# La CAF (CDF Analysis Farm)

**Massimo Casarsa**

*Sez. INFN di Trieste*

for the CDF Collaboration

**IFAE 2005**
Catania, 30 Marzo - 2 Aprile 2005

# Talk layout

| | |
|---|---|
| **1** | CAF motivation and goal |

| | |
|---|---|
| **2** | CAF overview from user standpoint |

| | |
|---|---|
| **3** | Structure and features of the CAF software |

| | |
|---|---|
| **4** | CAF evolution towards the GRID |

| | |
|---|---|
| **5** | Conclusion |

# CAF motivation and goal

## The CDF Collaboration

**Canada**
McGill Univ.
Univ. of Toronto

**USA**
Argonne National Laboratory,IL
Brandeis Univ.,MS
Univ. of Chicago,IL
Davis UC,CA
Duke Univ.,NC
FNAL,IL
Univ. of Florida,FL
Harvard Univ.,MA
Univ. of Illinois,IL
The Johns Hopkins Univ.,MD
LBNL,CA
MIT,MA
Michigan State Univ.,MI
Univ. of Michigan,MI
Univ. of New Mexico,NM
The Ohio State Univ.,OH
Univ. of Pennsylvania,PA
Univ. of Pittsburgh,PA
Purdue Univ.,IN
Univ. of Rochester,NY
Rockefeller Univ.,NY
Rutgers Univ.,NJ
Texas A&M Univ.,TX
Texas Tech Univ.,TX
Tufts Univ.,MA
UCLA,CA
Univ. of Wisconsin,WI
Yale Univ.,CT

**China**
Academia Sinica,
Taiwan

**Korea**
KHCL

**Russia**
JINR, Dubna
ITEP, Moscow

**Germany**
Univ. Karlsruhe

**Switzerland**
Univ. of Geneva

**UK**
Glasgow Univ.
Univ. of Liverpool
Univ. of Oxford
Univ. College London

**Italy**
Univ. of Bologna,INFN
Frascati, INFN
Univ. di Padova,INFN
Pisa, INFN
Univ. di Roma I,INFN
INFN-Trieste
Univ. di Udine

