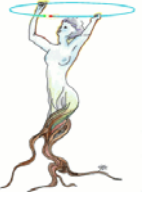# PROOF news

Gerardo Ganis / CERN

People working on the project:
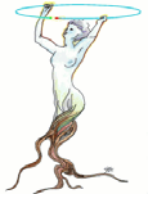B. Bellenot, M. Biskup, G. Kickinger, A. Peters, F. Rademakers / CERN
M. Ballintijn / MIT, D. Feitchinger / PSI

**2005 Intl. ROOT Users Workshop
28 – 30 Sept 2005, CERN, Switzerland**

# Outline

- PROOF overview
- New features and improvements
- Plans
- Summary

# PROOF – Parallel ROOT Facility

- System to export the ROOT analysis model on clusters of computers for interactive analysis of large data sets

- Flexible multi-tier architecture
  - in GRID contexts adapts to cluster of clusters or wide area *virtual clusters*

- Exploit inter-independence of entries in a tree or directory to achieve basic parallelism
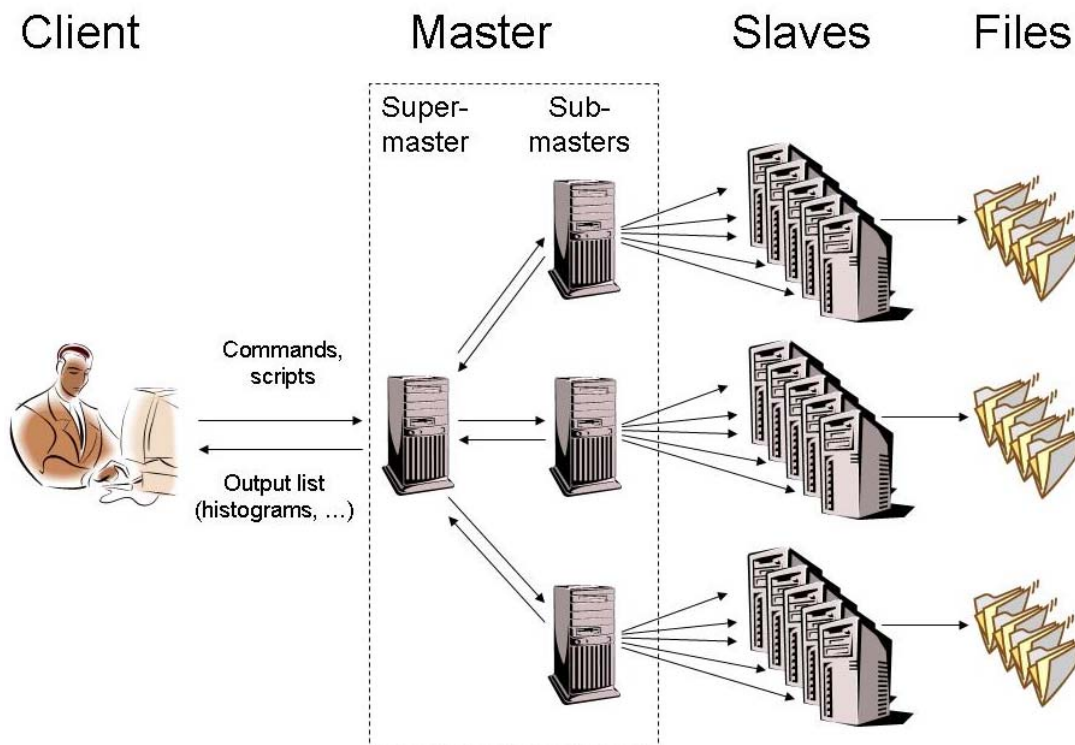  - data set split into *packets* assigned to worker nodes on demand

# PROOF – Design goals

- ## Transparency
  - Same syntax as in normal ROOT session
  - Input objects copied from client
  - Output objects merged, returned to client

- ## Scalability and Adaptability
  - Handle heterogeneous server performances
  - Vary packet size according to number of worker nodes and relative performance
  - Multi-Level-Master addresses cases of wide area clusters

