

Condor

A Distributed Job Scheduler

Mikkel Bystrup Stensgaard (s001434)
Matias Dons Dollerup (s991367)



Agenda

- Introduction
- Scheduling
- Universes
- Architecture
- Administration
- Conclusion



Introduction to Condor (1)



- Condor is a distributed job scheduler
- Provides environment for handling several jobs simultaneously
- Provides snapshot mechanism for restoration and migration of jobs

Introduction to Condor (2)



- A Condor pool is a set of grids, opportunistic workstations and dedicated workstations
- Multiple pools can be connected through Condor
- Condor takes advantage of opportunistic workstations to minimize idle time

Introduction to Condor (3)



- Example:

A Condor cluster at UW-Madison Department of Sciences containing more than a 1000 workstations including a 500 CPU Beowulf cluster, delivers around 650 CPU-days on a typical day.

Condor advertising (1)



- Condor operates through an advertising interface
- Both jobs and workstations advertise their presence in the pool
- A job advertisement presents the specifications for a job
- A workstation advertisement presents the capabilities of a workstations

Condor advertising (2)



- Jobs are specified in a Job Description File that Condor uses to create an advertisement
- The main fields in a description file is the requirements and the ranking expressions
- Expressions is given by a logical syntax, hence allowing very specific formulation

Condor advertising (3)



- Example: A job description file

```
universe = vanilla          # select runtime environment
executable = some_job
requirements = (Arch=="INTEL" && OpSys=="LINUX") #target
rank = (Memory * 10000) + KFlops                #target
arguments = -verbose
input = in.dat          # redirect to stdin
output = out.dat        # redirect to stdout
log = log.txt
queue                   # add job to queue
```

Job scheduling in Condor



- Jobs can be added from any terminal within the pool
- Condor matches job advertisements with machine advertisements and schedules execution

Universes in Condor



- **Condor provides scheduling for a variety of applications by using different execution environments called universes.**
 - Vanilla Universe: Sequential programs based on a shared file system.
 - MPI Universe: Message Parsing programs, based on process parallelism. These jobs are only run on Beowulf clusters – not on the opportunistic workstations.
 - Parallel Virtual Machine - PVM Universe: Supports jobs where machines can enter and/or leave during job execution.
 - Globus Universes: Grid computing environment.
 - Standard Universe: Special Condor environment optimized for effective job scheduling.

Standard Universe



- Application checkpointing enables process migration and process restoration
- Requires the application to be linked with Condor libraries
- Limitations: (running in user mode)
 - Multi-process jobs are not allowed including multiple kernel threads
 - Network communication must be brief
 - No interprocess communication (shared memory, pipes, semaphores)

Checkpointing (1)



- a "snapshot" of the jobs current state
- Freedom to preempt job
 - Rescheduling a higher priority job
 - User reclaims non-dedicated PC
 - Enables migration
- Increases fault tolerance

Checkpointing (2)



- Condor writes the following to a file/socket
 - CPU state, including registers
 - Stack
 - State of all open files
 - Signal handlers
 - Pending signals

Program must be linked with the Condor library!

Access to data



- Shared file system
 - Common on Beowulf cluster
 - Must be specified as a requirement expression
- Condor File Transfer mechanism
 - Condor automatically transfer specified files before starting the job. Also transfer output files back when job are finished of pre-empted.
- Remote I/O calls
 - Only possible in the Standard universe

