

Building A Science DMZ

Eli Dart, Network Engineer

ESnet Network Engineering Group

Joint Techs, Winter 2013

Honolulu, HI

January 13, 2013





Outline of the Afternoon

Eli Dart, ESnet

- Science DMZ architecture, security

Brian Tierney, ESnet

- Data transfer node, tools overview

Raj Kettimuthu, ANL and University of Chicago

- Globus Online

-Short break-

Jason Zurawski, Internet2

- perfSONAR

Guy Almes, Texas A&M University

- University case study



Motivation

Science data increasing both in volume and in value

- Higher instrument performance
- Increased capacity for discovery
- Analyses previously not possible

Lots of promise, but only if scientists can actually work with the data

- Data has to get to analysis resources
- Results have to get to people
- People have to share results

Common pain point – data mobility

- Movement of data between instruments, facilities, analysis systems, and scientists is a gating factor for much of data intensive science
- Data mobility is not the only part of data intensive science – not even the most important part
- However, without data mobility data intensive science is hard

We need to move data – how can we do it consistently well?

Motivation (2)



Networks play a crucial role

- The very structure of modern science assumes science networks exist – high performance, feature rich, global scope
- Networks enable key aspects of data intensive science
 - Data mobility, automated workflows
 - Access to facilities, data, analysis resources

Messing with the network is unpleasant for most scientists

- Not their area of expertise
- Not where the value is (no papers come from messing with the network)
- Data intensive science is about the science, not about the network
- However, it's a critical service – if the network breaks, everything stops

Therefore, infrastructure providers must cooperate to build consistent, reliable, high performance network services for data mobility

Here we describe one blueprint, the Science DMZ model, for addressing issues of high performance data mobility

TCP Background



Networks provide connectivity between hosts – how do hosts see the network?

- From an application's perspective, the interface to “the other end” is a socket
- Other similar constructs exist for non-IP protocols
- Communication is between applications – mostly over TCP

TCP – the fragile workhorse

- TCP is (for very good reasons) timid – packet loss is interpreted as congestion
- Packet loss in conjunction with latency is a performance killer
- Like it or not, TCP is used for the vast majority of data transfer applications

TCP Background (2)



It is far easier to architect the network to support TCP than it is to fix TCP

- People have been trying to fix TCP for years (with some success)
- However, here we are – packet loss is still the number one performance killer in long distance high performance environments

Pragmatically speaking, we must accommodate TCP

- Implications for equipment selection
 - Equipment must be able to accurately account for packets
 - Equipment must be able to provide loss-free service
- Implications for network architecture, deployment models
 - Infrastructure must be designed to allow easy troubleshooting
 - Test and measurement tools are critical – they have to be deployed



How Do We Accommodate TCP?

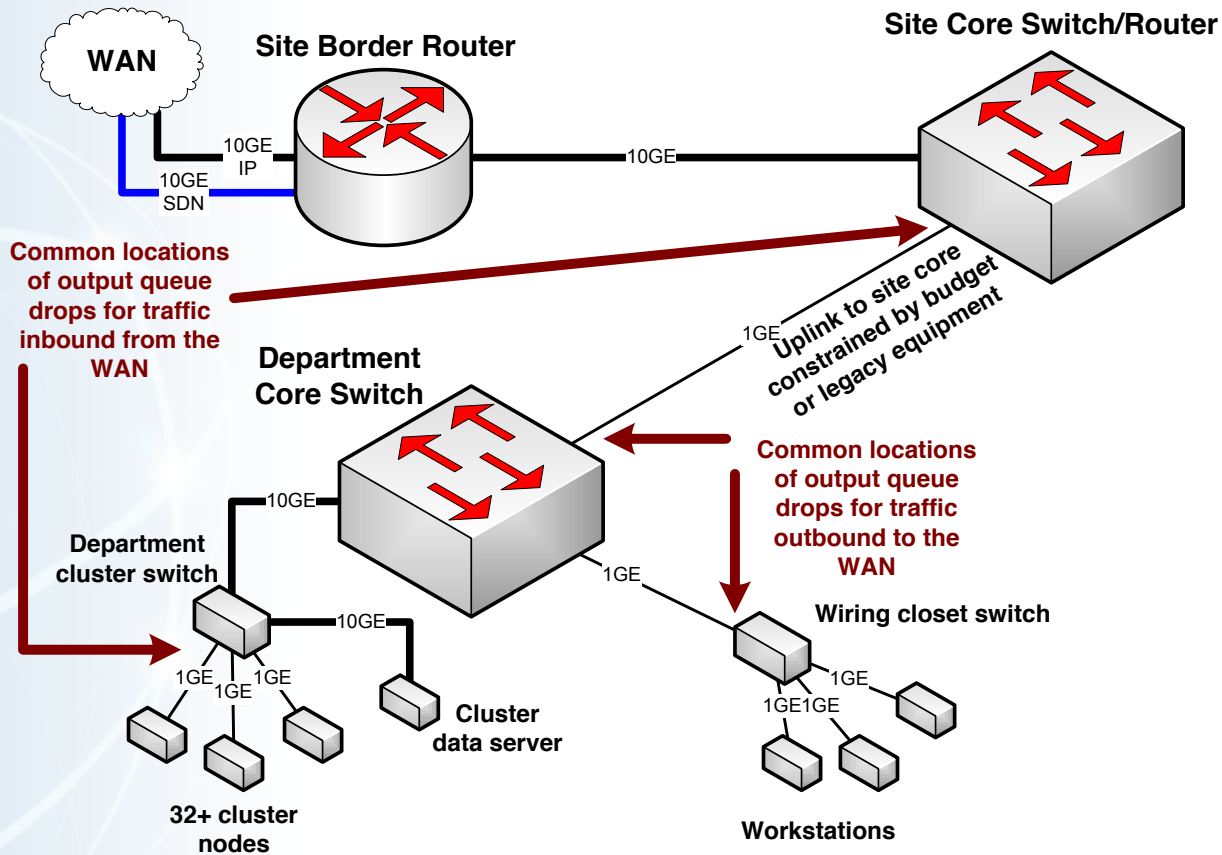
High-performance wide area TCP flows must get loss-free service

- Sufficient bandwidth to avoid congestion
- Deep enough buffers in routers and switches to handle bursts
 - Especially true for long-distance flows due to packet behavior
 - No, this isn't buffer bloat

Equally important – the infrastructure must be verifiable so that clean service can be provided

- Stuff breaks
 - Hardware, software, optics, bugs, ...
 - How do we deal with it in a production environment?
- Must be able to prove a network device or path is functioning correctly
 - Accurate counters must exist and be accessible
 - Need ability to run tests - perfSONAR
- Small footprint is a huge win – small number of devices so that problem isolation is tractable

Example Loss Locations





Services Overview – Wide Area

Data transfer takes advantage of wide area services:

High-performance routed IP with global connectivity

- Bread and butter
- Must be high-bandwidth, verifiably loss-free in the general case

Virtual circuit service

- Traffic isolation, traffic engineering
- Bandwidth and service guarantees
- Support for non-IP protocols

