

Results of HEP CG WP 3, GSI: Interactive Data Analysis with PROOF

P.Malzacher@gsi.de

Anna Kreshuk, Peter Malzacher, Anar Manafov, Victor Penso, Carsten
Preuss, Kilian Schwarz, Mykhaylo Zynovyev

16. June 2008

HEP CG Workshop

16.-17. June 2008

Dresden



Interactive Data Analysis with PROOF


PROOF Overview

Integration of the GSI AF into the
general purpose batch farm
(for Grid and local batch)

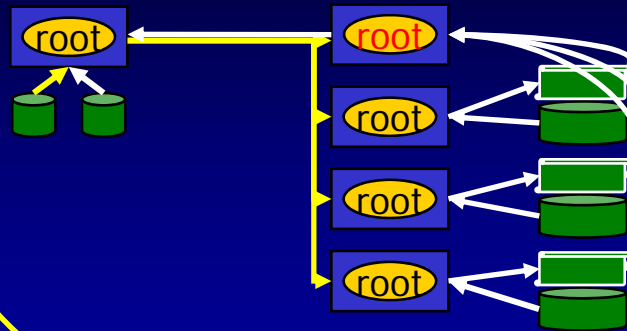
Extending PROOF

PROOF on the Grid

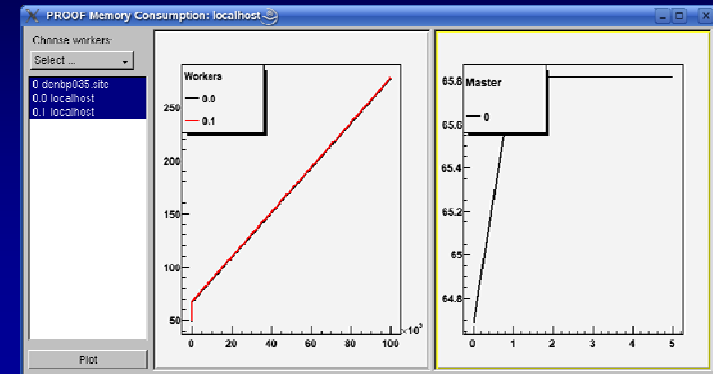
**GSI results of
HEP CG WP 3**



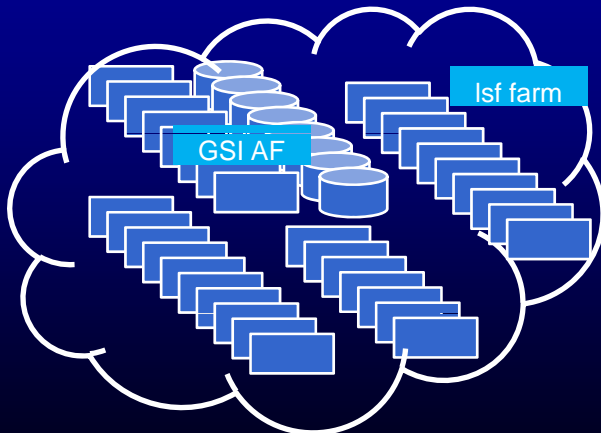
PROOF Overview



Extending PROOF



Integration of PROOF in the farm at GSI



PROOF on the Grid

| ID | Status |
|---|-----------|
| https://dgrid-rb.fzk.de:9000/YD5kzjA... | Running |
| https://dgrid-rb.fzk.de:9000/RyLF... | Scheduled |
| https://dgrid-rb.fzk.de:9000/VtD... | Running |
| https://dgrid-rb.fzk.de:9000/kbH... | Done |
| https://dgrid-rb.fzk.de:9000/q3Lj... | Running |
| https://dgrid-rb.fzk.de:9000/qEm... | Running |

PROOF: Parallel ROOT Facility

Interactive parallel analysis on a local cluster

Parallel processing of (local) data (trivial parallelism)

Fast Feedback

Output handling with direct visualization

Not a batch system, no Grid

The usage of PROOF is transparent

The same code can be run locally and in a PROOF system
(certain rules have to be followed)

~ 1997 : First Prototype

Fons Rademakers

2000...: Further developed by MIT Phobos group

Maarten Ballintijn, ...

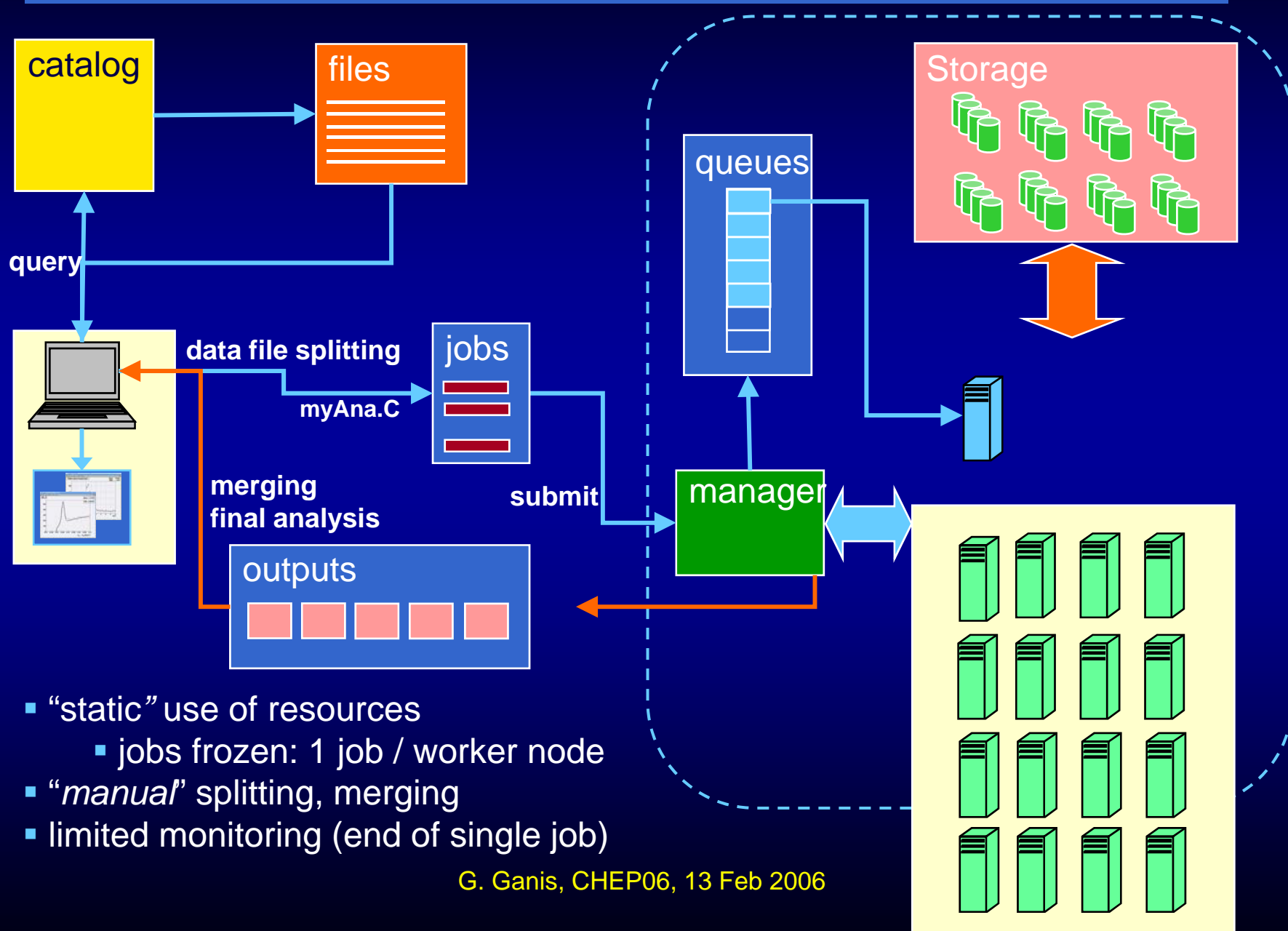
2005...: Alice sees PROOF as strategic tool

