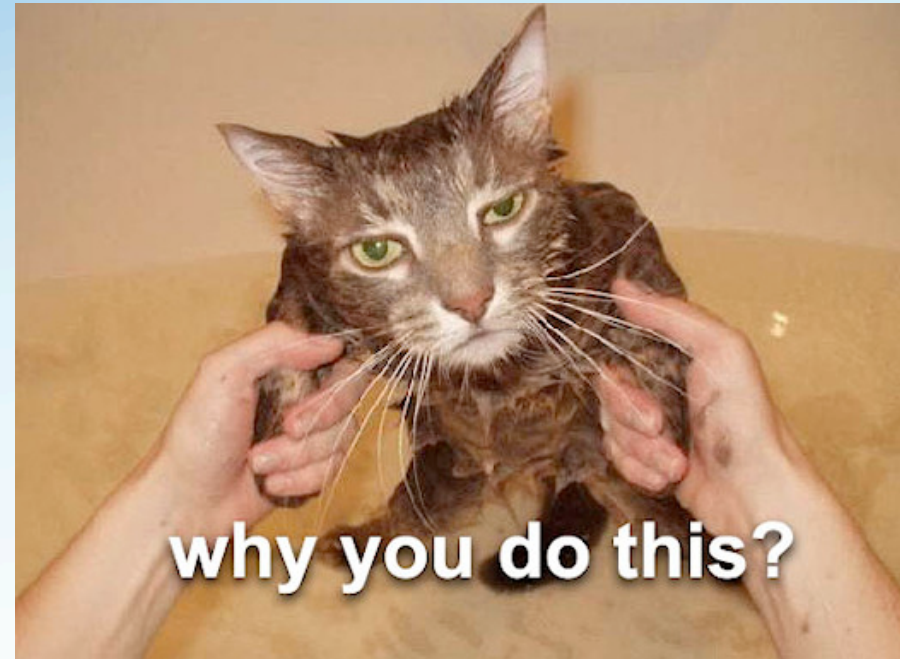
# A Jon Postel quote

- ***"In any large system, there's always something broken."***
- Networks are large and complex. There are multiple “layers” and we employ experts with knowledge of specific parts just to keep things running
  - Anything that can give an expert (or a layman) more insight into what is really going on is very valuable

# “Why?”

- Everyone should be familiar with what perfSONAR is about, this talk is not about that

- Mentioned in the “2013 NITRD Program Supplement to the President's Budget” (page 50) - <http://www.nitrd.gov/PUBS%5C2013supplement%5CFY13NITRDSupplement.pdf>



- If you are not running it, **this is not a sales pitch**
- Things I will highlight:
  - It is being used widely
  - It is finding problems

# perfSONAR-PS Status

- USATLAS
  - All Tier2s and Tier1 upgrading to new Dell R310/R610 (available as 'perfsonar node' in the portal)
  - Dashboard:  
<https://perfsonar.usatlas.bnl.gov:8443/exda/?page=25&cloudName=USATLAS>
  - Other non-US clouds (Canada, Japan, Italy) coming up as well
- CMS
  - All Tier2s (and Tier1) have monitoring in place. Should be testing to each other

# perfSONAR-PS Status – cont.

- LHCOPN

- All Tier1s and Tier 0 have machines in place with tests in place

- Dashboard:

<https://perfsonar.usatlas.bnl.gov:8443/exda/?page=25&cloudName=LHCOPN>

- LHCONE

- 16 Sites are being monitored as a part of the LHCONE Arch prototype phase. Some are fully configured, others are not (work in progress – Shawn is leading this).

- Dashboard:

<https://perfsonar.usatlas.bnl.gov:8443/exda/?page=25&cloudName=LHCONE>



# perfSONAR-PS Software

- Current release – 3.2.1.1
  - Expect a 3.2.2 in mid 2012
  - Bugfixes for the most part, no real ‘new’ features
  - <http://psps.perfsonar.net/toolkit>
- Items on the longer list:
  - Controlling an entire deployment instead of an individual island (N.B. some are exploring CFEngine and the like in this space)
  - Integrating the tools into a more portable dashboard (basing this heavily on the work by BNL)
  - Bottom line – lots to do, little time and resources to do it (but this isn't news)

# Agenda

- **Current Networking**
  - perfSONAR Status (ATLAS, CMS, LHCOPN,LHCONE)
  - **Reaching for the Brass Ring (why we monitor)**
- **Future Networking**
  - DYNES
  - LHCONE