Test and measurement – perfSONAR

- Enable testing, verification of performance, problem isolation
- Understand nominal conditions → what's normal, what's broken



Services Overview – Site/Campus

High performance routed IP

- Well-matched to wide area science service
- Verifiably loss-free

Circuit termination/endpoints

- DYNES, Tier1, RDMA, ...
- Remote filesystem mounts
- Non-IP protocols

Data sources and sinks

- Instruments and facilities
- Analysis resources
- Data systems

It is at the site or campus that it all comes together – scientists, instruments, data, analysis

The Data Transfer Trifecta: The “Science DMZ” Model



Dedicated
Systems for
Data Transfer

Network
Architecture

Performance
Testing &
Measurement

Data Transfer Node

- High performance
- Configured for data transfer
- Proper tools

Science DMZ

- Dedicated location for DTN
- Easy to deploy - no need to redesign the whole network
- Additional info:
<http://fasterdata.es.net/>

perfSONAR

- Enables fault isolation
- Verify correct operation
- Widely deployed in ESnet and other networks, as well as sites and facilities

Science DMZ Service Interaction



WAN entry

- How do wide area services enter the site?
- If they don't come to the Science DMZ first, there must be a clean path to the Science DMZ
- Clean wide area path for long-distance flows is key

Circuit services entry

- Virtual circuits support DYNES, LHC experiments, remote filesystem mounts, non-IP protocols, ...

Local resources

- Data Transfer Nodes
- Test and measurement (perfSONAR)

Security policy

- Separation of science and business traffic

Science DMZ Security



Goal – disentangle security policy and enforcement for science flows from that of business systems

Rationale

- Science flows are relatively simple from a security perspective
- Narrow application set on Science DMZ
 - Data transfer, data streaming packages
 - No printers, document readers, web browsers, building control systems, staff desktops, etc.
- Security controls that are typically implemented to protect business resources often cause performance problems
- Sizing security infrastructure on business networks for large science flows is expensive



The Ubiquitous Firewall

Remember the motivation: high performance data mobility in support of scientific applications

The workhorse device of network security – the firewall – has a poor track record in high-performance contexts

- Firewalls are typically designed to support a large number of users/devices, each with low throughput requirements
- Data intensive science typically generates a much smaller number of connections that are much higher throughput requirements
- This mismatch in behavior often violates the design parameters of the firewall, causing problems – for example:
 - Internal bottlenecks cause packet loss
 - State table timeouts cause file corruption or transfer failures



Firewall Capabilities and Science Traffic

Firewalls have an incredible amount of sophistication in an enterprise setting

- Application layer protocol analysis (HTTP, POP, MSRPC, etc)
- Built-in VPN servers
- User awareness
- Large number of concurrent connections supported

Data-intensive science flows typically look different on the wire

- Common case – data on filesystem A needs to be on filesystem Z
 - Data transfer tool verifies credentials over an encrypted channel
 - Then open a socket or set of sockets, and send data until done (1TB, 10TB, 100TB, ...)
- One workflow can use 50% or more of a network link

Do we have to use a firewall?



Firewalls As Access Lists

In cases where the workload is to push terabytes of data over a small number of sockets at high speed, advanced firewall features don't mean much

- We're not managing per-user-class access control through a VPN
- We don't need to do HTTP protocol analysis
- In fact, the enterprise features don't really play here at all

When you ask a firewall administrator to allow data transfers through the firewall, what do they ask for?

- IP address of your host
- IP address of the remote host
- Port range
- ***That looks like an ACL to me – I can do that on the router!***

Firewalls make expensive, low-performance ACLs – especially since ACL capabilities are typically built into the router



Security Without Firewalls

Data intensive science traffic interacts poorly with firewalls

- Most of the high-value features of firewalls do not apply
- Firewalls incur significant performance costs
- For high-performance data mobility workflows, firewalls work just like address/port filters

Does this mean we ignore security? **NO!**

- We **must** protect our systems
- We just need to find a way to do security that does not prevent us from getting the science done

Lots of options

- Intrusion detection (Bro, Snort, others), flow analysis, ...
- Tight ACLs reduce attack surface (possible in many but not all cases)
- ***Key point – performance is a mission requirement, and the security policies and mechanisms that protect the Science DMZ should be architected so that they serve the mission***

Security posture should be a product of collaboration between security people, network people, and science constituents

Science DMZ Takes Many Forms



There are a lot of ways to combine these things – it all depends on what you need to do

- Small installation for a project or two
- Facility inside a larger institution
- Institutional capability serving multiple departments/divisions
- Science capability that consumes a majority of the infrastructure

Some of these are straightforward, others are less obvious

Key focal point is the high-latency path for TCP

- The high-latency path for TCP is where loss has the worst impact
- Concentrate on minimizing complexity, maximizing instrumentation



Ad Hoc Deployment

This is often what gets tried first

Data transfer node deployed where the owner has space

- This is often the easiest thing to do at the time
- Straightforward to turn on, hard to achieve performance