2007...: Gerri Ganis, ..., CERN AF

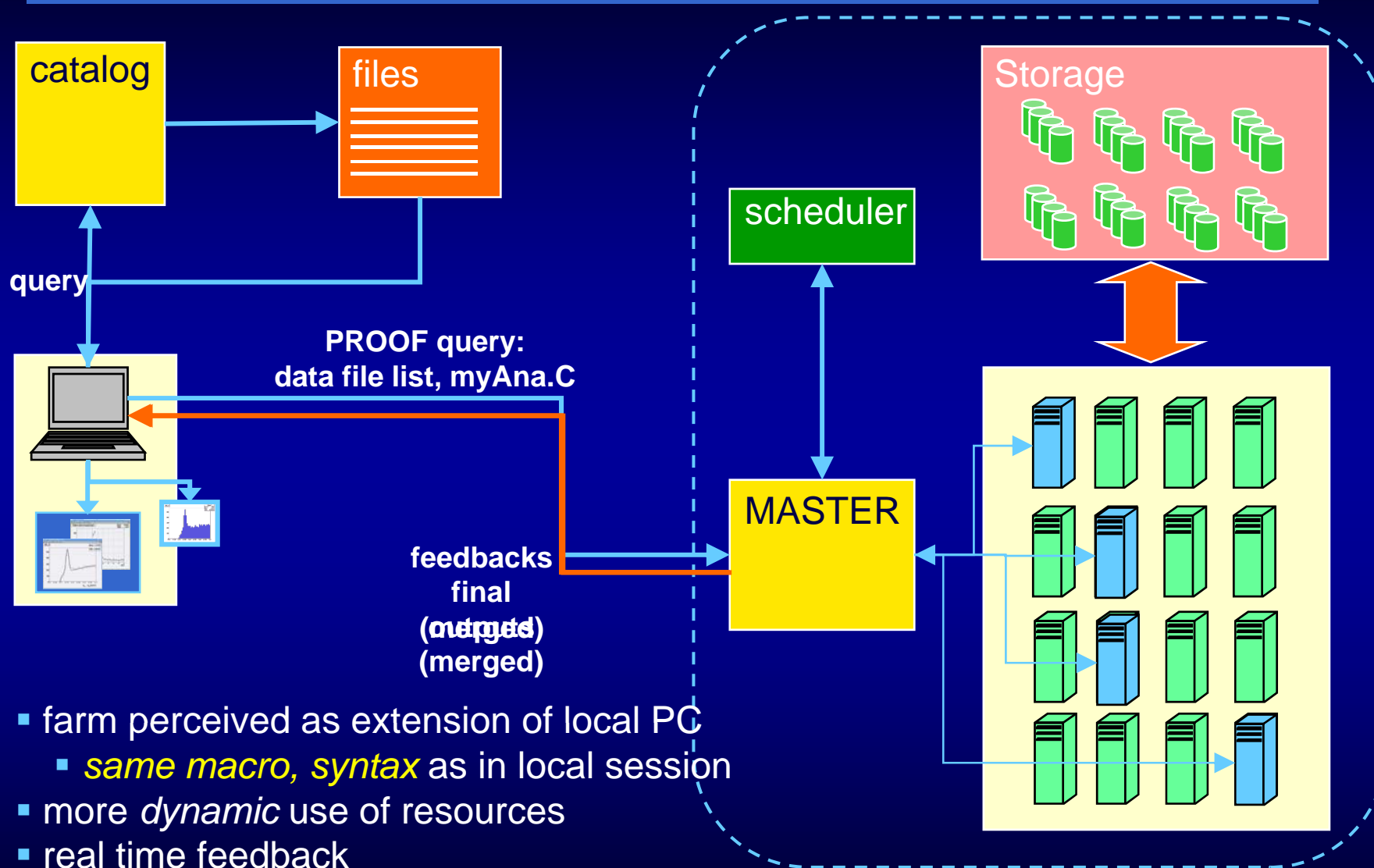
<http://root.cern.ch/root/PROOF2007/>

~ 60 participants, most from Alice, individuals from other exp.

Job Split approach



The PROOF approach



- farm perceived as extension of local PC
 - *same macro, syntax* as in local session
- more *dynamic* use of resources
- real time feedback
- automated *splitting* and merging

Classes TTree /TChain

A **tree** is a container for data storage

It consists of several *branches*

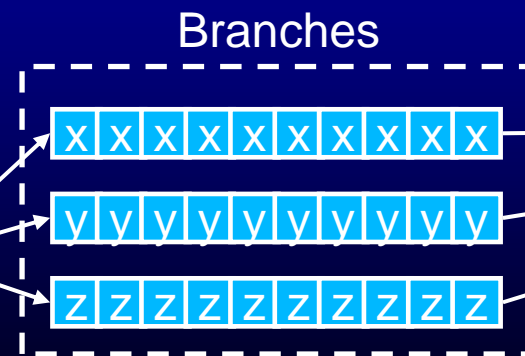
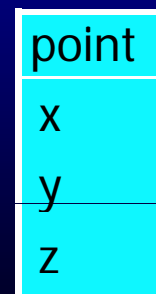
These can be in one or several files

Branches are stored contiguously (split mode)

When reading a tree, certain branches can be switched off → speed up of analysis when not all data is needed

Compressed

A chain is a list of trees (in several files)



Chain

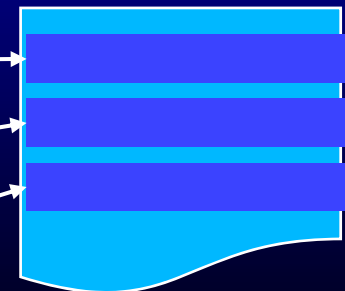
Tree1 (File1)

Tree2 (File2)

Tree3 (File3)

Tree4 (File4)

File



TSelector

Classes derived from TSelector can run locally and in PROOF

– Begin() once on your client

– SlaveBegin() once on each Slave

– Init(TTree* tree) for each tree

– Process(Long64_t entry) for each event

– SlaveTerminate()

– Terminate()

Progress dialog

PROOF Query Progress: jgrosseo@lxb6046.cern.ch

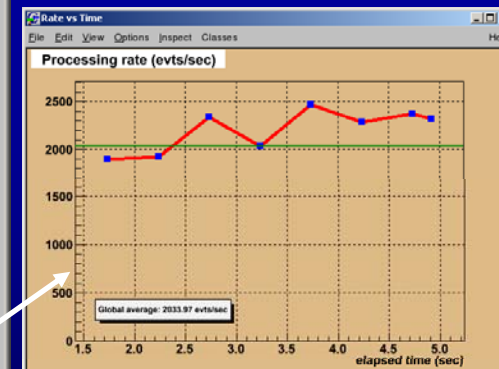
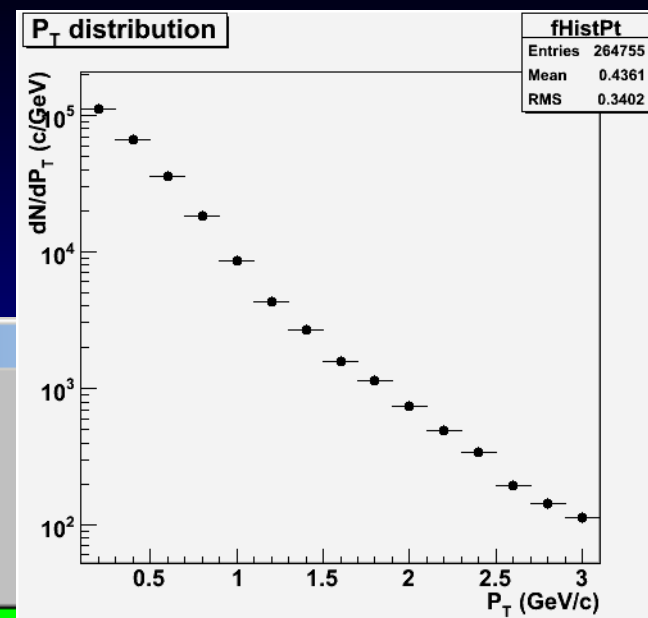
Executing on PROOF cluster "lxb6046.cern.ch" with 33 parallel workers:
Selector: TMySelector.cxx
100 files, number of events 10000, starting event 0

Initialization time: 0.9 secs
Processed: 10000 events (1790.14 MBs) in 4.9 sec
Processing rate: 2034.0 evts/sec (364.1 MBs/sec)

☐ Close dialog when processing is complete
☐ Show only logs from query last

Stop Cancel Close Show Logs Rate plot

Query statistics



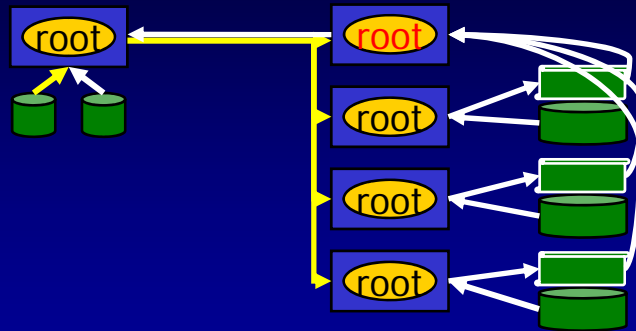
Abort query and
view results
up to now

Abort query and
discard results

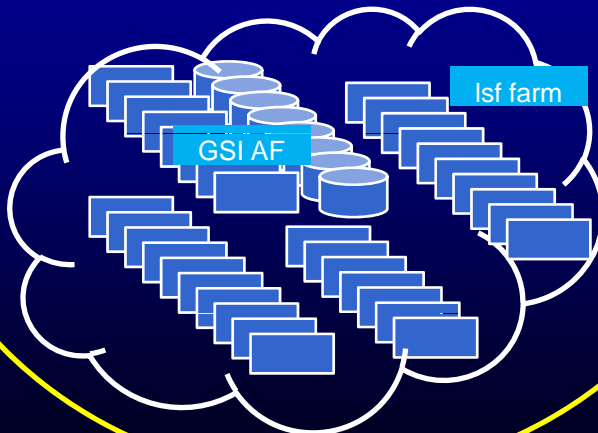
Show log
files

Show processing
rate

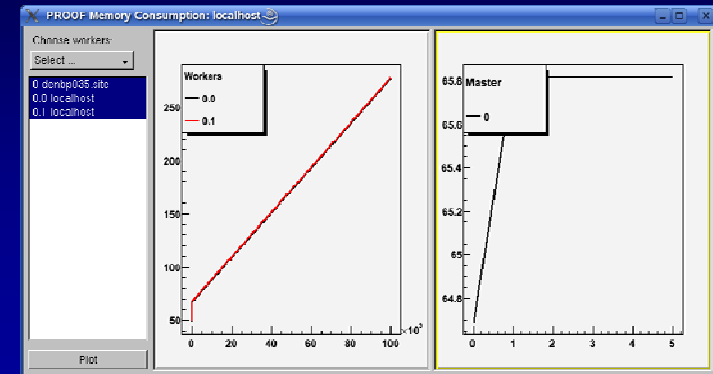
PROOF Overview



Integration of PROOF in the farm at GSI



Extending PROOF



PROOF on the Grid

| ID | Status |
|---|-----------|
| https://dgrid-rb.fzk.de:9000/YD5kzjA... | Running |
| https://dgrid-rb.fzk.de:9000/RyLF... | Scheduled |
| https://dgrid-rb.fzk.de:9000/VtD... | Running |
| https://dgrid-rb.fzk.de:9000/kbH... | Done |
| https://dgrid-rb.fzk.de:9000/q3Lj... | Running |
| https://dgrid-rb.fzk.de:9000/qEm... | Running |

Plans for the Alice Tier 2&3 at GSI: Size

| Year | 2007 | 2008 | 2009 | 2010 | 2011 |
|-------------|------|------|------|------|------|
| ramp-up | 0.4 | 1.0 | 1.3 | 1.7 | 2.2 |
| CPU (kSI2k) | 400 | 1000 | 1300 | 1700 | 2200 |
| Disk (TB) | 120 | 300 | 390 | 510 | 660 |
| WAN (Mb/s) | 100 | 1000 | 1000 | 1000 | ... |

| Germany, GSI, Darmstadt | Pledged | Planned to be pledged | | | | |
|-------------------------|---------|-----------------------|------|------|------|--|
| | 2006 | 2007 | 2008 | 2009 | 2010 | |
| CPU (kSI2K) | 100 | 260 | 660 | 860 | 1100 | |
| Disk (Tbytes) | 30 | 80 | 200 | 260 | 340 | |
| Nominal WAN (Mbits/sec) | 100 | 100 | 1000 | 1000 | 1000 | |

2/3 of that capacity is for the tier 2 (fixed via WLCG MoU)
1/3 for the tier 3

To support ALICE and to learn for FAIR computing.