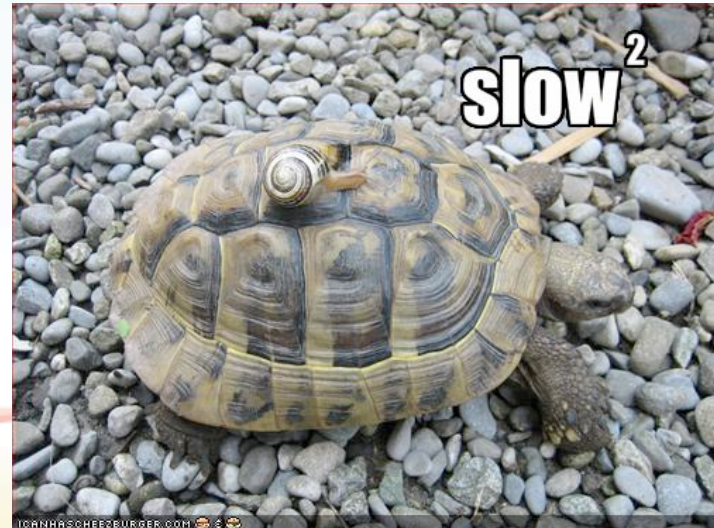


# “Why Care/Devote Resources?”

- Network monitoring *is*:
  - A way to pick out problems (packet loss, congestion, routing changes, low throughput)
  - Used by operators to find problems before the users (*you*) find them
  - Used by users (*you*) to keep the operators honest
- Network monitoring *isn't*:
  - An instant way to solve said problems. It will tell you ‘what’, it won’t tell you ‘how’ or ‘why’ without spending some time on the problem
  - Automatic. There is some work that needs to be put in by all levels (operators, VOs, etc.)

# Anatomy of a Problem

- “The Network is Slow”
  - Yes, its ok to say this. Don’t overdo it though (e.g. complaining at getting 8.5 Gbps when you got 9.3 Gbps yesterday), and try to evidence when you do say it (e.g. your graphs)
- Looking at the regular data (and alarming on it)
  - ATLAS, LHCOPN, etc. have the regular tests for this exact reason
- One off tests
  - Log on to the boxes (its easy, just like any other linux machine) and run some tests. Don’t know how? Ask!



# Anatomy of a Problem – cont.

- Escalation
  - You can escalate when you are in over your head. ESnet/Internet2 are here to help.
  - Also – talk to your local IT people so they are aware. They don't bite.
- Waiting (is the hardest part)
  - Debugging sucks.
  - It takes a long time.
  - It involves multiple parties (this is what makes it take longer)

# Current Problem (some of you know this)

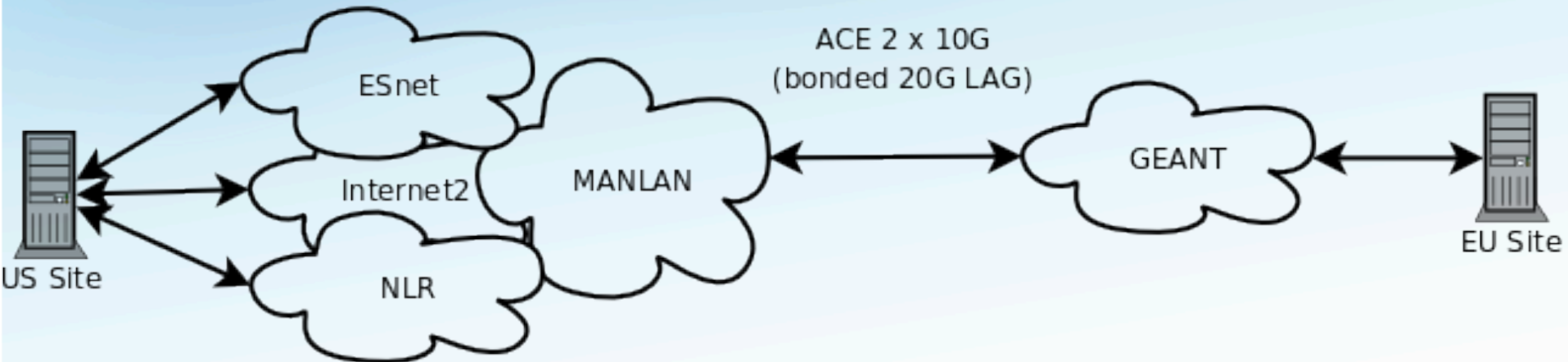
- 1 of the Transatlantic Link Pairs (New York to Amsterdam)
- Performance bad in one direction (from the EU to the US).
  - No problems seen in the other (US to EU) direction.
  - Common issue – downloaders (e.g. people not in your network) see a problem vs uploaders (people in your network).
- Depending on who/where you are, this may not be an issue for you:
  - US sites 'downloading' from EU may see this
  - EU sites that use the NLR routes to reach locations in the US will be affected (NLR uses AMS->NEWY route exclusively)
  - EU sites that use the Internet2/ESnet routes through Amsterdam to NY to reach US sites will be effected. If the EU site uses FRANK->WASH to reach US, there will be no problem.

# Why wasn't this caught?

- It was actually – GEANT, Internet2, and ESnet commissioned regular inter-domain testing between the networks in late 2011
- Reports came in in late 2011
- The hard part(s):
  - Debugging
  - Passive vs Active
  - LHCONe

