



# Software challenges for the next decade

DESY: October 29 2007

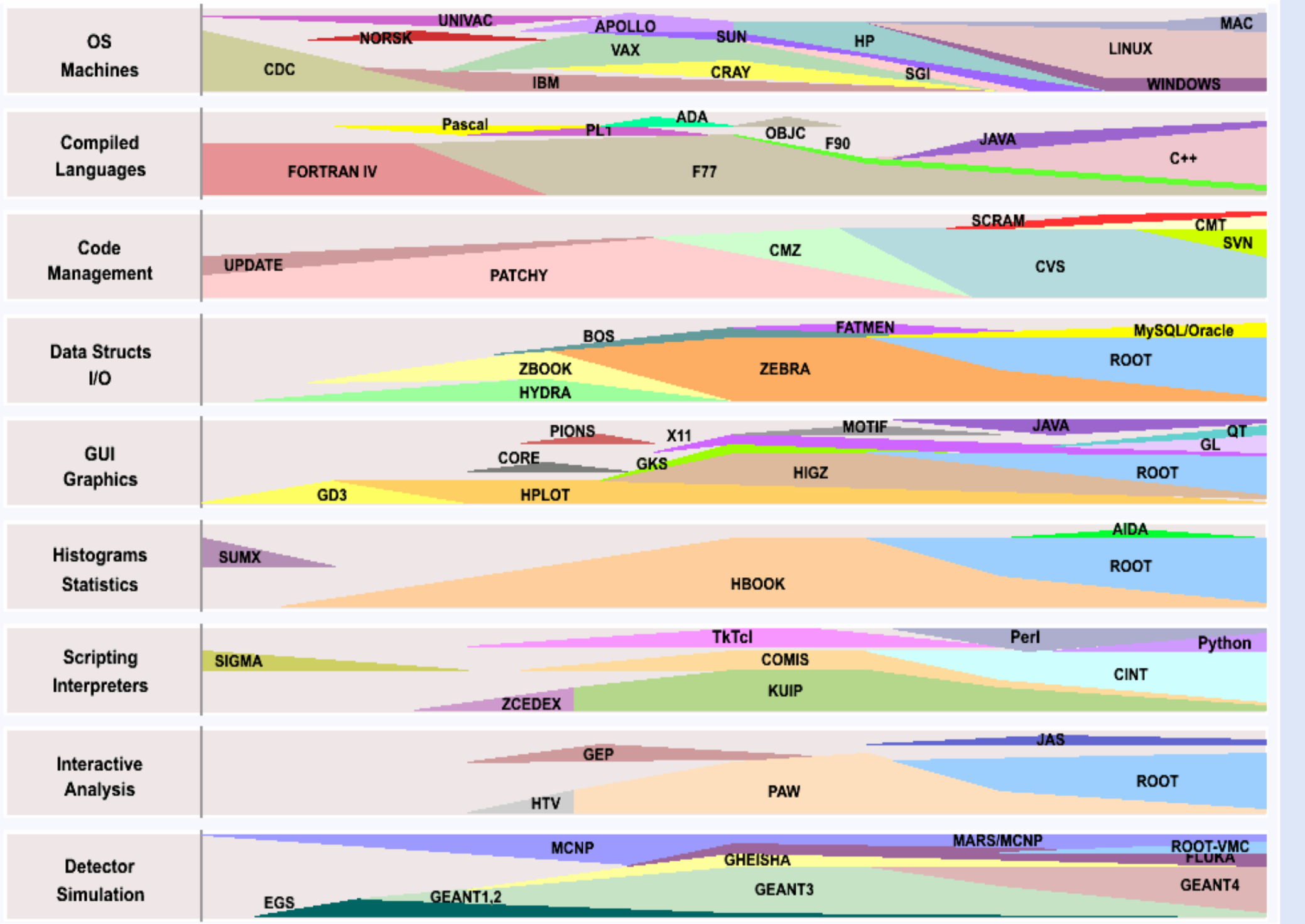
René Brun  
CERN



# My mini crystal ball



- Concentrate on HEP software only
- I am biased by the development of large software frameworks or libraries and their use in large experiments.
- Because of inertia or time to develop, there are things easy to predict.
- Technology can come with good or/and bad surprises. **See weather Forecast**



1970 1975 1980 1985 1990 1995 2000 2005 2010



# Time to develop: Main message



- Big latency in the process to collect requirements, make a design, implementation, release and effective use.
- It takes between 5 and 10 years to develop a large system like **ROOT** or **Geant4**.
- It takes a few years for users to get familiar with the systems.
- It takes a few years for systems to become de facto standards.
- For example, LHC expts are discovering now the advantages of split-mode I/O designed in 1995 and promoted in 1997.
- This trend will likely continue with the large collaborations.
- This has of course positive and negative aspects. Users like stability, BUT competition is mandatory



# The crystal ball in 1987

- **Fortran 90X** seems the obvious way to go
- OSI protocols to replace TCP/IP
- Processors: Vector or MPP machines
- PAW, Geant3, Bos, Zebra: Adapt them to F90X
- Methodology trend: **Entity Relationship Model**
- Parallelism: **vectorization or MPP (SIMD and MIMD)**
- BUT hard to anticipate that
  - The WEB will come less than 3 years later
  - The 1993/1994 revolution for languages and projects
  - The rapid grow in CPU power starting in 1994 (Pentium)



# Situation in 1997

- LHC projects moving to **C++**
  - **Several projects proposing to use Java**
  - Huge effort with OODBMS (ie Objectivity)
  - Investigate Commercial tools for data analysis
  - ROOT development not encouraged
  - Vast majority of users very sceptic.
- 
- RAM <256 MB
  - Program Size < 32 MB
  - <500 KLOcs
  - libs < 10
  - static linking
  - HSM: tape->Disk pool <1 TByte
  - Network 2MB/s



# The crystal ball in 1997

- **C++** now, **Java** in 2000
- Future is OODBMS (ie Objectivity)
- Central Event store accessed through the net
- Commercial tools for data analysis
  
- But fortunately a few people did not believe in this direction :☺
- First signs of problems with Babar
- FNAL RUN2 votes for ROOT in 1998
- GRID: an unknown word in 1997 :☺

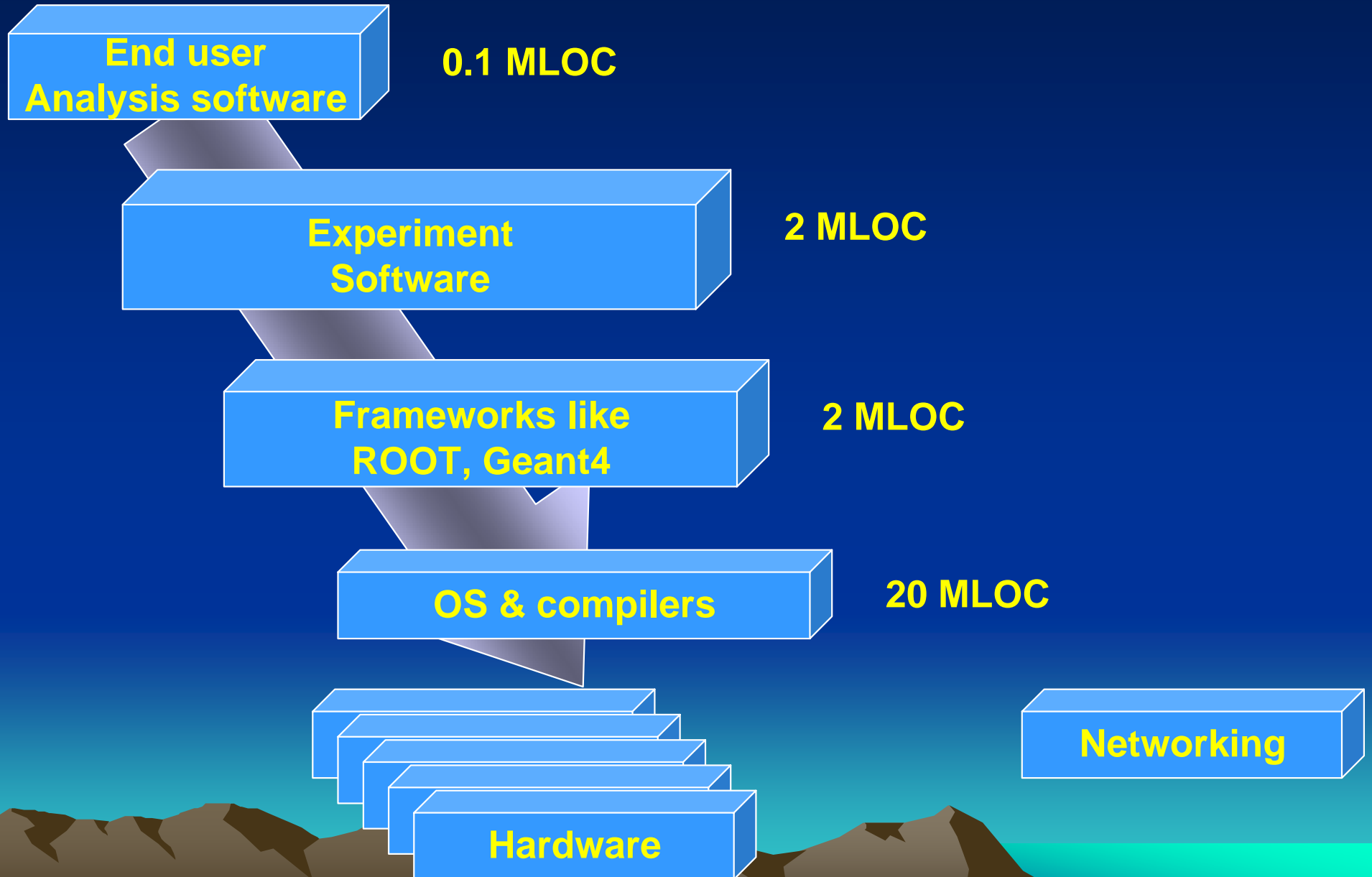


# Situation in 2007

- It took far more time than expected to move people to C++ and the new frameworks.
- ROOT de facto standard for I/O and interactive analysis.
- The GRID:
- Experiment frameworks are monsters



# Software Hierarchy





# Challenge 0

## Usability: Making things SIMPLER

- Guru view vs user view
- A normal user has to learn too many things before being able to do something useful.
- LHC frameworks becoming monsters
- fighting to work on 64 bits with <2 GBytes
- take for ever to start because too much code linked (shared libs with too many dependencies)
- fat classes vs too many classes
- It takes time to restructure large systems to take advantage of plug-in managers.