GSI Setup: ~40% = ALICE Tier2/3 usable via batch, grid and PROOF

~1400 cores

batch system Isf

debian sarge, etch32 & etch64

including

80 2*4core 2.67GHz Xeon with
4*500 GB internal disk

~15 used as PROOF cluster

= **GSI AF**

~ 500 TB in file server

3U 15*500GB SATA, RAID 5

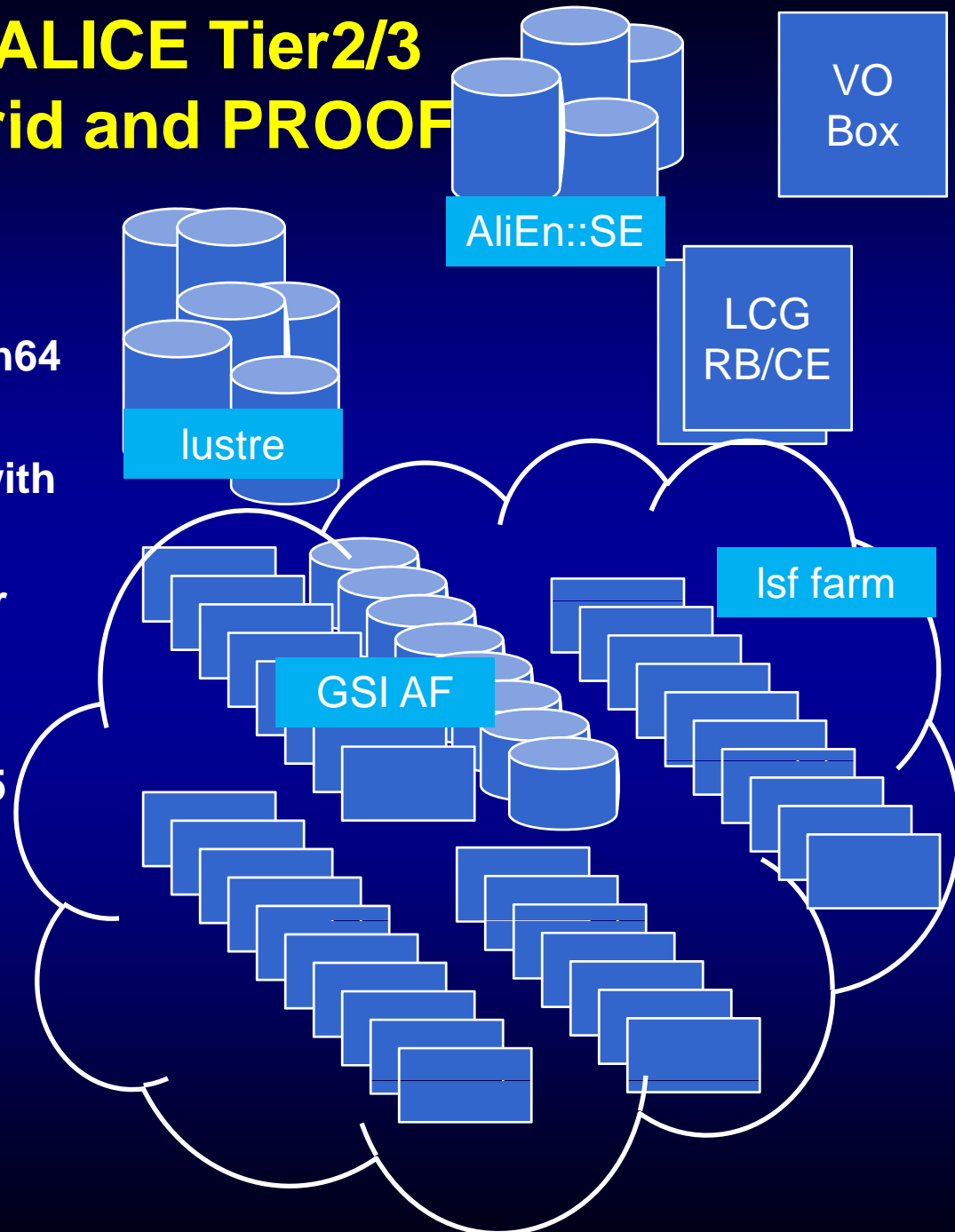
~ 50 AliEn storage element

~ 450 lustre as cluster file
system

data import via AliEn SE

movement to lustre or PROOF

via staging scripts





Result: Experience Report

Set-up of an integrated batch/grid/proof farm

Layout of the queues

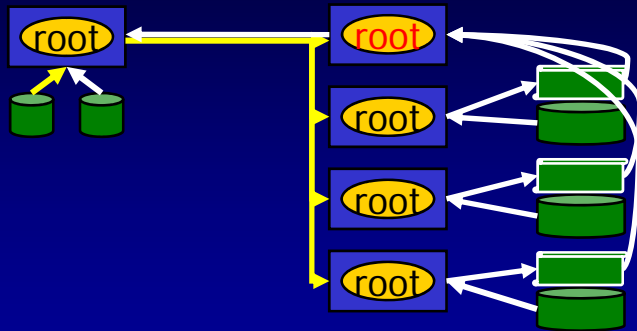
I/O: Mix of lustre and SE via xroot

GSI AF with local I/O

Suspend batch/grid jobs when proof is active

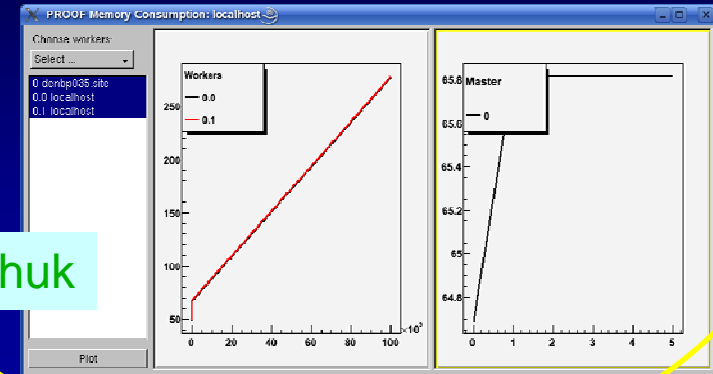
Staging to GSI AF vs usage from SE

PROOF Overview

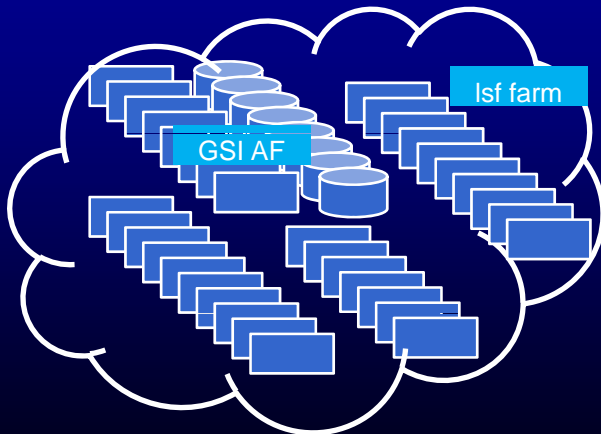


Extending PROOF

A. Kreshuk



Integration of PROOF in the farm at GSI



PROOF on the Grid

Grid

JDL File:

Endpoint:

worker(s)

100% submitted

Information about the last submitted job:

| ID | Status |
|---|-----------|
| https://dgrid-rb.fzk.de:9000/YD5kzjA... | Running |
| https://dgrid-rb.fzk.de:9000/RyLF... | Scheduled |
| https://dgrid-rb.fzk.de:9000/VtD... | Running |
| https://dgrid-rb.fzk.de:9000/kbH... | Done |
| https://dgrid-rb.fzk.de:9000/q3Lj... | Running |
| https://dgrid-rb.fzk.de:9000/qEm... | Running |

TEntryList: store lists of events



Run some analysis code on the data of myChain

Find interesting events

Keep a list of interesting events for further analysis

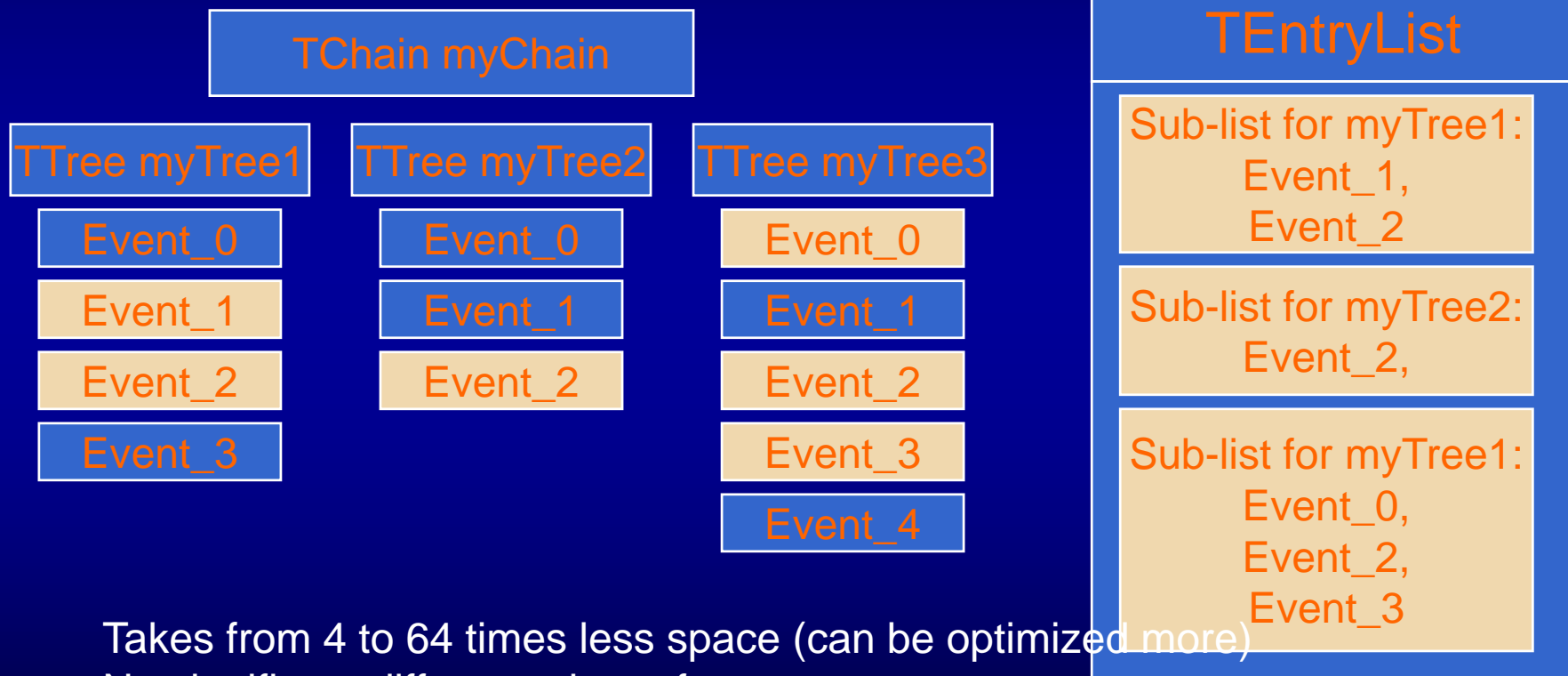
A "naive" implementation: TEventList

Basically, an array of longlong numbers



Because of global indexing in the TEventList, it's not very suitable for parallel processing – it's hard to extract sub-parts that belong to separate trees and process them independently