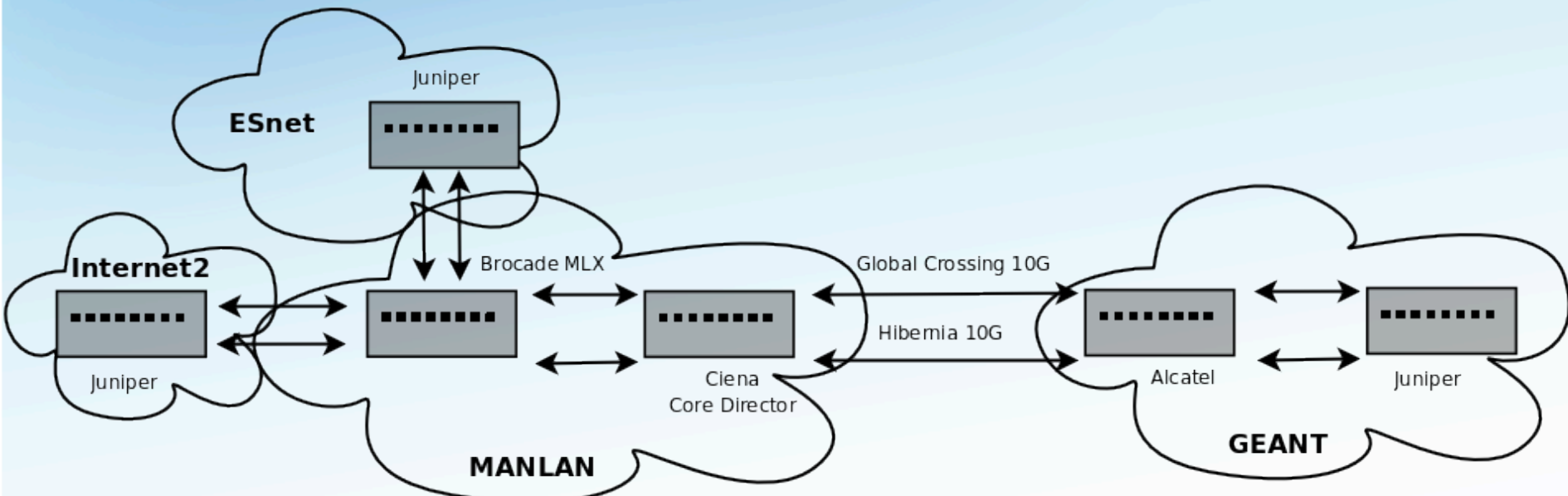


# A Basic Topology



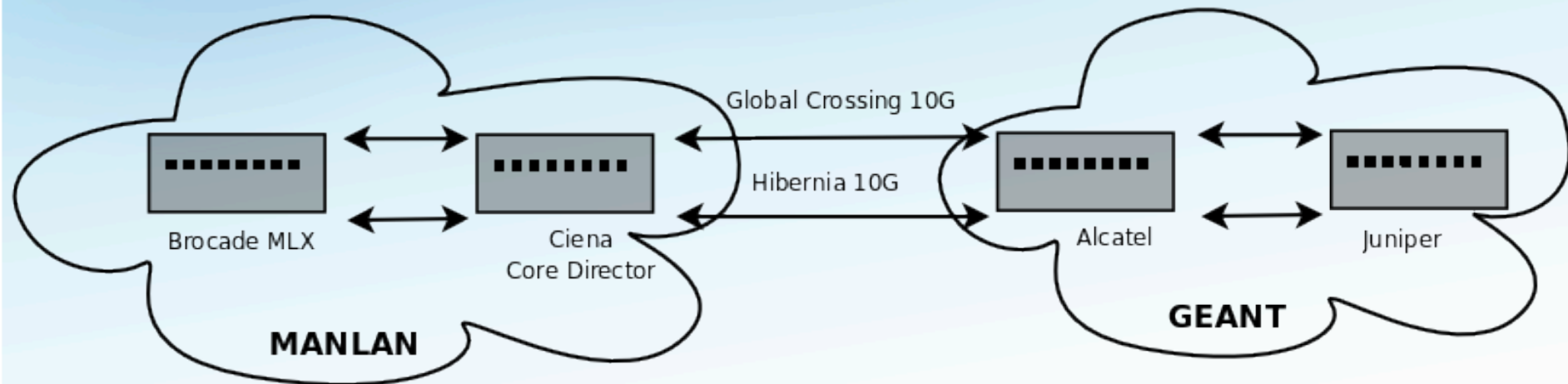
- All of the major networks show up at MANLAN XP
- Recent upgrade to switching fabric
- Major R&E Path to Europe is ACE (America Connects to Europe) IRNC Link
  - 2 x 10G LAGed Circuit
- GEANT Amsterdam Exchange feeds into other networks (GEANT, SURFnet, etc.)

# An even Better topology



- TA Circuits are SONET. Ciena CD and Alcatel terminate these on either end
- Switching/routing Fabric is connected to these two devices to support more connections (10G Ethernet for the most part)

