# U.S. ATLAS Tier-3 Computing Implementation Committee Report

September 22, 2014

*Doug Benjamin, Michael Ernst, Rob Gardner, Mark Neubauer (co-Chair), Jason Nielsen (co-Chair), Anyes Taffard, Erich Varnes, Torre Wenaus.*

*Ex-officio: U.S. ATLAS Operations Program: Jim Cochran, Srini Rajagopalan; U.S. ATLAS Institutional Board: Chip Brock, John Rutherfoord*

In 2013, the U.S. ATLAS Institutional Board convened a Tier-3 Task Force to survey the current use of Tier-3 computing and to identify future needs for LHC Run 2 physics. A set of six high-level recommendations from the Task Force, reproduced in Appendix A, affirmed the

need for local Tier-3 computing centers that can access expanded remote computing and storage resources on demand.[1]  The current Tier-3 systems are prized for the reliability of job submission, the high speed and low turnaround time for analysis jobs, and very few (if any) failed or lost jobs.  Most analyzers contacted in a related survey stated flatly that the local Tier 3 systems are essential for their analysis, going so far as to tie their success in achieving ATLAS physics goals to the resources provided by their Tier 3. Even though the Tier 3 systems were perfectly sized for Run 1 physics analysis, the larger Run 2 dataset has motivated a fresh look at the evolution of the Tier 3 systems as a whole.

The U.S. ATLAS Operations Program has charged the current committee with conducting a comprehensive study of requirements and technologies in order to propose a cost-effective plan to address the U.S. ATLAS IB Tier-3 Task Force recommendations.  This report summarizes the computing requirements for Run 2 physics analysis and introduces several possible plans to address the needs of U.S. ATLAS physicists.  These implementation plans leverage recent computing solutions and ongoing development within the U.S. ATLAS Operations Program, but full implementation will likely require additional resources from the operations and research programs.

# 1.  Current and Future Physics Analysis

The ATLAS experiment collected approximately 25 fb$^{-1}$ of 8 TeV data, and the current Tier-3 computing systems were designed to process this dataset. (The design and provisioning are described in the 2009 U.S. ATLAS Tier-3 Task Force report[2]).  Run 2 of the LHC is expected to deliver roughly 100 fb$^{-1}$ of 13 TeV collisions to ATLAS, and this increased dataset size requires a re-evaluation of the analysis model and use cases.  In this section, we describe the current analysis model and resource requirements, and we use the new ATLAS analysis model to extrapolate resource requirements into the future for Run 2.

## Current ATLAS Analysis
Under the current ATLAS analysis model, most analyzers process data in a two- or three-step process using a combination of remote and local computing resources.  In "step 1," D3PDs from central production are skimmed and slimmed on the grid (Tiers 1,2)  to produce reduced analysis ntuples.  These ntuples are downloaded in their entirety to Tier-3 systems, where they can be used to explore event selections and optimize the analysis.  In "step 2," the data are reduced further to include only the variables needed for the final analysis and in many cases, this can be accomplished using local Tier 3 resources.  Physics object calibrations and systematic variations are performed at this step in the workflow.  A rapid turnaround allows the analyzers to optimize selections and calculate permutations of systematic uncertainties.  In "step 3," the final plots and limit calculations are produced from the reduced ntuples.  The

---

[1] The report of the 2013 U.S. ATLAS IB Tier 3 Task Force is available at
http://www.pa.msu.edu/~brock/file_sharing/IB/USTier3Report2013.pdf.
[2] The report of the 2009 U.S. ATLAS Tier 3 Task Force is available at
http://www.pa.msu.edu/~brock/file_sharing/T3TaskForce/final/TierThree_v1_executiveFinal.pdf.

following three example analysis workflows demonstrate the range of resources needed to analyze the 25 fb$^{-1}$ dataset.

- The first example analyses, in the dilepton or multilepton final states, benefit from a relatively small dataset size.  Roughly 20 TB of group D3PD data are reduced in "step 1" to produce 1 TB of skimmed data ntuples.  This step requires 2000 jobs occupying 500 total CPU hours on a Tier 1 or Tier 2 facility.  Once the reduced dataset has been downloaded to a Tier-3 system, the "step 2" processing of dedicated ntuples requires a few hundred jobs occupying 100 total CPU hours.  Since the current Tier-3 systems have 50-200 cores, "step 2" can be accomplished in a few hours, allowing analyzers to develop the analysis rapidly.  The output of this step 2 is a much smaller final dataset of order 20 GB.  Many versions of these output datasets are stored, so that analyzers can track effects of systematic variations and perform statistical analysis.  "Step 3" involves the final plotting steps and statistical analysis, which range in wall clock time from a few minutes (simple plots) to several days (toy Monte Carlo experiments).  "Step 4"