TEntryList - new



Takes from 4 to 64 times less space (can be optimized more)

No significant difference in performance

Can be only partially loaded in memory

Sub-lists can be extracted and used to compose lists for other chains

Sub-lists can be extracted and processed independently

Improving Debugging: Memory consumption monitoring

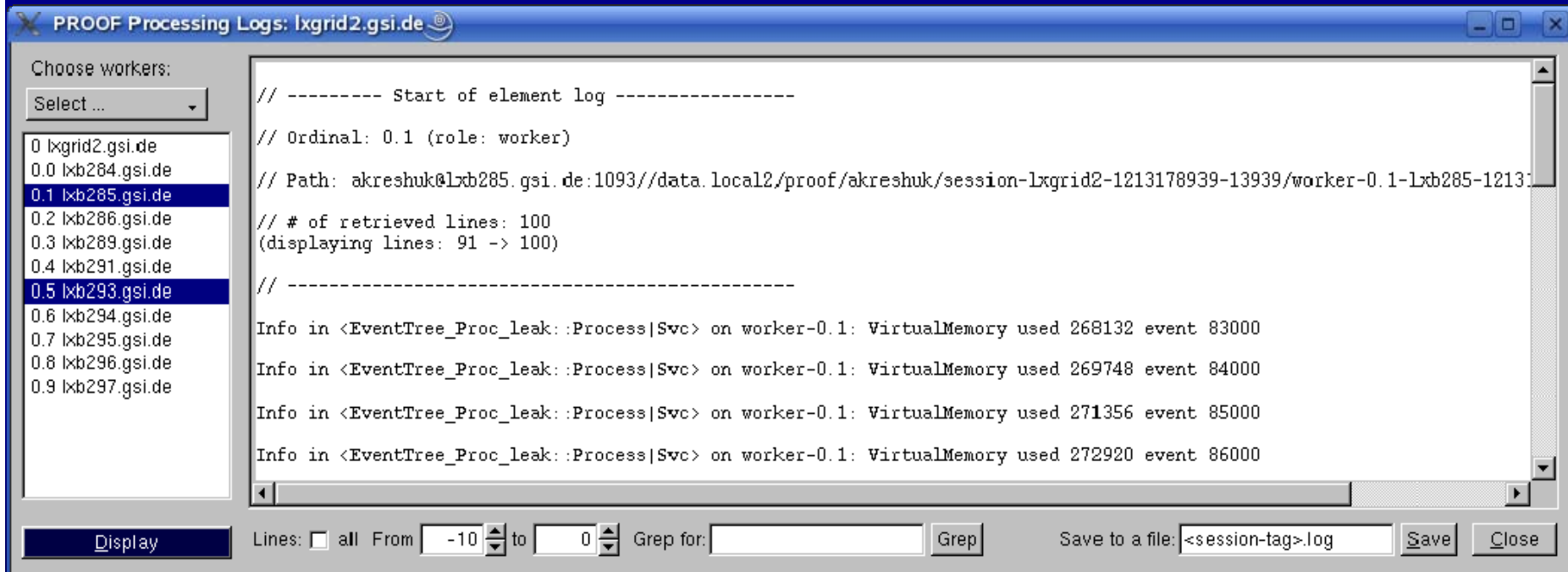
An easy way to access logfiles after a session crashed.

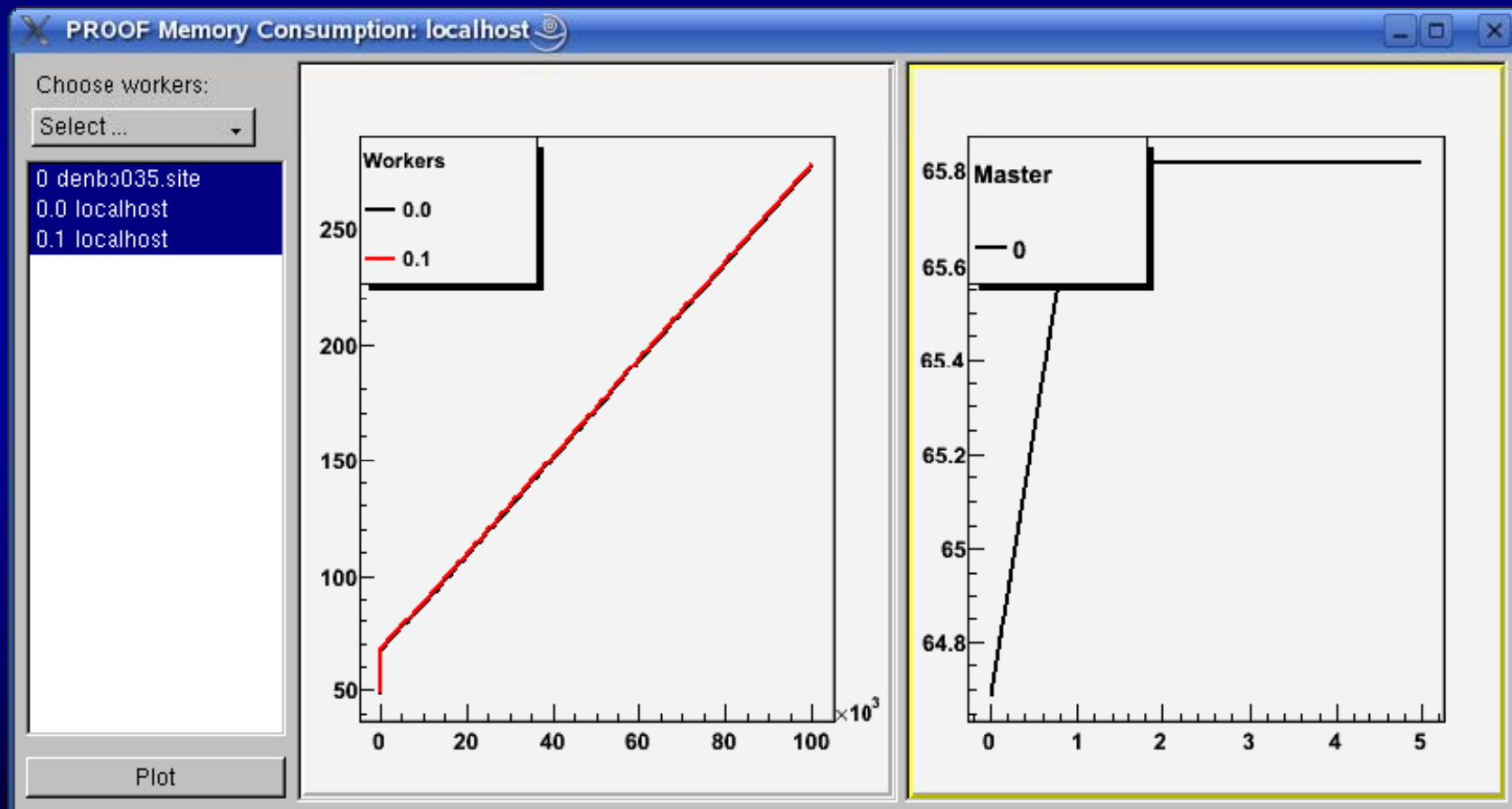
Workers monitor their memory usage and save info in the log file.

New button in the dialog box to display the evolution of memory usage per node in real time.

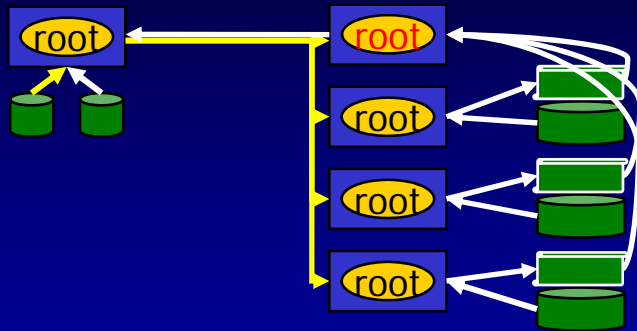
Client get warned of high usage:

The session may be eventually killed

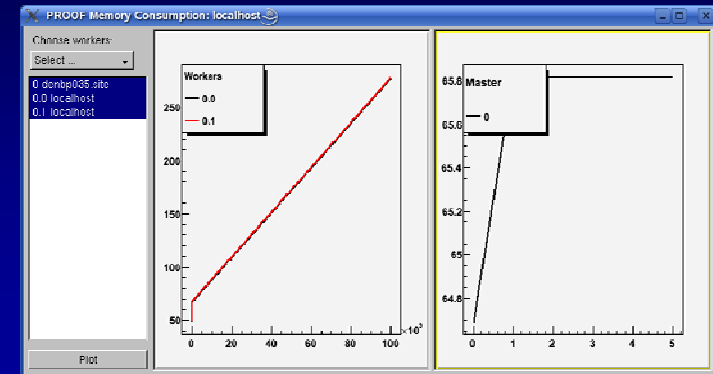




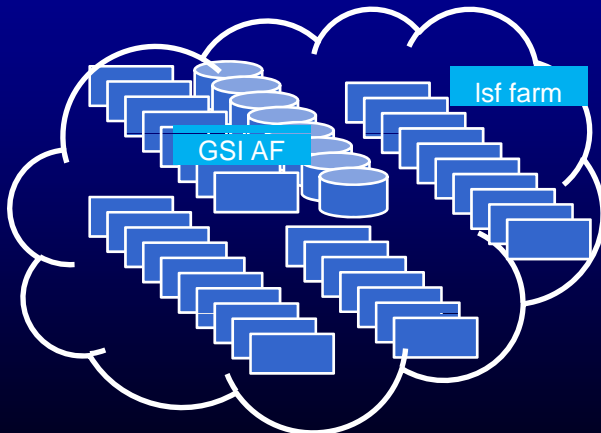
PROOF Overview



Extending PROOF

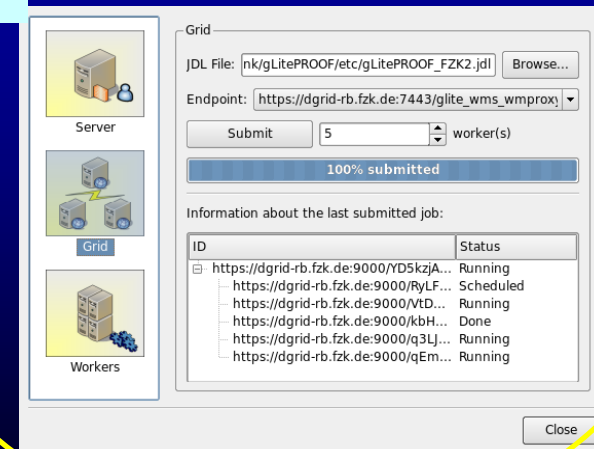


Integration of PROOF in the farm at GSI



PROOF on the Grid

A. Manafov



How to create a PROOF Cluster

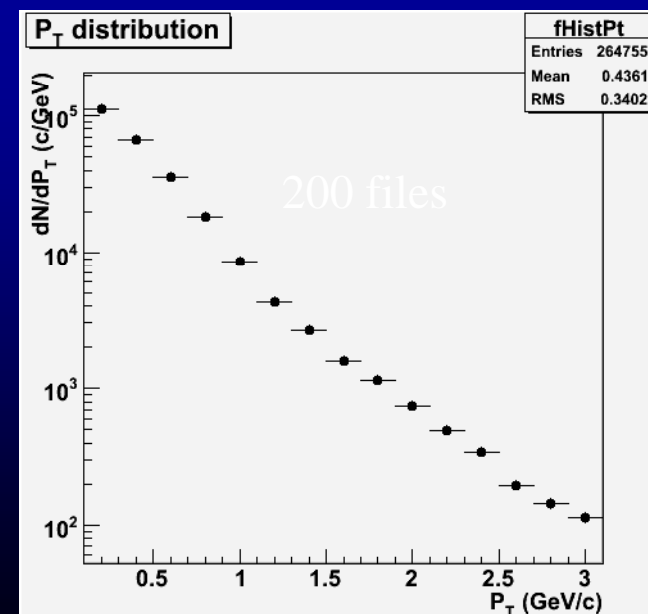
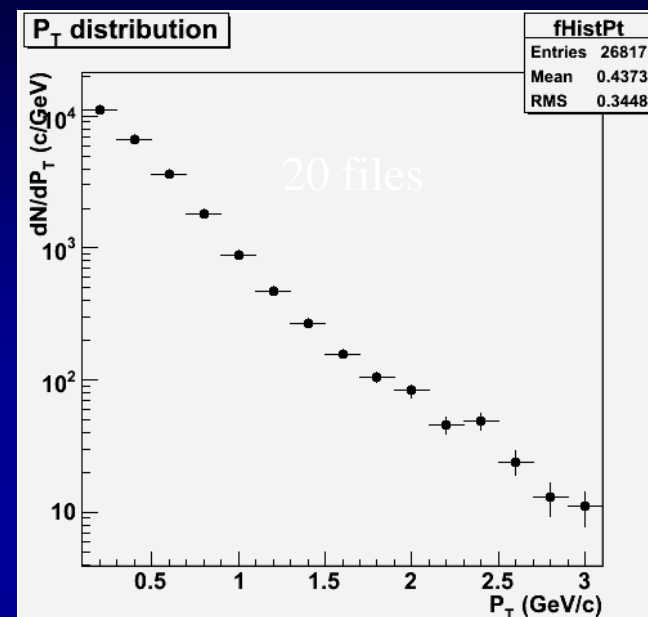
Add connecting to the cluster

➤ `TProof::Open("lxb6046")`

A PROOF Cluster is a set of demons waiting to start PROOF processes (master, or worker)

It can be setup

1. statically by the system administrator
e.g. CERNAF, GSI AF, ...
2. by the user
on machines where he can login
multiple processes on a multicore laptop
at GSI we have scripts for our batch system
3. via gLitePROOF on the GRID



gLitePROOF : **a gLite PROOF package**

A number of utilities and configuration files to implement a PROOF distributed data analysis on the gLite Grid.

Built on top of RGLite:

TGridXXX interface are implemented in RGLite for gLite MW.

ROOT team accepted our suggestions to TGridXXX interface.

gLitePROOF package

It setups "on-the-fly" a PROOF cluster on gLite Grid.

It works with mixed type of gLite worker nodes (x86_64, i686...)

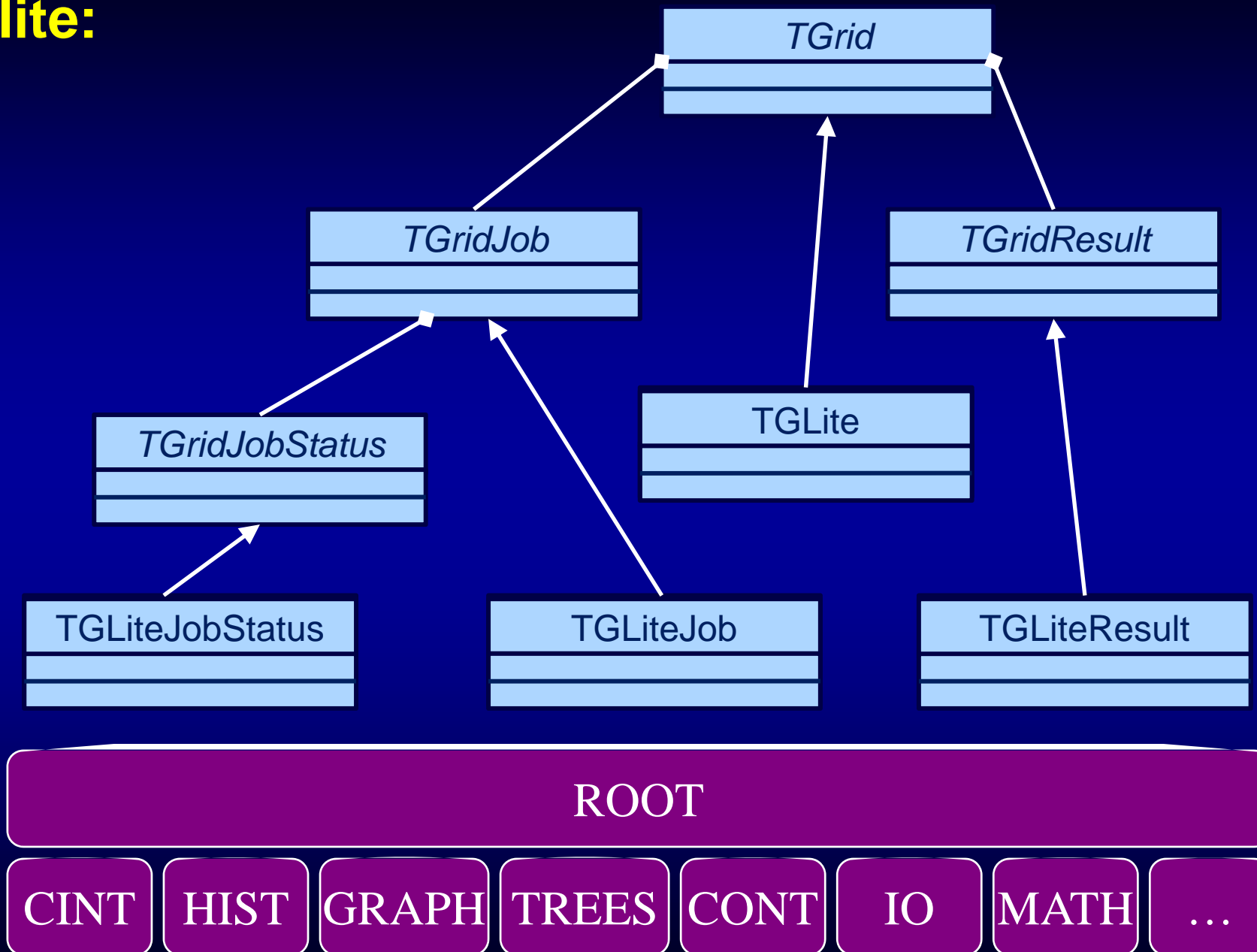
It supports reconnection.

<http://www-linux.gsi.de/~manafov/D-Grid/docz/>

RGlite:

g
L
i
t
e

A
L
I
E
N



RGLite example

```
// Initializing RGLite plug-in
TGrid::Connect("glite");
// Submitting a Job to gLite Grid
TGridJob *job = gGrid->Submit("JDLs/proofd.jdl");
// querying a Status of the Job
TGridJobStatus *status = job->GetJobStatus();
status->GetStatus();
// Getting a Job's output back to the user
job->GetOutputSandbox("/home/anar/");
```

- ❖ Job submission,
- ❖ status querying,
- ❖ output retrieving.

```
// Initializing RGLite plug-in
TGrid::Connect("glite");
// Changing current File Catalog directory to "dteam"
gGrid->Cd("dteam");
// Querying a list of files of the current FC directory
TGridResult* result = gGrid->Ls();
// Printing the list out
Int_t i=0;
while (result->GetFileName(i))
    cout << "File " << result->GetFileName(i++);
```

- ❖ Changing file catalog directory,
- ❖ querying lists of files.

ROOT Version 5.19/02 Release Notes

ROOT version 5.19/02 has been released March 15, 2008. In case you are upgrading from version 5.14, please read the releases notes of version 5.16 and version 5.18 in addition to these notes.