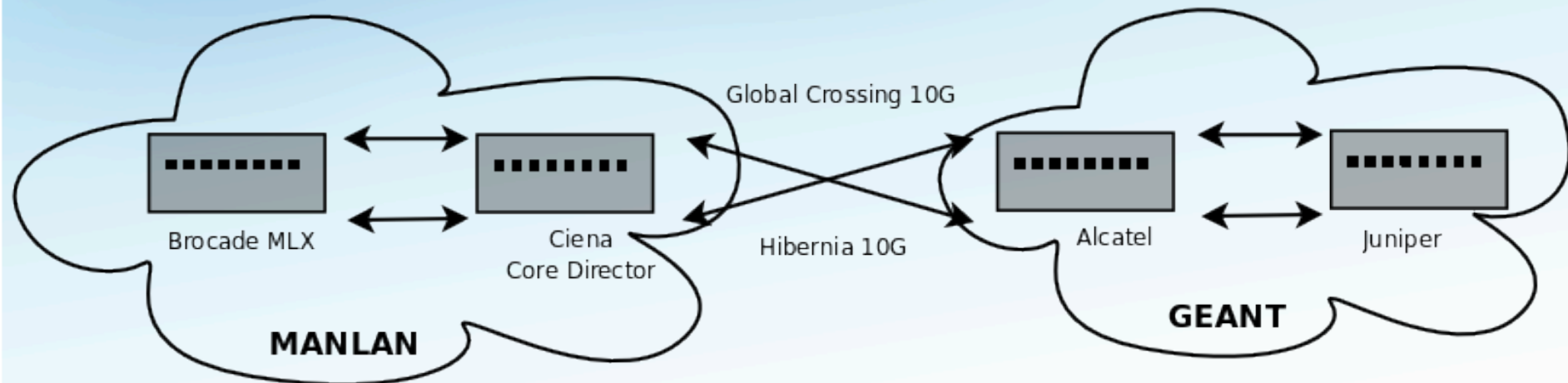
# Where we are spending time right now



- Narrowed the problem as much as possible. Testers on Internet2/GEANT are 1 hop off of the switching fabric on either end (and we still see loss)
- Is this a buffering issue? Is this a protocol issue? Is this an equipment fabric issue?

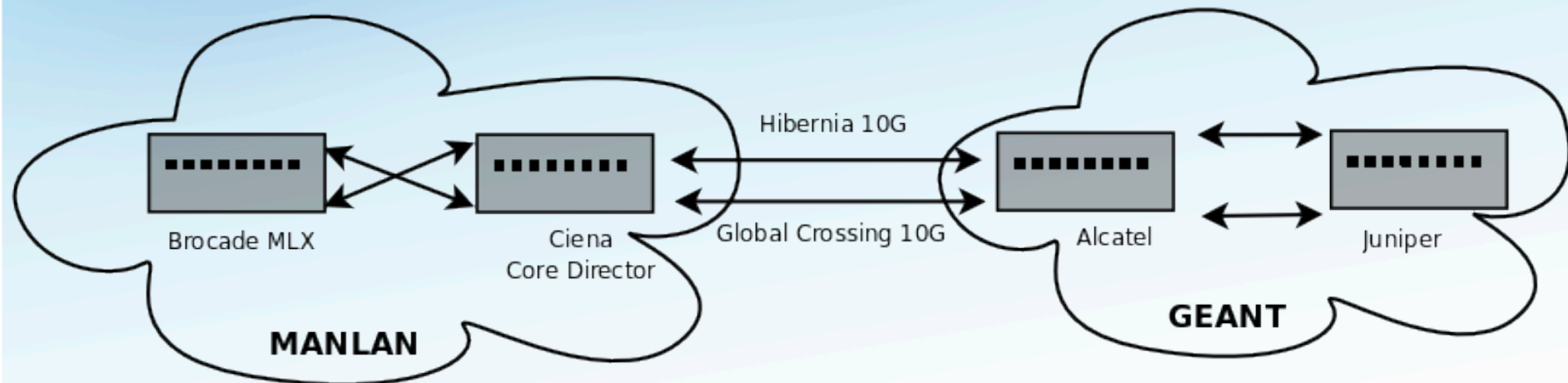


# 1<sup>st</sup> test – interface swapping @ MANLAN



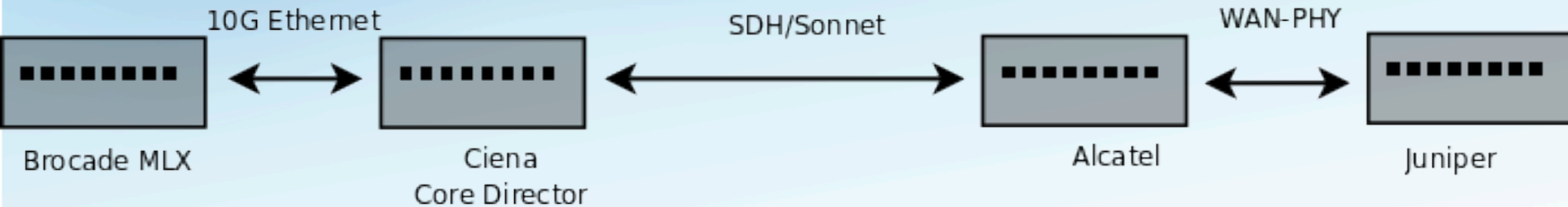
- Configuration change on Ciena (MANLAN) side to verify this device
- Blast through a set number of packets, make sure in and out packet counters agree
  - They did...

