# Efforts to Build a T3gs at Illinois

Mark Neubauer, Dave Lesny
University of Illinois at Urbana-Champaign

US ATLAS Distributed Facilities Meeting
LIGO Livingston Observatory

March 3, 2009

# Tier-3 Task Force and Site Structure

## Key recommendation of Tier-3 Task Force (T3TF):

(my paraphrasing)

• US ATLAS computing infrastructure should be deeper than the Tier1+2's and accommodate a set of 4 possible Tier-3 architectures: T3gs, T3g, T3w, and T3af.

  • The goal is to utilize significant, additional computing resources to produce a system which is flexible and nimble enough to weather the realities of ATLAS analysis once the data starts flowing from CERN.

  • The Tier-2's should be considered a finite and precious resource that must be used efficiently with as much end-user analysis as feasible done on the local Tier-3's

# What is a T3gs?

## T3gs: Tier-3 w/ Grid Services

- Cluster w/ software deployed for pAthena jobs, DQ2 services, and possibly LFC

- Expected to have "significant" CPU and storage resources:
  - ≈ several hundred CPU cores
  - ≈ several tens of TB

- Expected to a level of infrastructure (power, cooling), networking, and administrative support commensurate with available resources

Functionally-equivalent to Tier-2 from a services perspective

Distinct from Tier-2 in important ways:

- Size (not likely to be as large at Tier-2 in terms of available resources)
- Local policy control (resources funded at institutional-level)

# Why bother trying to deploy a T3gs?

## Value Added From a T3gs

1) Leverage University and Lab-based computing to increase overall capacity for production tasks like Data Format Production (e.g. $D^1PD \rightarrow D^2PD$) and Monte Carlo production

   - Given the specific T3 nature of T3gs, this is expected to be during intensive periods of critical need (e.g data reprocessing, MC before conference seasons)

   - T3gs with their T2-like configuration and capabilities in best position to pitch in when needed

2) Offload from Tier-2's computationally intensive jobs that do not need high-bandwidth access to data (e.g. ME calculations, MVAs, pseudo-expt generation for systematic evaluations). Generically true for any Tier-3 site (T3gs could also offload data-intensive jobs given its nature)

See the T3TF Document for some quantitative arguments regarding 1) and 2)

3) Spread around the wealth (of ATLAS computing knowledge)

My candid impression in trying to get our T3gs (T2-like) site together:

• Much of the documentation resides with a few experts (T2 site admins, etc). Some is contained within a disparate set of Twiki pages, NNs and disseminated when need arises through email, HN, etc.

-The experts have been very willing to help out when asked!

- To an outsider, this can appear to be a very large, very black box (I know that most don't need to know how this sausage is made..)

Main point: there is some value for outsiders like us to try to replicate a T2-like site, to help shed some light on how all of this works (widen knowledge base for future T3 installations)

# The IllinoisHEP Grid Site

We've been participating OSG site for ≈ 1 yr

We've been plugged into the Panda job mgmt on a similar time scale

• We've participated (not on purpose!) in M* reprocessing

• Done FDR data analysis on our site (on purpose)

Up until a month or so ago, we've been peer-configured to utilize the IU SE

• We've basically been running as a T3g before we came up with the name!

# IllinoisHEP in Action



Fully integrated into the pATHENA ATLAS job control infrastucture.

Analysis job running over the full FDR dataset to the IllinoisHEP Tier-3 site

# IllinoisHEP in Action (Ganglia Monitoring)

# Infrastructure & Support at Illinois

There is a great deal of infrastructure on campus for scientific HPC

• National Center for Supercomputing Applications (NCSA)

  • Enterprise-quality building in terms of space, power, cooling, networking
  • Blue Waters in 2011: World's 1st sustained Petaflop computer (likely of limited use for our particular applications)

• 10 Gbps connectivity to ICCN, Internet-2, … w/ access from our Physics Building (Loomis Lab) and NCSA via a "Campus Research Network" designed to avoid bottlenecks in the campus firewalls for data-intensive research (we're the poster child for this)

• Loomis Lab has sufficient space, power, cooling, networking for a few full racks of compute nodes and storage devices

• Senior Research Physicist (Dave Lesny): The one who put all of this together

• I have some past experience with scientific computing in putting together and operating the CDF CAF

# Deployment Strategy

Note: We are <span style="color:red">not</span> a large site at this point in time
- 16 CPU cores and 20 TB of disk

Our focus has <span style="color:red">not</span> been on deploying a large amount of hardware, but rather deploying the <span style="color:red">full set of services</span> required to make our site a functional T3gs
- Some of our "utility" hardware needs to be upgraded/replaced
- We've deployed extensive monitoring (via Ganglia) to discover bottlenecks
- Leave scaling of CPU cores and storage until later

We've chosen to deploy our IllinoisHEP site in Loomis Lab rather than NCSA in order to have "full" control over what we do in the development stage

Once our T3gs is in production and scaled to a reasonable level, we would like to consider replicating this site to NCSA within a partnership that we are trying to foster (Neubauer NCSA faculty fellow)

In other words, we would like to house a future "substantial" procurement in NCSA
- This depends upon a number of factors like a willing partnership with NCSA, demonstrated utility of the IllinoisT3gs, and external $$