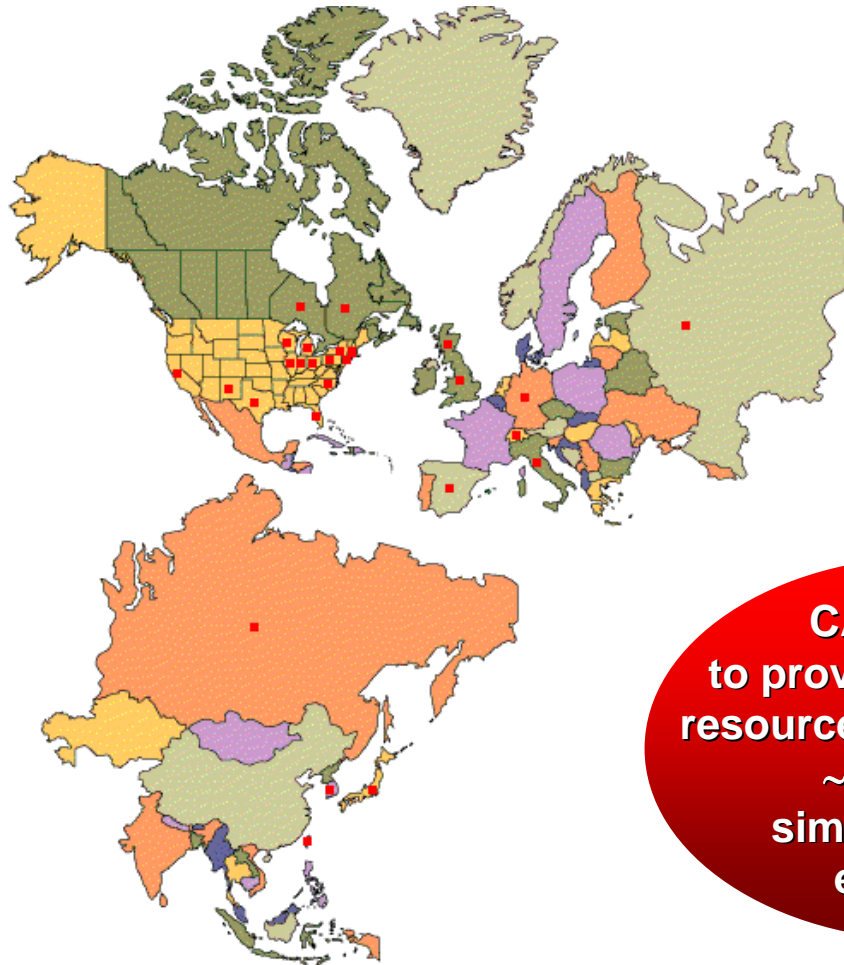**Spain**
Univ. of Cantabria

**Japan**
Hiroshima Univ.
KEK
Osaka City Univ.
Univ. of Tsukuba
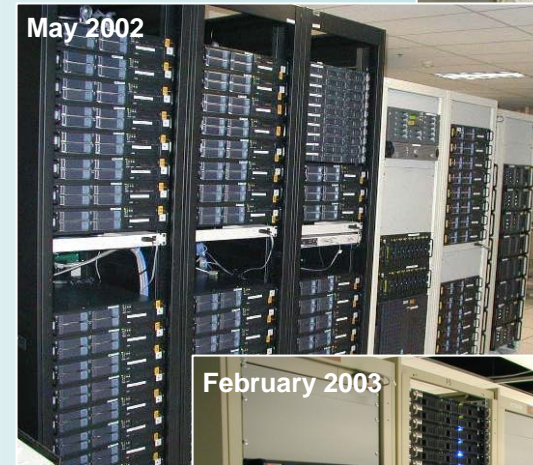Waseda Univ., Tokyo

► ~600 Physicists ◄

► 56 Institutions ◄

► 11 Countries ◄

**CAF goal is
to provide computing
resources for analysis to
~200 users
simultaneously
every day**

# CAF milestones

- Assembled in 2002 to meet the Collaboration needs for computing resources:

    - ✓ <u>data analysis</u>: CDF produces datasets of 100s of TBs whose processing takes several days (0.1-0.5 s/event).

    - ✓ <u>MC production</u>: detector simulation is heavy CPU consuming (~1 s/event).

- Born as a farm localized at FNAL with the FBSNG batch manager, then migrated to Condor.