## 2<sup>nd</sup> – interface swapping @ MANLAN



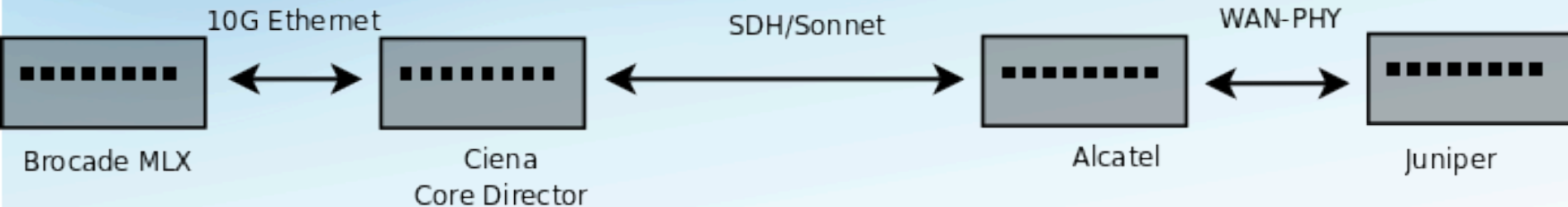
- Configuration change on Brocade (MANLAN) side to verify this device
- Blast through a set number of packets, make sure in and out packet counters agree
  - They did...

# 3<sup>rd</sup> – It's the buffering, stupid



- All of these devices are are functioning at 10Gbps line rates
- Ethernet, SONET, and WAN-PHY do have minor speed differences
  - A burst of packets on an input could overdrive an output.
  - There needs to be enough buffering to cover these cases
  - Input vs output have different queues
- Buffering was increased to the max – around 32K (yes, this doesn't sound like a lot, and its not. Enough to handle a couple of frames only...
  - It did reduce the loss percentage

# New testing (~ 1 week from now)



- Protocol encapsulation is tricky
  - Ethernet frame is shoved into a SONET frame for transit
  - WAN-PHY (a form of Ethernet w/ extra encapsulation) would be in the same boat
  - Is the translation getting garbled? Note that some devices will happily pass a bad packet on a given layer and as it gets handed back up error correction will reject it.
- Testing these theories are a bit invasive, so its taking a little time to schedule

# Where it worked, where it isn't working

- Test Coverage - **B+**
  - Internet2, ESnet, GEANT, and the experiments all have testers available
  - Some of the GEANT testers are limited in functionality
- “Reportability” - **D**
  - I took the role of ‘user’ this time. My ticket was closed **3** (!) times:
    - The day after I opened it, because there were no counters reporting loss. It was re-opened after I complained they had to “try harder”
    - 1 week later, after testing in MANLAN revealed no issues (I was told to “go ask someone else”). It was re-opened after I noted the problem is not solved from a “user” perspective
    - 1 week after that, when I was told “open tickets count against the engineer assigned” [*maybe they are not fed that day?*]. I let it be closed this time, and dealt with my tickets in other systems
  - This is something that needs to be fixed



# Where it worked, where it isn't working

- NOC to Customer Interactions - **C-**
  - NOC treated report with skepticism. Calling it 'my' packet loss (e.g. they don't trust the measurement tools, and look to the passive counters as the law of the land)
  - I had to escalate this into management to keep things 'open'. Strong desire to close tickets that are viewed as 'not my problem'. There is no home for the homeless...
- NOC to NOC Interactions - **B-**
  - NOCs coordinate resources well, but timelines to find a fix are slow. A downtime of 5 minutes is scheduled a full 2 weeks out, and only after approval at high levels
- Getting a resolution - **Incomplete**
  - More testing is needed/is expected.
  - This is a very challenging problem, and the time it has taken to solve reflects this (e.g. no clear sign of packet loss on **INTERNET** devices, but applications react poorly).

# Actions

- Jason
  - Still trying to update ATLAS/CMS when I hear news
  - Staying on top of them to get this fixed (there are still some that deny this exists)
- USATLAS Throughput Group
  - Thinking about the process to recommend for the end scientist/site to report issues in a trackable manner
  - “Customers” to the networks, use that relationship when possible
- Networks
  - Do a better job of coordinating resources and responding to problems

# Agenda

- Current Networking
  - perfSONAR Status (ATLAS, CMS, LHCOPN,LHCONE)
  - Reaching for the Brass Ring (why we monitor)
- Future Networking
  - DYNES
  - LHCONE



# Updates on DYNES

- What is it – read the content here if you need to:  
<http://www.internet2.edu/dynes>
- Basic idea:
  - Provide hardware and OpenSource software to address data intensive science on campuses
    - Switch, data movement server, controller PC for hardware
    - FDT, OSCARS, and perfSONAR for software
    - Goal is to encourage campuses to create a research grade network (e.g. the ‘science dmz’ -  
<http://fasterdata.es.net/fasterdata/science-dmz/> )
  - Can’t provide raw capacity, but is a tool to manage existing capacity
    - Layer 2 networking (e.g. dynamic capacity – possibility of bandwidth guarantees)
    - End to end ‘circuit’ capabilities (e.g. protected VLANs)

# Campus Nets – Clogging Ur Tubes



“Internets”

# Campus Nets – What about Science?



“Internets”