# PROOF – Multi-tier Architecture

- Optimize for data locality
- If not possible, remote data via (x)rootd, rfiod, dCache …
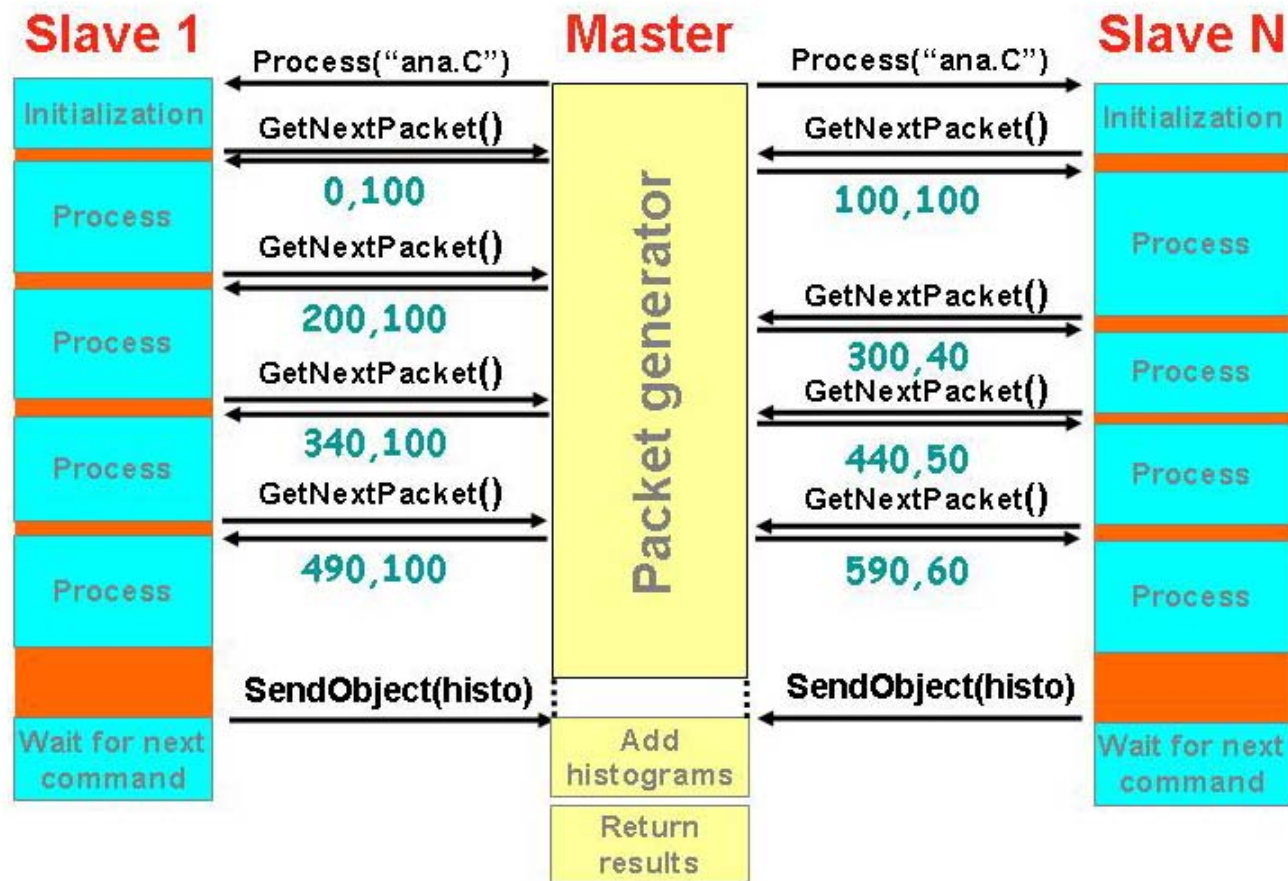
G. Ganis, ROOT05, 29 Sept 2005

# PROOF – Workflow

- ## Pull architecture
  - *dynamic load balancing* naturally achieved

# PROOF – User Sandbox

- User's have their own sandbox on each worker node

- File transfers minimized
    - cache packages, selector
    - File integrity: MD5 checksums, timestamps

- Package manager to upload files or packages
    - binary or source
    - PAR (PROOF Archive, like Java jar)
        - provides ROOT-INF directory, BUILD.sh, SETUP.C to control setup in each worker
    - `TProof` API to handle all this

# PROOF – Additional issues
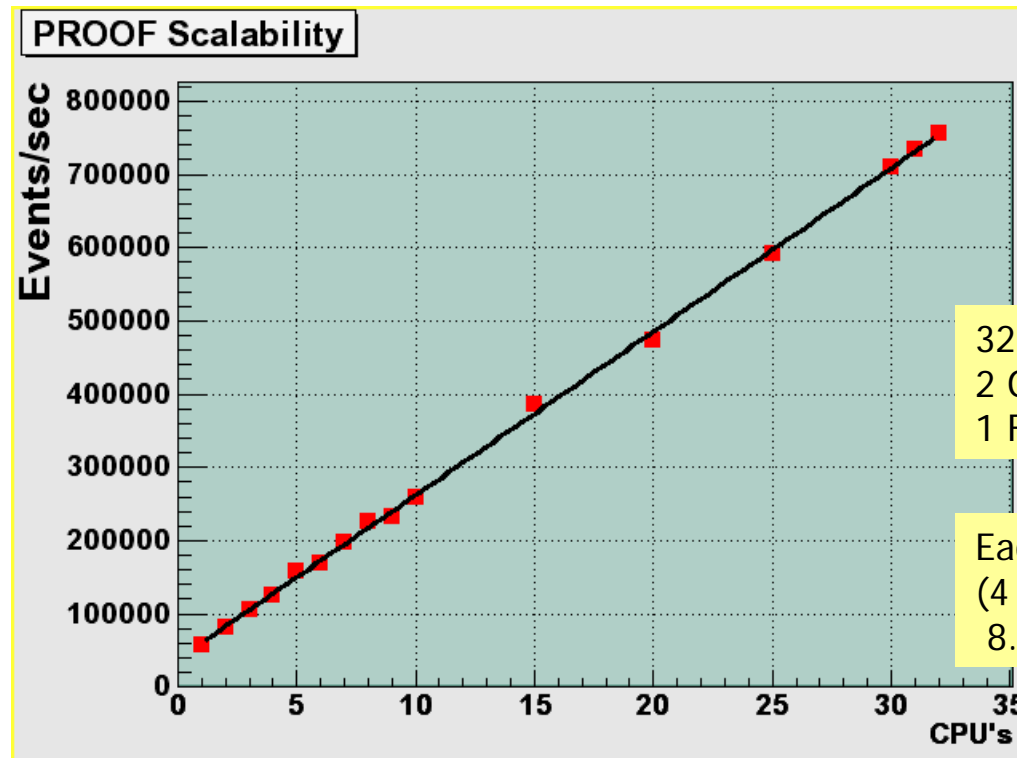
- ## Error handling
  - Death of master: fatal
  - Death of worker: redirect its work to other workers
  - Ctrl-C interrupt: send as OOB message (C->M, M->W)

- ## Authentication
  - Supported: password-based protocols, globus/GSI, Krb5, fast identification
  - Credential forwarding
  - Mixed configuration possible
    - e.g. pwd (C->M), fast ID internally (M->W)