- The second example is an analysis of single-lepton + jets signatures, performed on a larger dataset collected with single-lepton triggers.  In this analysis, "step 1" starts with 80 TB of group-produced ntuples and skims the dataset to 20 TB.  This step is performed centrally as another group production.  With all systematic variations included, the step 1 CPU requirement is 600 CPUs for one week, for a total of 100K CPU-hours on the grid.  In "step 2," the final analysis selections are made, leading to an output of 5-10 GB of small ntuples.  For most cases, this step takes about 1-3 hours on a local farm of 50 processors.  "Step 3" includes plotting, calculating statistical interpretations, and training multivariate selection algorithms.  These take no more than 30 minutes on the local resource, and the final step results in output of hundreds of MB.

- The third example is an analysis of the jets + missing energy final state, performed on a very large dataset collected with jet triggers.  In this analysis workflow, "step 1" reduces 150 TB of group D3PDs to 5 TB on the Grid.  This processing takes 1-2 weeks.  "Step 2," which is run locally, results in a 10-20 GB dataset.  This processing can take up to a week, if all systematic variations are computed.  The final "step 3" calculations take about 12 hours in the nominal fit configuration.

**To summarize, we have found that each analysis workflow in the current Run 1 ATLAS analysis model requires roughly 5 TB of user-defined storage to store the reduced datasets from "step 1" and roughly 100 cores to achieve a "step 2" turnaround of no more than 2 hours wall-clock time.**  At the moment, these resources are provisioned locally in the Tier-3 computing systems, where the CPU resources, but not storage, can be shared between analyses.  It is interesting to compare these numbers with the 2009 U.S. ATLAS

Tier-3 Task Force report, which predicted the need for a typical "T3g" system having 80 CPU cores and 20 TB of storage, which would be used to perform several analyses.

In the discussion above, we have deliberately ignored the event generation, private Monte Carlo production, and statistical interpretation use cases. Even though the Tier-3 systems have proven to be very valuable for these activities, the use of private Monte Carlo generation varies widely from analysis to analysis and is difficult to extrapolate. We assume that sufficient resources are available at Tier 1, Tier 2, and other shared resources like HPC's and clouds to carry out Monte Carlo generation in the official production framework.

## Future ATLAS Analysis

The new ATLAS analysis model allows analyzers to produce results directly from a common analysis format (xAOD) that can be read outside of the Athena framework. It also provides a centralized skimming production service to be used by physics groups. As shown in Figure 1, the skimmed/slimmed common analysis dataset can be further reduced in steps to a final ntuple. These reductions correspond to reduction steps 1 and 2 as defined above.
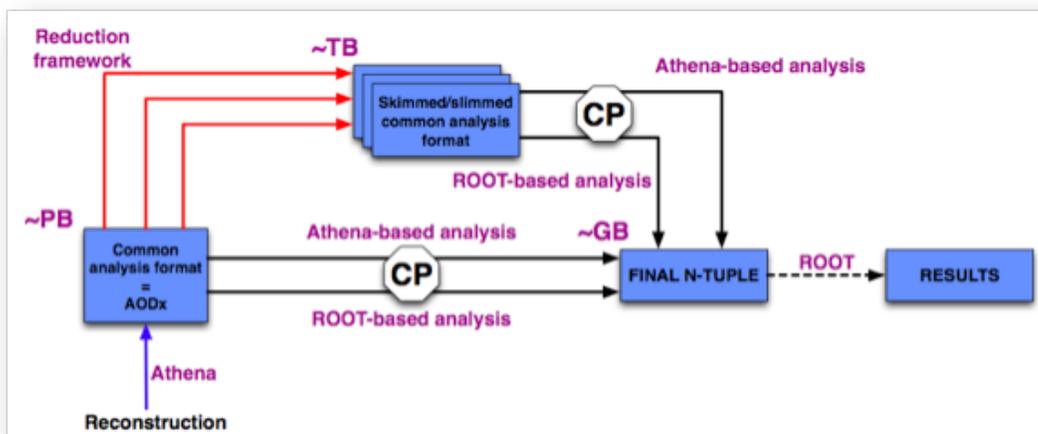


*Figure 1: Schematic overview of the new ATLAS analysis model for Run 2 physics analysis.*

Because all group production will use the common xAOD analysis format, the number of distinct skimmed datasets is expected to be smaller in Run 2 than in Run 1. (The large number of D3PD formats in the physics groups will be deprecated.) We expect that the xAOD event size will be similar to the current group D3PD event size, but it may be possible to share skimmed datasets between analyses.

The main reduction step may include additional calibrations and corrections to physics objects, so the CPU requirements per event for step 1 are not reduced substantially with respect to Run 1. The skimmed/slimmed output dataset might be transferred to local computing resources or staged to a central LOCALGROUPDISK. The Athena- and

ROOT-based analysis activity that follows the creation of the skimmed dataset again requires a fast turnaround for event selection development and systematic studies.