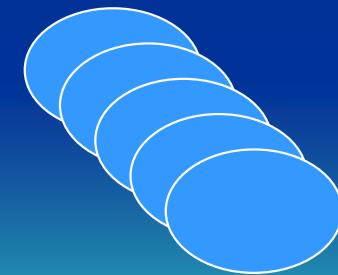
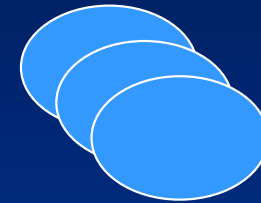
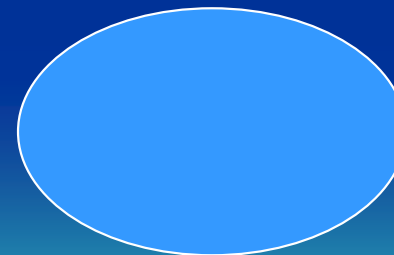
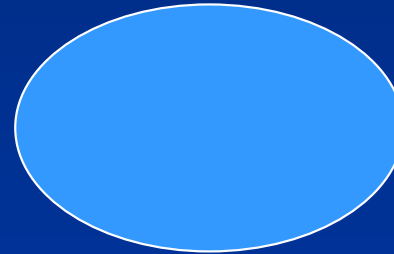
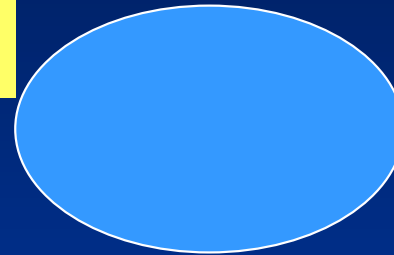


# Challenge 1

## Problem decomposition

Will have to deal with many shared libs

Only a small fraction of code used





# Some Facts

10 shared libs

200 classes

**ROOT  
In  
1995**

PAW model

100 shared libs

1800 classes

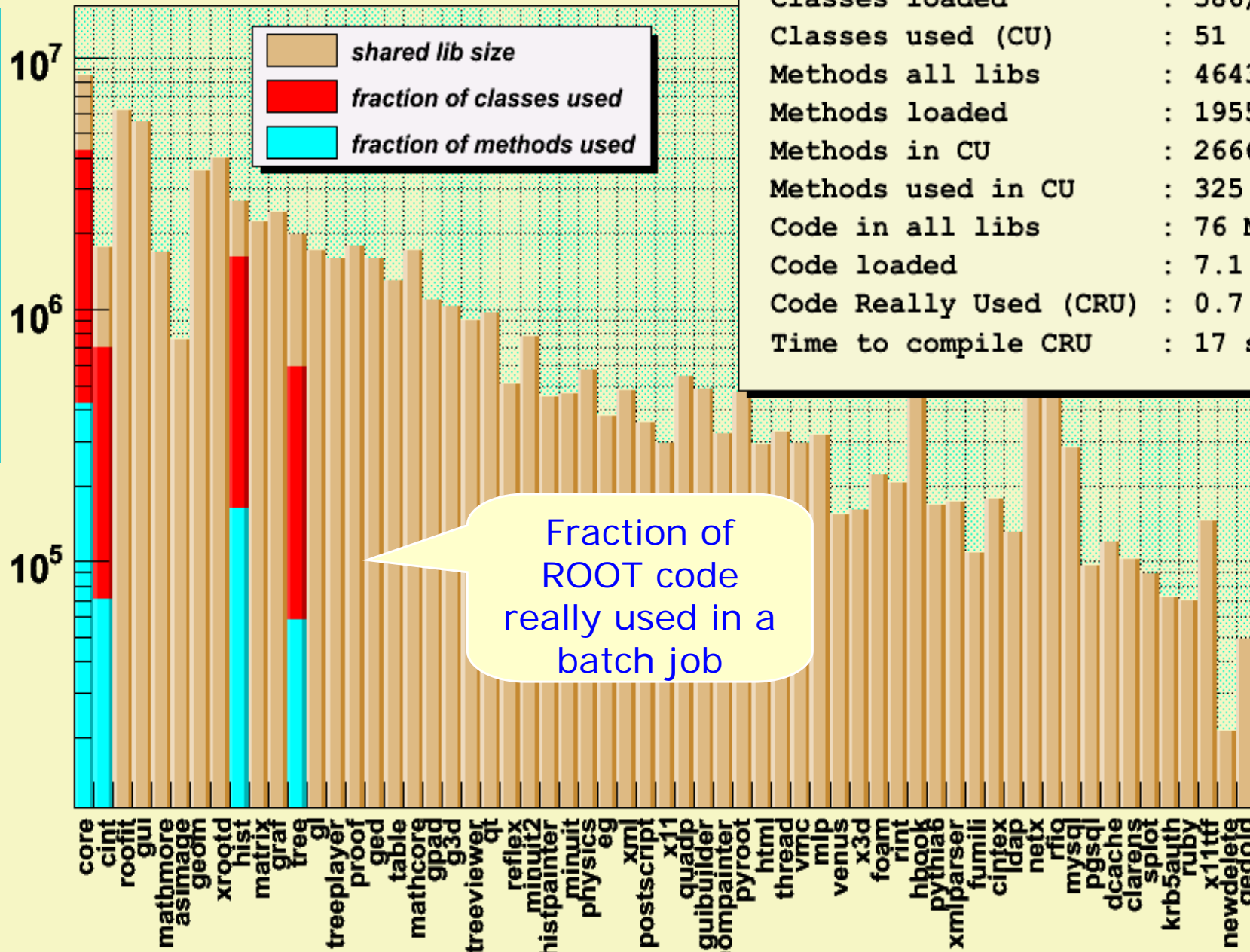
**ROOT  
In  
2007**

Plug-in manager



# code used in a batch use case

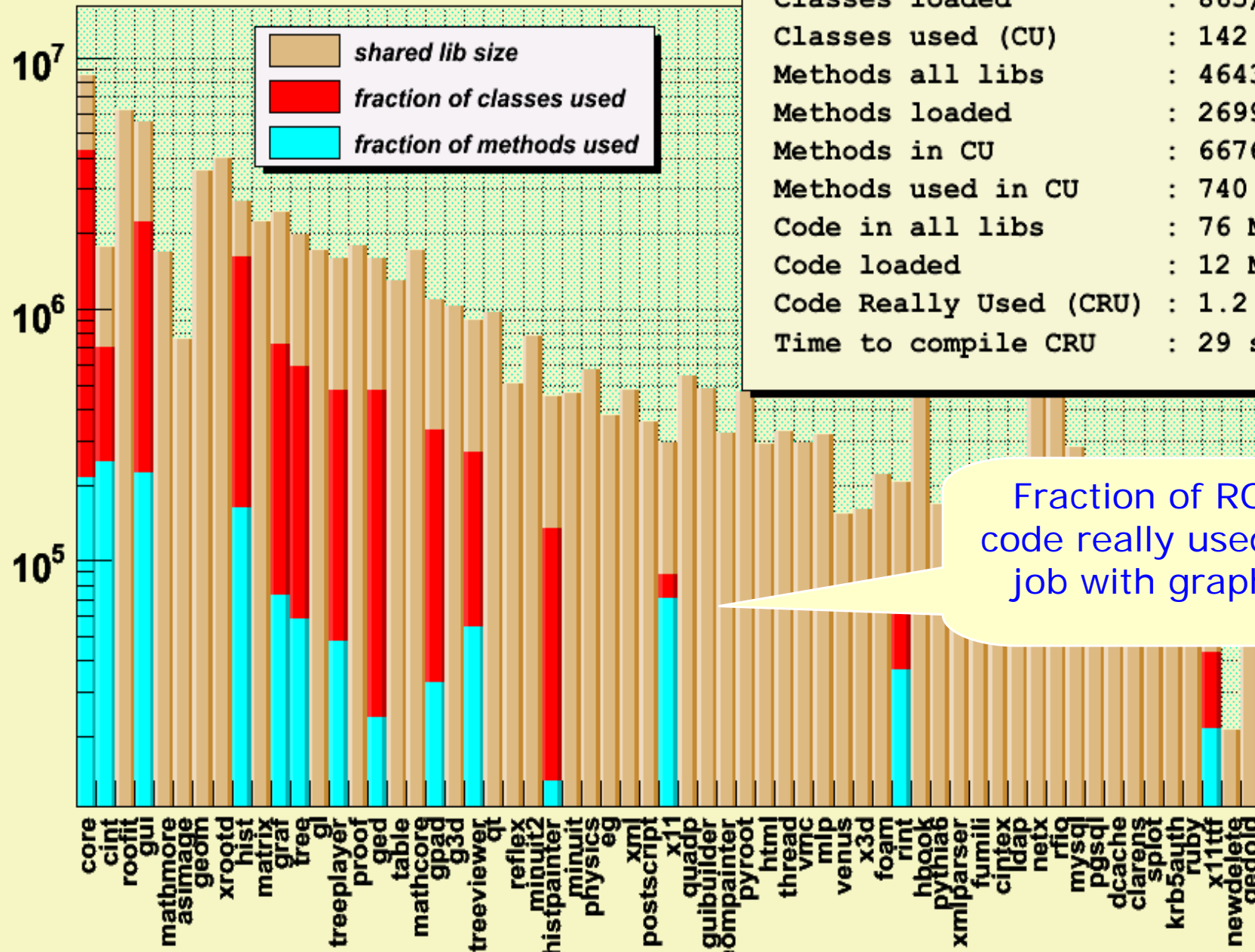
Shared lib size in bytes



Libs used	: 4/86
Classes loaded	: 586/1459
Classes used (CU)	: 51
Methods all libs	: 46438
Methods loaded	: 19550
Methods in CU	: 2666
Methods used in CU	: 325
Code in all libs	: 76 Mb
Code loaded	: 7.1 Mb
Code Really Used (CRU)	: 0.7 Mb
Time to compile CRU	: 17 s