# PROOF – Scalability



8.8GB, 128 files
1 node: 325 s
32 nodes in parallel: 12 s

32 nodes: dual Itanium II 1 GHz CPU's,
2 GB RAM, 2x75 GB 15K SCSI disk,
1 Fast Eth, 1 GB Eth nic (not used)

Each node has one copy of the data set
(4 files, total of 277 MB), 32 nodes:
 8.8 Gbyte in 128 files, 9 million events

# PROOF – data analysis

## Normal ROOT

- **TChain**: collection of **TTree**
- **TSelector**: frame for **Begin(), Process(), Terminate()**

```
TChain a("h42");
{// Define the data set
    a.Add("root://oplapro62.cern.ch//tmp/dstarmb.root");
    a.Add("root://oplapro62.cern.ch//tmp/dstarp1a.root");
    a.Add("root://oplapro62.cern.ch//tmp/dstarp1b.root");
    a.Add("root://oplapro62.cern.ch//tmp/dstarp2.root");
    // Process the selector
    a.Process("h1analysis.C");
}
```
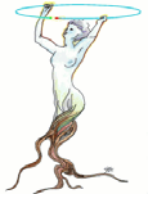
Local processing

## PROOF

- Same chain, same selector

```
{// Open PROOF
    TProof *proof = new TProof("master")
    // Process the selector
    a.Process("h1analysis.C");
}
```

Remote processing

# PROOF on wide area clusters

- Use "file catalogs" for file location
- Use "resource brokers" to identify best worker nodes
- Examples: interface with Condor COD, PEAC
  - see M. Ballintijn at ROOT04
- In GRID context
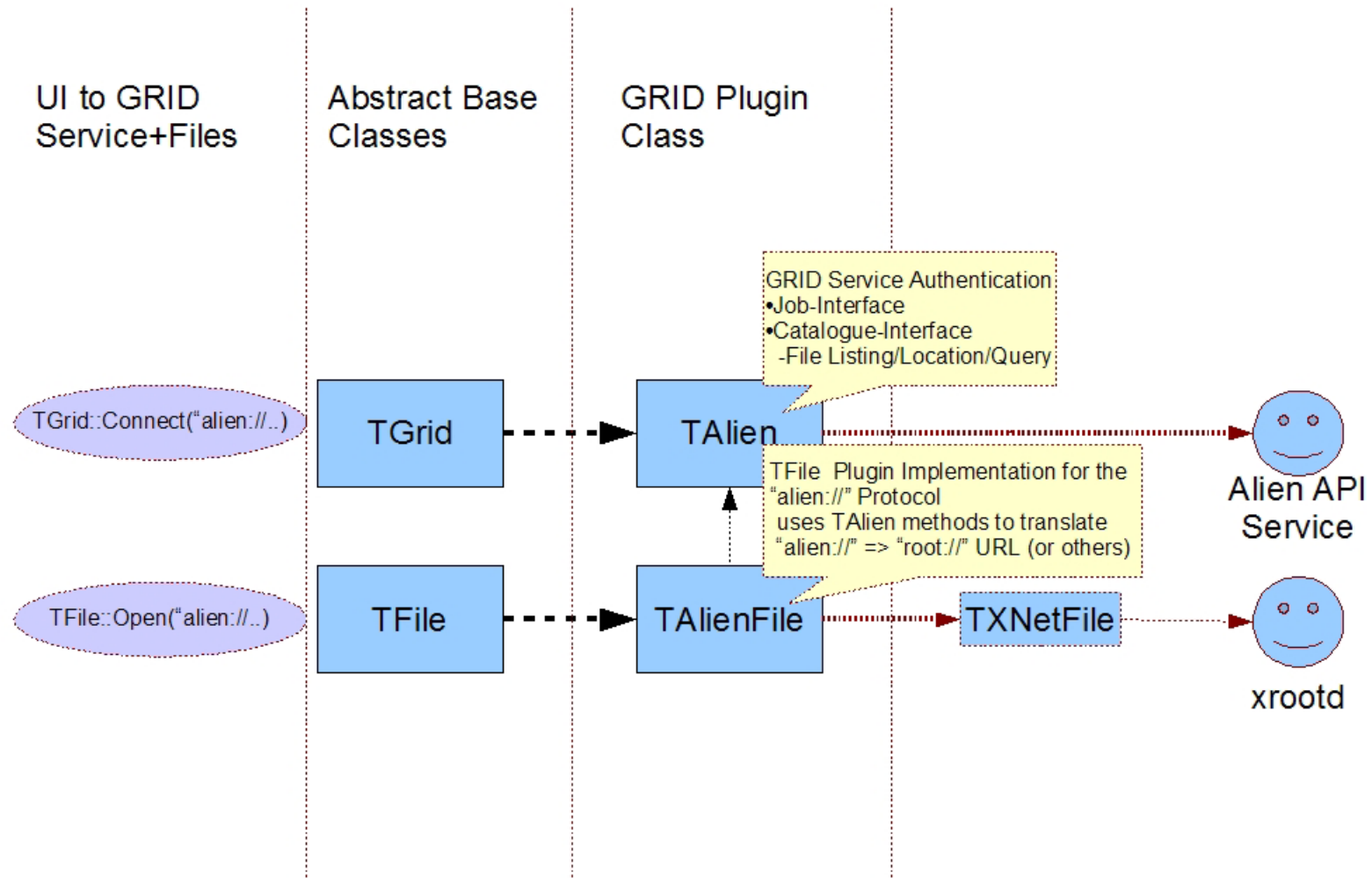  - **TGrid**: abstract interface for all services

```
class TGrid : public TObject {
public:
   …
   virtual TGridResult *Query (const char *query) = 0;

   static TGrid *Connect (const char *grid, const char *uid = 0,
                          const char *pw = 0);


   ClassDef(TGrid,0)  // ABC defining interface to GRID services
};
```