The U.S. ATLAS Tier-3 Task Force estimated that "the requirements to store and process 100 fb-1 of Run 2 data will entail about 3-4 times the resources used to analyze the Run 1 data sample.  This estimate is based on the expected increase in peak luminosity, the increased rate of trigger system output, and other such factors which are known to scale over Run 1."  Our committee concurs with this estimate.  We have no doubt that there may be some savings due to "general resourcefulness on the part of ATLAS physicists," but we use the naive scaling factor of 4 as a baseline estimate.  These estimates are consistent with the responses from a U.S. ATLAS Analysis Support survey conducted in late 2012.  At that time, the median expectations were that Run 2 analysis would require 2-3 times the resources used in Run 1.

## 2.    Computing Resource Estimates for Run 2 Physics Analysis

The following key constraints drive our physics analysis computing resource requirements for Run 2:
- There must be sufficient CPU resources to perform the initial skim ("step 1") from group-produced xAOD datasets to analysis ntuples.
- There must be sufficient CPU resources to ensure a fast turnaround time for final ntuple analysis ("step 2").
- There must be sufficient user-defined storage resources to store results of steps 1 and 2 temporarily and to archive results of step 3.
- There must be sufficient network bandwidth to allow efficient use of the CPU and storage resources. While the network requirements are not quantified in the Tier 3 implementation document, they are necessary and required for successful analysis of Run 2 data.

We foresee that a large fraction of the group production datasets will be stored on LOCALGROUPDISK storage at Tier 1 and Tier 2 sites in the U.S. facilities.  There is sufficient storage there to hold the output of "step 1" jobs, even if more than 1 version is stored.

The 2013 US ATLAS T3TF2 estimated that the requirements to analyze 100 fb -1 of Run 2 data will entail approximately 3-4 times the resources used to analyze the Run 1 data sample. This estimate was based on the expected increase in peak luminosity, the increased rate of trigger system output, and other such factors which are known to scale over Run 1. ***Our Committee concurs with this estimate.*** These estimates are consistent with results from a US Analysis Support survey conducted in 2012 where the median expectation from users was that Run 2 analysis would require 2-3 times the Run 1 resources.

In *isolation*, we expect a **typical Run 2 analysis** to need approximately **20 TB of user-defined storage and 400 CPUs** to satisfy these constraints.

# 3.   Tools and Technologies

This section describes tools and technologies that are relevant to the U.S. ATLAS Tier-3 implementation plans described in the next section.  Many of these tools follow the commercial trend of leveraging high-speed networking to enable transparent access to remote resources.  The first four tools were initiated by the U.S. ATLAS operations program, and the last one has been actively utilized within the program.

### 3.1.   Federated Data Access (FAX)

FAX, already in use by the central production system PanDA, **provides direct read access to storage in the US ATLAS computing facility**.  It can be used from Tier-3 computers.  A global namespace is used to search for files requested by a user job.

- Status as R&D project: in production.
- Status as a tool for Tier 3s: in production.

### 3.2.   ATLAS Connect

The deep and responsive batch queues of the ATLAS Tier-3 clusters were cited as a key feature by the Tier-3 Task Force.  ATLAS Connect (http://connect.usatlas.org/) is a set of computing services designed to augment local systems by focusing on simple batch-like analysis processing.  ATLAS Connect provides a login host to submit jobs directly to the U.S. ATLAS Computing Facility as if it were a local batch queue, and to other resources such as a shared cluster on a university campus.  It also **provides a utility to virtually extend the batch capacity of an existing Tier-3 cluster by sending jobs to Tier-2 sites if additional job slots are needed**.  ATLAS Connect also provides a storage service called FAXbox for staging user job input and output datasets.

- Status as R&D project: in early production.
- Status as a tool for Tier 3s: several Tier-3's already configured to flock into ATLAS Connect: UC, IU, UIUC, UTexas

### 3.3.   PanDA on Tier3

PanDA (Production and Distributed Analysis System) is the distributed workload management system developed in the U.S. and used by all collaborators in ATLAS. While PanDA is the workhorse for all central production, group production, and user analysis performed on distributed Tier 1 and Tier 2 resources worldwide, it has also been deployed at many Tier 3 sites (most often in Europe) as PanDA on Tier 3. For local batch processing at U.S. Tier 3 sites, the overhead of PanDA is usually unnecessary and has not been worth the extra effort to setup and maintain.  A project to provide a simple PanDA-in-a-box is underway, partially funded by DOE-ASCR and non-U.S. funding sources. This project will provide a simple-to-install PanDA client for local Tier 3 resources so that they can connect to the U.S. ATLAS PanDA server. **Deploying PanDA at a Tier 3 will provide seamless integration**

**with resources accessible through PanDA at Tier 1, Tier 2, and other distributed resources.**

- Status as R&D project: in production
- Status as a tool for Tier 3s: PanDA on Tier 3 is in production at 47 Tier 3 sites worldwide (6 in U.S.). A simplified PanDA-in-a-box is under development outside of U.S. ATLAS operations.