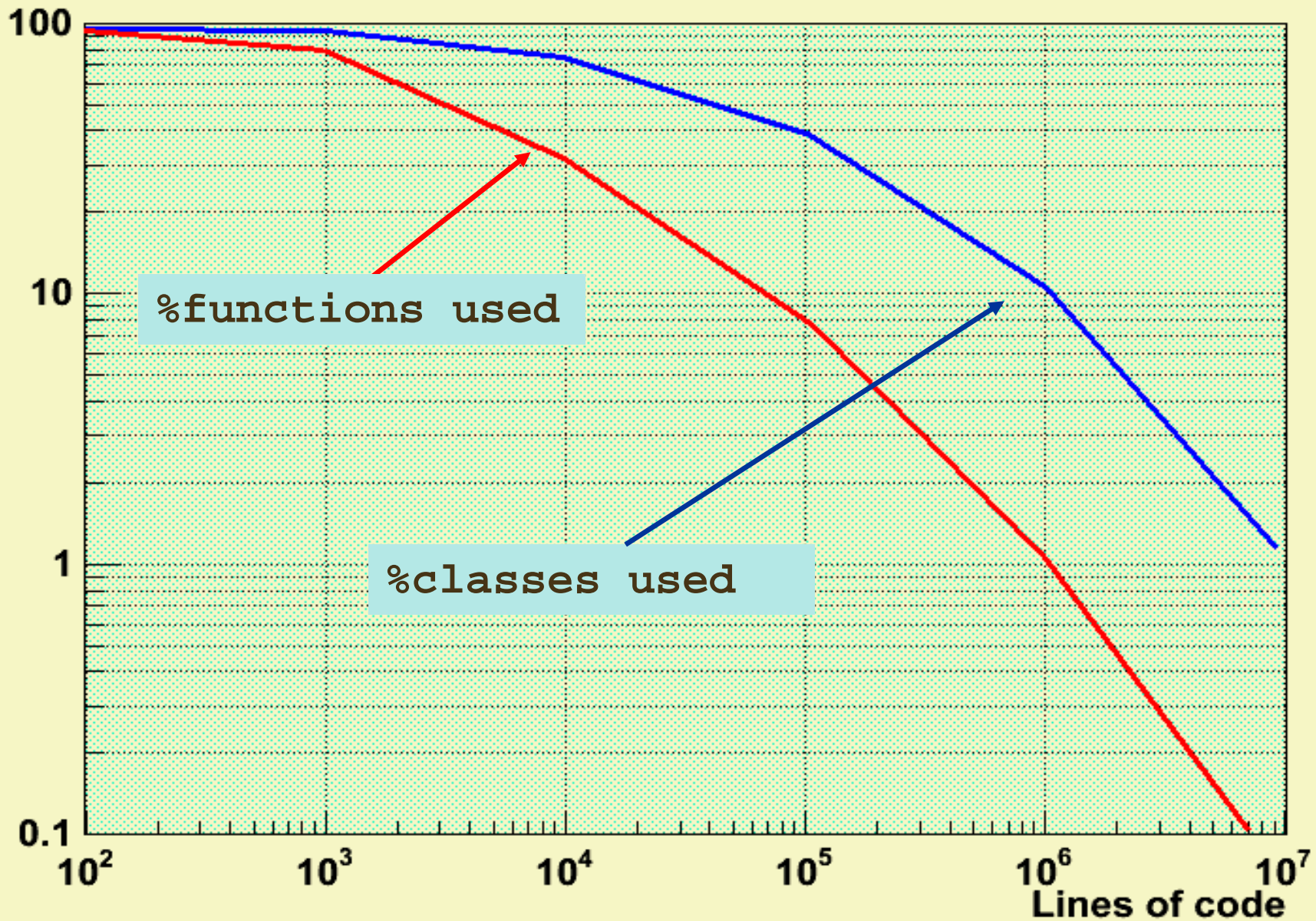
# code used in a graphics use case



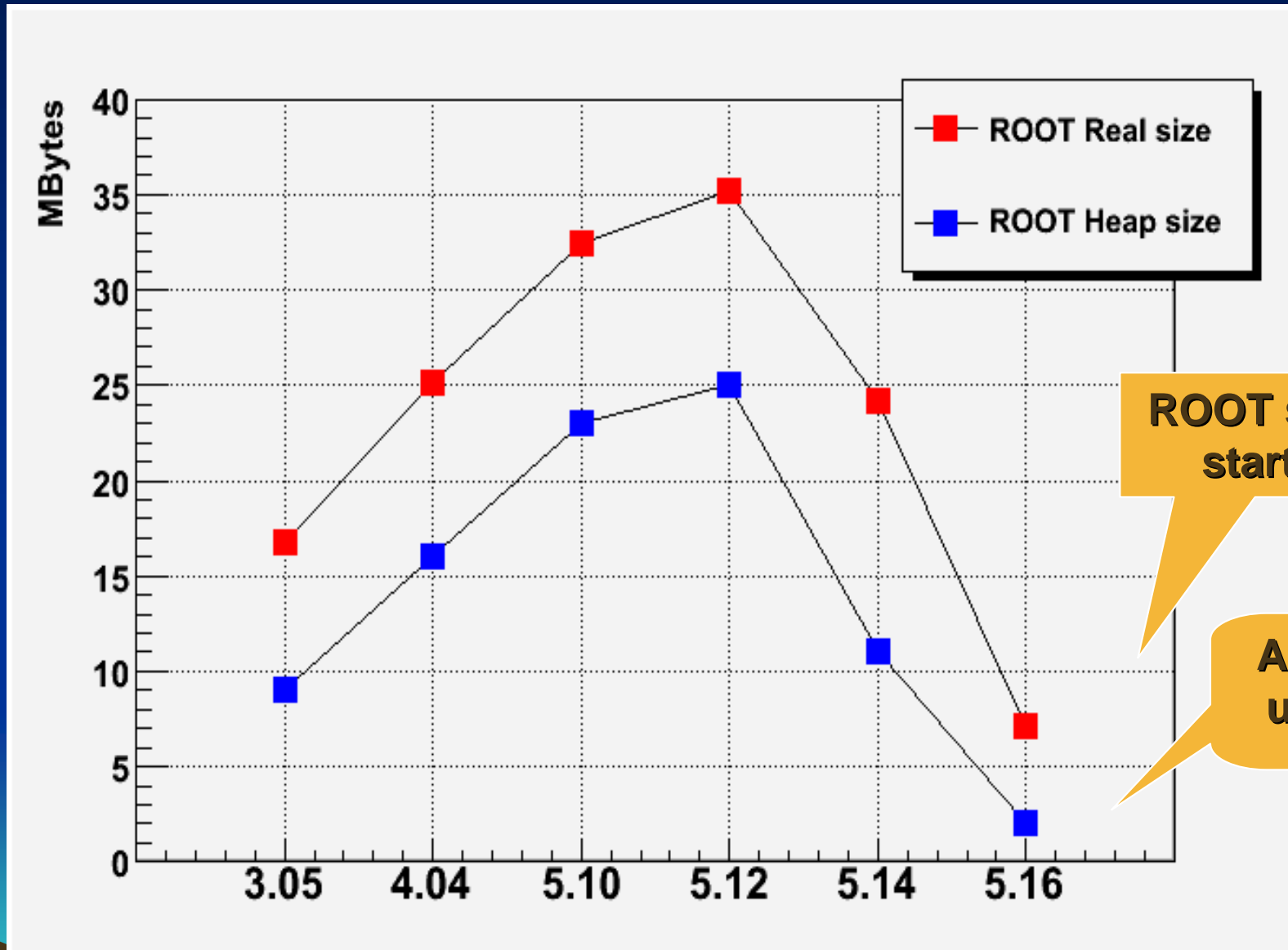
Libs used	: 14/86
Classes loaded	: 865/1459
Classes used (CU)	: 142
Methods all libs	: 46438
Methods loaded	: 26996
Methods in CU	: 6676
Methods used in CU	: 740
Code in all libs	: 76 Mb
Code loaded	: 12 Mb
Code Really Used (CRU)	: 1.2 Mb
Time to compile CRU	: 29 s

Fraction of ROOT code really used in a job with graphics

## Per cent of code used



# Large Heap Size Reduction



**ROOT size at start-up**

**Also speed-up start-up time**



# Challenge 2

## Hardware will force parallelism

- Multi-Core (2-8)
- Many-Core (32-256)
- Mixture CPU + GPU-like
- Virtualization
- May be a new technology?
- Parallelism: a must



# Challenge 3

## Design for Parallelism

- The GRID is a parallel engine. However you will not use the GRID software on your 32-core laptop.
- Minimize globals and make tasks as independent as possible.
- Be thread-safe and better thread-capable
- Think Top->Down and Bottom->Up



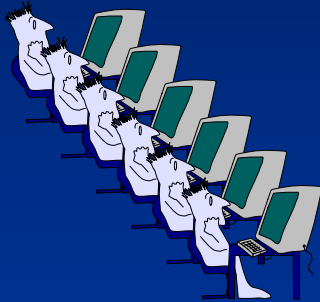


# Parallelism: Where ?



**Multi-Core CPU laptop/desktop**

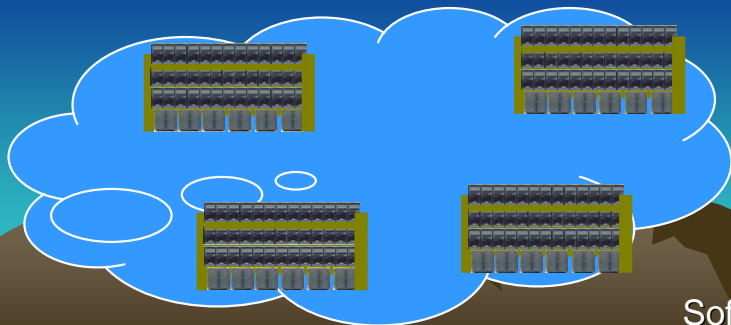
**2(2007) 32(2012?)**



**Network of desktops**



**Local Cluster  
with multi-core CPUs**



**GRID(s)**



# Challenge 4

## Design for Client-Server

- The majority of today's applications are client-server (xrootd, Dcache, sql, etc).
- This trend will increase.
- Be able to stream objects or objects collections.
- Server logic robust against client changes.
- Server able to execute dynamic plug-ins.
- Must be robust against client or network crash



# Challenge 5

## Sophisticated Plug-in Managers

- When using a large software base distributed in hundred of shared libs, it is essential to discover automatically where to find a class.
- The interpreters must be able to autoloading the corresponding libraries



## Challenge 6

# The Language Reflexion System

- Develop a robust dictionary system that can be migrated smoothly to the reflexion system to be introduced in C++ in a few years.
- Meanwhile reduce the size of dictionaries by doing more things at run time.
- Replace generated code by objects stored in ROOT files.
- Direct calls to compiled code from the interpreter instead of function stubs. This is compiler dependent (mangling/de-mangling symbols).



# Problem with Dictionaries



Today `cint/reflex` dictionaries are machine dependent. They represent a very substantial fraction of the total code

We are now working to reduce this size by at least a factor 3!

	<code>*.o</code>	<code>G_*.o</code>	Dict %
<code>mathcore</code>	2674520	2509880	93.8%
<code>mathmore</code>	598040	451520	75.5%
<code>base</code>	6920485	4975700	71.8%
<code>physics</code>	786700	558412	71.0%
<code>treeplayer</code>	2142848	1495320	69.8%
<code>geom</code>	4685652	3096172	66.1%
<code>tree</code>	2696032	1592332	59.1%
<code>g3d</code>	1555196	908176	58.4%
<code>geompainter</code>	339612	196588	57.9%
<code>graf</code>	2945432	1610356	54.7%
<code>matrix</code>	3756632	2020388	53.8%
<code>meta</code>	1775888	909036	51.2%
<code>hist</code>	3765540	1914012	50.8%
<code>gl</code>	2313720	1126580	48.7%
<code>gpad</code>	1871020	781792	41.8%
<code>histpainter</code>	538212	204192	37.9%
<code>minuit</code>	581724	196496	33.8%



# Challenge 7

## Opportunistic Use of Interpreters

- Use interpreted code only for:
  - External and thin layer (task organizer)
  - Slots execution in GUI signal/slots
  - Dynamic GUI builder in programs like event displays.
- Instead optimize the compiler/linker interface (eg **TACLIC**) to have
  - Very fast compilation/linking when performance is not an issue
  - Slower compilation but faster execution for the key algorithms
- ie use ONE single language for 99% of your code and the interpreter of your choice for the layer between shell programming and program orchestration.