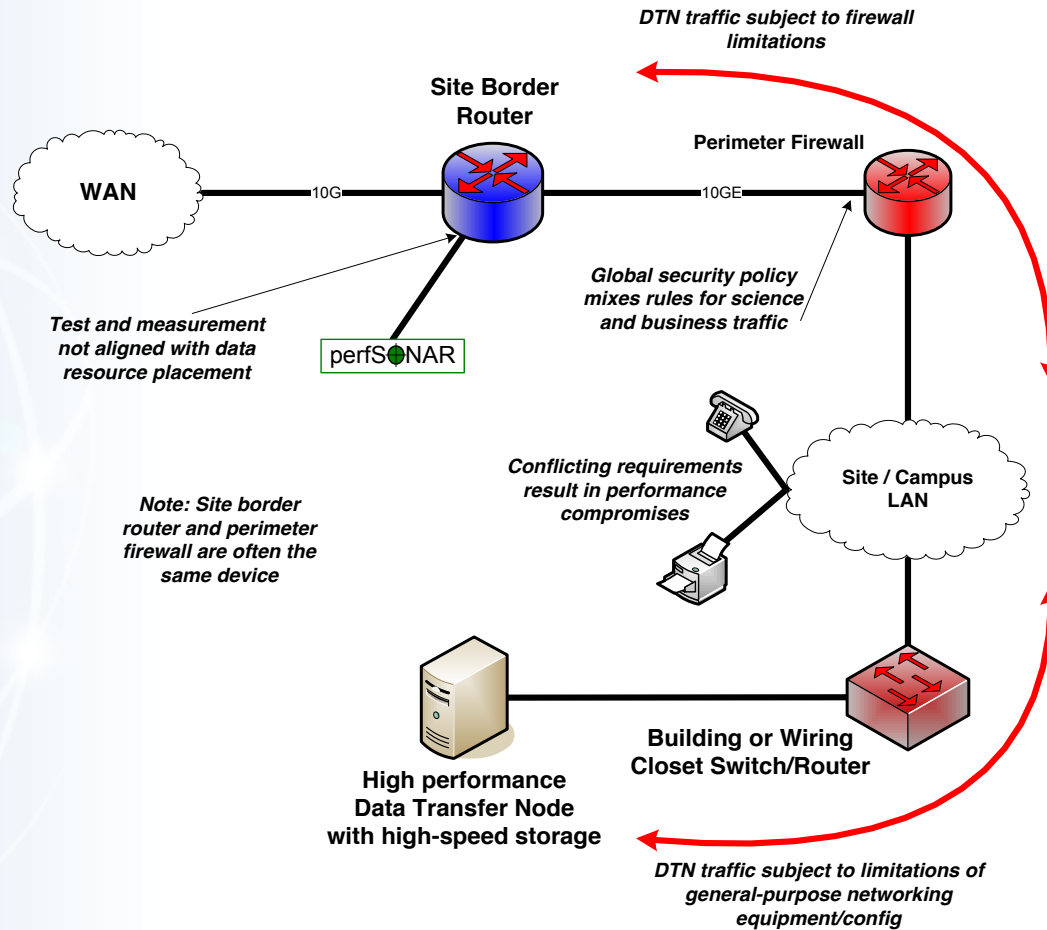
perfSONAR at the border

- This is a good start
- Need a second one next to the DTN

Entire LAN path has to be sized for data flows

Entire LAN path is part of any troubleshooting exercise

Ad Hoc DTN Deployment





Small-scale or Prototype Deployment

Add-on to existing network infrastructure

- All that is required is a port on the border router
- Small footprint, pre-production commitment

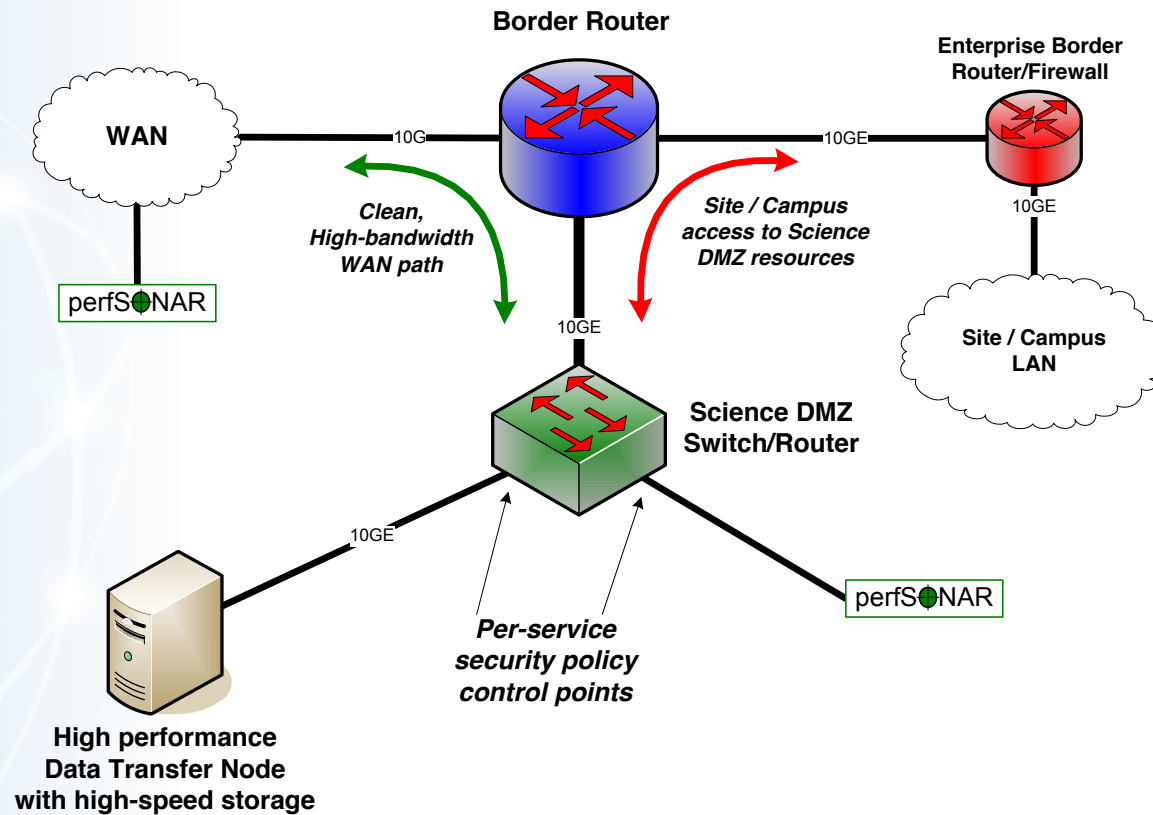
Easy to experiment with components and technologies

- DTN prototyping
- perfSONAR testing

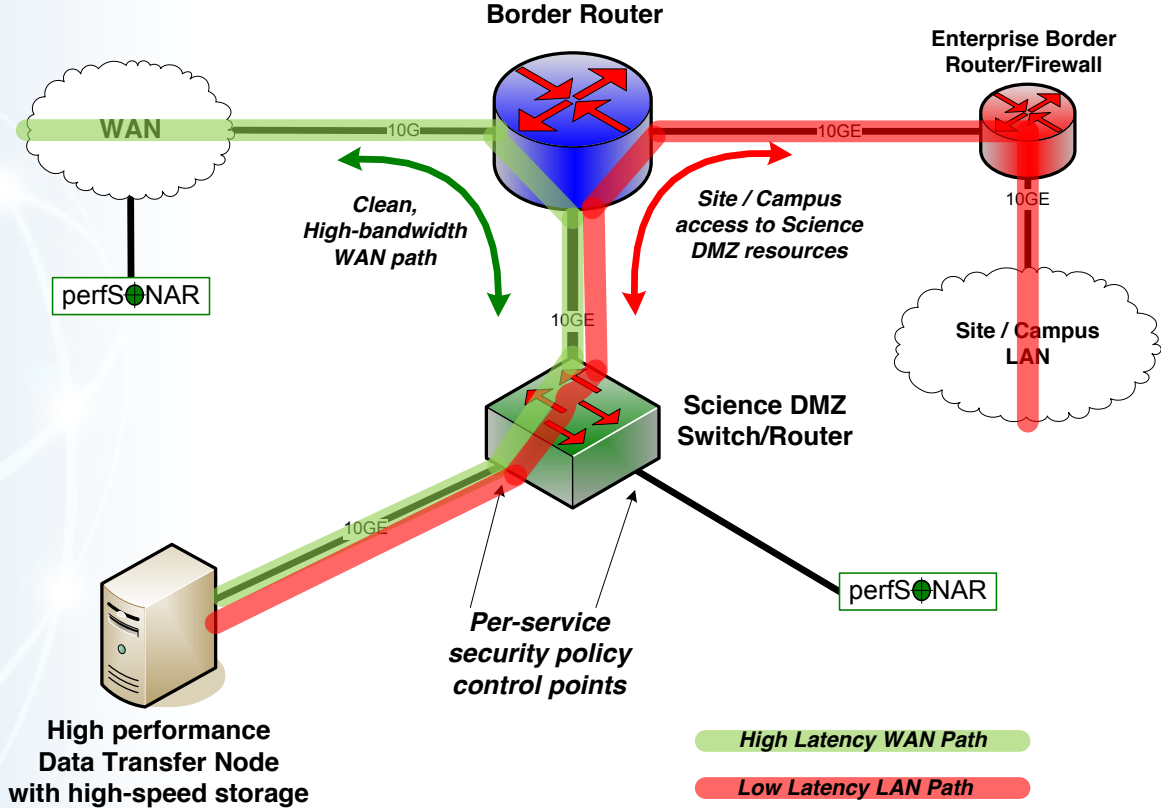
Limited scope makes security policy exceptions easy