### 3.4. Rucio Cache

The new ATLAS distributed data management system being deployed is called Rucio. Several U.S. ATLAS Tier-3 sites have data servers connected to the current DQ2 system so they can receive data automatically via the gridftp protocol.  Currently, data are sent either through an explicit request for data transfer using Web tools (DaTRI) or as output written directly by a user job running on the WLCG computing resources to their local storage resources.  For those sites who have these resources, they have proved to be very helpful.  A replacement is needed under the new Rucio data management system.

The Rucio data management system is built on the concept of a Rucio Storage Element (RSE), the smallest unit of storage space addressable within Rucio.  **A Rucio storage cache may be used as part of a local Tier-3 storage system to receive output datasets from remote computing jobs and write them to a locally-managed RSE.**  The technical details of the system are being developed in conjunction with the Rucio development team in Europe.

- Status as R&D project: under development with U.S. ATLAS personnel and CERN developers
- Status as a tool for Tier 3s: see above.  Tier 3 sites are the focus of the development.

### 3.5. Agile Infrastructure

Agile Infrastructure refers to an extensible computing infrastructure enabled by virtualization and cloud computing.  In the Infrastructure-as-a-Service model (IaaS), the hardware resources (servers, storage, and networking) are provided as a commodity service, while the applications and data are actively managed by administrators.  **This Agile Infrastructure approach allows for increased flexibility in provisioning computing resources in public and private clouds**, and it offers a path toward providing on-demand computing resources for physics analysis.

Pilot programs in U.S. ATLAS are currently implementing and testing interfaces for the submission of physics analysis jobs to cloud computing instances.  For the future, one could imagine the responsive allocation of resources as compute nodes, interactive login nodes, or virtual Tier-3 nodes.

- Status as R&D project: in production

● Status as a tool for Tier 3s: under development

# 4. Candidate Plans for Tier-3 Evolution

In this section, we describe several specific plans under consideration for evolution of the Tier-3 Run 1 system to Run 2. There are two parts to this section: **distinct plans,** which are different in essence and therefore amenable to comparison with one another, and **other recommended actions** that we strongly suggest U.S. ATLAS undertake because they enhance the analysis capabilities of U.S. ATLAS physicists.

## 4.1. Distinct Plans

In this section, we present several possible plans for Tier-3 evolution. These plans differ in their essential elements, particularly in how they would need to be funded and in what ways they make use of available resources, such that they are amenable to direct comparison for cost effectiveness. The plans are driven primarily by the following two observations:

1. It is cheaper to have the majority of the batch horsepower in support of U.S. ATLAS physics located in centrally-managed facilities (primarily the Tier-2 complex), where there are economies of scale in terms of hardware procurement, configuration, and support.

2. However, the majority of U.S. ATLAS groups strongly desire some computing resources which are under their complete control. By "control" we mean that the resources are "immediately" available for interactive or batch use by members of the local group *when* they want to use them, and the group determines *how and by whom* they are used.

These plans were chosen to have an element of realism from both funding and implementation perspectives and be sufficiently broad to span the range of possibilities. More granularity in the plans and additional plans-of-action are possible, of course, but we aim to keep things as simple as possible in this report, but no simpler.

Further discussion and comparison of the plans is included in Section 4.1.4.

### 4.1.1. Plan 1

This plan would execute an *approximate refresh* of the *current* Tier-3 system and utilize *existing* tools for Tier-3 analysis. At most sites, the Tier-3 computing hardware is either out of warranty or very nearly in that state. New funds would primarily be spent to replace the aging Tier-3 hardware at the existing sites. In a climate of constrained budgets, this plan would require new funds at approximately the same level as the past funding used to procure the present Tier-3 system -- primarily the DOE ARRA (~$650k) + NSF MRI ($620k) awards.  (A similar or greater amount was contributed by universities through on-campus funds or in-kind

contributions.) This plan would involve no new U.S. ATLAS operations program investment in tools and technologies to utilize Tier-3 or Tier-2 facility resources for U.S. physics analysis.

We note that hardware bought at the same dollar amount in 2015 and beyond comes with substantially more CPU and storage capability as compared to hardware purchased for the existing Tier-3 systems. We therefore anticipate that a sizeable fraction of the factor of 3-4 in extra resources needed for Tier-3 computing in Run 2 as compared with Run 1 would be covered in this way.