# Challenge 8

## LAN and WAN I/O caches

- Must be able to work very efficiently across fat pipes but with high latencies.
- Must be able to cache portions or full files on a local cache.
- This requires changes in data servers (Castor, Dcache, xrootd). These tools will have to **interoperate**.
- The ROOT file info must be given to these systems for optimum performance. See TTreeCache improvements.



# Disk cache improvements with high latency networks



- The file is on a CERN machine connected to the CERN LAN at at 100MB/s.
- The client **A** is on the same machine as the file (local read)
- The client **F** is connected via ADSL with a bandwidth of 8Mbits/s and a latency of 70 milliseconds (Mac Intel Coreduo 2Ghz).
- The client **G** is connected via a 10Gbits/s to a CERN machine via Caltech latency 240 ms.
- The times reported in the table are realtime seconds

One query to  
a 280 MB Tree  
I/O = 16.6 MB

client	latency(ms)	cache size=0	cache size=64KB	cache size=10MB
A	0.0	3.4	3.4	3.4
F	72.0	743.7	48.3	28.0
G	240.0	>1800s	125.4s	9.9s

We expect to reach 4.5 s



# I/O: More

- -Efficient access via a LAN AND WAN
- -Caching
- -Better schema evolution
- -More support for complex event models
- -zip/unzip improvements (separate threads)
- -More work with SQL data bases



# Challenge 9

## Code Performance

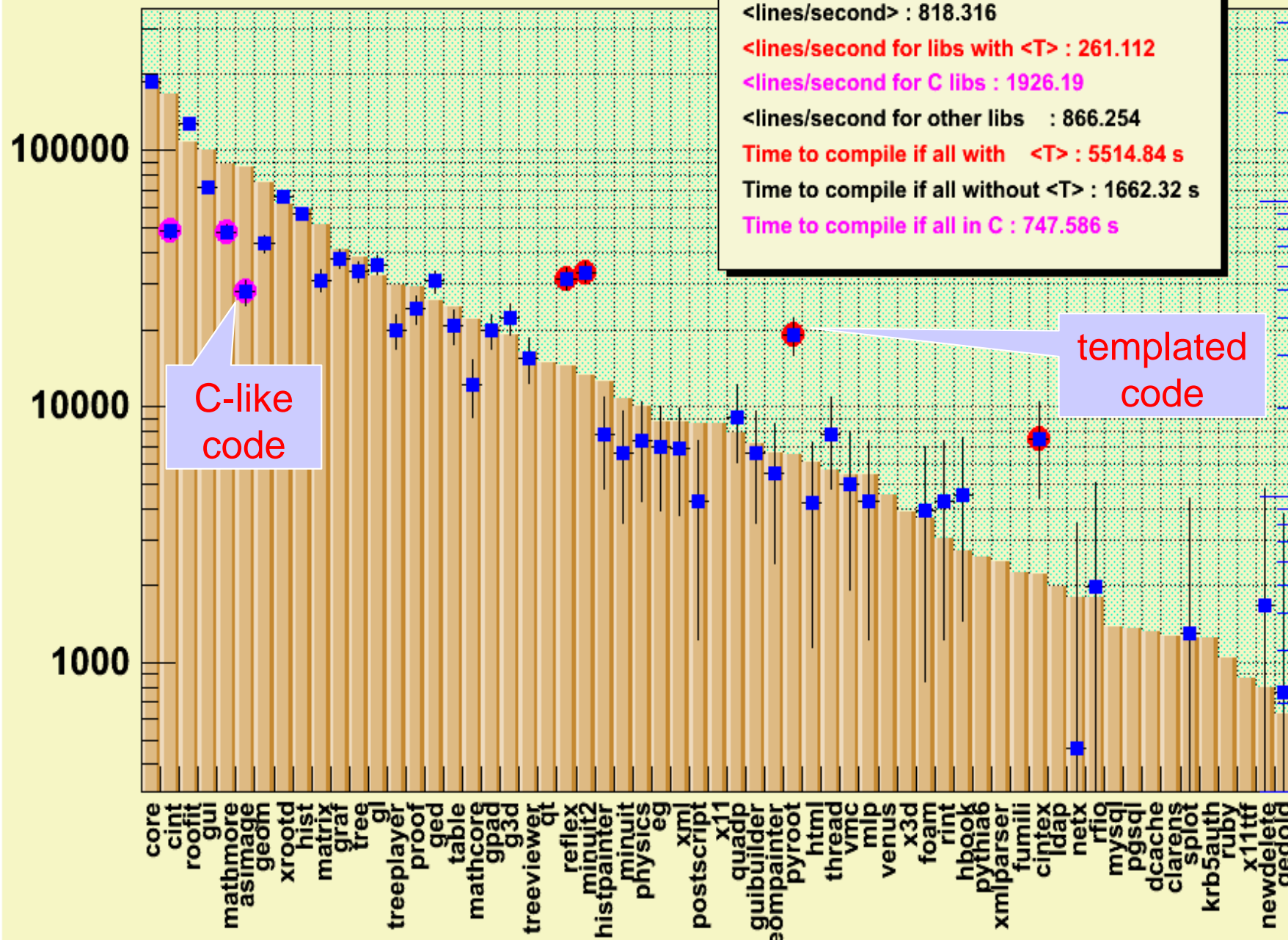
- HEP code does not exploit hardware (see [S.Jarp](#) talk at CHEP)
- Large data structures spread over >100 Megabytes
- templated code pitfall
  - STL code duplication
  - good perf improvement when testing with a toy.
  - disaster when running real programs.
- `std::string` passed by value
- abuse of `new/delete` for small objects or stack objects
- linear searches vs hash tables or binary search
- abuse of inheritance hierarchy
- code with no vectors -> do not use the pipeline



# Compilation Time



## Lines of code per library



Time to compile in seconds (VC++7.1)

C-like code

templated code



# LHC software



	Alice	Atlas	CMS	ROOT
number of lines in header files	102282	698208	104923	153775
classes total	1815	8910	???	1500
classes in dict	1669	>4120 2140	835	1422
lines in dict	479849	455705	103057	698000
classes c++ lines	577882	1524866	277923	857390
total lines Classes+dict	1057731	???	380980	1553390
total f77 lines	736751	928574	???	3000
directories	540	19522	<500	958
comp time	25'	750'	90'	30'
lines compiled/s	1196	50 (70)	71	863



# Challenge 10

## Towards Task-oriented programming

The image displays three sequential screenshots of the ROOT Object Browser interface, illustrating the transition from a flat file system to a task-oriented structure.

- Left Screenshot:** Shows a standard file browser view of the file system. A red arrow points to the directory tree with the label **OS files**.
- Middle Screenshot:** Shows a hierarchical tree structure of objects. A yellow box labeled **Browsing** is positioned above the tree. A red arrow points to the tree structure with the label **Data hierarchy**.
- Right Screenshot:** Shows a task-oriented structure where objects are organized into tasks. A red arrow points to the task-based hierarchy with the label **Dynamic tasks**.

At the bottom of each screenshot, the filter is set to "All Files (\*.\*)".



# Challenge 11

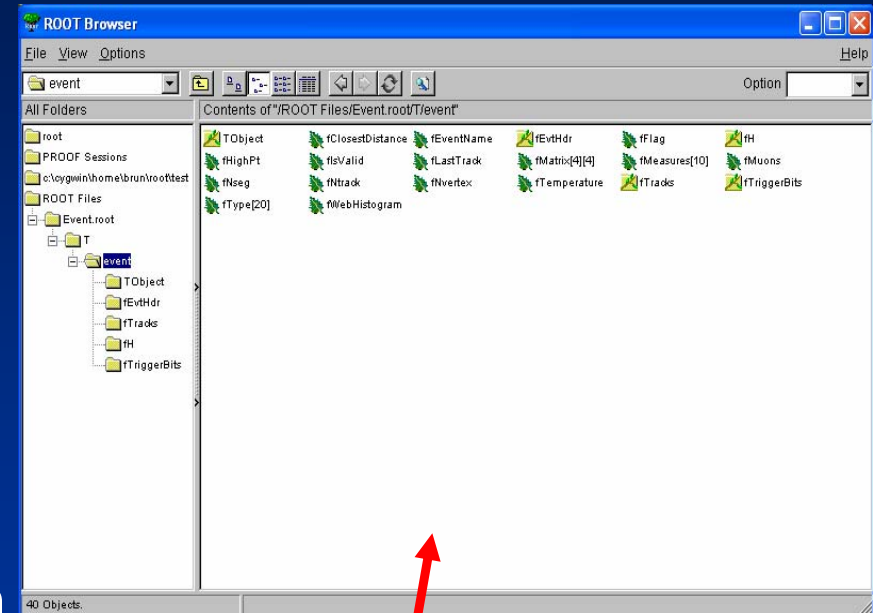
## Customizable and Dynamic GUIs

- From a standard browser (eg ROOT **TBrowser**) one must be able to include user-defined GUIs.
- The GUIs should not require any pre-processor.
- They can be **executed/loaded/modified** in the same session.



# Browser Improvements

- The browser (**TBrowser** and derivatives) is an essential component (from beginners to advanced applications).
- It is currently restricted to the browsing of ROOT files or Trees.
- We are extending **TBrowser** such that it could be the central interface and the manager for any GUI application (editors, web browsers, event displays, etc).