- Only allow traffic from partners
- Add-on to production infrastructure – lower risk

Prototype Science DMZ



Prototype Science DMZ Data Path





Prototype With Virtual Circuits

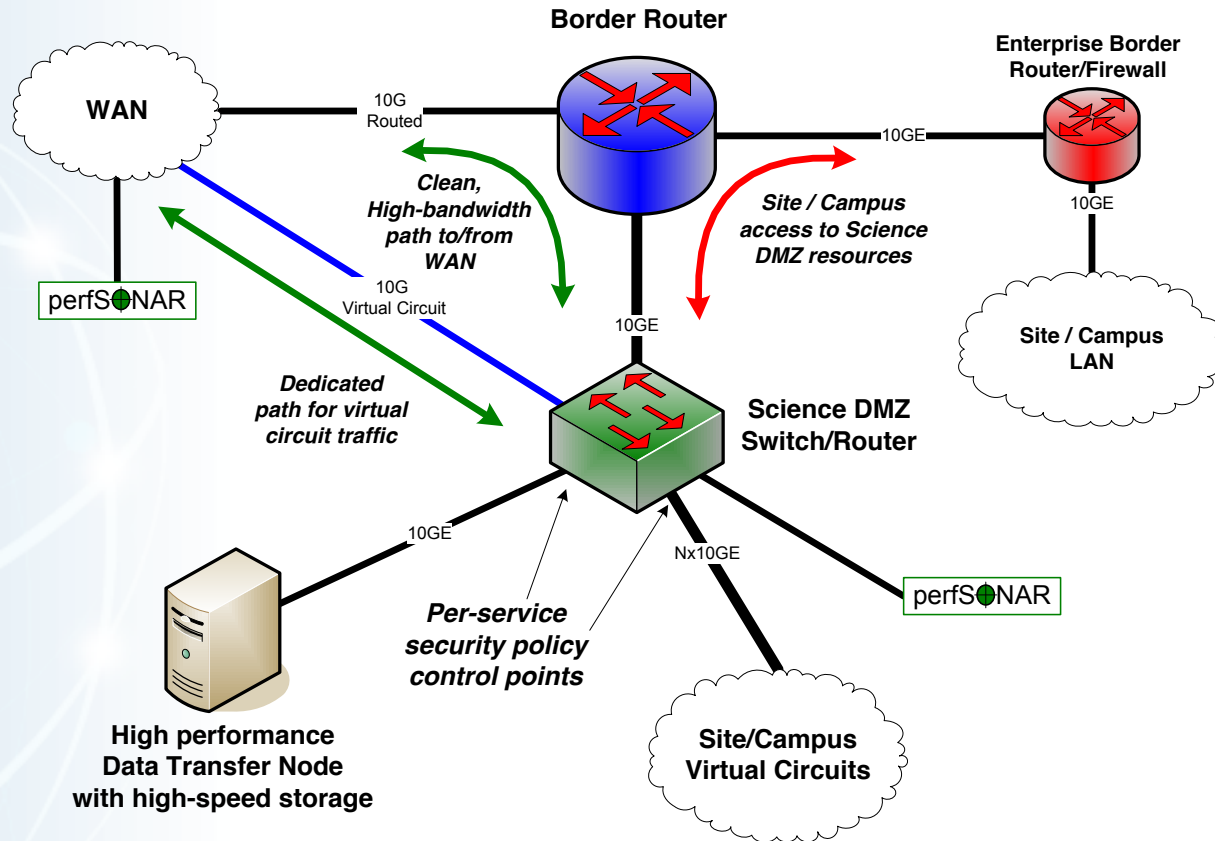
Small virtual circuit prototype can be done in a small Science DMZ

- Perfect example is a DYNES deployment
- Virtual circuit connection may or may not traverse border router

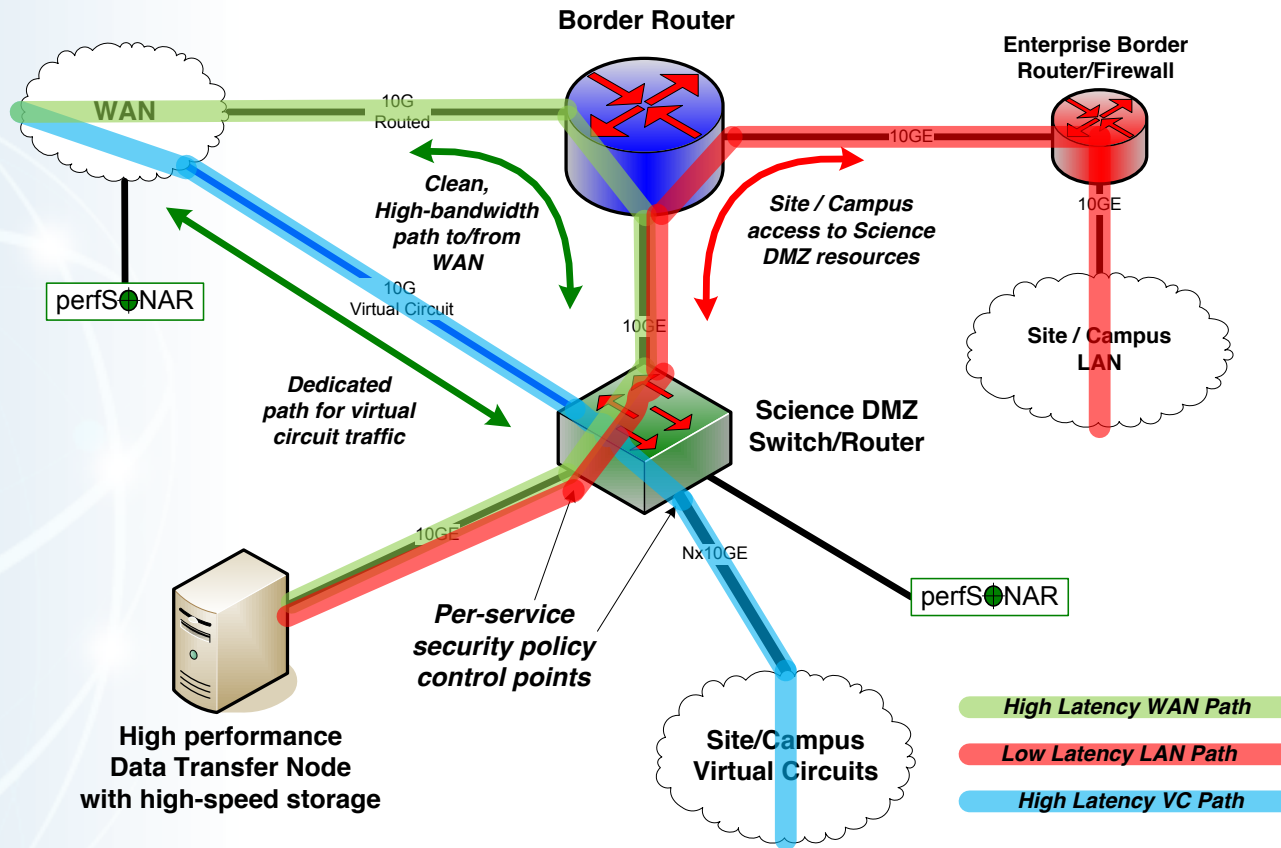
As with any Science DMZ deployment, this can be expanded as need grows