### 4.1.2. Plan 2

This plan would repurpose existing Tier-3 hardware at each site to provide a more optimal mix of interactive and batch computing within the Tier-3 complex. Sites with Tier-3 resources would use them for interactive and small-scale batch jobs characteristic of the last few steps of analysis workflows. Heavier analysis workloads that cannot be run on the local resource in a reasonable amount of time would be submitted to the Tier-2 facilities *using existing tools* and use the beyond-pledge CPU and storage resources there for U.S. physicists.

Note that this plan does not require any appreciable level of new funding for Tier-3 hardware or U.S. ATLAS operations program investment in new tools and technologies, but it does require that the U.S. ATLAS computing facilities continue to receive the appropriate level of funding to provide beyond-pledge resources for use by U.S. physicists.

### 4.1.3. Plan 3

In this plan, each institution deploys a *modest* Tier-3 resource locally using the latest hardware. The local resource is to be used primarily for interactive and small-scale batch jobs characteristic of the last few steps of analysis workflows (as in Plan 2). We propose to deploy technologies described in Section 3 that transparently expand the Tier-3 resources into the beyond-pledge resources on the facilities when the workload is such that it would not complete in a reasonable amount of time on the local Tier-3 resource. In this way, the local resources are used first for workflows in support of analysis but are augmented as needed by resources available on the U.S. facilities for physics. Finally, we suggest that Tier-3 sites consider repurposing their existing Tier-3 hardware to provide a more optimal mix of interactive and batch computing (as in Plan 2).

This plan requires new funds to deploy a modest-but-modern resource locally at each institution. Given the use cases discussed in Section 1 and the Run 2 resource estimates presented in Section 2, the size of this new local resource is a few hundred CPU cores and several tens of TB of storage.

We refer to the building block of this resource as a ***local resource unit* (LRU)**. From the Run 2 estimates presented in Section 2, one might naively conclude that the LRU is 4 times the resource needs for an isolated analysis. However, this would neglect that (1) people at an institution collaborate on analyses leading to sharing of resources and (2) not all Tier-3

computing is local in this plan. With these two considerations in mind, we define the disk LRU as **10 TB** of user-managed storage and the CPU LRU as **100 CPU** cores. The LRU is the minimum new Tier-3 resource for a single active analyzer at an institution. Some example hardware configurations that reasonably map onto the LRU concept at the time of writing are presented in Appendix A.

The exact new Tier-3 resource needs are determined by the group size and scope of analysis work (similar to the 2009 T3TF report). We recommend that these needs be represented by a number of LRUs, with the number of LRUs at a given institution being scaled by the number of active analyzers $A$ at the institution in the following way:

$$CPU = \ln(A) * 100 \text{ cores}$$

$$disk = A/2 * 10 \text{ TB}$$

An active analyzer is a faculty member, postdoc, or graduate student involved in data analysis on ATLAS. Teaching faculty are weighted by a factor ½ relative to the other categories, due to the level of responsibilities outside of ATLAS analysis. The number of active analyzers is self-reported by each institution.

This plan also requires U.S. ATLAS operations program investment to develop, deploy, and support technologies that facilitate transparent use of beyond-pledge resources on the facilities to augment Tier-3 resources.

### 4.1.4. Comparison of the Candidate Plans

In this section, we include further discussion of the plans just described and compare them relative to each other for cost-effectiveness.

Plan 1 is certainly the simplest of the plans since it uses the current Tier-3 system as a template and deploys new hardware at a level which is both at the approximate scale needed for Run 2 and probably not impossible to fund even in a climate of constrained budgets. It addresses the issue of aging hardware in the current system directly and thoroughly. It does not address issues of non-optimal distribution of resources across U.S. groups and the mix of batch and interactive computing that exists at some Tier-3s, as is done explicitly in Plans 2 and 3. In Plan 1 and Plan 2, no new investment in technical improvements to use the beyond-pledge resources for U.S. physicists is made, unlike in Plan 3. Use of beyond-pledge resources on the facilities for analysis benefits from economy of scale in both hardware and technical support.

Plan 2 addresses issues of non-optimal distribution of resources across U.S. groups and the mix of batch and interactive computing that exists at some Tier-3s. It also utilizes the beyond-pledge resources but requires new investment in technical improvements in order to use those resources as efficiently for analysis as in Plan 3. However, without additional investment in Tier-3 hardware at institutions, Plan 2 has the problem of aging hardware that becomes increasingly unreliable (requiring more admin support) and obsolete over time, does

not increase campus presence in ATLAS computing (which could potentially allow for opportunistic use of more CPU for U.S. ATLAS overall), and does not incentivize future in-kind contributions from institutions (e.g., hardware, people, expertise), which have been very substantial in Run 1. Given the resource requirements outlined above, there is a strong consensus in the committee that U.S. physicists will fall behind in analysis productivity if Plan 2 is chosen, thereby limiting the U.S. contributions in achieving the physics goals of Run 2.