Old/current browser



# Hist Browser + stdin/stdout



ROOT Browser

Framework File Edit View Options Inspect Classes Help

Files Classes

Editor HTML Canvas

Test random numbers

h1f	
Entries	10000
Mean	3.646
RMS	1.84

Command

Command (local):

```
root [1] .ls
TFile**      fillrandom.root
TFile*       fillrandom.root
KEY: TFormula form1;1 abs(sin(x)/x)
KEY: TF1      sqrt;1      x*gaus(0)+[3]*form1
KEY: TH1F     h1f;1      Test random numbers
```



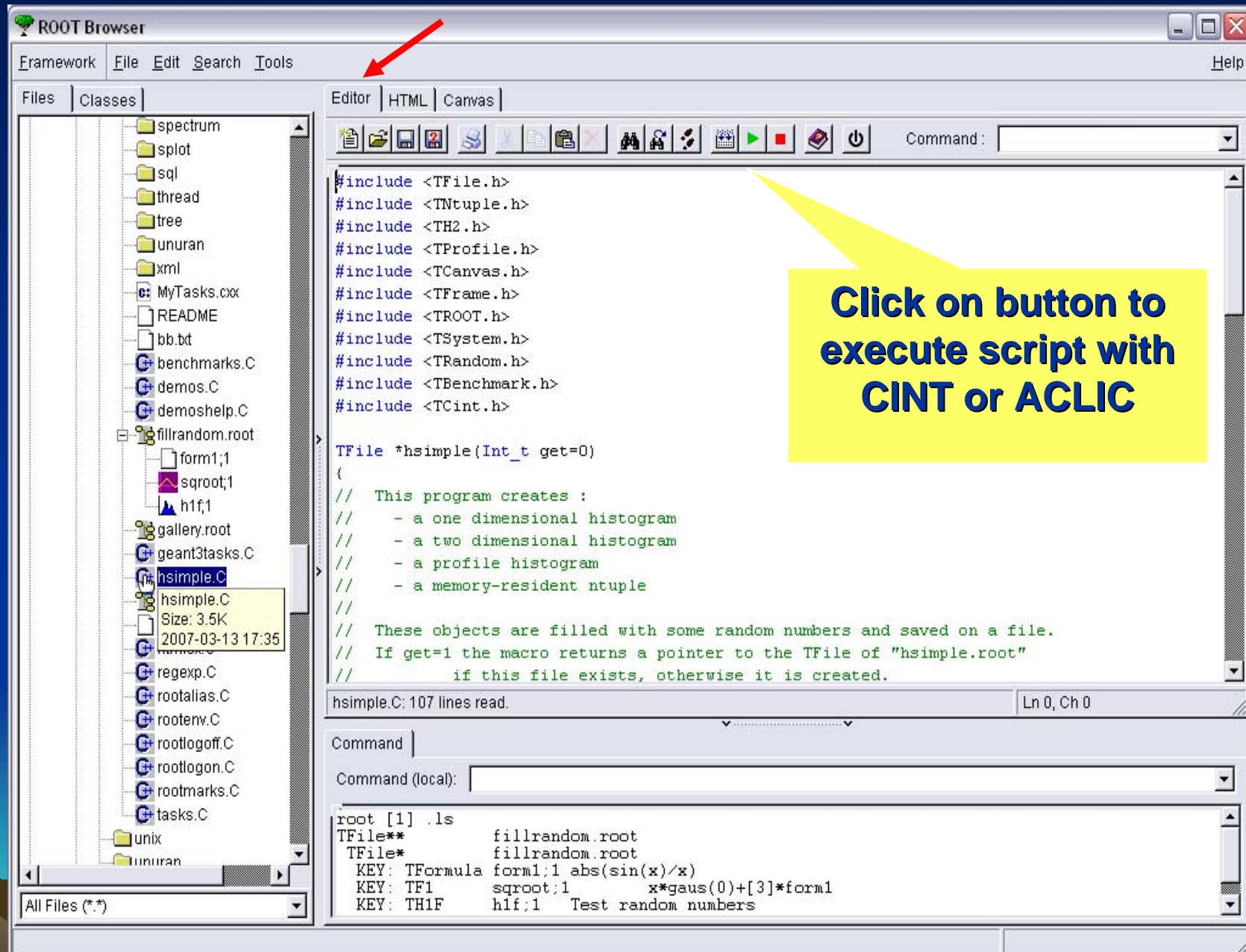
# TGhtml web browser plug-in



The screenshot shows the ROOT Browser interface. A red arrow points to the 'HTML' tab in the top menu. The address bar contains the URL `http://root.cern.ch/files/`, which is circled in red. Below the address bar is a table titled 'Index of /files' with columns for Name, Last modified, Size, and Description. The table lists various files and directories, including 'Parent Directory', 'hidisplay.root', 'aleph.root', 'aleph.all.C', 'alice.root', 'alice.all.C', 'alice.itsd.C', 'alice.itsv.C', 'ams.C', 'ams.all.C', 'archive.zip', and 'atlas.root'. A yellow callout box points to the 'aleph.root' entry with the text 'You can browse a root file'. Another yellow callout box points to the 'ams.C' entry with the text 'You can execute a script'. At the bottom, the Command window shows the output of a `ls` command:

```
root [1] .ls
TFile**      fillrandom.root
TFile*       fillrandom.root
KEY: TFormula form1:1 abs(sin(x)/x)
KEY: TF1      sqrt:1      x*gaus(0)+[3]*form1
KEY: TH1F     hf:1       Test random numbers
```

# Macro Manager/Editor plug-in



The screenshot shows the ROOT Browser window with the Macro Manager/Editor plug-in. The left pane displays a file tree with 'hsimple.C' selected. The main editor pane shows the C++ code for 'hsimple.C', which includes various ROOT headers and defines a 'TFile \*hsimple' function. A yellow callout box with the text 'Click on button to execute script with CINT or ACLIC' points to the 'Execute' button (a green play icon) in the toolbar. A red arrow points to the 'Tools' menu in the top menu bar. The bottom pane shows the command prompt output for the 'root [1] .ls' command, listing the contents of the 'fillrandom.root' directory.

```
#include <TFile.h>
#include <TNTuple.h>
#include <TH2.h>
#include <TProfile.h>
#include <TCanvas.h>
#include <TFrame.h>
#include <TROOT.h>
#include <TSystem.h>
#include <TRandom.h>
#include <TBenchmark.h>
#include <TCint.h>

TFile *hsimple(Int_t get=0)
{
  // This program creates :
  //   - a one dimensional histogram
  //   - a two dimensional histogram
  //   - a profile histogram
  //   - a memory-resident ntuple
  //
  // These objects are filled with some random numbers and saved on a file.
  // If get=1 the macro returns a pointer to the TFile of "hsimple.root"
  // if this file exists, otherwise it is created.
}

hsimple.C: 107 lines read. Ln 0, Ch 0

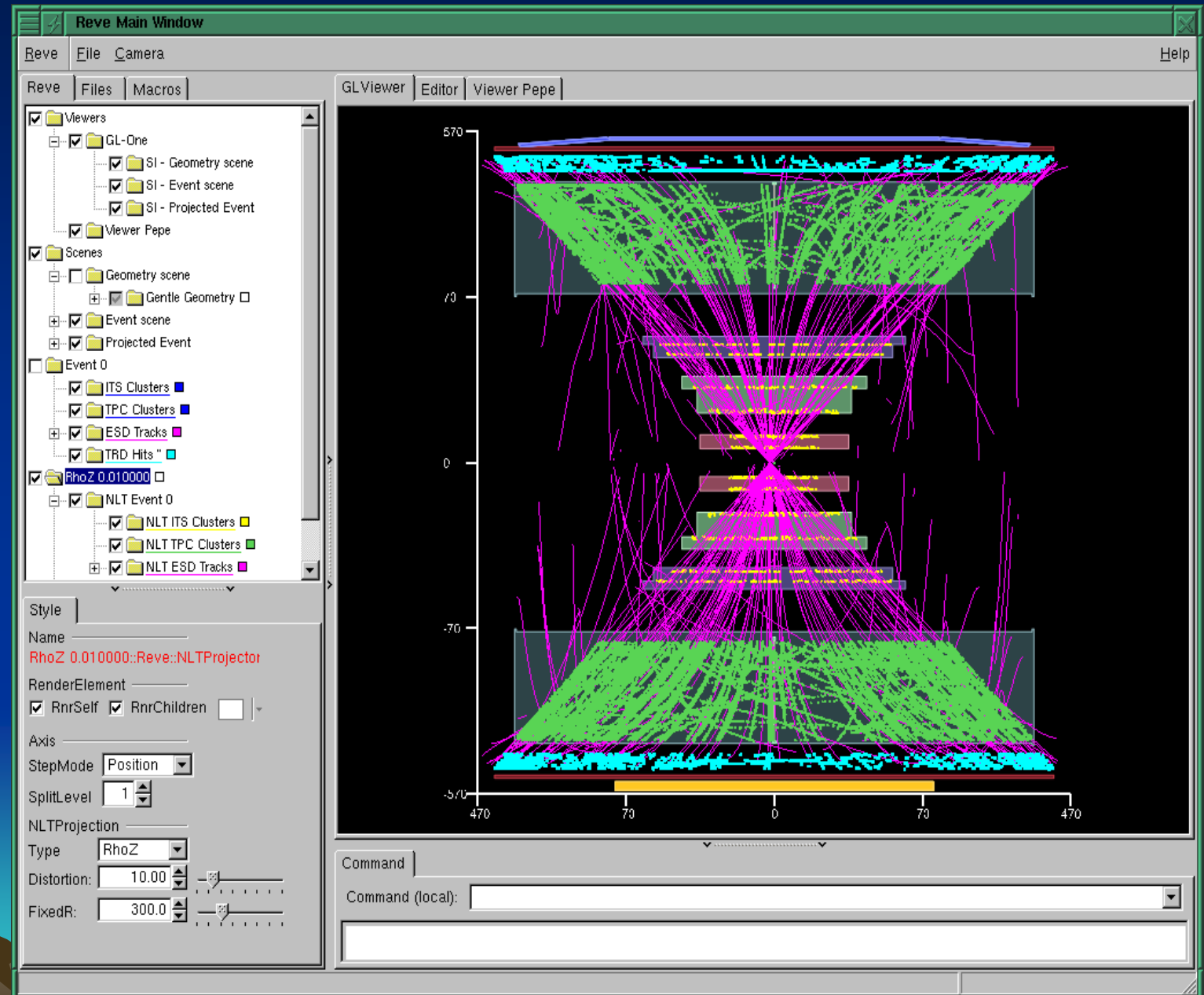
Command
Command (local):
root [1] .ls
TFile**      fillrandom.root
TFile*       fillrandom.root
KEY: TFormula form1:1 abs(sin(x)/x)
KEY: TF1      sqrt:1      x*gaus(0)+[3]*form1
KEY: TH1F     h1f:1      Test random numbers
```



# GL Viewer plug-in



Alice event display prototype using the new browser





# Challenge 12

## Executing Anywhere from Anywhere

- One should be able to start an application from any web browser.
- The local UI and GUI can execute transparently on a remote process.
- The resulting objects are streamed to the local session for fast visualization.
- Prototype in latest ROOT using ssh technology.

```
root > .R lxplus.cern.ch  
lxplus > .x doSomething.C  
lxplus > .R  
root > //edit the local canvas
```