In this particular diagram, Science DMZ hosts can use either the routed or the circuit connection

Virtual Circuit Prototype Deployment



Virtual Circuit Prototype Data Path





Support For Multiple Projects

Science DMZ architecture allows multiple projects to put DTNs in place

- Modular architecture
- Centralized location for data servers

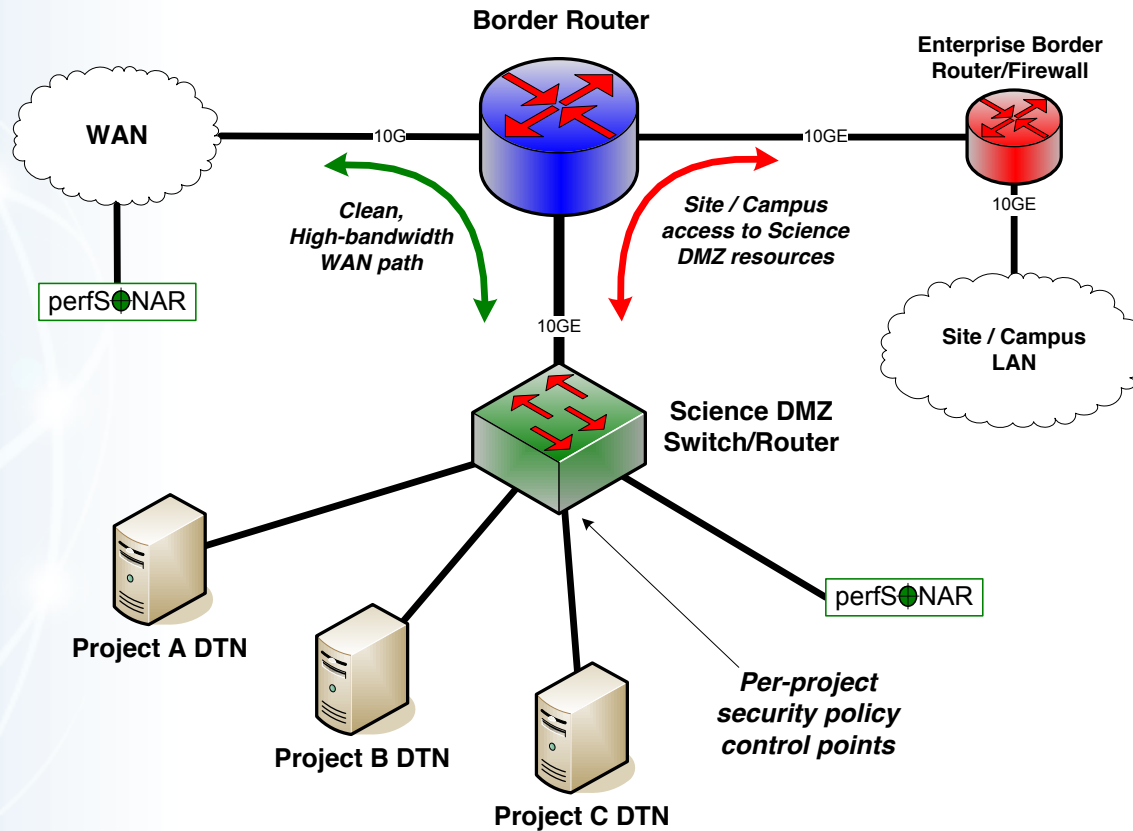
This may or may not work well depending on institutional politics

- Issues such as physical security can make this a non-starter
- On the other hand, some shops already have service models in place

On balance, this can provide a cost savings – it depends

- Central support for data servers vs. carrying data flows
- How far do the data flows have to go?

Multiple Projects



Supercomputer Center Deployment



High-performance networking is assumed in this environment

- Data flows between systems, between systems and storage, wide area, etc.
- Global filesystem often ties resources together
 - Portions of this may not run over Ethernet (e.g. IB)
 - Implications for Data Transfer Nodes

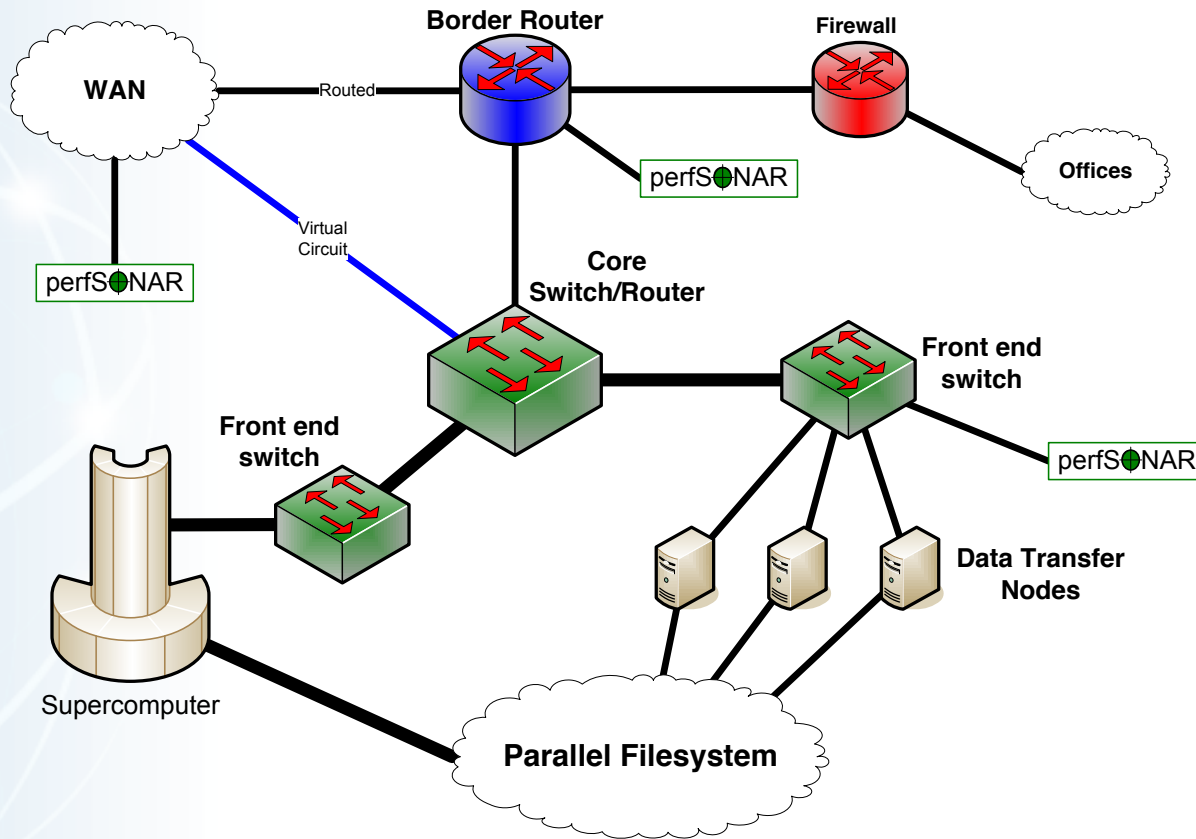
“Science DMZ” may not look like a discrete entity here