# TGrid: AliEn implementation

UI to GRID Service+Files | Abstract Base Classes | GRID Plugin Class

GRID Service Authentication:
• Job-Interface
• Catalogue-Interface
  -File Listing/Location/Query

TGrid::Connect("alien://..")  →  TGrid  ⇢  TAlien  →  Alien API Service

TFile Plugin Implementation for the "alien://" Protocol
uses TAlien methods to translate "alien://" => "root://" URL (or others)

TFile::Open("alien://..")  →  TFile  ⇢  TAlienFile  →  TXNetFile  →  xrootd

# TGrid: AliEn implementation (cnt'd)

```cpp
// Connect
TGrid *alien = TGrid::Connect("alien://");

// Query
TGridResult *res =
    alien->Query("/alice/cern.ch/user/p/peters/analysis/miniesd/",
                 "*.root");
// List of files
TList *listf = res->GetFileInfoList();

// Create chain
TChain *chain = new TChain("Events", "session");
chain->AddFileInfoList(res);

// Start PROOF
TProof *proof = new TProof("remote");
chain->SetProof();                              DEMO

// Process your query
chain->Process("selector.C");
```
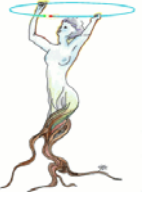
# New developments

- Towards "interactive batch"
    - Support for asynchronous running mode
    - Support for multi sessions
    - Query result management

- GUI controller

# Interactive batch

## Why? Analysis session: example

AQ1: 1s query produces a local histogram

AQ2: a 10mn query submitted to PROOF1

AQ3->AQ7: short queries

AQ8: a 10h query submitted to PROOF2

**Monday at 10h15 ROOT session on my laptop**

BQ1: browse results of AQ2
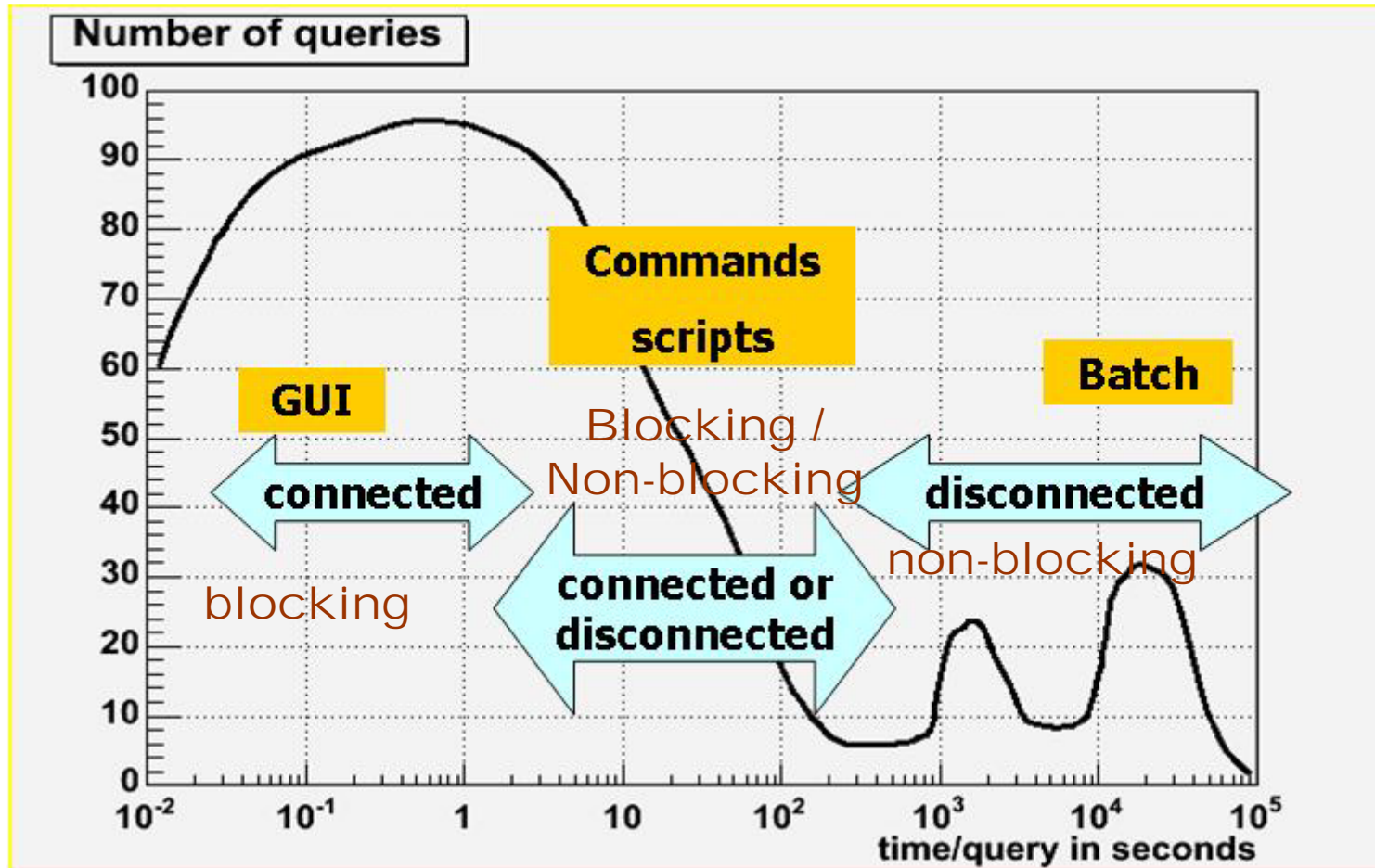
BQ2: browse temporary results of AQ8