- CAF model exported and farm decentralized to many sites around the world (DCAFs): at present ~50% of CDF computing power outside Fermilab.
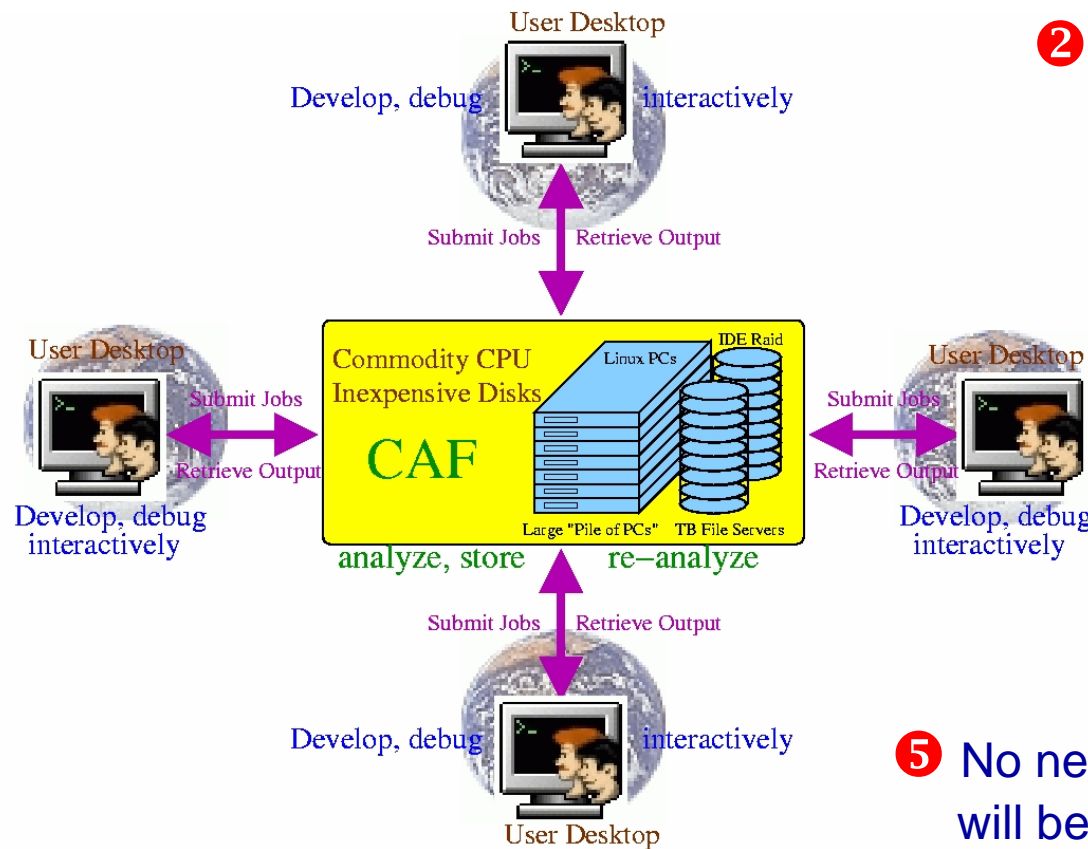
February 2002

May 2002

February 2003

# DCAFs resources around the world

| Cluster name | Location | CPU [GHz] | Disk space [TB] |
|---|---|---|---|
| Original FNAL CAF | FNAL | 1200 | 300 |
| FNAL CondorCAF | | 2000 | |
| **CNAFCAF** | **Bologna (Italy)** | **300** | **22** |
| KORCAF | KNU (South Korea) | 120 | 0.6 |
| ASCAF | Acc. Sinica (Taiwan) | 134 | 3.0 |
| SDSCCAF | San Diego (USA) | 280 | 4.0 |
| HEXCAF | Rutgers (USA) | 100 | 4.0 |
| TORCAF2 | Toronto (Canada) | 576 | 10 |
| JPCAF | Tsukuba (Japan) | 152 | 5.0 |
| CANCAF | Cantabria (Spain) | 52 | 1.5 |
| MITCAF | Boston (USA) | 110 | 2.0 |
| **TOTAL** | | **5024** | **352.1** |

- Data reside at FNAL, DCAFs mostly used for MC production, recently dataset replicas allow to run also analysis jobs.

# CAF philosophy

❶ Develop and debug analysis code on personal desktop or laptop.



❷ Submit jobs to CAF from anywhere in the world.

❸ Have the output delivered wherever you like.

❹ Interact with running jobs as they were local.

❺ No need to stay connected, will be notified at job completion.

# CAF user interface (CafGui)

**Fill in appropriate fields:**

- ▶ **farm**
- ▶ **data access method and dataset**
- ▶ **process type**
- ▶ **group**
- ▶ **command**
- ▶ **original directory**
- ▶ **output location**
- ▶ **e-mail address**

**submit**



CDF Run II CAF GUI

fcdfhead3.fnal.gov:8000

Analysis Farm: caf

Data Access: Method: DFC
Dataset: xbhd0d

Process Type: long

Group: italy

**section range**

Initial Command: ./run.sh $   1   800

Original Directory: /cdf/home/casarsa/analysis/mixing/cafjobs/   Browse...