- By the time you get through interconnecting all the resources, you end up with most of the network in the Science DMZ
- This is as it should be – the point is appropriate deployment of tools, configuration, policy control, etc.

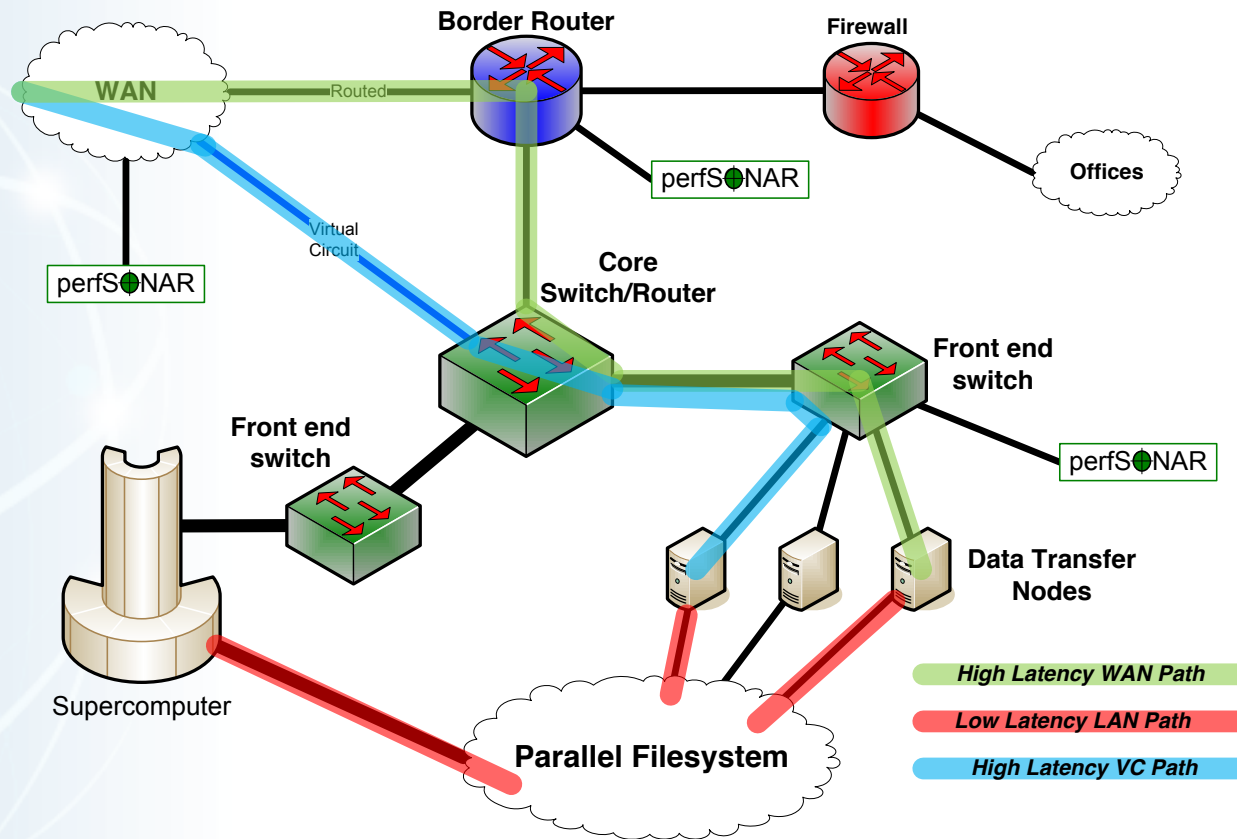
Office networks can look like an afterthought, but they aren't

- Deployed with appropriate security controls
- Office infrastructure need not be sized for science traffic

Supercomputer Center



Supercomputer Center Data Path





Major Data Site Deployment

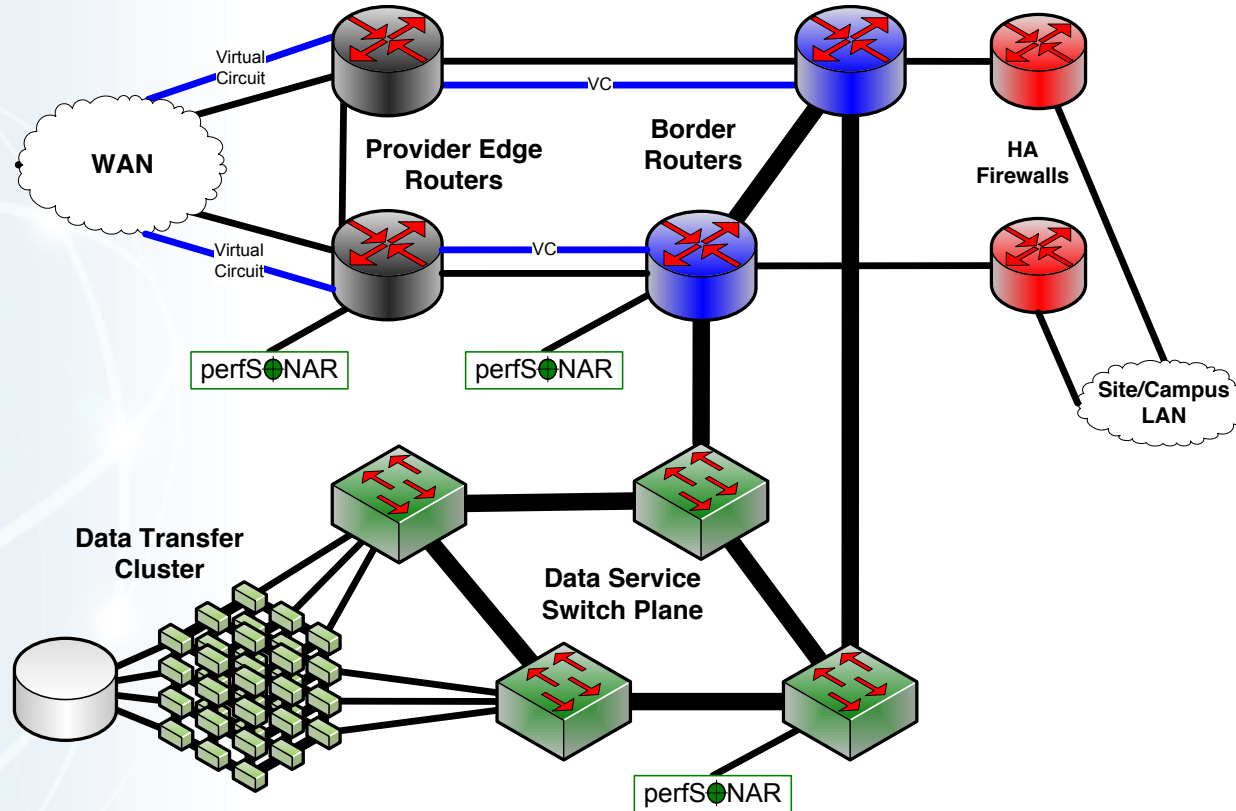
In some cases, large scale data service is the major driver