BQ3->BQ6: submit 4 10mn queries to PROOF1

**Monday at 16h25 ROOT session on my laptop**

CQ1: Browse results of AQ8, BQ3->BQ6

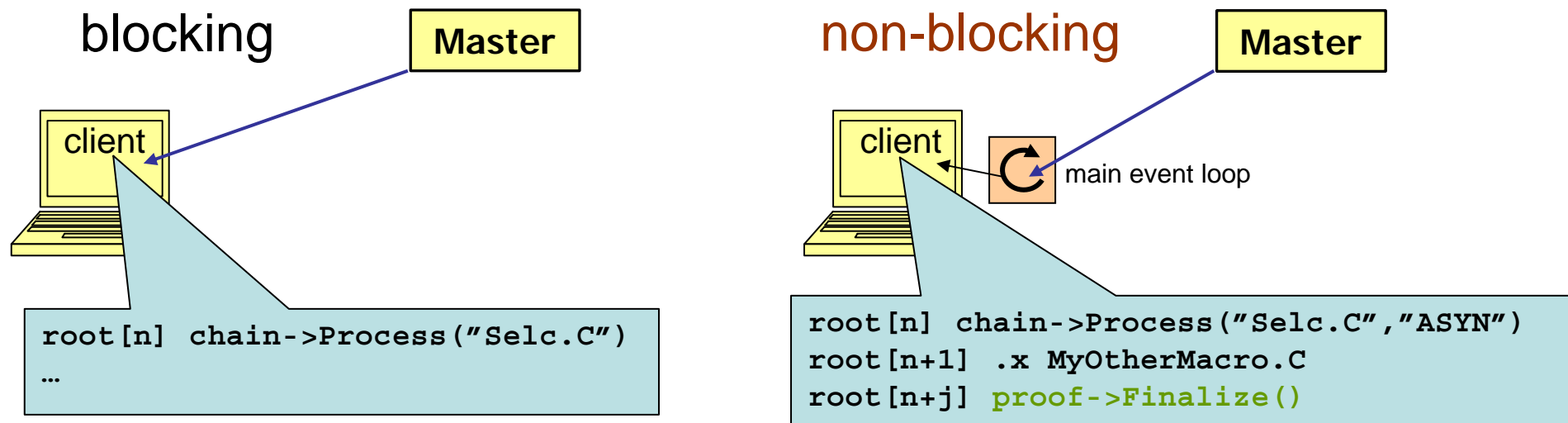**Wednesday at 8h40 Carrot session on any web browser**

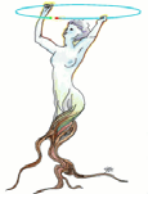# Typical query-time distribution
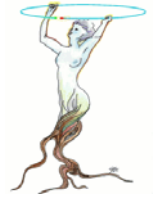
## Asynchronous mode

- PROOF initially designed for short queries:
  - **`TProof::Process()`** blocks, as in a local session
- Inconvenient for longer queries:
  - *idle* client session waiting for master reply
- Non-blocking processing achieved handling input from master via the main event loop

blocking

| Master |

| client |

```
root[n] chain->Process("Selc.C")
…
```

non-blocking

| Master |

| client |

main event loop

```
root[n] chain->Process("Selc.C","ASYN")
root[n+1] .x MyOtherMacro.C
root[n+j] proof->Finalize()
```
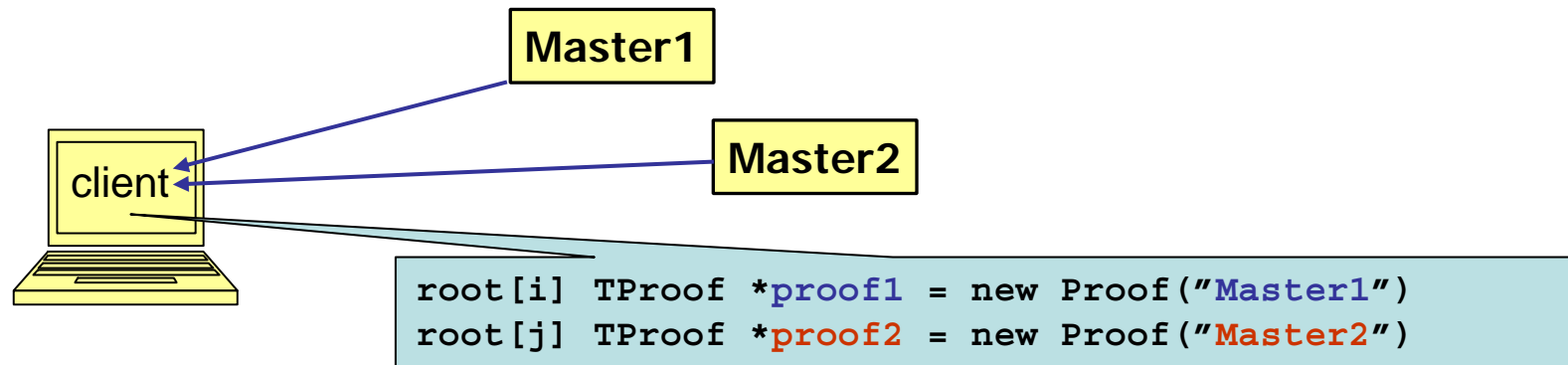
## Asynchronous mode (cnt'd)

- End-of-Processing signaled by reception of a **TQueryResult** object (see below)

- **TProof::Finalize()**
  - Re-initialize **TSelector**, if needed
  - Merge outputs
  - Run **TSelector::Terminate()**

- *Submitted* queries are added to the waiting list on the master and processed sequentially

# Multi sessions

- ## Start PROOF sessions to more than one cluster



```
root[i] TProof *proof1 = new Proof("Master1")
root[j] TProof *proof2 = new Proof("Master2")
```

- ## List of open sessions in **TROOT**

```
root[n] gROOT->GetListOfProofs()->ls()
```

