



# Matchmaker Policies: Users and Groups

HTCondor Week, Madison 2017

**Jaime Frey ([jfrey@cs.wisc.edu](mailto:jfrey@cs.wisc.edu))**

**Center for High Throughput Computing  
Department of Computer Sciences  
University of Wisconsin-Madison**

# HTCondor scheduling policy

- › So you have some resources...  
... how does HTCondor decide which job to run?
- › The admin needs to define a policy that controls the relative priorities
- › What defines a “good” or “fair” policy?

# First Things First

- › HTCondor does not share the same model of, for example, PBS, where jobs are placed into a first-in-first-out queue
- › It instead is based around a concept called “Fair Share”
  - Assumes users are competing for resources
  - Aims for long-term fairness

# Spinning Pie

- › Available compute resources are “The Pie”
- › Users, with their relative priorities, are each trying to get their “Pie Slice”
- › But it’s more complicated: Both users and machines can specify preferences.
- › Basic questions need to be answered, such as “do you ever want to preempt a running job for a new job if it’s a better match”? (For some definition of “better”)

# Spinning Pie

- › First, the Matchmaker takes some jobs from each user and finds resources for them.
- › After all users have got their initial “Pie Slice”, if there are still more jobs and resources, we continue “spinning the pie” and handing out resources until everything is matched.

# Relative Priorities

- › If two users have the same relative priority, then over time the pool will be divided equally among them.
- › Over time?
- › Yes! By default, HTCondor tracks usage and has a formula for determining priority based on both current demand and prior usage
- › However, prior usage “decays” over time

# Pseudo-Example

- › Example: (A pool of 100 cores)
- › User 'A' submits 100,000 jobs and 100 of them begin running, using the entire pool.
- › After 8 hours, user 'B' submits 100,000 jobs
  
- › What happens?

# Pseudo-Example

- › Example: (A pool of 100 cores)
- › User 'A' submits 100,000 jobs and 100 of them begin running, using the entire pool.
- › After 8 hours, user 'B' submits 100,000 jobs
- › The scheduler will now allocate MORE than 50 cores to user 'B' because user 'A' has accumulated a lot of recent usage
- › Over time, each will end up with 50 cores.



# Overview of Condor Architecture

Schedd A

Greg Job1  
Greg Job2  
Greg Job3  
Ann Job1  
Ann Job2  
Ann Job3

Schedd B

Greg Job4  
Greg Job5  
Greg Job6  
Ann Job7  
Ann Job8  
Joe Job1  
Joe Job2  
Joe Job3

Central  
Manager

Usage  
History

worker

worker

worker

worker

worker

worker

# Negotiator metric: User Priority

- › Negotiator computes, stores the user prio
- › View with `condor_userprio` tool
- › Inversely related to machines allocated (lower number is better priority)
  - A user with priority of 10 will be able to claim twice as many machines as a user with priority 20

# What's a user?

- › Bob in schedd1 same as Bob in schedd2?
- › If have same UID\_DOMAIN, they are.
- › We'll talk later about other user definitions.
- › Map files can define the local user name

# User Priority

- › (Effective) User Priority is determined by multiplying two components
- › Real Priority \* Priority Factor

# Real Priority

- › Based on actual usage
- › Starts at 0.5
- › Approaches actual number of machines used over time
  - Configuration setting **PRIORITY\_HALFLIFE**
  - If **PRIORITY\_HALFLIFE** = +Inf, no history
  - Default one day (in seconds)
- › Asymptotically grows/shrinks to current usage

# Priority Factor

- › Assigned by administrator
  - Set/viewed with `condor_userprio`
  - Persistently stored in CM
- › Defaults to 1000 (`DEFAULT_PRIO_FACTOR`)
- › Allows admins to give unequal prio to different users
- › “Nice user”s have Prio Factors of 10,000,000,000

# condor\_userprio

## > Command usage:

condor\_userprio

User Name	Effective Priority	Priority Factor	In Use (wghted-hrs)	Last Usage	
lmichael@submit-3.chtc.wisc.edu	5.00	10.00	0	16.37	0+23:46
blin@osghost.chtc.wisc.edu	7.71	10.00	0	5412.38	0+01:05
osgtest@osghost.chtc.wisc.edu	90.57	10.00	47	45505.99	<now>
cxiong36@submit-3.chtc.wisc.edu	500.00	1000.00	0	0.29	0+00:09
ojalvo@hep.wisc.edu	500.00	1000.00	0	398148.56	0+05:37
wjiang4@submit-3.chtc.wisc.edu	500.00	1000.00	0	0.22	0+21:25
cxiong36@submit.chtc.wisc.edu	500.00	1000.00	0	63.38	0+21:42

# Accounting Groups (2 kinds)

- › Manage priorities across groups of users and jobs
- › Can guarantee maximum numbers of computers for groups (quotas)
- › Supports hierarchies
- › Anyone can join any group (well...)