# Campus Nets w/ DYNES Vision

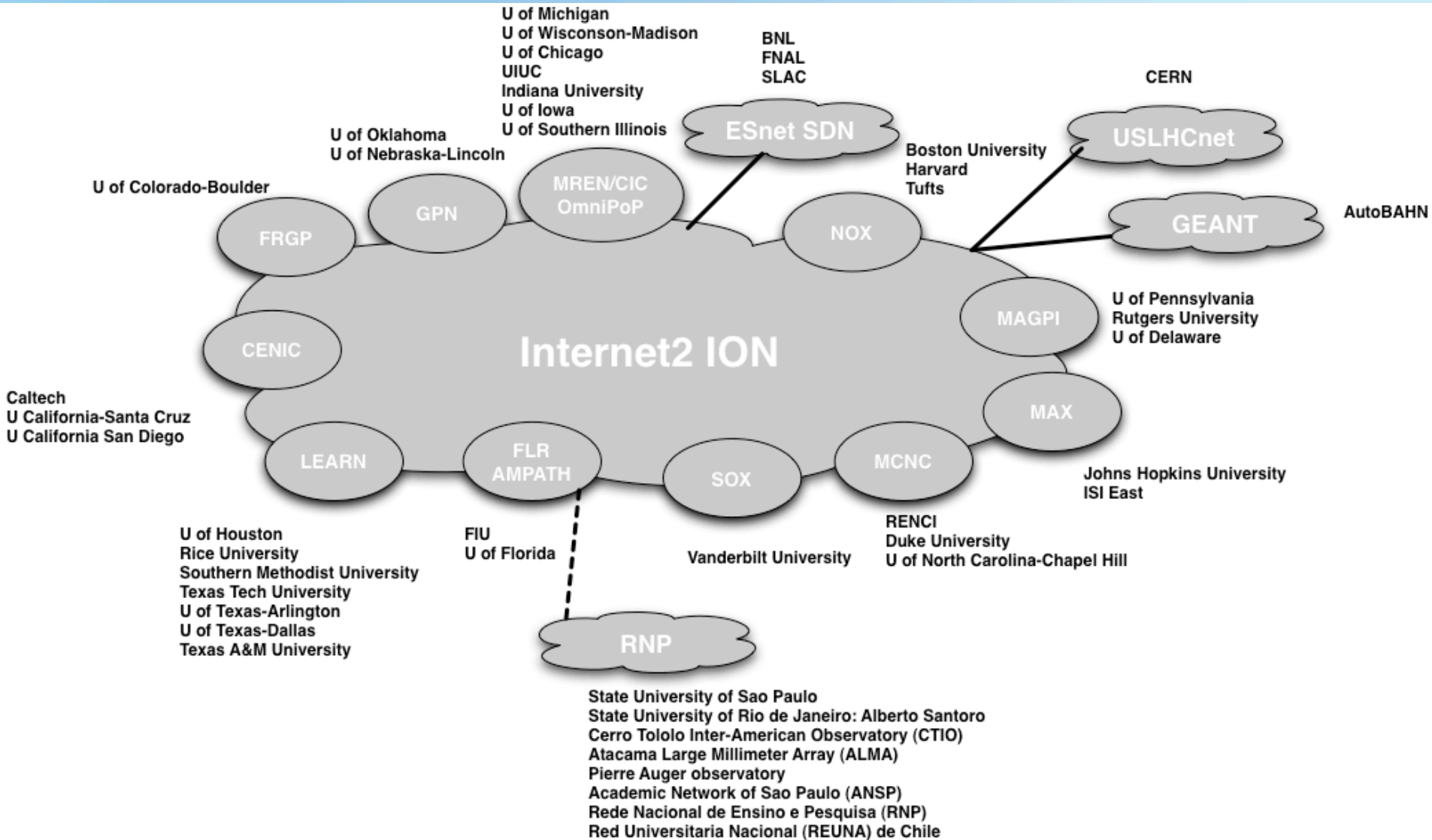


“Internets”

Encapsulated  
Layer 2 (MPLS)

INTERNET  
*2*

# Updates on DYNES – cont.



# Updates on DYNES – cont.

- Status (see web for more details):
  - Group A (~9 sites), deployed and working
  - Group B (~11 sites), deployed, and starting to come online
  - Group C (~14 sites), ordered and being configured, deployment in the next month
  - **We have funding left if you are not connected, and are still interested**
- Related Work:
  - Working w/ AMPATH and RNP in Brazil to connect OSCARS circuits to research facilities (e.g. SPRACE). Demos were done last year and were successful.
  - Early talks with LSST (telescope in Chile) to support management of data flows approaching 80Gbps in 2020
  - Early talks with GlobusOnline to integrate support into this tool to reach DYNES sites using OSCARS and traditional IP networking

# DYNES Open Questions/Next Steps

- Where do we go from here?
- Applications
  - FDT is integrated and can use the APIs to use Layer 2 technologies (OSCARS/ION + maybe someday soon 'OpenFlow')
  - What about PhEDEX/DQ2 directly?
  - FTS (since this is the scheduling bit under the data movers)
  - What about the underlying OSG tools?
    - Which ones make sense, SRM? Others?
  - Why integrate an application?
    - Layer 2 technologies are 'HOT/FAST Lane' compared to campus IP. Can give you a direct path to the Campus WAN and through the regional network (congestion free)
    - IP connectivity may 'work', but its hard to manage end to end (especially for TCP)
    - Data movers that can take advantage of this are more likely to get resources in constrained environments

INTERNET

# DYNES Open Questions/Next Steps

- Network
  - LHCONE (see next) will have support for Layer 2 services
  - Regionals/Campuses in the US are being invited to participate in Layer 2 networks
    - DYNES via Internet2 ION/ESnet SDN, etc.
    - OpenFlow is gaining a lot of traction
- Vision (being implemented by some already) intelligent applications that make the choice for the user.
  - Don't have to care about the network on the bottom, things just 'work'
  - Let the scientists be scientists, not engineers