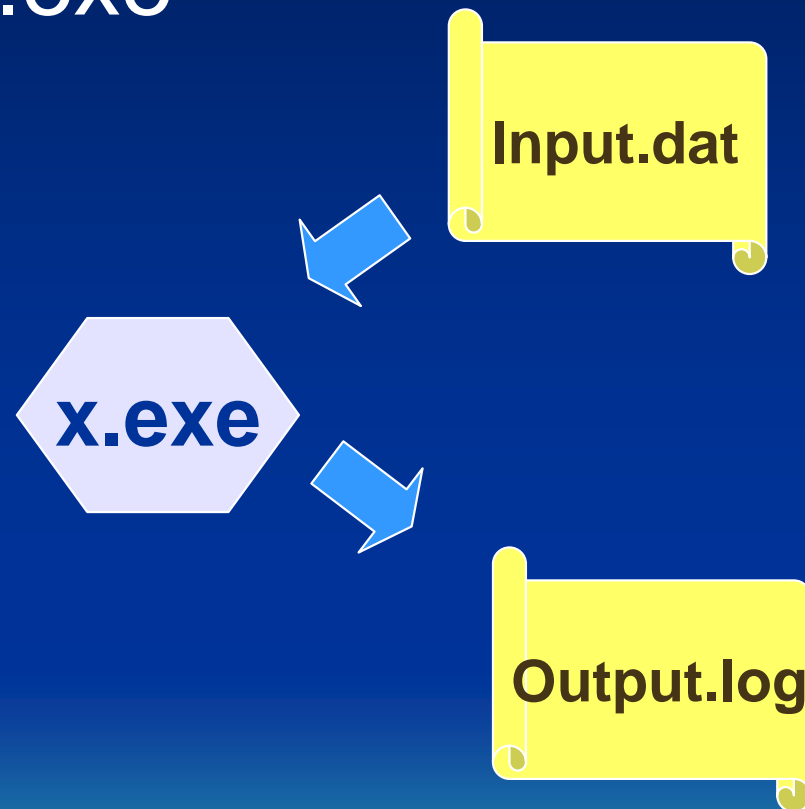
# Challenge 13

## Evolution of the Execution Model

- From stand alone modules
- To shared libs
- To plug-in managers
- To distributed computing
- To distributed and parallel computing