Binaries for all supported platforms are available at:

<http://root.cern.ch/root/Version519.html>

Versions for AFS have also been updated. See the list of supported platforms:

<http://root.cern.ch/Welcome.html>

For more information, see:

<http://root.cern.ch>

RGLITE: A ROOT GRID interface

RGLite plug-in - a ROOT plug-in module, which implements the ROOT Grid interface and offers to ROOT users possibilities to perform a number of operations using gLite middleware from within ROOT. Supported features:

- Workload Management System operations:
 - job submission – normal, DAG and parametric jobs (gLite WMPProxy API),
 - smart look-up algorithm for WMP-Endpoints,
 - job status querying (gLite LB API),
 - job output retrieving (Globus GridFTP).
- File Catalog operations (gLite/LCG LFC API):
 - smart session manager,
 - set/query the current working catalog directory,
 - list files, directories and their stats,
 - add/remove files in a catalog namespace,
 - add/remove directories,
 - add/remove replicas from a given file.
- An executive logging.
- Support of an external XML configuration file with according XML schema.

gLitePROOF components:

PROOFAgent – a lightweight, standalone C++ application. Acts as a multifunctional proxy client/server and helps to use proof/xrootd on the Grid worker nodes behind a firewall.

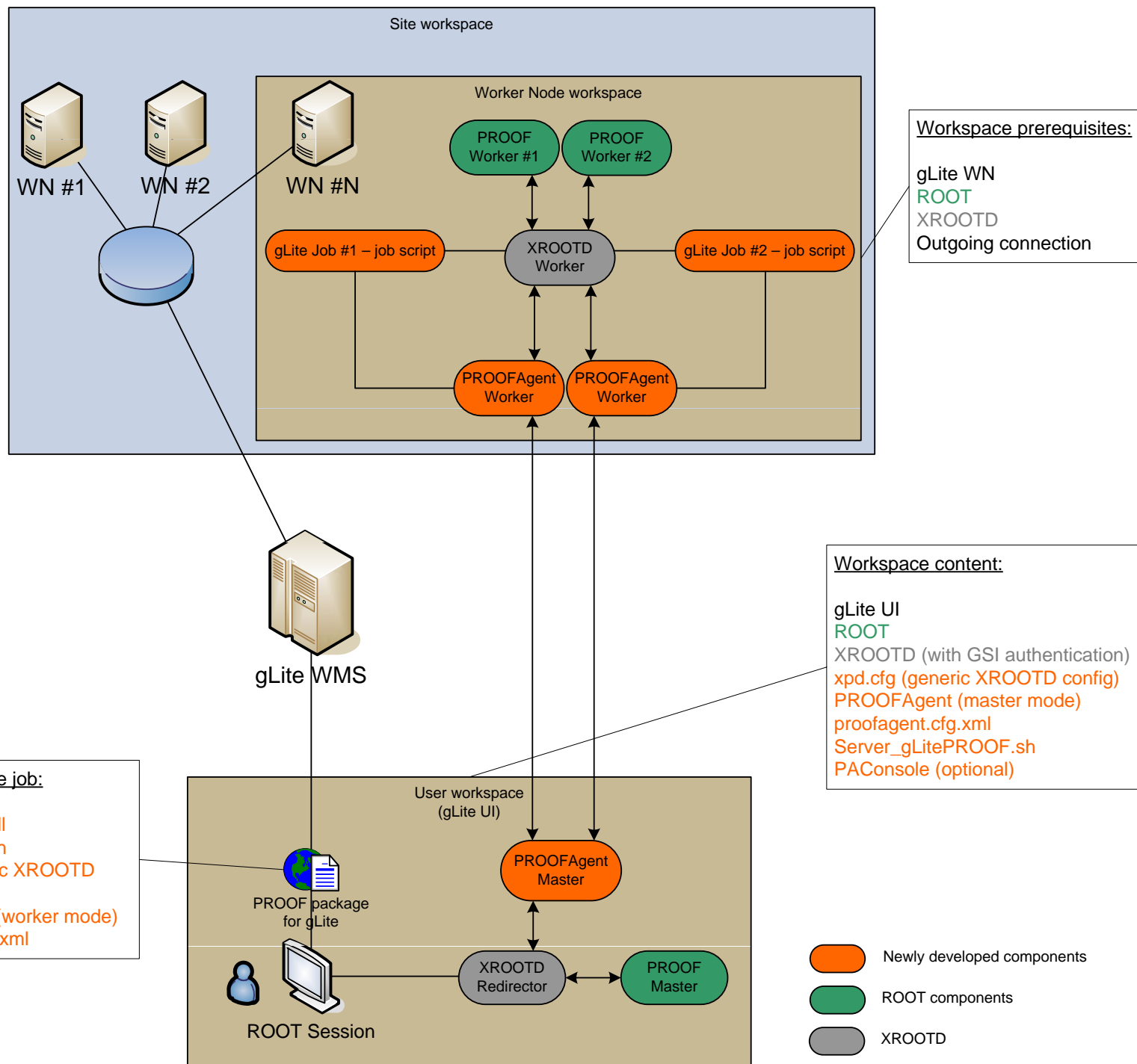
PAConsole – a standalone C++ application, provides a GUI and aims to simplify the usage of PROOFAgent and gLitePROOF configuration files. PAConsole uses GAW to perform gLite job submissions. Users can control jobs directly using ROOT and RGLite plug-in instead of using PAConsole.

xpd.cfg – a generic XROOTD configuration file (configures redirector and remote Grid workers)

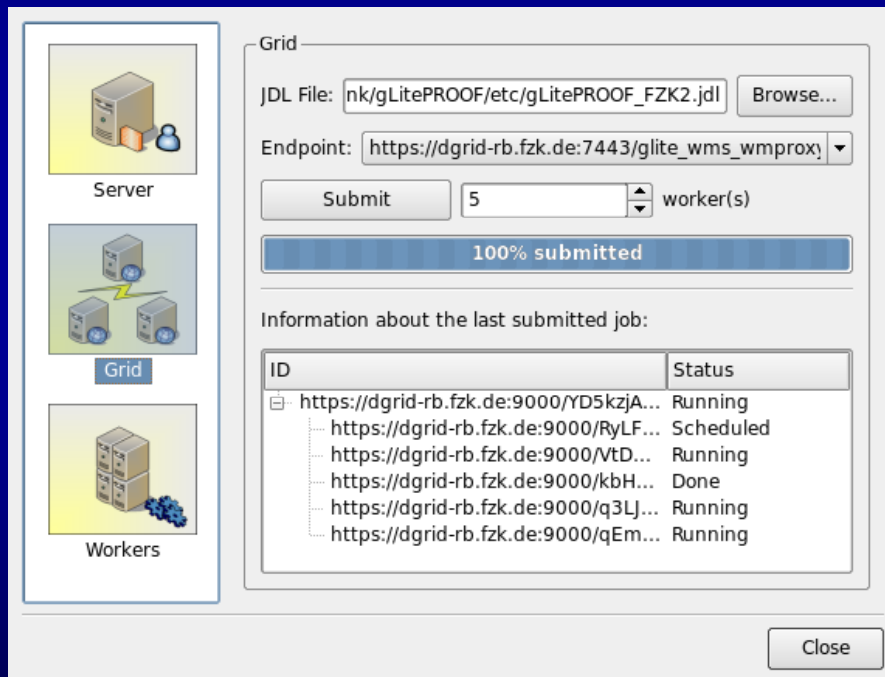
Server_gLitePROOF.sh – a server side script. Helps to start/stop services of gLitePROOF. Could be used via command line or PAConsole GUI.

gLitePROOF.jdl – a JDL file, describes a generic, parametric Grid job, which is submitted to gLite and aims to execute gLitePROOF workers on Grid worker nodes.

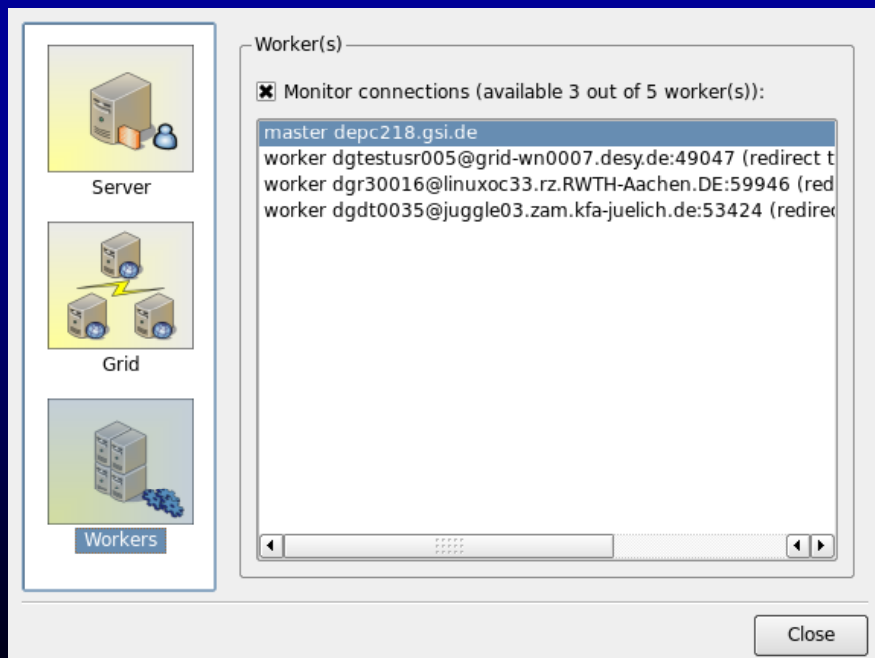
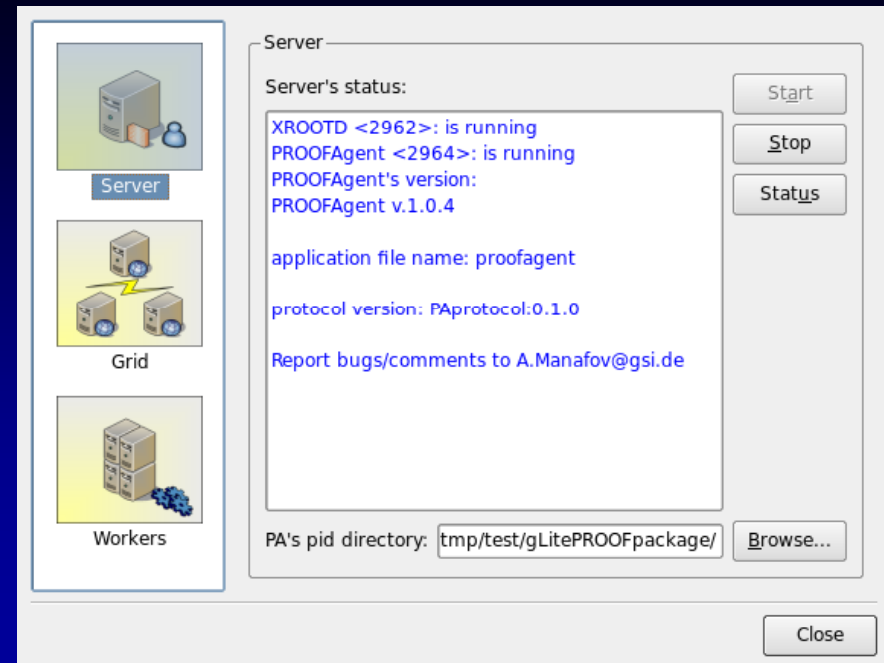
gLitePROOF.sh – a job script. Executed by LRMS on remote workers. Script makes environment recon, uploads necessary packages and starts gLitePROOF services.