- ## Results of queries saved in lists owned by **TProof**

```
root[m] proof1->GetQueryResults()->ls()
```

# Query result management

- Asynchronous finalization needs selector, output list
- Keep track of what produced a given set of results

- **TQueryResult:** everything about a query
  - unique identifier (PROOF session + sequential #)
  - selector files (in compressed form, **TMacro**)
  - data set definition, list of loaded packages
  - processing info:
    - logs, start/end time, performance parameters
  - list of output objects
- Global list of **TQueryResult** in **TProofPlayer**

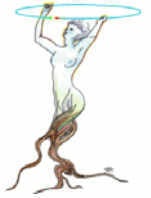- Will also be used for local sessions

# Query result management (cnt'd)

- **`TQueryResult`** object created by the master
  - Copy saved and continuously updated in sandbox

```
> ls <SandBox>/queries
./  ../  session-0-pcepsft43-1234678932-23456/
> ls <SandBox>/queries/session-0-pcepsft43-1234678932-23456
1/ 2/ 3/
> ls <SandBox>/queries/session-0-pcepsft43-1234678932-23456/2
query-result.root
```

- Master has access to all queries in the sandbox
  - can retrieve content at any time, from any session
  - can archive to any mass storage
  - can cleanup (if not in use)

# Query result management (cnt'd)

- **`TProof`** API to handle query results:
  - **`ShowQueries("A")`**
    - show list of queries, either local or known to master
  - **`ShowLog(Int_t query)`**
    - show logs from processing the query
  - **`Finalize(Int_t query)`**
    - run **`Terminate()`** of the output list (only once)
  - **`Retrieve(Int_t query, const char *path)`**
    - retrieve the result of processing from the master
  - **`Remove(Int_t query)`**
    - remove a query from the lists
  - **`TList *GetQueryResults()`**
    - get the list of **`TQueryResult`** already retrieved
  - …

# New features at work - 1



Start new session

**gProof** is the current session

ShowQueries

Unique ID

Selector name

Status

Event range

```
Shell - Konsole <5>
Session Edit View Settings Help
pcepsft43:~/local/root/test/proof/demo>
pcepsft43:~/local/root/test/proof/demo>root -l
root [0]
root [0]
root [0] TProof *proof1 = new TProof("localhost")
Starting master: opening connection ...
Starting master: OK
Opening connections to workers: OK (4 workers)
Setting up worker servers: OK (4 workers)
PROOF set to parallel mode (4 workers)
root [1]
root [1]
root [1] gProof == proof1
(int)1
root [2]
```

```
Shell - Konsole <5>
View Settings Help
root [2]
root [2] proof1->ShowQueries("A")
+++
+++ Queries processed during other sessions: 3
+++ #:1 ref:"session-0-pcepsft43-1127858035-16151:q1" sel:SimpleSel stopped   evts:0-1078
+++ #:2 ref:"session-0-pcepsft43-1127858035-16151:q2" sel:SimpleSel completed evts:0-1349
+++ #:3 ref:"session-0-pcepsft43-1127858689-16282:q2" sel:SimpleSel completed evts:0-1349
+++
+++ Queries processed during this session: selector: 0, draw: 0
+++
root [3]
```

# New features at work - 2

Define the data set

```
Shell - Konsole <6>
Session Edit View Settings Help

root [3]
root [3]
root [3]   TDSet* chain = new TDSet("TTree","bg_filtered")
root [4]   chain->Add("$ROOTSYS/tutorials/mlpHiggs.root")
(Bool_t)1
root [5] chain->Print("a")
OBJ: TDSet        type TTree        bg_filtered     in /     elements 1
TDSetElement file="$ROOTSYS/tutorials/ml            g_filtered" first=0 num=-1 msd=""
root [6] █
```

Draw query:
histograms filled on
the cluster

```
Shell - Konsole <5>
Session Edit View Settings Help

root [6]
root [6]
root [6]  chain->Draw("nch","")
<TCanvas::MakeDefCanvas>: created default TCanvas with name c1
(Int_t)1350
root [7]
root [7]
```

# New features at work - 3

Selector query

```
5>
                                    Settings  Help
+++ Queries processed during this session: selec
+++
root [9]
root [9]
root [9]  chain->Process("SimpleSel.C","ASYN")
(Int_t)2
root [10]
```

**PROOF Query Progress: localhost**

Executing on PROOF cluster "localhost" with 4 parallel slaves:

Selector: SimpleSel.C

1 files, number of events 1350, starting event 0

Processed:        1350 events in 10.4 sec

Processing rate:    127.4 events/sec

☐ Close dialog when processing is complete

☐ Show only logs from query | last

Show Logs

**PROOF Processing Logs: localhost**

```
Info in <TProofServ::SetQueryRunning> on master0: starting query: 1

Info in <TProofPlayerRemote::Process>: starting new query
Info in <TSelector::Begin>: 0x95f0ad0:     * Option = "ASYN"

Info in <TProofServ::SetQueryRunning> on master0: starting query: 2
Info in <TSelector::SlaveBegin> on slave 0.2:     * Option = "ASYN"
Info in <TSelector::SlaveBegin> on slave 0.2: histo: 0x927efd8 created and added to output
Info in <TSelector::SlaveBegin> on slave 0.2: tree: (nil)
Info in <TSelector::Process> on slave 0.2: entry: 41,  nch: 0.000000
Info in <TSelector::Process> on slave 0.2: entry: 101,  nch: 0.000000
Info in <TSelector::Process> on slave 0.2: entry: 111,  nch: 0.000000
```

Close

**Shell**

Session

```
root [10
root [10
+++
+++ Quer
+++ #:1                                           0-1078
+++ #:2                                           0-1349
+++ #:3                                           0-1349
+++
+++ Quer
+++ #:4                                           0-1349
+++
root [11]
root [11]
```