Ouput File Location: icaf:ntu_BsDsKstarK_$.tgz

Email?   Email Address: casarsa@fnal.gov

Submit   Quit   **Ready**

```
(2005-03-11 12:58:08) DFC data access method selected
(2005-03-11 12:58:15) long process type selected
(2005-03-11 12:58:20) group italy  selected
(2005-03-11 13:00:00) Email sent to casarsa@fnal.gov upon job completion
(2005-03-11 13:00:15) Submission canceled by user
(2005-03-11 13:00:51) Email sent to casarsa@fnal.gov upon job completion
(2005-03-11 13:00:52) Continuing with submission
(2005-03-11 13:00:52) /bin/tar -chvzf /cdf/scratch/casarsa/casarsa68364.tgz *
(2005-03-11 13:01:21) Remove /cdf/scratch/casarsa/casarsa68364.tgz
(2005-03-11 13:01:21) Job Submission is successful, JID: 562882
(2005-03-11 13:02:15) Sending of email disabled
(2005-03-11 13:02:16) Sending of email enabled
```

# User interactive tools

CafMon allows users to interact with remote jobs as they were running locally:

- **job management tasks:**

  - ✓ kill            ⟹  kill job/section;
  - ✓ hold/release    ⟹  hold/release jobs;
  - ✓ chprio/chgroup  ⟹  change the priority/group of a job.

- **jobs and remote system interactive monitoring:**

  - ✓ jobs           ⟹  list submitted jobs and sections;
  - ✓ log            ⟹  print out log file of a specific section;
  - ✓ top            ⟹  show status of the node running a specific section;
  - ✓ ps             ⟹  show the processes of a specific section;
  - ✓ dir            ⟹  show the working dir of a specific section;
  - ✓ tail/head/cat  ⟹  inspect any text file in section working directory.

- **debug:**

  - ✓ debug          ⟹  attach a debugging session to a remote running process.

# User web monitoring

# User web monitoring: worker node status

# Farm implementation

- **Hardware**:

  need lots of CPUs $\Rightarrow$ commodity CPUs $\Rightarrow$ dual Intel/AMD;
  need lots of disk $\Rightarrow$ cheap disk $\Rightarrow$ IDE RAID 50 arrays.

- **Nodes software**:

  ✓ operating system: Linux;
  ✓ have access to CDF software.

- **Batch manager**: Condor

  ✓ six virtual-machines (VM) per node;
  ✓ higher priority for groups on institutions' proprietary hardware;
  ✓ process type for job length:

  | | |
  |---|---|
  | test: | 2 h (max 4 sections), |
  | short: | 6 h, |
  | medium: | 12 h, |
  | long: | 72 h. |

Condor configuration:

headnode

schedd → negotiator ← collector
schedd ← negotiator

claim the match

matched to schedd

status

startd

starterd

compute node    starterd

# CAF software overview

- <u>Design goal</u>:

  - ✓ Give user access to CAF resources from anywhere in the world.

- <u>Design constraints/desirables</u>:

  - ✓ Fermilab computing security policy $\rightarrow$ Kerberos;

  - ✓ administrative ease $\rightarrow$ no user accounts
    $\rightarrow$ non-interactive batch, jobs run as generic users (one for each VM);

  - ✓ user identity $\rightarrow$ unique privileges for batch jobs, disk space;

  - ✓ large scale parallelization with single submission (Condor dagman).

- <u>Result</u>:

  - ✓ very user-friendly software;

  - ✓ user provides a shell script and an executable + all needed files;

  - ✓ everything is tarred up and sent to CAF;

  - ✓ user is notified when his job is completed.

# CAF software: user desktop

**User desktop**

CafMon

Web browser

CafGui/ CafSubmit

*kerberos*

*kerberos*

**Head node**

monitor

mailer

submitter

CoD

startd

CafExe

r script

er Exe

DFC/SAM

**Data Handling system**

User interactive tool to manage and monitor running jobs:
- ▶ communicates with head node via a Kerberos authenticated connection.