- Huge volumes of data – ingest, export
- Large number of external hosts accessing/submitting data

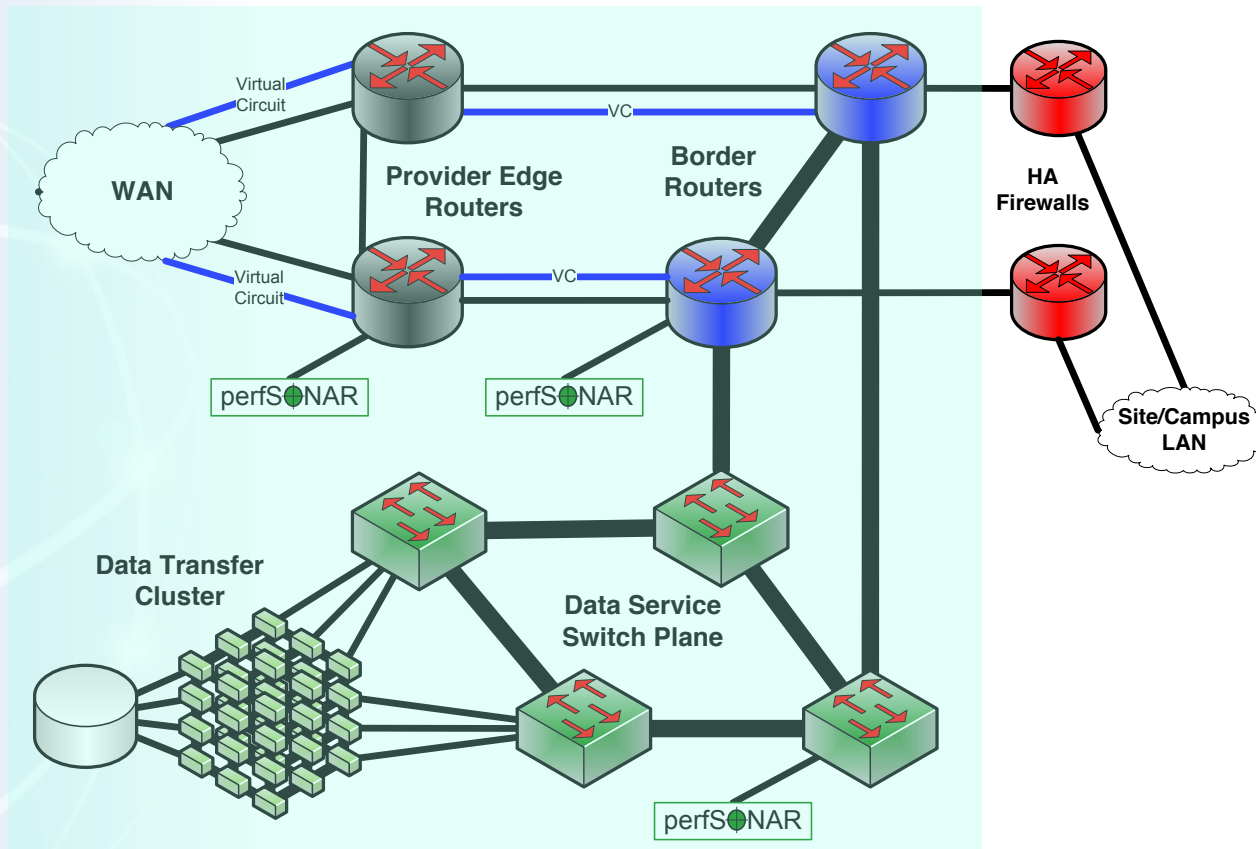
Single-pipe deployments don't work

- Everything is parallel
 - Networks (Nx10G LAGs, soon to be Nx100G)
 - Hosts – data transfer clusters, no individual DTNs
 - WAN connections – multiple entry, redundant equipment
- Choke points (e.g. firewalls) cause problems

Data Site – Architecture



Data Site – Data Path





Distributed Science DMZ

Fiber-rich environment enables distributed Science DMZ

- No need to accommodate all equipment in one location
- Allows the deployment of institutional science service

WAN services arrive at the site in the normal way

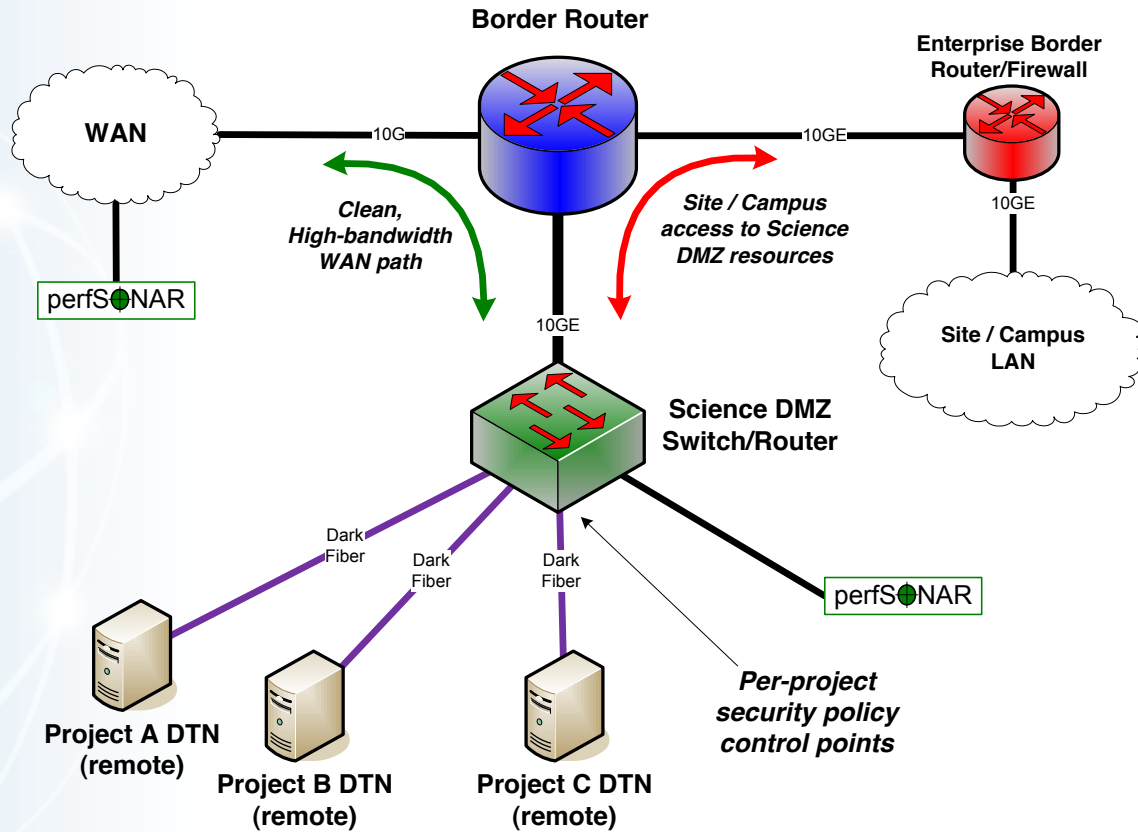
Dark fiber distributes connectivity to Science DMZ services throughout the site