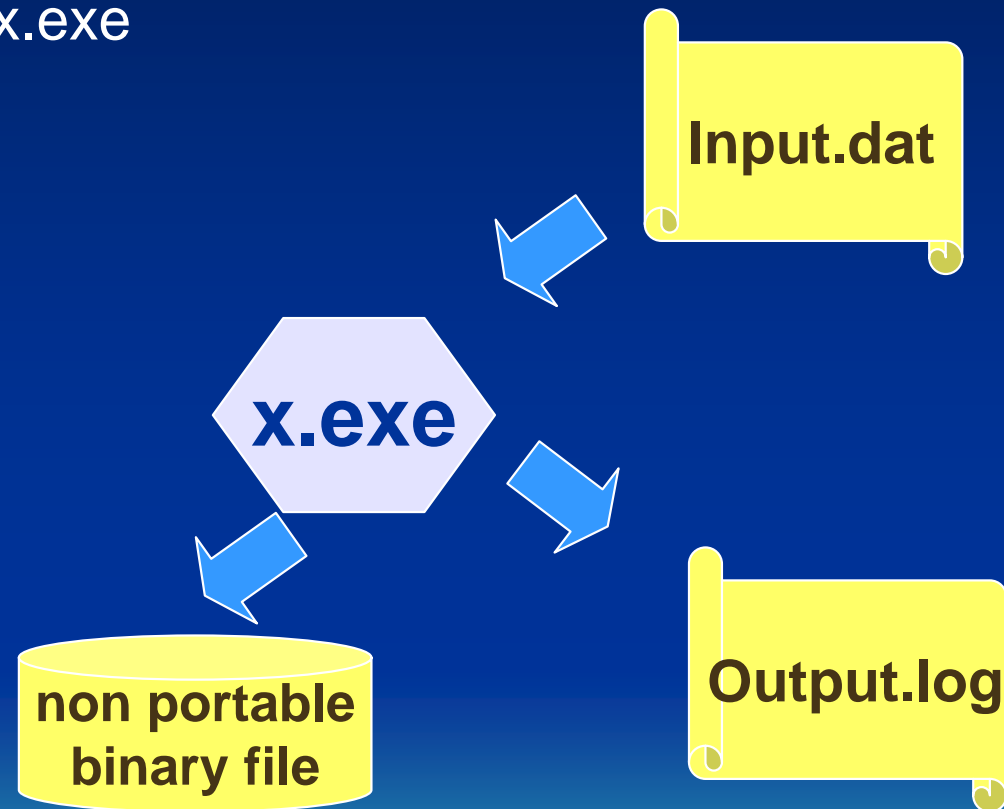
# Executable module in 1967

- `x.f` -> `x.o` -> `x.exe`



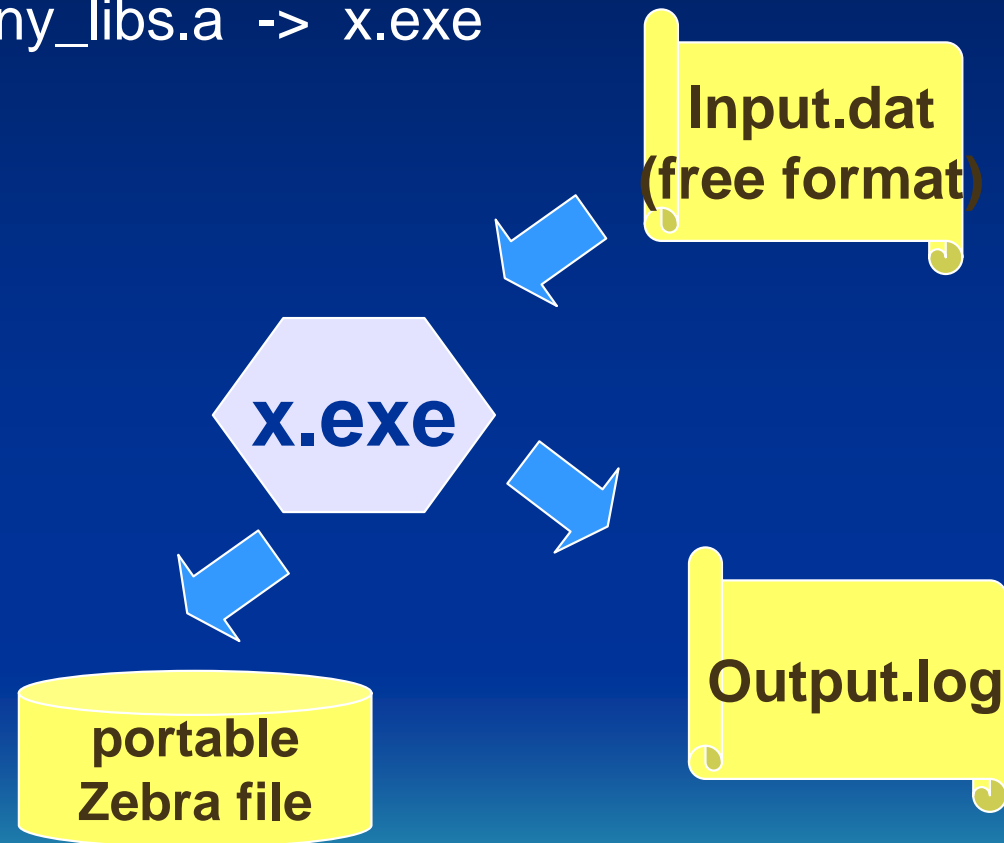
# Executable module in 1977

- `x.f` -> `x.o`
- `x.o + libs.a` -> `x.exe`



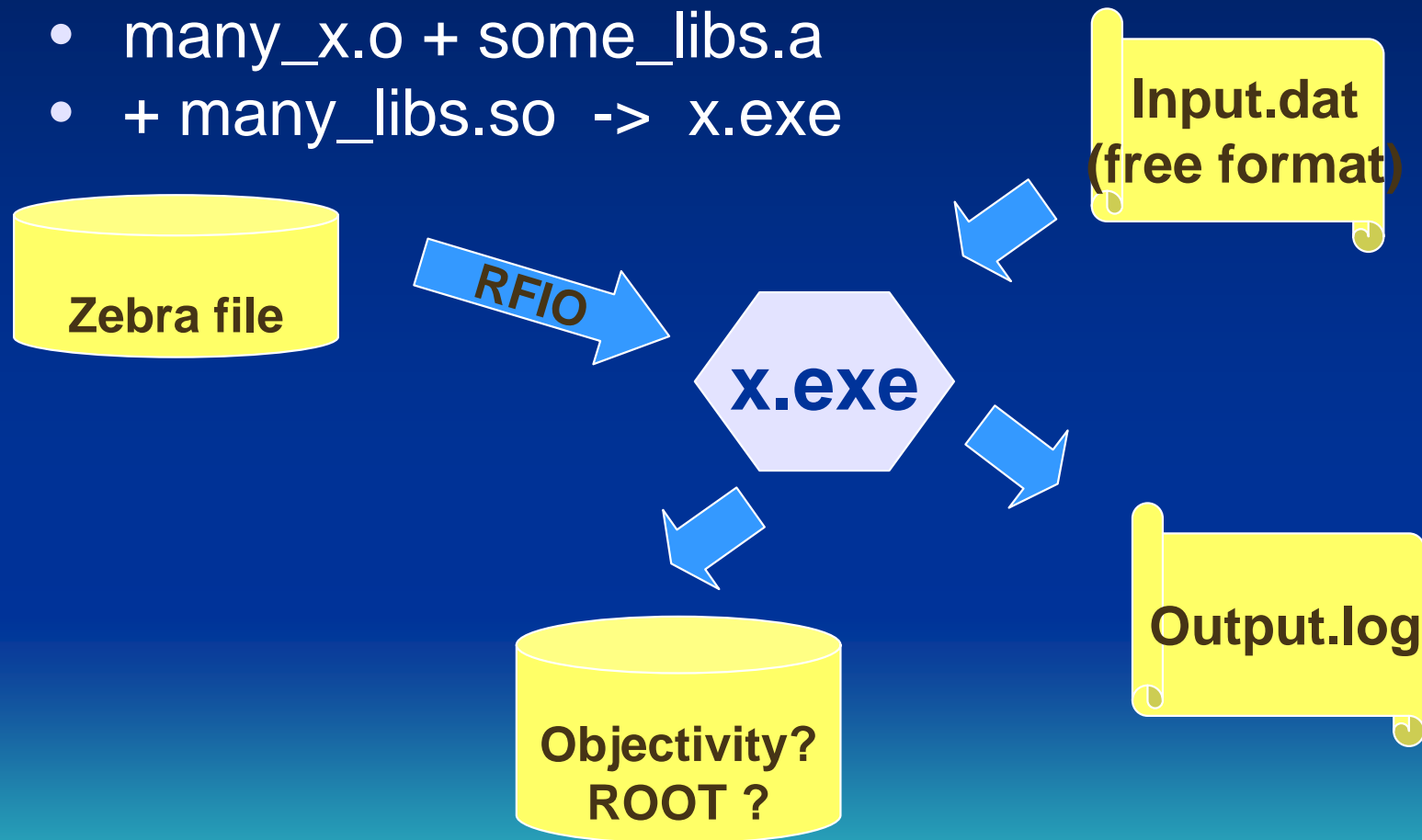
# Executable module in 1987

- many\_x.f -> many\_x.o
- many\_x.o + many\_libs.a -> x.exe



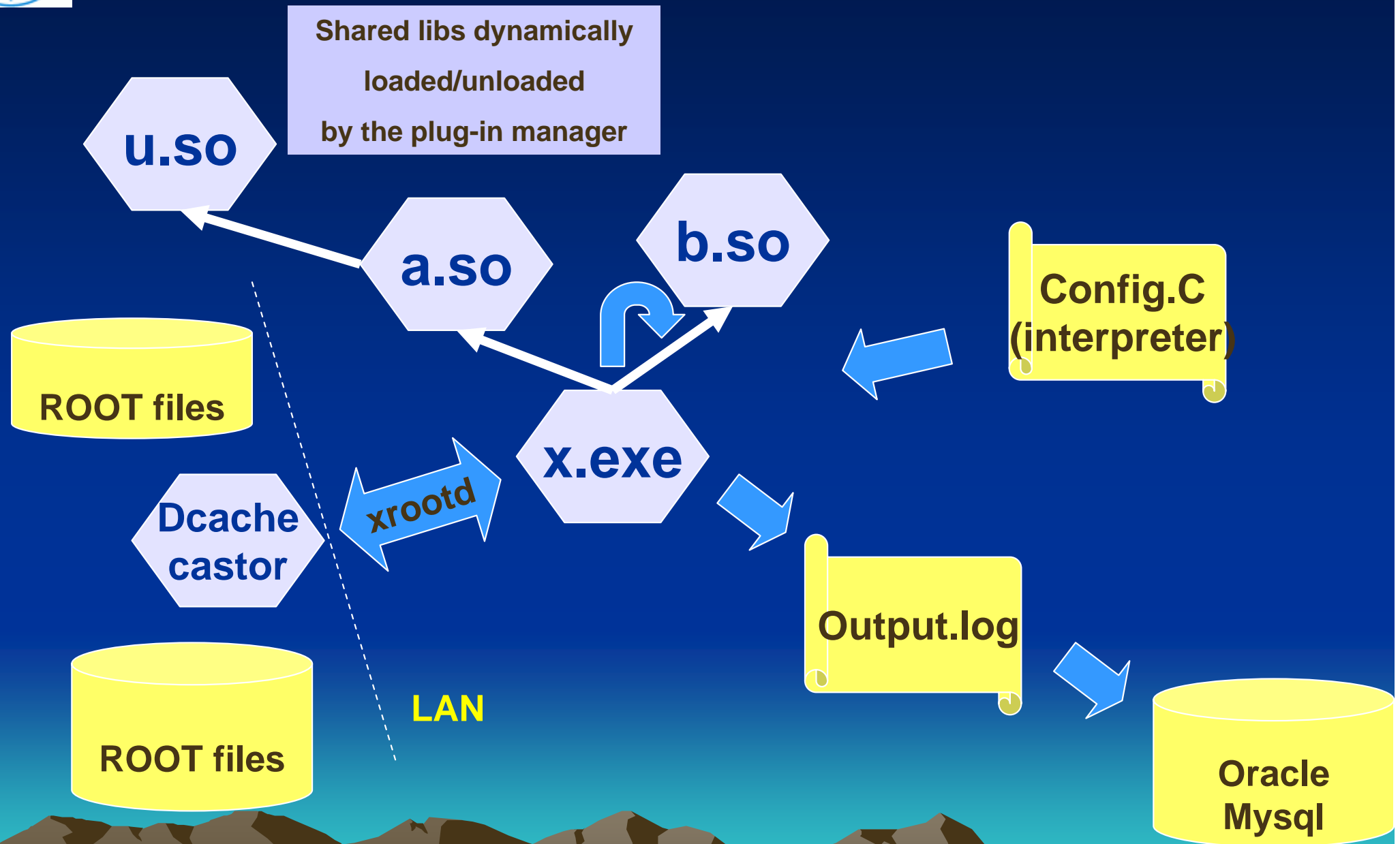
# Executable module in 1997

- many\_x.f -> many\_x.o
- many\_x.o + some\_libs.a
- + many\_libs.so -> x.exe



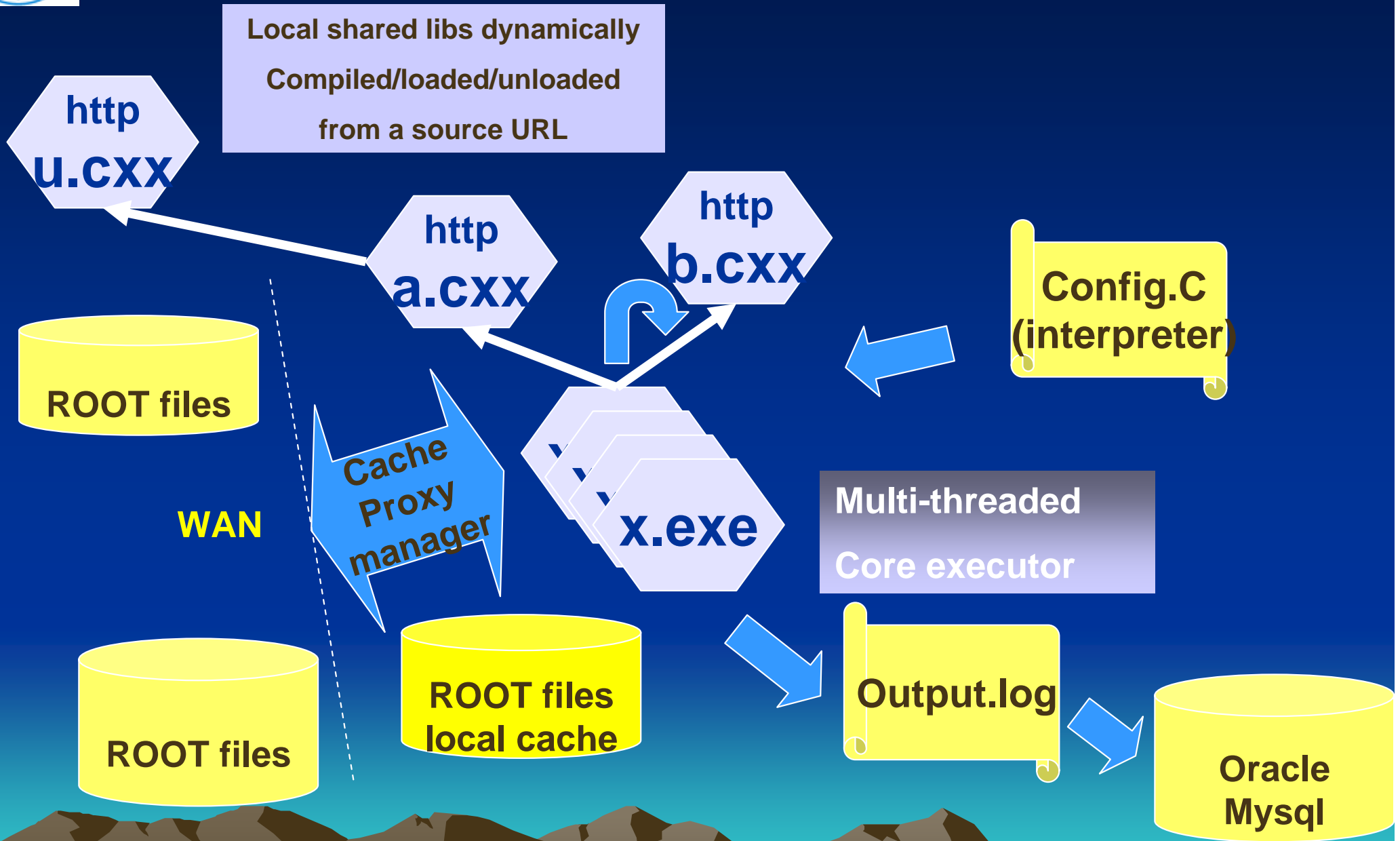


# Executable module in 2007





# Executable module in 2017 ?





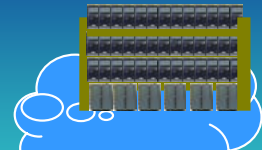
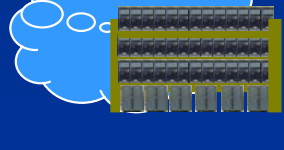
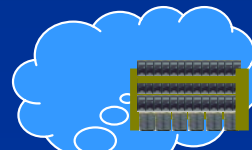
# Challenge 14

## Data Analysis on the GRID

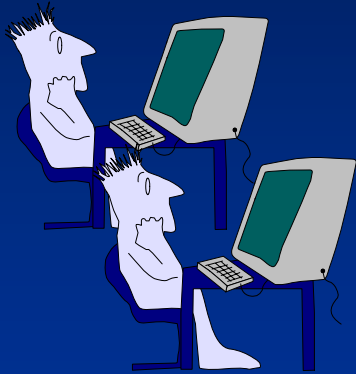


5,000 physicists  
in 1000 locations

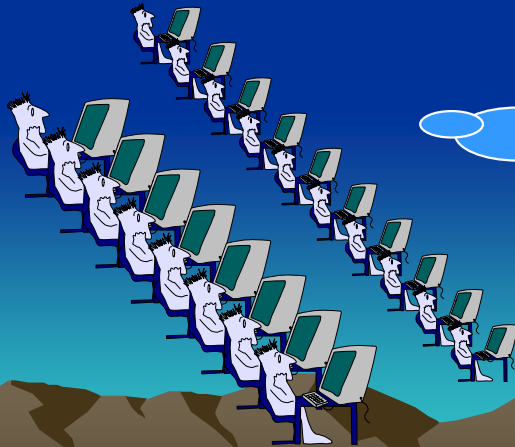
100,000 computers  
in 1000 locations



# GRID: Users profile



Few big users submitting  
many long jobs (Monte Carlo,  
reconstruction)  
**They want to run many jobs in  
one month**



Many users submitting  
many short jobs (physics  
analysis)  
**They want to run many jobs  
in one hour or less**



# Big but few Users

- Monte Carlo jobs (one hour → one day)
  - Each job generates one file (1 GigaByte)
- Reconstruction job (10 minutes -> one hour)
  - Input from the MC job or copied from a storage centre
  - Output (< input) is staged back to a storage centre
- Success rate (90%). If the job fails you resubmit it.
- For several years, GRID projects focused effort on big users only.