User front-end for job submission:
- ▶ communicates with head node via a Kerberos authenticated connection;
- ▶ sends user submission parameters to head node.

Web browser
- ▶ give access to the CAF web monitoring pages.

le serv

rootd

- – – – monitoring path
- ———— job request/exe path
- ———— data server path
- ———— job output path

# CAF software: head node

**Monitoring daemon :**

▶ provides monitoring information;

▶ manages user interactive requests.

**CAF router:**

▶ allows communication between the monitor daemon and CafExe.

**Head node**

- monitor
- CafRout
- mailer
- CoD
- submitter → Condor

**Worker node**

- startd
- CafExe
- CoD

Kerberos

Kerberos

root

**Notification daemon:**

▶ sends summary mail at the end of job:
  - exit codes of all sections,
  - waiting, CPU, real times,
  - input/output data summary.

▶ cleans up Condor files.

**Submission request manager:**

▶ grants Kerberos authentication to user;

▶ transfers CAF tar-ball from user desktop;

▶ creates Condor submit files;

▶ submits CafExe to Condor.

**Data Handling system**

SAM

- – – – monitoring path
- ——— job request/exe path
- ——— data server path
- ——— job output path

# CAF software: worker node

**User desktop**

CafMon

Web browser

**Head node**

monitor

CafRout

*kerberos*

CAF job wrapper:

▶ gets user Kerberos credentials;

▶ untars user tar-ball;

▶ runs the user initial command (usually a shell script);

▶ at the end tars up the working directory and transfers the tar-ball to the output location;

▶ monitoring tasks:
  ▪ create job summary file,
  ▪ serve as interactive CafMon callback.