# New features at work - 4

# New features at work - 5



Finalize at any time, any order

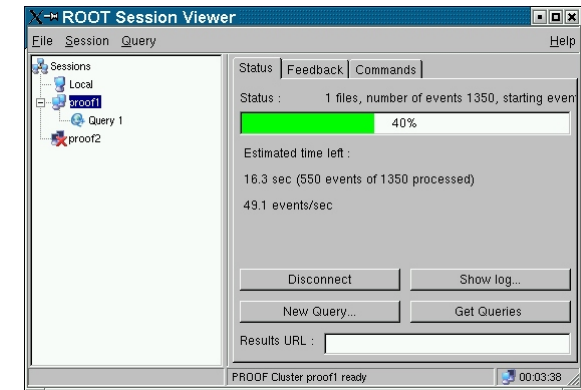# GUI controller

## Allows full *on-click* control on everything



- define a new session
- submit a query, execute a command
- query editor
  - execute macro to define or pick up a **TChain**
  - browse directories with selectors
- online monitoring of feedback histograms
- browse folders with results of query
- retrieve, delete, archive functionality
- start viewer for fast **TChain** browsing

**Improvements**

- **Session startup**
  - Optimized sequential startup
    - separate `proofd` phase from server setup
    - `execv("proofserv")` run in parallel
    - significant speed-up for medium/slow machines
  - Parallel startup
    - Start workers in dedicated threads
    - almost full parallelism
      - protect authentication, lists updates
    - `Proof.ParallelStartup: yes` (on the master)
  - Startup status and progress bar
    - Psychological speed-up …

## Improvements (cnt'd)

- New packetizer (`TPacketizerProgressive`)
  - optimized data file opening
- Support for `TTree` friends in PROOF
- Full draw functionality via the master  M.Biskup's talk

# Plans

- **Consolidate what in place**
  - Error handling, fault tolerance
- **Complete wiring, tuning of the GUI**
- **Complete implementation of "interactive batch"**
  - Stateless connection (see next)
- **Cluster configuration**
  - dynamic master-worker setup
  - allow workers come and go
- **Documentation**

# PROOF – Connection layer



slave 1

fork()
execv()

proofd

proofd

proofslave

· · ·

slave n

fork()
execv()

proofd

proofd

proofslave

master

execv()

proofd

fork()

proofd

proofserv

client

○ *parent* proofd (always running)

○ *child* proofd (executing proofserv / proofslave)

○ proofserv / proofslave : `TProofServ` instances (`TApplication`)

G. Ganis, ROOT05, 29 Sept 2005

# Stateless connection
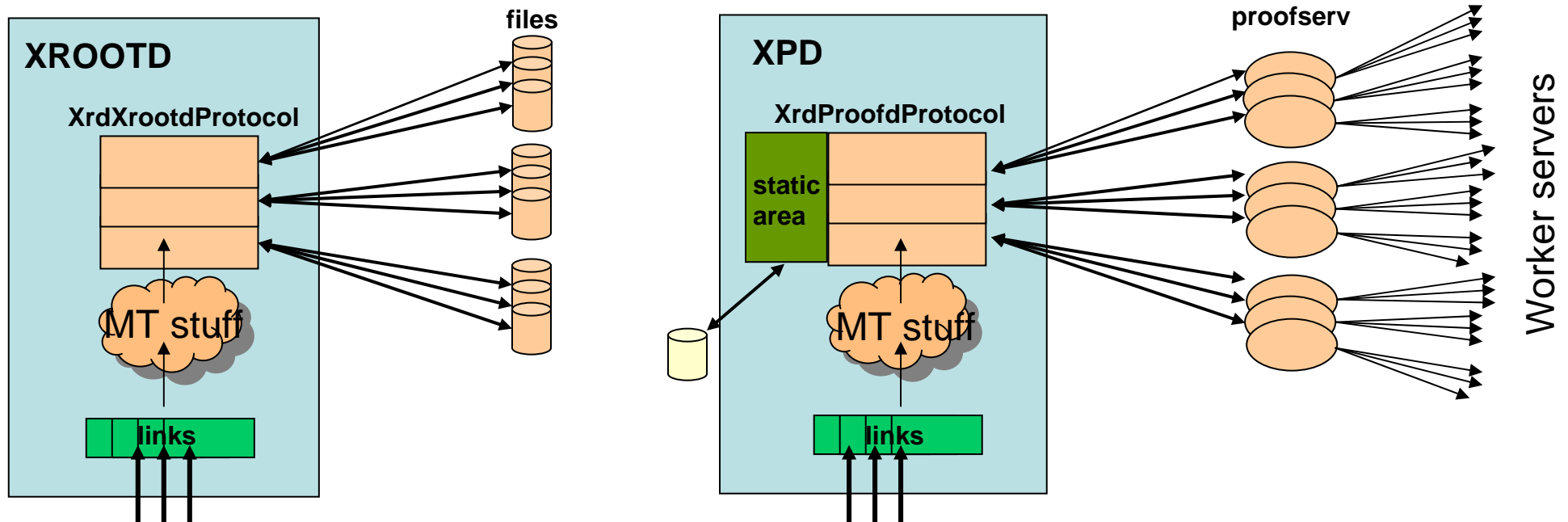
- Introduce a coordinator on the server side able to keep the **proofserv** processes alive when the client goes away
- **proofserv** must be separate processes
  - occasional crashes should not touch other clients

- Main component of **xrootd** (*xrd*) handles networking, security, generic work dispatching