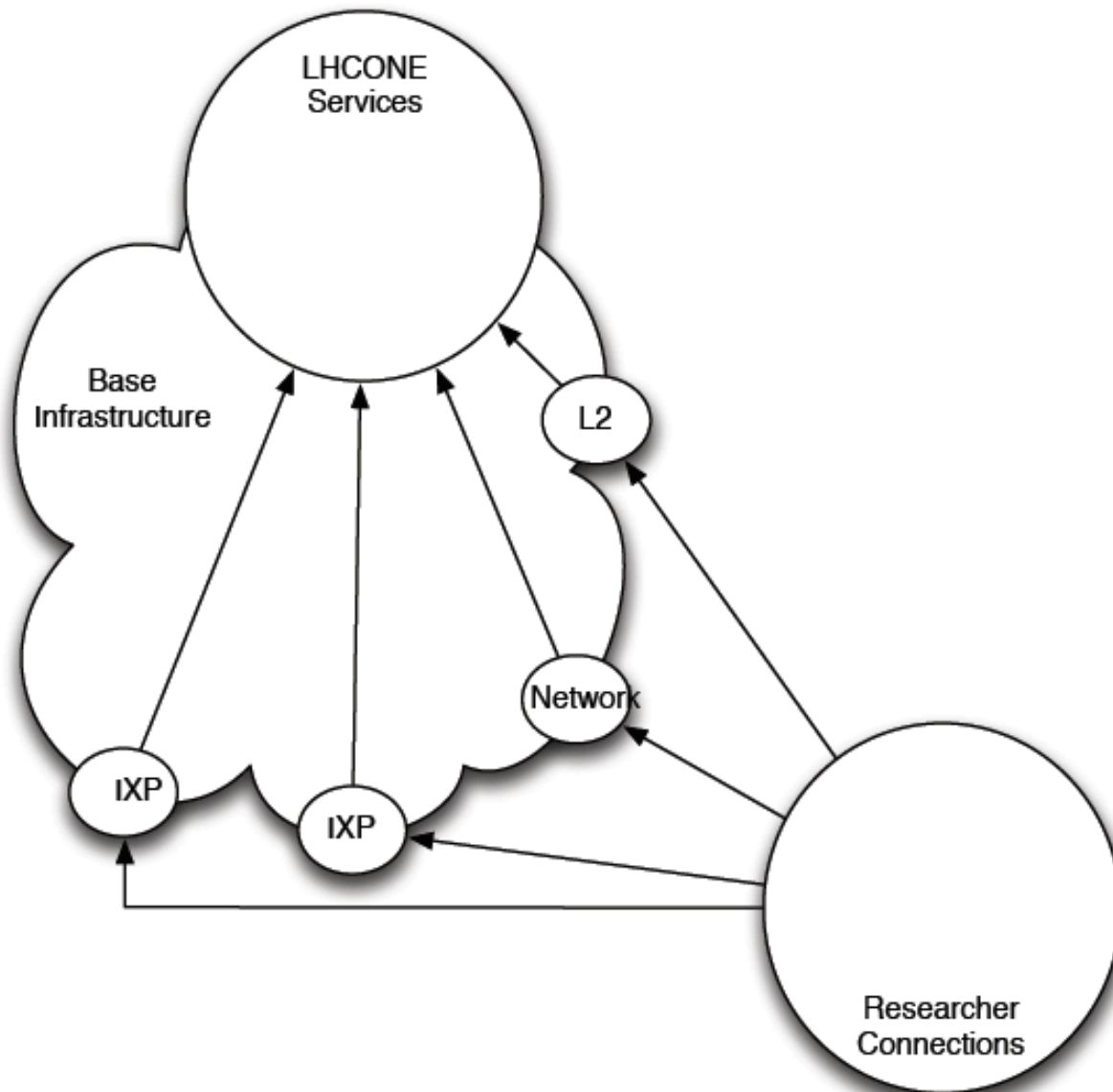
# Agenda

- Current Networking
  - perfSONAR Status (ATLAS, CMS, LHCOPN,LHCONE)
  - Reaching for the Brass Ring (why we monitor)
- Future Networking
  - DYNES
  - LHCONE

# LHCONE

- In case you missed it...
  - No more global VLAN (not scalable, too much of a pain)
  - Direct L2 circuits (e.g. through OSCARS or similar technologies) still being explored
  - Current work is on Islands of L3 VPNs
    - VRF – Virtual [VPN] Routing and Forwarding is being used
- Purpose?
  - Allows participants to move traffic between one another as needed.
  - Built using available components of the R&E networking infrastructure (e.g. ESnet, GEANT, Internet2, USLHCnet, ACE, CERNLIGHT Starlight, MANLAN, etc.)