Plan 3 optimizes the use of existing Tier-3 hardware, deploys new hardware at each institution at the appropriate level given the resource needs for Run 2 and the scale of analysis activities at each institution, and invests in new tools and technologies to make better use of the beyond-pledge-resources for analysis. The tools and technologies follow the commercial trend of levering and high-speed networking to enable transparent access to remote resources for analysis and are highly-leveraged from existing software, facilities, and ongoing efforts with the US ATLAS Physics Support, Software and Computing program.

|  | **Plan 1** (full hardware refresh) | **Plan 2** (no new resources) | **Plan 3** (hybrid model) |
|---|---|---|---|
| **Description** | Refresh of current systems. Utilize existing tools. | Repurpose existing T3 hardware. No new investment in T3 hardware or tools to expand beyond local T3 resource | Provision LRUs at each institution. Deploy Tier-3 flocking extensions and FAX-accessible user scratch storage |
| **Pros** | Straightforward. Less risk. Requires development of no new technologies. Does not require additional support personnel. | Simple. Requires development of no new technologies. | Modern T3 hardware provisioned at the appropriate level. Retain flexibility and control of local T3. Transparently expand beyond local T3s into facilities, when workflow demands. |
| **Cons** | Require largest amount of new computer hardware. Increased local infrastructure required (space, power, cooling). Continued local support for a clustered | Aging hardware becomes increasingly unreliable (requiring more admin support). Loss of local T3s → crucial loss of flexibility and control | Requires significant additional support beyond Plans 1 and 2 to scale existing system and support. Increased local infrastructure and continued local support for a |

| | environment which would be more than a simple set of LRU hosts. | | clustered environment. |
|---|---|---|---|
| **Assessment** | | | Recommended |
| **Rank** | 2 | 3 | 1 |

After evaluating the pros and cons of the various plans, we ranked Plan 3 (local/remote hybrid model) highest, just ahead of Plan 1 (equipment refresh).  It was felt that Plan 2, with no new Tier 3 resources, would jeopardize the ability of U.S. ATLAS physicists to complete data analyses.  The core functional components of Plan 3 are the following:

- Modest local computing and storage, scaled to the LRU needs for each institute
- Mechanisms to send the overflow of locally-submitted jobs to remote resources via ATLAS Connect or PanDA
- Remotely-accessible central storage of large datasets using FAXbox or ATLAS Distributed Data Management (DDM)

## 4.2.  Other Recommended Actions

In this section, we describe action that we strongly recommend U.S. ATLAS undertake because they enhance the analysis capabilities of U.S. ATLAS physicists.  These actions are complementary to the distinct plans described above, and they should be pursued regardless of which implementation plan is followed.

### 4.2.1.  Develop and Support Large Shared Tier-3 Systems

Several U.S. ATLAS groups have jointly provisioned large shared Tier-3 systems, pooling resources to locate equipment at a single host institution.  Such large systems will also benefit from the elastic resource allocation described above.  They will have access to remote storage elements in the U.S. facilities, and they will be able to send jobs transparently to remote computing resources.  We expect that very large systems will be administered locally as part of the host institution's commitment.

There are two models for developing and supporting the large shared Tier-3 systems.  The first model is that of the existing "Tier-3 consortia," which enjoy economies of scale due to pooled resources and which foster close collaboration on physics analysis.  The success of this model in at least one site was highlighted in the 2013 Tier-3 Task Force report.  The second model is one or more national or regional shared facilities selected in a competitive process.  The U.S. ATLAS Operations Program would be responsible for developing a solicitation for proposals that rewards institutional buy-in or matching as a cost-savings measure.  Just as with the Tier-3 consortia, the host institution would be expected to

administer the shared facility for U.S. ATLAS, with some enhanced resource priority for the host institution's users.

The challenge for such facilities is to provide some measure of flexible resource control on shared resources.  Local resource control to ensure low batch queue latency is one key feature of the current Tier-3 systems, and some work will be needed to develop similar control mechanisms in shared facilities.

# 5.    Management and Budget Considerations

In this section, we discuss how the Tier-3 computing resources and personnel would be managed under the preferred implementation Plan 3. We also identify what can be accomplished within the current U.S. ATLAS Operations Program budget guidance and prioritize other possible requests in the event that supplemental funding from the Operations Program or other sources materializes.

The U.S. ATLAS Operations Program already provides for provisioning and support of beyond-pledge resources in the Facilities, as well as Tier-3 hardware installation and operating system maintenance assistance.  Continuing these efforts can be accomplished within the current program guidance.  The additional requests required to implement Plan 3 are as follows.