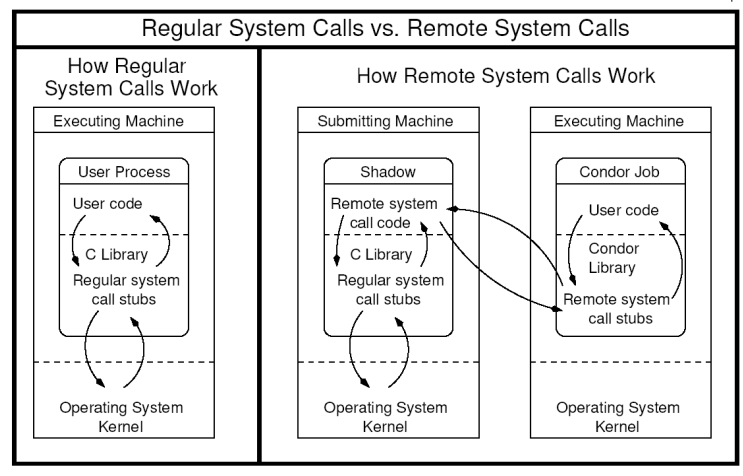


Remote I/O operations

- Shared file system
 - Bandwidth to central server limits possible number of nodes. (or performance)
 - Requires account on file server
- Condor Remote System operations
 - No special account. (easier administration)
 - Open/Read/write operation are handled on submitter machine.
 - Implemented by replacing call stub



Remote I/O operations

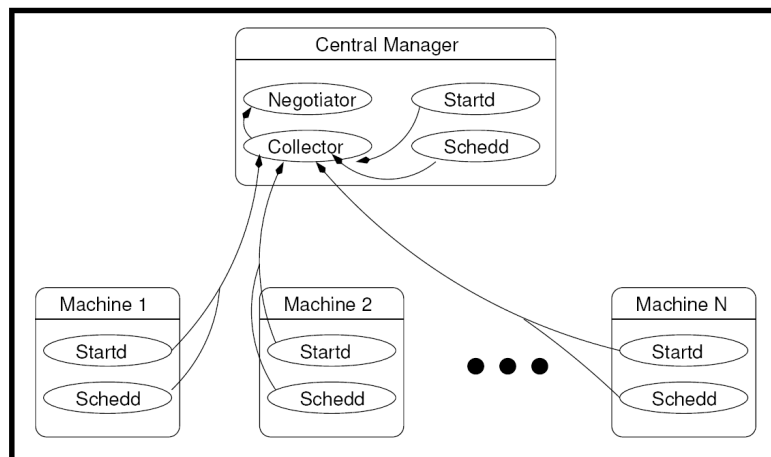


Architecture

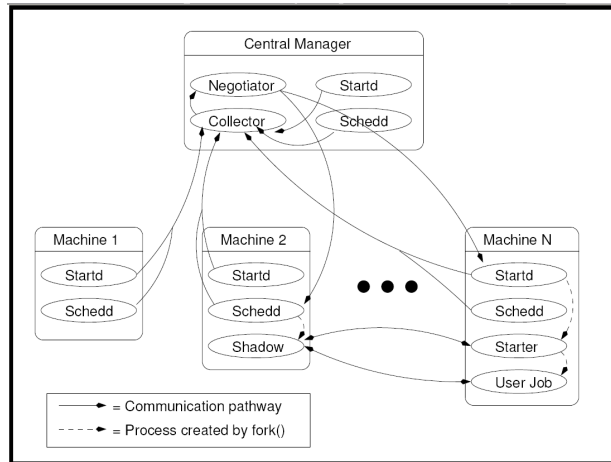


- One central manager
 - Collects Class-adds
 - Match jobs to machines
- Nodes
 - Contains local job-queue
 - Negotiates with central manager

Deamons, Idle



Deamons, job



Deamons (1)



- *"condor_master"* process runs on each node
 - Can restart other daemons
 - Enable network update of binary files
 - Inform administrator of problems
 - Enable remote start/stop of daemons

Deamons (2)

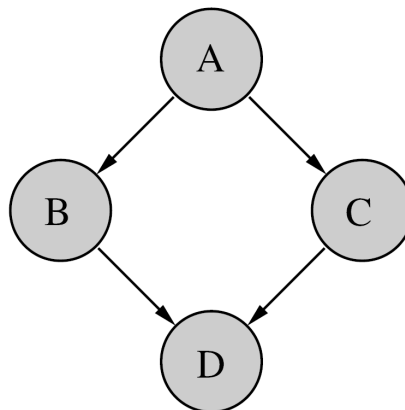


- “*start_d*”
 - Enables the node to excecute jobs
- “*condor_schedd*”
 - Manages local job-queue
- “*condor_shadow*”
 - *Manages file-transfers, logs etc.*