- Departments with their own networking groups can manage their own local Science DMZ infrastructure
- Facilities or buildings can be served without building up the business network to support those flows

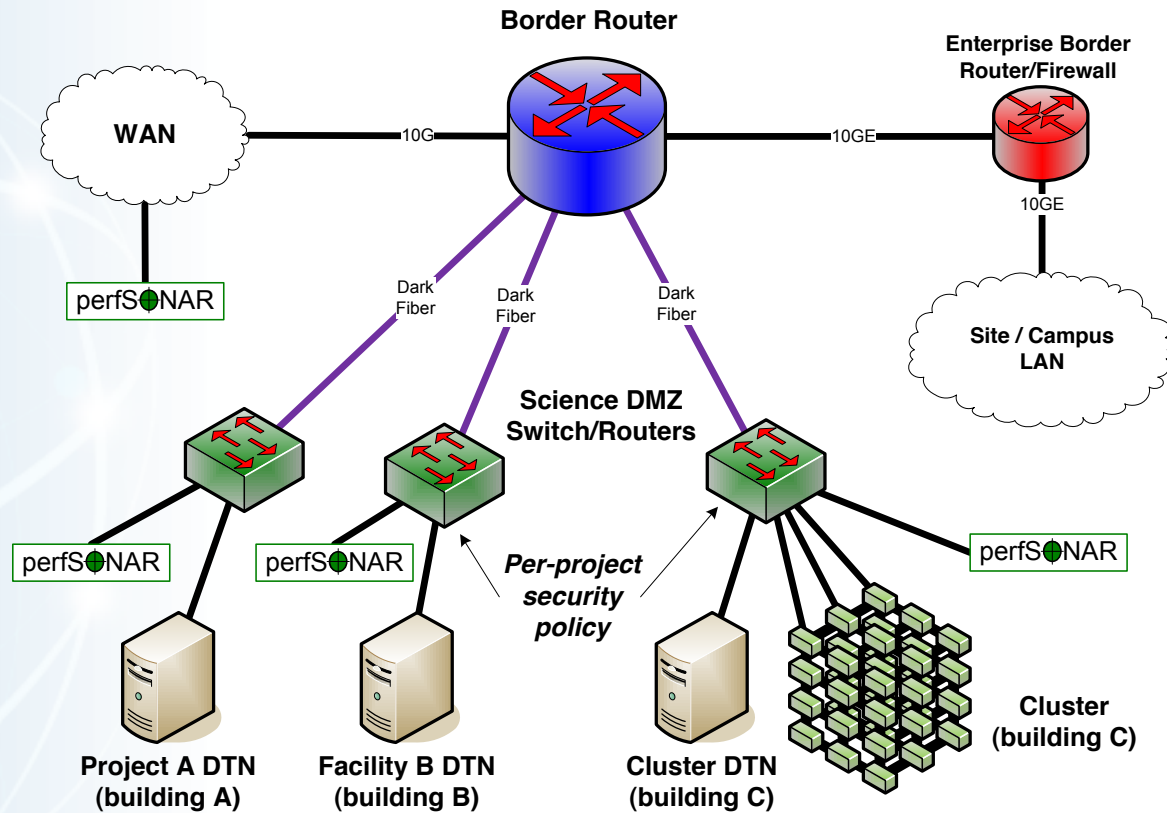
Security is made more complex

- Remote infrastructure must be monitored
- Several technical remedies exist (arpwatch, no DHCP, separate address space, etc.)
- Solutions depend on relationships with security groups

Distributed Science DMZ – Dark Fiber



Multiple Science DMZs – Dark Fiber





Development Environment

One thing that often happens is that an early power user of the Science DMZ is the network engineering group that builds it

- Service prototyping
- Upgrade planning for production Science DMZ
- Deployment of test applications for other user groups to demonstrate value

The production Science DMZ is just that – production

- Once users are on it, you can't take it down to try something new
- Stuff that works tends to attract workload

Take-home message: plan for multiple Science DMZs from the beginning – at the very least you're going to need one for yourself



Common Threads

Two common threads exist in all these examples

Accommodation of TCP

- Wide area portion of data transfers traverses purpose-built path
- High performance devices that don't drop packets

Ability to test and verify

- When problems arise (and they always will), they can be solved if the infrastructure is built correctly
- Small device count makes it easier to find issues
- Multiple test and measurement hosts provide multiple views of the data path
 - perfSONAR nodes at the site and in the WAN
 - perfSONAR nodes at the remote site



Science DMZ Benefits

Better access to remote facilities by local users

Local facilities provide better service to remote users

Ability to support science that might otherwise be impossible

Metcalf's Law – value increases as the square of connected devices

- Communication between institutions with functional Science DMZs is greatly facilitated
- Increased ability to collaborate in a data-intensive world

Cost/Effort benefits also

- Shorter time to fix performance problems – less staff effort
- Appropriate implementation of security policy – lower risk
- No need to drag high-speed flows across business network → lower IT infrastructure costs



Questions?

Thanks!

Eli Dart - dart@es.net

<http://www.es.net/>

<http://fasterdata.es.net/>



U.S. DEPARTMENT OF
ENERGY
Office of Science





Outline of the Afternoon

~~Eli Dart, ESnet~~

- ~~• Science DMZ architecture, security~~

Brian Tierney, ESnet

- Data transfer node, tools overview

Raj Kettimuthu, ANL and University of Chicago

- Globus Online

-Short break-

Jason Zurawski, Internet2

- perfSONAR

Guy Almes, Texas A&M University

- University case study