- To support the new ATLAS Connect system in production, 1.0 FTE in Facilities Integration and Operations will be needed for the following:
    1. Assisting Tier3 sites with the software installation and proper system configuration into the distributed computing system:
        a. Provide a "thinly-provisioned" Tier-3 configuration system so as to reduce the needed local technical expertise.
        b. Provide on-going development support for tools which extend the local batch schedulers to the US ATLAS Computing Facility (the local scheduler might be HTCondor, PBS, SLURM, SGE, etc.)
        c. Support for configurations that allow joint use of both local Tier-3 clusters and any local, institutional computing resource such as a campus research computing center.
        d. Assistance / guidance for network troubleshooting to optimize performance for jobs which run remote to the local Tier-3
    2. System build and configuration of ATLAS Connect central services:
        a. Central login host
        b. Flocking host (routing Tier-3 jobs to external resources)
        c. FAXbox backend storage system
        d. Web servers (job monitor, accounting, and a continuous validation service)

e. Support for interactive user accounts, home directory backups, access to ATLAS software and grid tools, and local access to FAXbox storage for the US ATLAS Collaboration
f. Support for groups (by institution or physics working group) for resource access and accounting mechanisms.

- To support Tier-3 analyzers using both local and distributed systems, 0.5 FTE in Physics Support will be needed for the following:
    - Providing templates and wrappers for specific analysis tasks using ATLAS Connect, PanDA-in-a-box, and Agile infrastructure.
    - Troubleshooting and supporting user code for analysis jobs which access data over the wide area network.
    - Supporting local Rucio data caches on Tier-3 sites if this technology becomes available.

## 6. Postscript: A Physicist Uses the ATLAS Tier-3 in 2018

The year is 2018, and Jill is a U.S. ATLAS physicist working on a hot topic at the LHC: the search for evidence of new high-mass particles decaying to top quark pairs, specifically in the lepton+jets final state signature.

Jill's data analysis begins with approximately 300 TB of centrally-produced single-lepton-trigger xAOD files, which have been placed by U.S. ATLAS operations on the LOCALGROUPDISKs at the Tier 2 facilities. Jill needs to apply a new set of lepton calibrations and calculate uncertainties due to systematic variations before making her event selection, so she defines and submits a PanDA job from her Tier 3 system. This "Step 1" stage of the analysis runs entirely in the analysis queues on the Tier 1 and Tier 2 facilities, taking advantage of the LOCALGROUPDISK storage there. Jill can expect access to thousands of CPUs during the week it takes her "Step 1" job to finish. Luckily this step is not repeated very often, perhaps once or twice per month. The output of this step is a much smaller xAOD file, perhaps 10 TB, that is stored in the LOCALGROUPDISK under the ATLAS DDM system.

Jill copies the output dataset of 10 TB back to her local Tier 3 system, leaving the original at the Tier 2 facility on LOCALGROUPDISK. Now she is ready to tune her event selection by applying cuts to the dataset in a "step 2" analysis stage. To ensure that her reduction step will finish in a few hours, Jill defines a task with 400 subjobs running in parallel and submits them to her Tier 3 cluster. Because her Tier 3 cluster has only 150 job slots free, the remaining 250 jobs are sent to available queues in the Tier 2 facilities. The remote jobs access the input dataset on LOCALGROUPDISK via the FAX service. After 3 hours, Jill's jobs have finished, and the output ntuples of 50 GB have been delivered from the remote sites back to her Tier 3 disks.

Now Jill can plot the results of her event selection using lightweight analysis tools on the small local dataset.  Since she is fitting the background contributions from the data distributions, she has implemented a complicated likelihood fit that takes about 1 hour to run on the local CPU. She also calculates the expected discovery reach for her search using similar tools that generate a large number of pseudoexperiments.  These tools take advantage of the significant available local CPU on Jill's Tier 3 cluster.

# APPENDICES
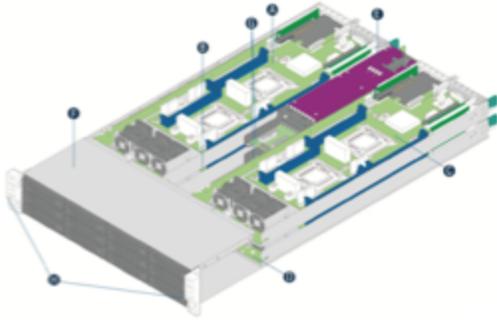
## A. Example LRU Hardware Configurations

In Section 2 on projected needs for Run 2 Tier-3-like workflows, it was argued that, in *isolation*, we expect a *typical* Run 2 analysis to need approximately 20 TB of user-defined storage and 400 CPUs. In the recommended plan (Plan 3), only a fraction of that need is to be satisfied by (new) local Tier-3 resources given the anticipated level of CPU and storage resource sharing and availability of beyond-pledge resources at the Tier-2 facilities for analysis workflows. In Section 4.1.3, we refer to the building block of this resource as a **local resource unit (LRU)**. We define the disk LRU as **10 TB** of user-managed storage and the CPU LRU as **100 CPU** cores. The LRU is the minimum new Tier-3 resource for a single active analyzer at an institution. In reality, there is no single-analyzer institution, so the actual recommended resource at each institution is scale up from the LRU just mentioned.