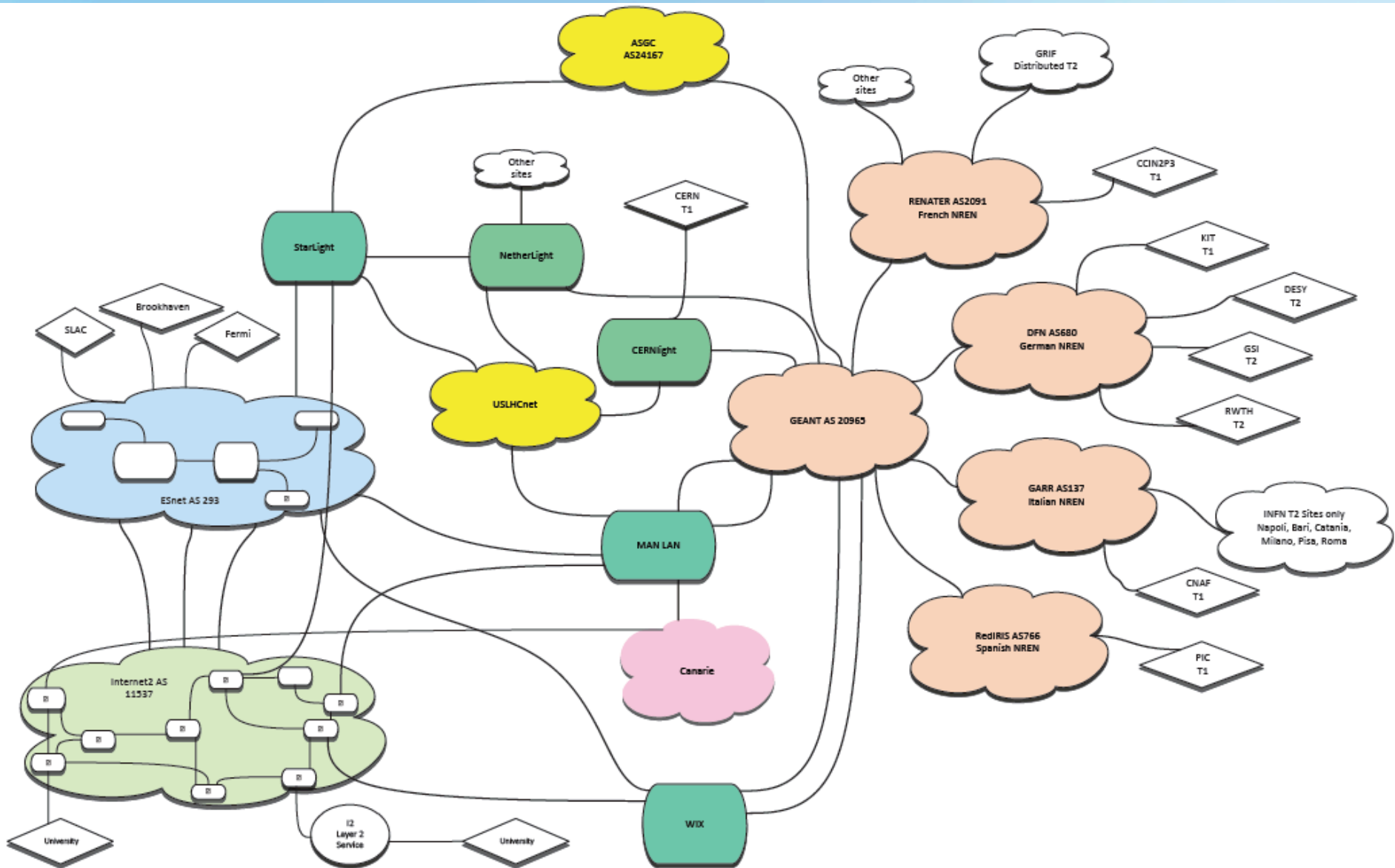
# LHCONE – The Idea



# LHCONE Guts

- How is this done?
  - it is possible to implement a shared broadcast domain using a specific IP prefix or it can be implemented via a VRF
    - Virtual routers (Internet2) vs dedicated resources (e.g. Starlight Cisco)
  - Difference between this and shared VLAN:
    - There are routed boundaries between portions of the shared structures
    - There is a requirements for the exchange of routing information across those boundaries.
      - This information will be exchanged using BGP.

# LHCONE – More Exact



# LHCONE – Whats Next?

- Hard for me to answer this – I am not the user 😊
- As “users”, you all have some important things to do:
  - Do your science as before
  - Can you reach the places you need to reach? Are things any better or worse than before?
  - Is your life measurably better with LHCONE vs without (don't answer this now, have a cookie or something first)
    - Since this is ‘just the network’ you may not even notice (unless its not working)
- All kidding aside – the next steps for this lie with the stakeholders, and it is anticipated that you will ‘vote’ with your opinions as well as funding dollars.

# Closing Thoughts

- Monitoring
  - Monitoring is not a sexy topic, it's a means to an end
  - We (networks, as well as VOs) need it to make sure that things are working so that users (all of you) aren't sad
- L2 & Advanced Networking
  - Lots of opportunity to use new technologies
  - Hard sale to add features into applications
  - We (network providers) can 'help' with adaptations, but we don't have the manpower/funding to lead in this area.



**Addressing the “things that go bump in the net” –  
perfSONAR/DYNES/LHCONE**

March 20<sup>th</sup> 2012, OSG/ATLAS/CMS

Jason Zurawski, Internet2 Research Liaison

[zurawski@internet2.edu](mailto:zurawski@internet2.edu)

For more information, visit <http://www.internet2.edu/research/>