



ISTITUTO NAZIONALE DI FISICA NUCLEARE

CNAF Bologna

INFN/TC-99/20

11 Ottobre 1999

**CONDOR ON WAN
IMPLEMENTATION PROPOSAL**

A. Ghiselli, P. Mazzanti,
E. Becchetti, D. Bortolotti, C. Bulfon, G. Cabras, P. Capiluppi, R. Elia, E. Fasanelli,
T. Ferrari, L. Fonti, U. Gasparini, R. Gomezel, E. Leonardi, P. Mastroserio, D. Menasce,
M. Michelotto, S. Parlati, F. Prelz, A. Rappoldi, F. Semeria, L. Servoli, M. Sgaravatto,
C. Vistoli, F. Taurino

INFN-CNAF, V.le Berti Pichat 6/2, I-40127 Bologna, Italy

Abstract

This deliverable is part of the project 'Condor on WAN' and describes the Condor pool implementation as computing resource for INFN. The Condor pool specifications have been defined taking into account the results of the test phase, terminated by the end of 98. The structure of the operational service is also described

Keywords: HTC, WAN, Checkpoint domain, Sub-pool.

*Published by SIS-Pubblicazioni
Laboratori Nazionali di Frascati*

1 INTRODUCTION

This proposal describes the configuration of a wide area Condor batch system at INFN to provide the INFN users with a global computing resource.

High Throughput Computing (HTC) on WAN can be deployed to satisfy the computing needs of INFN by accessing the huge CPU capacity distributed in all INFN sites: thank to the increasing CPU power of PCs and WSs and the decrease in cost the total computing capacity of the Institute has been increasing substantially.

Since the Condor philosophy aims at using only *idle* CPU cycles, Condor has been identified as the right candidate to satisfy the computing needs at INFN: Condor optimizes the usage of existing computing capacity. Local users still access local machines with higher priority. Resource sharing is controlled by local policies: for example subsets of machines (sub pools) can be defined. In each subset (e.g. all the machine belonging to some research group) specific jobs can have absolute priority on other jobs, that may be eventually vacated when high priority jobs need to execute.

The implementation of the Condor System on the INFN wide area network as described in this document is a result of the first part of the project “Condor on WAN”, developed in collaboration with the Condor-Team from the Computer Science Department of the Wisconsin-Madison University.

2 TEST PHASE

During the first part of the project an experimental Condor pool has been set up in order to check the reliability and efficiency of the system on WAN, and its suitability to INFN computing needs. Tests of CPU intensive jobs gave good results in terms of very good workload, whereas jobs with an high frequency I/O were less efficient of CPU usage. In latter case performance improves if jobs run in a uniform file system: this means that with appropriate configuration even with I/O intensive jobs can run efficiently. Future mechanisms like caching or dedicated file systems will be investigated to improve the efficiency of I/O intensive programs. Other findings are the need of flexible policies in CPU usage of the machines and the importance of an adequate location of checkpoint servers.

3 IMPLEMENTATION PHASE

The choice to have only one ‘pool’ (i.e. a pool with only one Central Manager) was made in order to optimize the CPU usage of all the INFN hosts available for Condor. The need of providing guarantees to local jobs in CPU usage can be satisfied by configuring sub-pools, while the overall efficiency of the system can still be achieved through a suitable setting of a set of checkpoint servers.

3.1 SUB-POOL

A sub-pool is a collection of machines configured in order to give higher priorities to jobs belonging to local users or to research group users. A sub-pool can be local to one INFN site or distributed between different sites connected through WAN. Sub-pool policies must be defined by local management in agreement with the responsible for the research groups.

3.2 CHECKPOINT TOPOLOGY

The need of an appropriate checkpoint server topology stems from the decision to limit the impact of checkpoint file transfers especially when the number of machines will increase up to several hundreds. The optimal checkpointing policies the following:

- Checkpointing a big size file should be accomplished in short time in order to let the owner access its machine without a visible delay.
- Sub-pools may make use of a dedicated checkpoint server.
- The definition of the “best” checkpoint server should be network adaptive.

Obviously checkpoint should not limit the overall computing throughput.

The solution adopted will be implemented in two steps:

- 1) Configuration of *checkpoint domains*
- 2) Implementing a *distributed (dynamic) checkpointing* as a new feature of the Condor System.

The most important characteristic of the solution adopted is that the *network* has been defined as *resource* of Condor: network bandwidth between checkpoints server and execution machines is a machine ClassAds attribute, dynamically updated, and used by checkpointing for the better choice between execution machines and checkpoint servers.