In this section, we provide some possible configurations that map reasonably well onto the LRU using candidate hardware at the time of this report. As compute and storage technologies (and pricing) evolve, so too will the best option for procuring new hardware locally at institutions.

### A.1 Colfax 2U Rackmount Server

Below is a possible configuration based on the the Colfax CX22850i-X5 2U Rackmount Server (inspired by CMS's "Frankiac").

- Colfax CX22850i-X5 2U Rackmount Server
- Four compute nodes in one 2U chassis
  - Dual Xeon E5-2650v2 processors (8 cores / 16 threads) per node
    - 15 C / 32 T / node → 48C / 96T / chassis
  - 12 x 4 TB drives total gives a total of 48 TB
  - Cost ~$20k total at present

## A.2 Dell C6220 2U Server

We have initiated discussions with Dell representatives regarding the appropriate configuration and quotes.

# B. Recap of U.S. ATLAS IB Tier 3 Task Force 2 Recommendations

The U.S. ATLAS IB Tier 3 Task Force 2 recommendations, released in 2013, highlighted specific needs for U.S. ATLAS data analysis. The implementation plans considered in this report satisfy those recommendations. In this appendix, we reproduce the T3TF2 recommendations and show how they are addressed in the implementation plans.

*T3TF2 Recommendation 1:*
*US ATLAS should continue to endorse and support Local Tier 3 Centers (LT3Cs), in the form of either single-institution LT3Cs or LT3Cs that are shared by consortia of Universities, as essential for physics analysis, detector studies and upgrade simulation for US ATLAS groups, and make a strong case to the funding agencies for continuing the support of such centers. Given the significant local funding that the Universities were able to raise to build and operate the existing LT3Cs, LT3Cs represent a very cost effective way for the funding agencies to support the US ATLAS physics program.*

Implementation plans 1 and 3 endorse the concept of LT3Cs, specifically in the shape of the Local Resource Unit. The U.S. ATLAS operations program continues to support LT3Cs through the Physics Support efforts.

*T3TF2 Recommendation 2:*
*US ATLAS should work to extend the current Tier 3 system by providing mechanisms for Tier 3 jobs to expand onto resources outside of the LT3Cs when the LT3C resources are fully utilized. US ATLAS should also continue to invest in technologies that give Tier 3 functionality to institutions that do not have an LT3C so they may continue to contribute to the ATLAS physics program.*

The Atlas Connect and/or PanDA services can provide access to remote CPU resources so that Tier 3 jobs can expand onto participating facilities. The central ATLAS Connect Login service would provide Tier 3 functionality to any users who request accounts.

***T3TF2 Recommendation 3:***
*US ATLAS should support the use of wide-area data access mechanisms to provide effort is needed to provide a robust solution for Tier 3 users.*

The FAX service provides remote access to datasets from any configured Tier-3 system.

***T3TF2 Recommendation 4:***
*US ATLAS should provide all users a sufficient amount of guaranteed storage space. This storage, assigned to each user, should be located at the Tier 1/Tier 2 centers where the batch computing resources are located.*

All users will have sufficient storage space for data analysis step 1 (on group disks in central facilities) and steps 2-4 (on the Local Resource Unit).

***T3TF2 Recommendation 5:***
*The US ATLAS Analysis Support Team should work with the ATLAS Distributed Data Management (DDM) and ATLAS Distributed Computing (ADC) developer teams to provide a solution that will allow users to direct their output from Grid jobs back to their local storage. The solution should allow the user to specify where the output is sent when the job is submitted and the job output must be returned with minimal delay immediately after the job is completed.*

The Rucio Cache service is being developed within U.S. ATLAS Physics Support, in collaboration with Rucio developers at CERN. Existing data management tools in Rucio can also be used.

***T3TF2 Recommendation 6:***
*To promote the adoption of the new components of the US ATLAS analysis environment such as the wide-area data access mechanism, the US ATLAS Analysis Support Group should provide documentation and organize comprehensive tutorials to train interested users on how to best take advantage of the new resources.*

The Analysis Support group has a good track record of supporting users through training and documentation. Continued and expanded user support is a key part of the implementation plan.