**Worker node**

startd

CafExe

User script

User Exe

*CoD*

*kerberos*

*kerberos*

*rootd*

*DFC/SAM*

**Data Handling system**

- - - monitoring path
——— job request/exe path
——— data server path
——— job output path
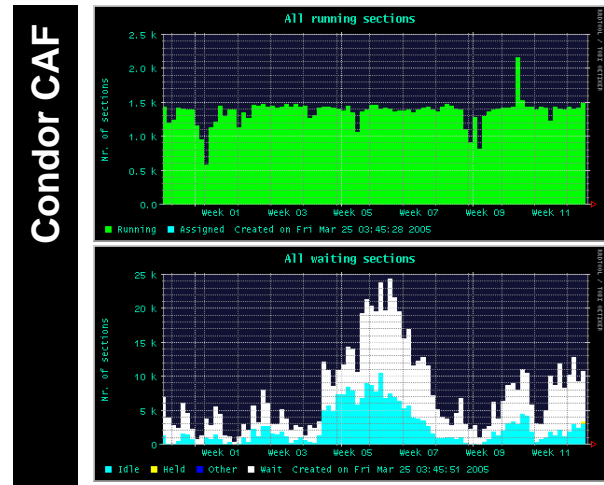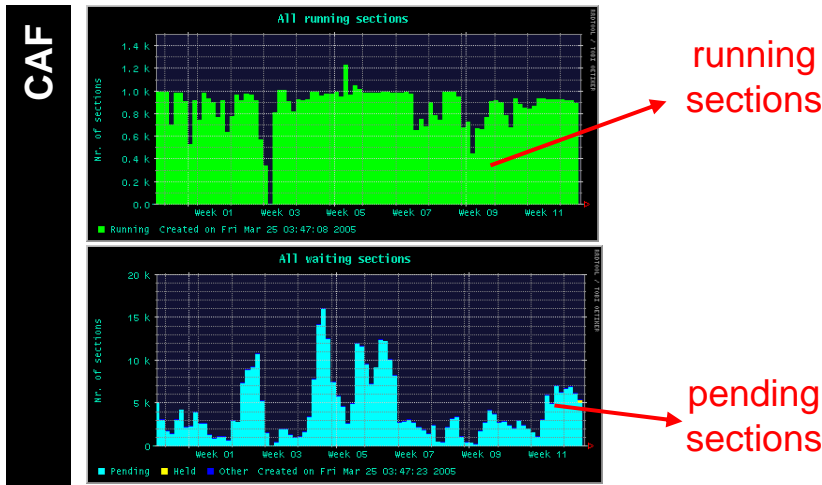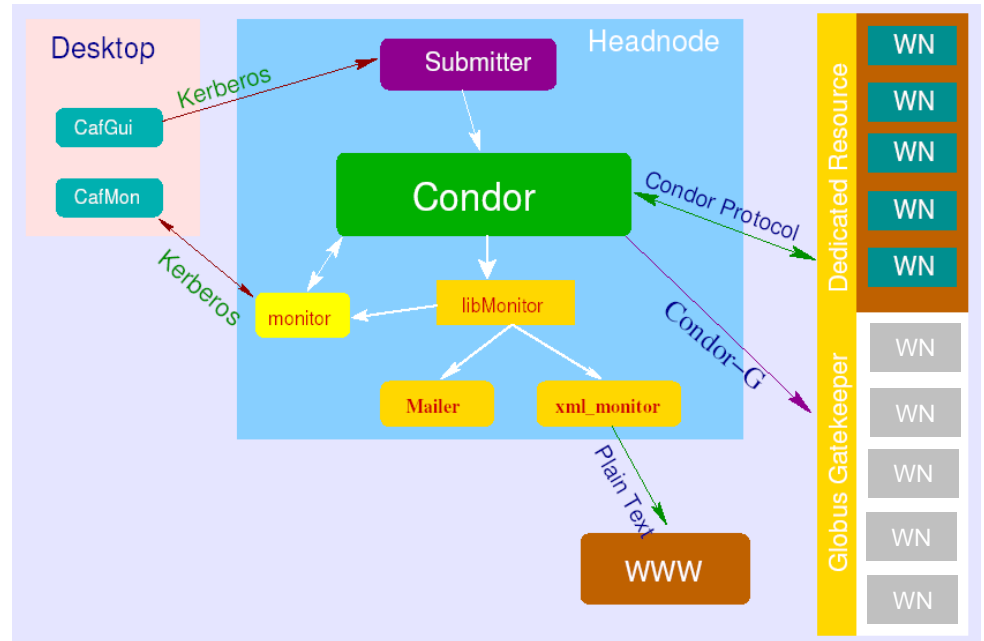
# CAF and DCAF utilization



- The CDF increasing volume of data pushes for more and more computer power in the next future $\Rightarrow$ GRID may offer plenty of resources, at least until LHC turns on.

# Toward CafGRID

- <u>First step</u>: **GlideCAF at CNAF Tier1**
  - ✓ based on condor_glidin;
  - ✓ first working prototype.

condor_glideins are Condor-G jobs which bypass the Resource Broker and reach the GRID site Gatekeeper directly:

- ▶ the GateKeeper distributes the jobs to the WNs;
- ▶ the WNs install/run Condor on the fly once the Condor-G jobs start;
- ▶ the WNs become a part of the Condor pool when the Condor daemons start.



- <u>Final goal</u>:
  - ✓ modify CAF software in order to allow jobs submission to GRID;
  - ✓ integrate DCAFs into the GRID as Computer Elements is desirable, but not yet designed.

# Conclusion

- The need for computing resources has been steadily urging the CDF Central Analysis Facility to evolve from an old-fashioned farm, originally localized at FNAL, to a GRID oriented structure, distributed worldwide (~50% of CDF computing power already outside FNAL).

- GRID represents an abundant reservoir of computing resources to fulfill CDF analysis needs.

- First step has been done to integrate CDF DCAFs into GRID, at CNAF a working prototype already exploits Tier1 resources:

  - ✓ users may submit MC jobs,
  - ✓ preliminary tests on data analysis jobs.