Dependencies



- Dependencies are specified as a Directed acyclic graph (DAG)
- Submitted via the DAGMan scheduler



SMPS in condor



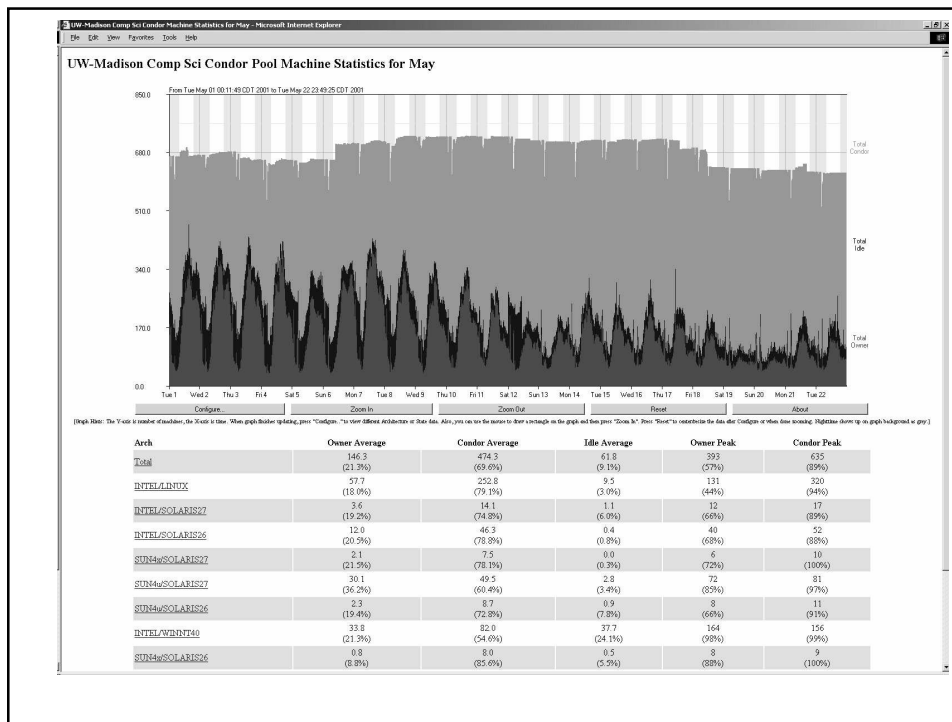
- Represented as multiple virtual machines
- Possibility to manually specify resources for each VM
 - Cpu's
 - Memory
 - Swap
 - Disk



Administration



- Rich set of tools for administration
- Includes
 - Status of pool, including detailed information about each node. (current load, memory etc)
 - Display job stats, including finished jobs
 - Changing scheduling policy on-the-fly



Conclusion

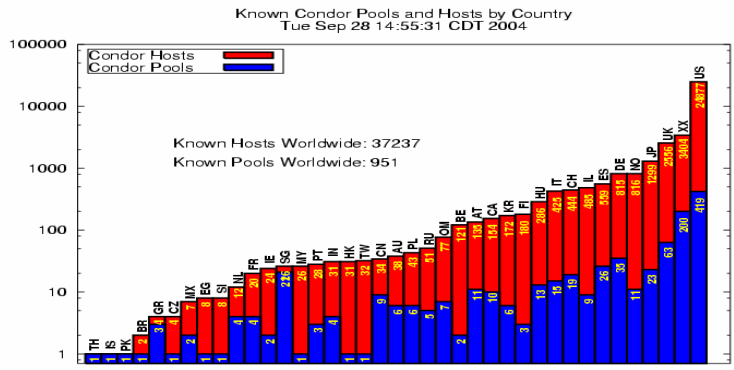
- Powerfull tool for scheduling jobs across multiple platforms
- Works with deticated clusters such as beowolf
- Utilizes resources on non-dedicated PCs
- Includes checkpointing and job migration.

More information

<http://www.cs.wisc.edu/condor/>



Known Condor Pools



Questions?



42