# Small but many Users

- **Scenario 1:** submit one batch job to the GRID. It runs somewhere with varying response times.
- **Scenario 2:** Use a splitter to submit many batch jobs to process many data sets (eg CRAB, Ganga, Alien). Output data sets are merged automatically. Success rate  $< 90\%$ . You see the final results only when the last job has been received and all results merged.
- **Scenario 3:** Use PROOF (automatic splitter and merger). Success rate close to  $100\%$ . You can see intermediate feedback objects like histograms. You run from an interactive ROOT session.



# GRID & Parallelism: 1

- The user application splits the problem in  $N$  subtasks. Each task is submitted to a GRID node (minimal input, minimal output).
- The GRID task can run synchronously or asynchronously. If the task fails or time-out, it can be resubmitted to another node.
- One of the first and simplest use of the GRID, but not many applications in HEP.
- Examples are **SETI**, **BOINC**



# GRID & Parallelism: 2

- The typical case of Monte Carlo or reconstruction in HEP.
- It requires massive data transfers between the main computing centres.
- This activity has concentrated so far a very large fraction of the GRID projects and budgets.
- It has been an essential step to foster coordination between hundreds of sites, improve international network bandwidths and robustness.



# GRID & Parallelism: 3

- Distributed data analysis will be a major challenge for the coming experiments.
- This is the area with thousands of people running many different styles of queries, most of the time in a chaotic way.
- The main challenges:
  - Access to millions of data sets (eg 500 TeraBytes)
  - Best match between execution and data location
  - Distributing/compiling/linking users code(a few thousand lines) with experiment large libraries (a few million lines of code).
  - Simplicity of use
  - Real time response
  - Robustness.



# GRID & Parallelism: 3a

- **Currently 2 different & competing directions for distributed data analysis.**
  - **Batch** solution using the existing GRID infrastructure for Monte Carlo and reconstruction programs. A front-end program partitions the problem to analyze ND data sets on NP processors.
  - **Interactive** solution PROOF. Each query is parallelized with an optimum match of execution and data location.

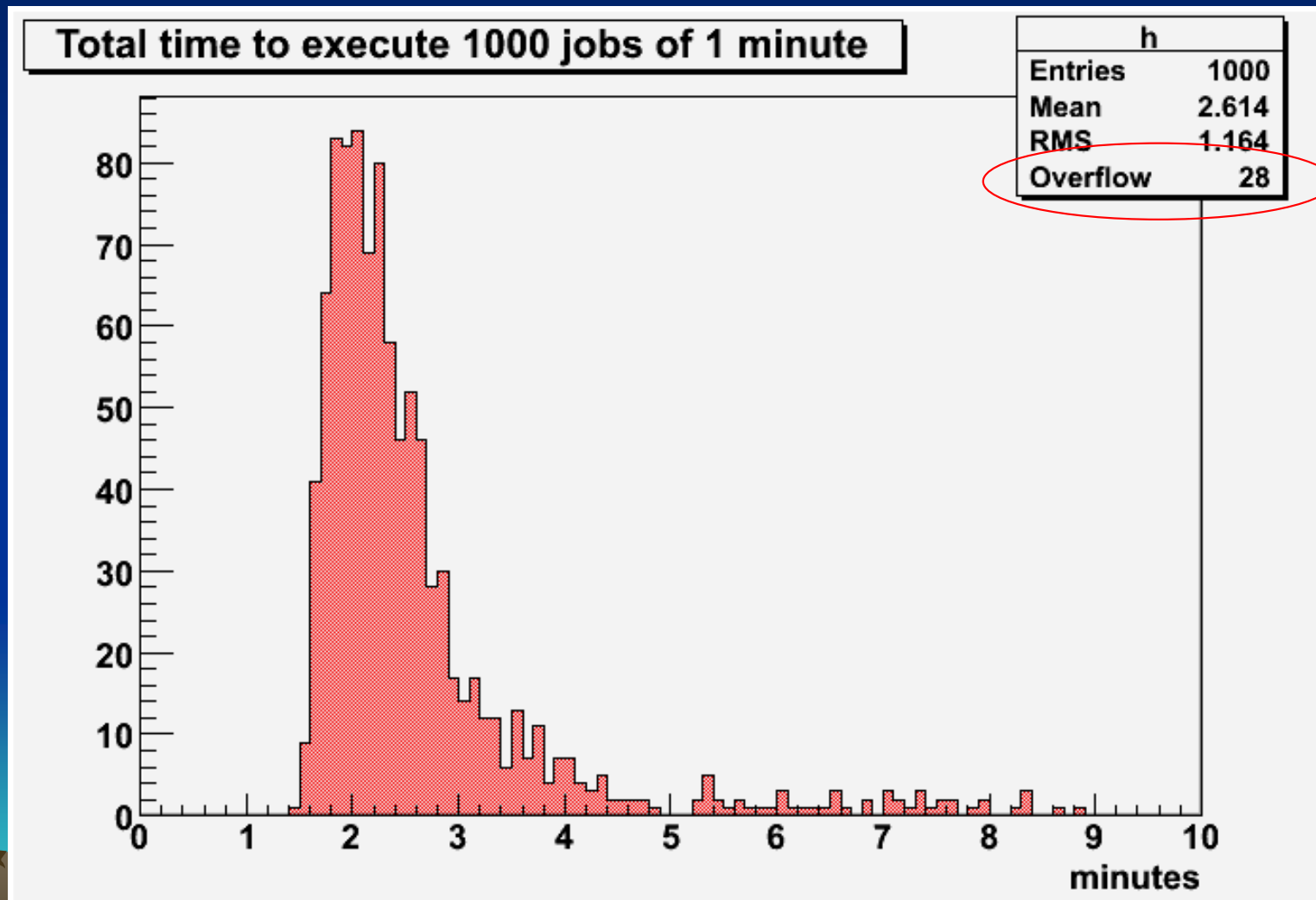


# Scenario 1 & 2: PROS

- Job level parallelism. Conventional model. Nothing to change in user program.
  - Initialisation phase
  - Loop on events
  - Termination
- Same job can run on laptop or on a GRID job.

# Scenario 1 & 2: CONS(1)

- Long tail in the jobs wall clock time distribution.





# Scenario 1 & 2: CONS(2)

- Can only merge output after a time cut.
- More data movement (input & output)
- Cannot have interactive feedback
- Two consecutive queries will produce different results (problem with rare events)
- Will use only one core on a multi core laptop or GRID node.
- Hard to control priorities and quotas.



# Scenario 3: PROS

- Predictive response. Event level parallelism. Workers terminate at the same time
- Process moved to data as much as possible, otherwise use network.
- Interactive feedback
- Can view running queries in different ROOT sessions.
- Can take advantage of multi core cpus



# Scenario 3: CONS

- Event level parallelism. User code must follow a template: the **TSelector** API.
- Good for a local area cluster. More difficult to put in place in a GRID collection of local area clusters.
- Interactive schedulers, priority managers must be put in place.
- Debugging a problem slightly more difficult.



# Challenge 15

## Languages

- C++ clear winner in our field and also other fields
- see, eg a recent compilation at <http://www.lextrait.com/vincent/implementations.html>
- From simple C++ to complex templated code
- Unlike Java, no reflexion system. This is essential for I/O and interpreters.
- C++2009: better thread support, Aspect-oriented
- C++2014: first reflexion system?



# Challenge 16

## Software Development Tools

- better integration with Xcode, VisualStudio or like
- fast memory checkers
- faster **valgrind**
- faster **profilers**
- Better tools to debug parallel applications
- **Code checkers** and **smell detection**
- Better html page generators



# Challenge 17

## Distributed Code Management

- patchy, cmz -> cvs
- cvs -> **svn**
- cmt? scram? (managing dependencies)
- automatic project creation from cvs/svn to VisualStudio or Xcode and vice-versa



# Challenge 18

## Simplification of Software Distribution

- tar files
- source + make
- install from `http://source`
- install from `http://binary proxy`
- install on demand via plugin manager, autoloader
- automatic updates
- time to install
- fraction of code used

See **BOOT**  
Project  
First release  
In June 08



# Challenge 19

## Software Correctness

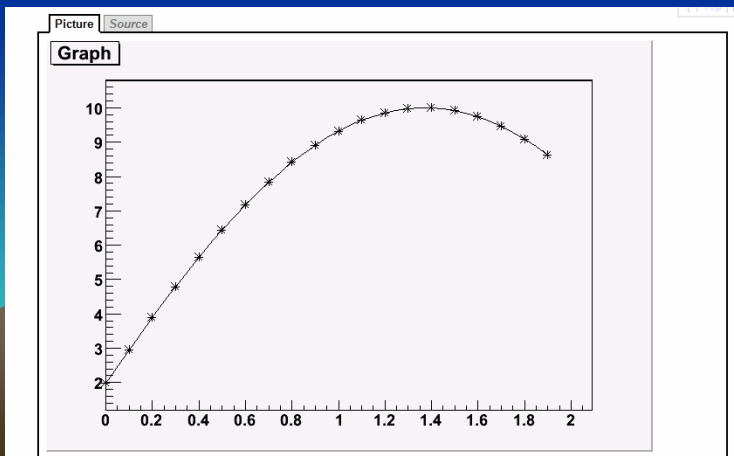
- -big concern with multi million lines of code
- -validation suite
- -unit test
- -combinatorial test
- -nightly builds (code + validation suite)



# Challenge 20

## Scalable Software Documentation

- Legacy **Doxygen**
- Need for something more dynamic, understanding Math, Latex, 2-D and 3-D graphics, interactive tutorials.
- See results of new **THtml** at:
  - <http://root.cern.ch/root/html/TGraph.html>



```
{
  TCanvas *c1 = new TCanvas("c1","A Simple Graph Example",200,10,700,500);
  Double_t x[100], y[100];
  Int_t n = 20;
  for (Int_t i=0;i<n;i++) {
    x[i] = i*0.1;
    y[i] = 10*sin(x[i]+0.2);
  }
  gr = new TGraph(n,x,y);
  gr->Draw("AC*");
  return c1;
}
```



# Challenge 21

## Education

- Training must be a continuous effort
- Core Software guys often desperate with newcomers.
- Software Engineering and discipline required to participate to large international projects is absent in University programs.



# Summary

- A large fraction of the software for the next decade already in place or shaping up.
- Long time between design and effective use.
- Core Software requires Open Source and international cooperation to guarantee stability and smooth evolution.
- Parallelism will become a key parameter
- More effort must be invested in software quality, training and education.



# Summary-2

- But the MAIN challenge will be to deliver **scalable systems**:
- **Simple to use** for beginners with some very basic rules and tools.
- Browsing (understanding) an ever growing dynamic code and data will be a must.



# Summary-3

- Building large software systems is like building a large city. One needs to standardize on the communication pipes for input and output and setup a basic set of rules to extend the system or navigate inside.