Initially each execution machine will have a fixed checkpoint server associated with it, then the association between execution machine and checkpoint server will be dynamically decided according to the network load.

3.3 INFN CONDOR TOPOLOGY

Most of the INFN sites have several machines in the WAN Condor pool and they are connected to the research network GARR-B with access speeds ranging from 2Mbps up to 8Mbps. The logical topology is described below.

The checkpoint server domain is defined according to the following guidelines:

- Presence of a sufficiently large CPU capacity
- Presence of a set of machines with an efficient network connectivity
- Set of site policies (eg. jobs have to run only locally)

The initial topology will have at least 10 checkpoint servers and the idea is to increase the number of the machines in order to have one checkpoint server in each site.

3.4 IMPLEMENTATION SCHEDULE

1. Phase 1 June 99

According to the guidelines above the first domains will be configured in Bologna, Napoli, Padova, cnaf and Milano. The Condor pool will have one ‘Central Manager’ located at CNAF and its backup located in Rome or Milan.

Some tools are under development in order to provide user friendly interfaces to submit and monitor jobs to the Condor system.

2. Phase 2 July 99

This phase is characterized by the dynamic definition of checkpoint domains, i.e. the association between execution machines and checkpoint servers will be done by the ‘Network Manager’ of the pool.

3.5 MANAGEMENT

3.5.1 *Central management*

The Admin Group will act as central management group and has to provide:

- Configuration, tuning and overall maintenance of the INFN Condor Wan pool
- management tools
- activity reports
- Condor resource usage statistics (CPU, Network, Ckpt-server)
- Which Condor release has to be installed
- Help desk for users and local administrators.
- Interface to condor support in Madison.

The required man-power should be 1.5FTE

Admin Group composition:

5 people: D. Bortolotti (30%), M. Sgaravatto (30%), E. Querzola (30%), P. Mastroserio and F. Taurino (30%), [1.2FTE in total]

3.5.2 *Local management*

Local management has to provide:

- release installation in agreement with the central management
- local condor usage policies (e.g. sub-pools)

man-power: 5% max per site.

Local management group:

Torino (L. Gaido), Milano (F. Prelz), Udine (G. Cabras), Trieste (R. Gomezel), Padova (M. Michelotto, M. Sgaravatto), Pavia (A. Rappoldi), Genova (A. Brunengo, C. Salvo), CNAF (E. Querzola), Bologna (D. Bortolotti), Firenze (R. Checchin), Pisa (F. Donno), Perugia (E. Becchetti, L. Servoli), Roma1 (C. Bulfon, E. Leonardi), Roma2 (R. Elia), LNGS (S. Parlati), Bari (P. Amendola, B. Tataranni), Lecce (E. Fasanelli), Napoli (P. Mastroserio, F. Taurino), Catania (G. Andronico, R. Barbera), Cagliari (A. Fara).

3.5.3 *Steering Committee*

The Steering committee should:

- consider the status of the condor system and suggest when upgrade the software
- interact with the Condor Team and suggest possible modifications of the system
- define the general policy of the condor pool
- organize meeting for condor administrators (and users)

Steering Committee:

A. Ghiselli, T. Ferrari, P. Mazzanti, F. Prelz, F. Semeria, C. Vistoli, M. Sgaravatto

4 EXTERNAL COLLABORATIONS AND SUPPORT

Part of the project is developed in collaboration with the Condor Team of the Computer Science Department of Wisconsin–Madison University. The collaboration is motivated by the Condor Team interest to investigate the Condor behavior in a WAN scenario.

Project documentation is available in *HTTP://www.mi.infn.it/condor*.

5 CONCLUSIONS

The Condor WAN pool test layout, where machines are distributed over 20 INFN sites and connected through the national research network GARR, gave the possibility to prove the reliability and robustness of the system and to study the most suitable ‘checkpoint domain’ topology in order to optimize checkpoint operations and to limit geographic network traffic as much as possible. This goal can be achieved by considering the ‘network’ as Condor resource. A ‘Network Manager’ for Condor has been developed in collaboration with the authors of Condor.

Furthermore in each site machines belonging to the pool can be configured in order to give absolute priority to the local research group jobs as if they were in a dedicated ‘sub-pool’.

The choice to have only one ‘pool’ does not represent a single point of failure because the last releases of Condor make the definition of central manager backups possible for higher stability of the WAN pool.