IllinoisHEP T3gs

Service
Compute Element
Storage Element
DDM
Worker

ICCN, Internet2, …

10 Gbps

Foundry FastIron SuperX

Campus Research Network (public)

1 Gbps

GUMS — osggums

NAT Routers — osgnx1 osgnx0

LFC Server — osgx3

DQ2 Server — osgx4

dCache admin, dCache door SRM — osgx1

dCache door — osgx2

Globus Gatekeeper GridFTP Server — osgx0

HP 2900 — Internal Network (private)

Worker nodes (16 CPU cores)

wn00
wn01

/home/osg* — fx01 fx02

2 × 1 Gbps

Condor Master — condor

System Monitoring — ganglia

mysql00 — LFC/DQ2 MySQL DBs

se01 — dCache PoolMgr

se00 — pnfs server

pn00 pn01
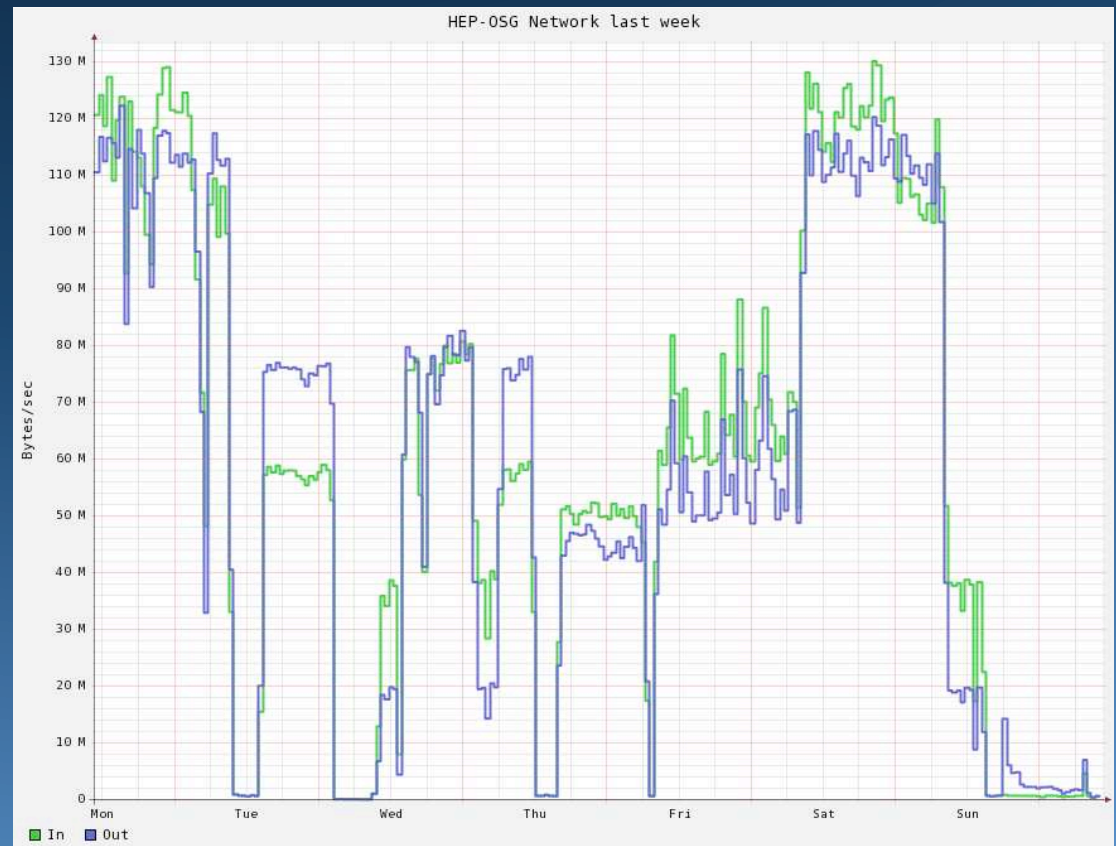
dCache Pools (20 TB)

# Some Baby Pictures

# SE and LFC Tests

Dave Lesny has been performing stress tests of our Storage Element and LFC by using the SRM interfaces to write/read files and register the files into our LFC

Hiro has also performed tests of our DDM configuration by doing simple FTS transfers between BNL and IllinoisHEP

# Current IllinoisHEP T3gs Status/Plans

As far as we know, we are fully configured as a T3gs site with the required services (Panda-aware, DQ2 site services, LFC, etc)

There are a (small?) number of configuration questions still to be ironed out like our DQ2 site name (IllinoisHEP or UIUC) and whether or not the Grid can handle more combinations of UI, IU, UC, …

At this point, I think we are ready for more high-level tests (e.g. DQ2 subscriptions and full chain w/ Panda jobs)

There are a number of policy questions that need to be answered about how the T3gs and the other sites that the T3TF recommend should be utilized within US ATLAS. Hopefully some will get answered at this workshop

It would be fantastic for US ATLAS if we had Tier-3 sites such that their resources could be dynamically allocated to accommodate local analysis needs and pitch in to help during times of intensive need for processing power within ATLAS.

## Happy sqrt() day to all!