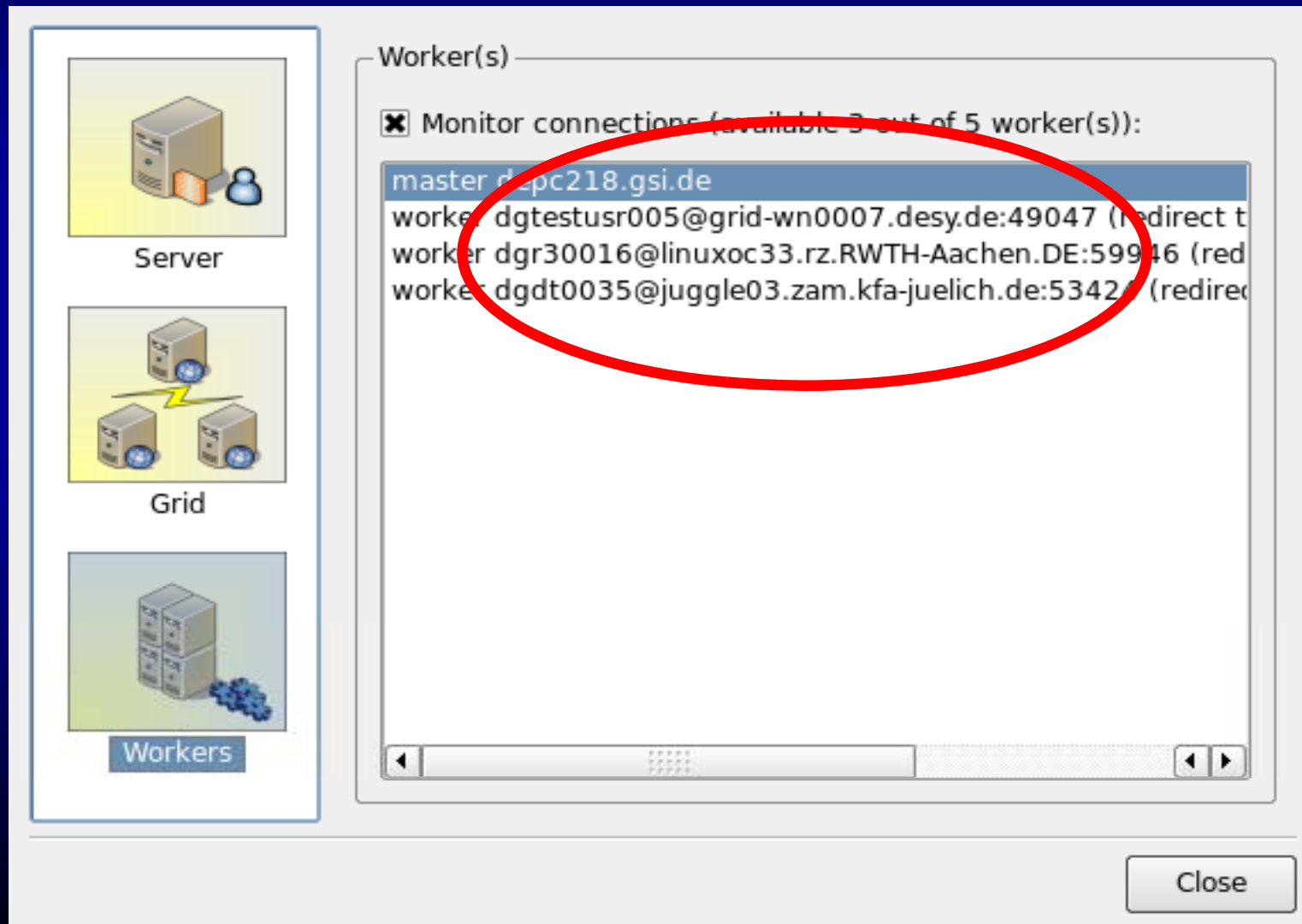
PAConsole: a GUI to setup a PROOF Cluster on demand



newest feature: supports reconnection



Workers on different sites



<https://subversion.gsi.de/trac/dgrid/>

 [Login](#) [Settings](#) [Help/Guide](#) [About Trac](#)

| | | | | | | | | |
|--|-------------|----------|---------|---------------|--------------|--------|-------------|-------|
| | Wiki | Timeline | Roadmap | Browse Source | View Tickets | Search | Doc.Portals | Build |
|--|-------------|----------|---------|---------------|--------------|--------|-------------|-------|

[Start Page](#) | [Index by Title](#) | [Index by Date](#) | [Last Change](#)

Welcome to RGLite and gLitePROOF

Projects

- [glite-api-wrapper](#) - a library, which wraps some parts of gLite API and adds automation and helpers to simplify access to the API.
- [RGLite plug-in](#) - a ROOT plug-in module, which implements the ROOT Grid interface and offers to ROOT users a possibility to use gLite middleware from within ROOT.
- [PROOFAgent](#) - a multi-functional client/server application. It helps to use proofd on the Grid's worker nodes which are behind a firewall, also PROOFAgent has number of additional useful functionality which helps to process PROOF interactive analysis on the Grid.
- [PAConsole](#) - a GUI application, helps to manage the PROOFAgent daemon and gLite_PROOF package.
- [gLitePROOF](#) - implementation of the PROOF distributed data analysis on the gLite MW.

Documentation

- Source code documentation you can find by pressing "Source Documentation" button in the main menu bar above or by the following [link](#).
- [Developers Area](#)

Support

Interactive Data Analysis with PROOF: Summary

ALICE sees PROOF as strategic tool

**Integration of the GSI AF into the general purpose
batch farm (for Grid and local batch)**

To be done: Experience Report

Extending PROOF:

all results go into the ROOT distr.

PROOF on the Grid:

RGLite is in the ROOT distr.

for the other packages see the project wiki:

<http://wiki.gsi.de/Grid/RGLiteAndGAW>





RGLite plug-in and gLitePROOF package



*Anar Manafov, GSI,
HEPCG Workshop, Jun 2008*



GAW - RGLite - gLitePROOF

Three modules, developed by GSI (Darmstadt) in terms of D-Grid Project AP3

gLite Middleware

gLite API

WMProxy

gLite LB

Globus

LFC

glite-api-wrapper library (GAW)

a library, which wraps parts of gLite API and adds automation and helpers to simplify access to it

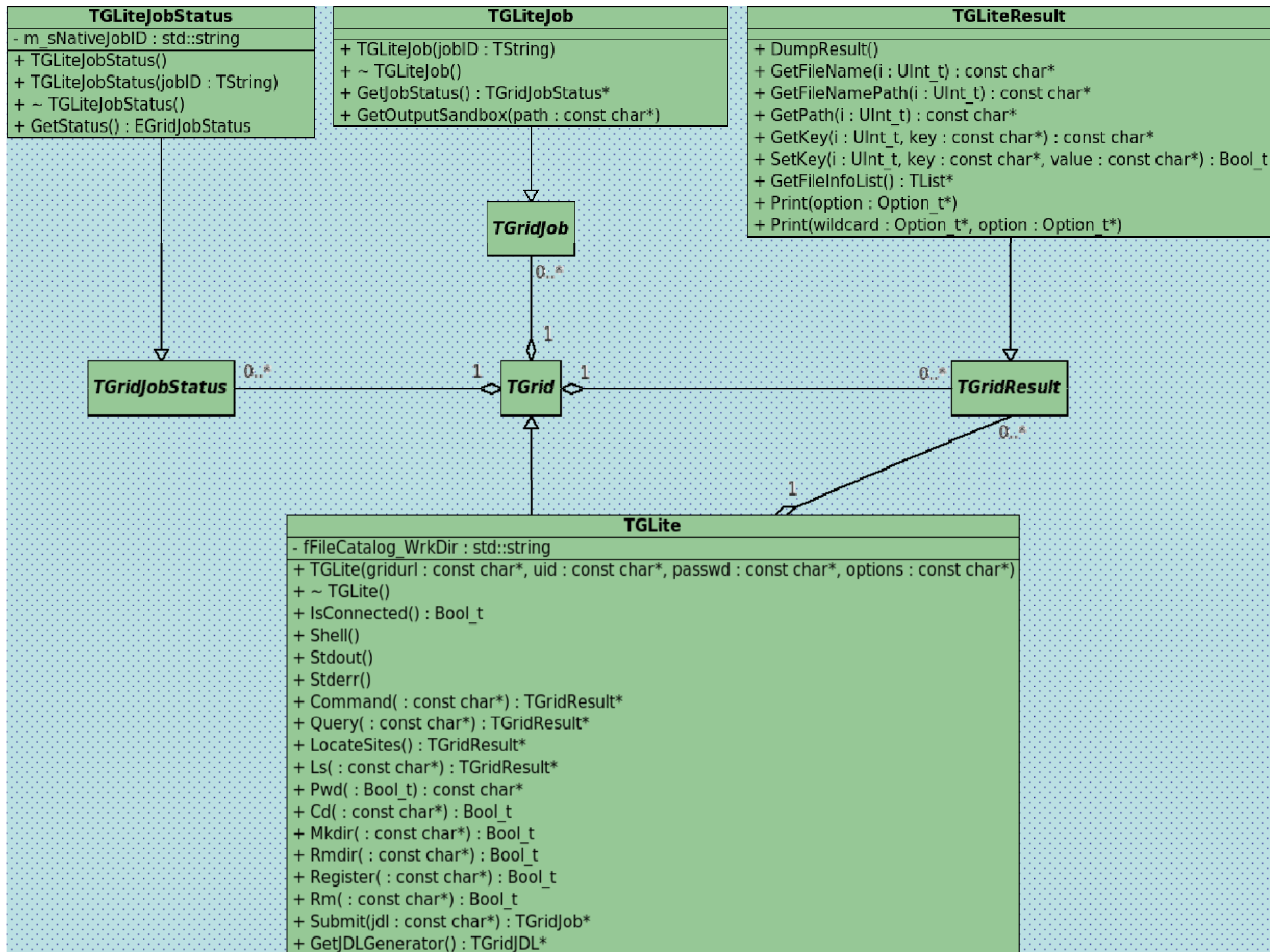
RGLite plug-in

a ROOT plug-in module, which implements the ROOT Grid interface and offers to ROOT users a possibility to use gLite middleware from within ROOT