A. Hanushevsky's talk

- Good candidate: already in ROOT, thoroughly tested

# XPD
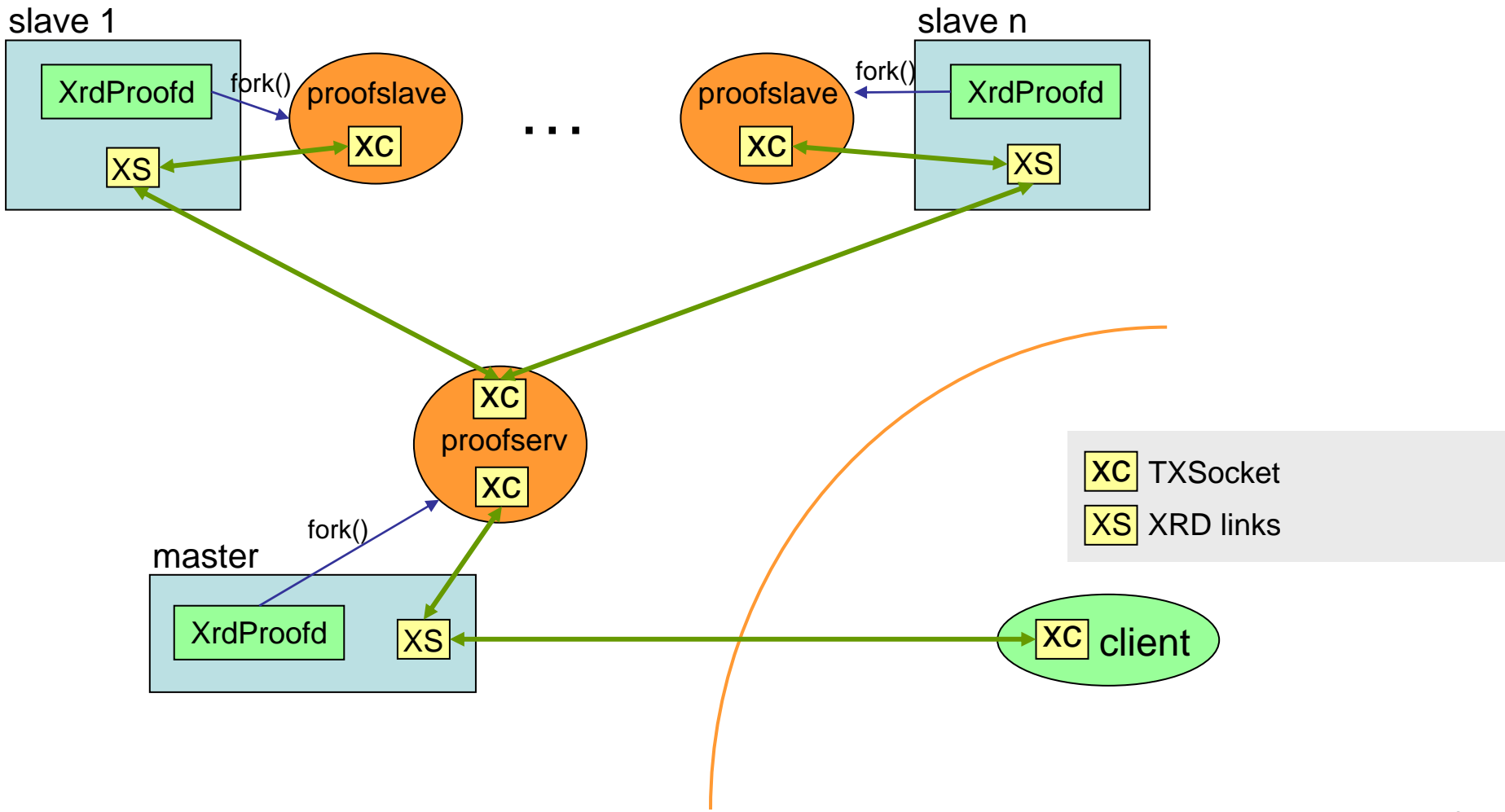
- Prototype based on XROOTD



- **XrdProofdProtocol: client gateway to proofserv**
- static area for all client information and its activities

# XPD communication layer



slave 1

XrdProofd — fork() → proofslave **XC**

XS

. . .

slave n

proofslave **XC** ← fork() — XrdProofd

XS

**XC** proofserv **XC**

master

XrdProofd — fork() → XS

**XC** TXSocket

**XS** XRD links

**XC** client

# XPD remarks

- First impressions for the prototype positive
    - Disconnect / reconnect can be implemented naturally
    - Asynchronous reading allows to setup a control interrupt network independent of OOB
    - Cleaner security system

- Main issue is to understand the impact on performances of the additional layer

# Demo

- CERN testbed:
  - 32 dual Pentium III 800 MHz / 512 MB memory
  - 100 MBit Ethernet
  - 600 GB hard disk

- CMS simulation data
  - 20 files, 1.4 GB total

# Demo: CMS selector

```cpp
Bool_t TCjets::Process (Long64_t entry)
{
   b_CaloJetCollection_obj_->GetEntry(entry);
   if (CaloJetCollection_obj_ < 11) return kFALSE;
   fChain->GetEntry(entry);

   for (Int_t i = 0; i < CaloTowerCollection_obj_; i++) {

      Float_t e   = CaloTowerCollection_obj_e[i];
      Float_t et  = CaloTowerCollection_obj_et[i];
      Float_t eta = CaloTowerCollection_obj_eta[i];
      fHe  ->Fill(e);
      fHet ->Fill(et);
      fHeta->Fill(eta);
      fHete->Fill(e,et);
   }
   return kTRUE;
}
```

# Summary

- Task force working on PROOF work-package since ROOT restructuring
- Several new features introduced
    - "interactive batch" mode
    - GUI controller
- Many other expected during coming months
- Next step is to get the experiments start using PROOF
- Test-bed being setup at CERN
    - experiments welcome to send data samples and complex selectors to better tune the system