ROOT framework

gLitePROOF

A PROOF distributed data analysis on the gLite Grid. A PROOF cluster on demand



RGLite example

```
// Initializing RGLite plug-in
TGrid::Connect("glite");
// Submitting a Job to gLite Grid
TGridJob *job = gGrid->Submit("JDLs/proofd.jdl");
// querying a Status of the Job
TGridJobStatus *status = job->GetJobStatus();
status->GetStatus();
// Getting a Job's output back to the user
job->GetOutputSandbox("/home/anar/");
```

- Job submission,
- status querying,
- output retrieving.

```
// Initializing RGLite plug-in
TGrid::Connect("glite");
// Changing current File Catalog directory to "dteam"
gGrid->Cd("dteam");
// Querying a list of files of the current FC directory
TGridResult* result = gGrid->Ls();
// Printing the list out
Int_t i=0;
while (result->GetFileName(i))\
> printf("File %s\n",result->GetFileName(i++));
```

- Changing file catalog directory,
- querying lists of files.

RGLite features

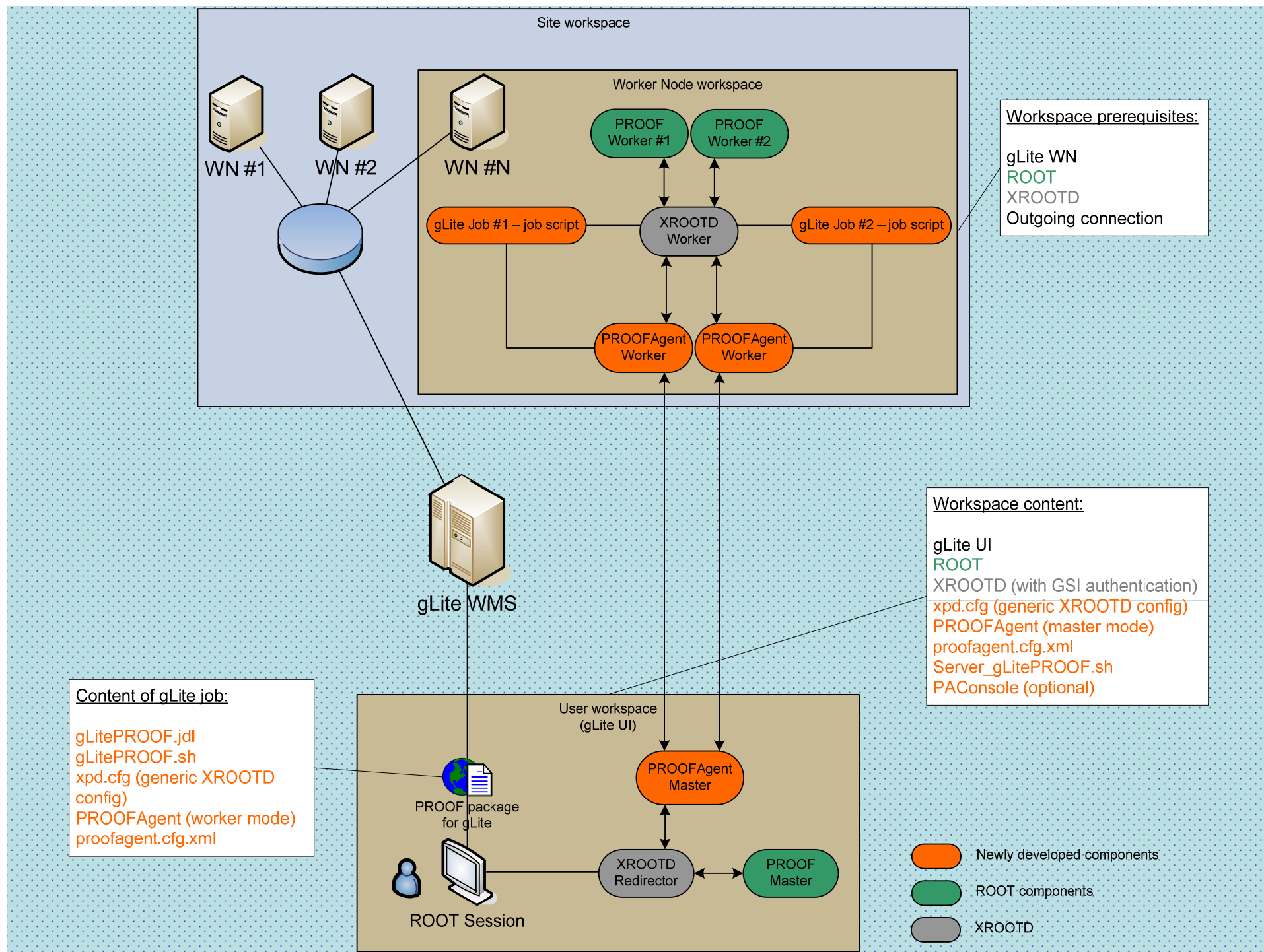


An official part of the ROOT distribution since ROOT v5.19.


- **Workload Management System operations:**
 - job submission – normal, DAG and parametric jobs (gLite WMPProxy API),
 - smart look-up algorithm for WMP-Endpoints,
 - job status querying (gLite LB API),
 - job output retrieving (Globus GridFTP).
- **File Catalog operations (gLite/LCG LFC API):**
 - smart session manager,
 - set/query the current working catalog directory,
 - list files, directories and their stats,
 - add/remove files in a catalog namespace,
 - add/remove directories,
 - add/remove replicas from a given file.
- **An executive logging.**
- **Support of an external xml configuration file.**

gLitePROOF – a gLite PROOF package.


It is a number of utilities and configuration files, developed at GSI in terms of the D-Grid project and aims to implement a PROOF distributed data analysis on the gLite Grid.




PAConsole - a GUI of gLitePROOF



Server



Grid



Workers

Server

Server's status:

```


XROOTD <31191>: is running
PROOFAgent <31200>: is running
PROOFAgent's version:
PROOFAgent v.1.0.4

application file name: proofagent

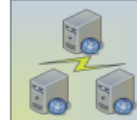
protocol version: PAprotocol:0.1.0

Report bugs/comments to A.Manafov@gsi.de
        
```

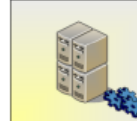
PA's pid directory:



Server



Grid



Workers

Grid

JDL File:

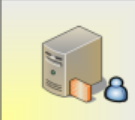
Endpoint:

worker(s)


100% submitted

Information about the last submitted job:


| ID | Status |
|---|-----------|
| https://dgrid-rb.fzk.de:9000/YD5kzjA... | Running |
| https://dgrid-rb.fzk.de:9000/RyLF... | Scheduled |
| https://dgrid-rb.fzk.de:9000/VtD... | Running |
| https://dgrid-rb.fzk.de:9000/kbH... | Done |
| https://dgrid-rb.fzk.de:9000/q3LJ... | Running |
| https://dgrid-rb.fzk.de:9000/qEm... | Running |



Server



Grid




Workers

Worker(s)

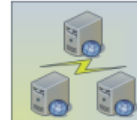
☒ Monitor connections (available 3 out of 3 worker(s)):

```

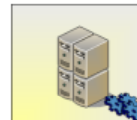
master depc218.gsi.de
worker dech001@grid44.gsi.de:34391 (redirect through local
worker dech001@grid4.gsi.de:53679 (redirect through localh
worker dech001@grid17.gsi.de:56334 (redirect through local
        
```



Server



Grid



Workers

Grid

JDL File:

Endpoint:

worker(s)

100% submitted

Information about the last submitted job:

| ID | Status |
|---|---------|
| https://grid25.gsi.de:7443/glite_wms_wmpro... | Running |
| https://grid25.gsi.de:7443/glite_wms_wmpro... | Waiting |
| https://grid25.gsi.de:7443/glite_wms_wmpro... | Waiting |
| https://grid25.gsi.de:7443/glite_wms_wmpro... | Waiting |



Ctrl+C
 Ctrl+O
 Ctrl+E

gLitePROOF summary

➤ Requirements

- Client-Side: gLite UI 3.1, ROOT 5.18+,
- Remote-Side: WMP Endpoint, open out. traffic (Globus ports on WNs).

➤ Easy in use

-  ➤ one-click-installation,
-  ➤ user's manual (<http://www-linux.gsi.de/~manafov/D-Grid/docz/>)
- user friendly GUI,
- works out of the box.

➤ Easy to extend

➤ Combines resources of the Grid and advantages of PROOF

- transparency (local ROOT analysis <-> PROOF <-> gLitePROOF),
- scalability (The basic architecture should not put any implicit limitations on the number of workers).

➤ Works on heterogeneous machines (gLite WNs)

➤ Uses xrootd and Grid methods to access data

Summary

➤ RGLite

- TGridXXX interface are implemented in RGLite for gLite MW.
- ROOT team accepted our suggestions to TGridXXX interface.
- uses WMP job submission.
- compatible with gLite UI 3.1.

NEW

- a part of the ROOT distributive (since ROOT v5.19).

➤ gLitePROOF package

NEW

- New stable release.

NEW

- Two officially registered users from ATLAS (CPPM and LAL, France).
- setups “on-the-fly” a PROOF cluster on gLite Grid.
- works with mixed type of gLite worker nodes (x86_64, i686...).
- supports reconnections.
- provides GUI.

➤ Use of Agile methods of software development:

- Continuous integration (automated builds on SLC3, SL4, and F8, nightly builds).
- Unit tests.
- Projects metrics.
- Task tickets and sprints...

Trac: <https://subversion.gsi.de/trac/dgrid>

Wiki: <http://wiki.gsi.de/cgi-bin/view/Grid/RGLiteAndGAW>

GEFÖRDERT VOM



Bundesministerium
für Bildung
und Forschung

16-Jun-2008