# Accounting Groups as Alias

- › In submit file
  - `Accounting_Group = group1`
- › Treats all users as the same for priority
- › Accounting groups not pre-defined
- › Admin can enforce group membership
  - Submit transforms and submit requirements
- › `condor_userprio` replaces user with group

# Prio factors with groups

```
condor_userprio -setfactor 10 group1@wisc.edu  
condor_userprio -setfactor 20 group2@wisc.edu
```

Note that you must get UID\_DOMAIN correct

Gives group1 members twice as many resources as group2

# Accounting Groups w/ Quota

- › Must be predefined in cm's config file:

```
GROUP_NAMES = a, b, c
```

```
GROUP_QUOTA_a = 10
```

```
GROUP_QUOTA_b = 20
```

- › And in submit file:

```
Accounting_Group = a
```

```
Accounting_User = gthain
```

# Group Quotas

- › “a” limited to 10

- › “b” to 20

- › Even if idle machines

- › What is the unit?

- Slot weight.

- › With fair share for users within group

- › Can create a hierarchy of groups, quotas

- › Must be predefined in cm’s config file:

```
GROUP_NAMES = a, b, c
```

```
GROUP_QUOTA_a = 10
```

```
GROUP_QUOTA_b = 20
```

- › And in submit file:

```
Accounting_Group = a
```

```
Accounting_User = gthain
```

# GROUP\_AUTOREGROUP

- › Also allows groups to go over quota if idle machines.
- › “Last chance” round, with every submitter for themselves.

# Rebalancing the Pool

- › Match between schedd and startd can be reused to run many jobs
- › May need to create opportunities to rebalance how machines are allocated
  - New user
  - Jobs with special requirements (GPUs, high memory)

# How to Rematch

- › Have startds return frequently to negotiator for rematching
  - CLAIM\_WORKLIFE
  - Draining
  - More load on system, may not be necessary
- › Have negotiator proactively rematch a machine
  - Preempt running job to replace with better job
  - MaxJobRetirementTime can minimize killing of jobs

# Two Types of Preemption

## › Startd Rank

- Startd prefers new job
  - New job has larger startd Rank value

## › User Priority

- New job's user has better priority (deserves increased share of the pool)
  - New job has lower user prio value

## › No preemption by default

- Must opt-in



# Negotiation Cycle

- › Gets all the machine ads
- › Updates user prio info for all users
- › Computes pie slice for each user
- › For each user, finds the schedd
  - For each job (until pie slice consumed)
    - Finds all matching machines for the job
    - Sorts the machines
    - Gives the best sorted machine to the job
- › If machines and jobs left, spins pie again

# Sorting Slots: Sort Levels

## › Single sort on a five-value key

- NEGOTIATOR\_PRE\_JOB\_RANK
- **Job Rank**
- NEGOTIATOR\_POST\_JOB\_RANK
- No preemption > Startd Rank preemption > User priority preemption
- PREEMPTION\_RANK

# Negotiator Expression Conventions

- › Evaluated as if in the machine ad
  - › MY.FOO : Foo in machine ad
  - › TARGET.FOO : Foo in job ad
  - › FOO : check machine ad, then job ad for Foo
- › Use MY or TARGET if attribute could appear in either ad

# Accounting Attributes

- › Negotiator adds attributes about pool usage of job owners
- › Info about job being matched
  - `SubmitterUserPrio`
  - `SubmitterUserResourcesInUse`
- › Info about running job that would be preempted
  - `RemoteUserPrio`
  - `RemoteUserResourcesInUse`

# Group Accounting Attributes

- › More attributes when using groups
  - `SubmitterNegotiatingGroup`
  - `SubmitterAutoregroup`
  - `SubmitterGroup`
  - `SubmitterGroupResourcesInUse`
  - `SubmitterGroupQuota`
  - `RemoteGroup`
  - `RemoteGroupResourcesInUse`
  - `RemoteGroupQuota`

# If Matched machine claimed, extra checks required

## > **PREEMPTION\_REQUIREMENTS**

- Evaluated when replacing a running job with a better priority job
- If False, don't preempt

## > **PREEMPTION\_RANK**

- Of machines negotiator is willing to preempt, which one to prefer

# No-Preemption Optimization

- › `NEGOTIATOR_CONSIDER_PREEMPTION = False`
- › Negotiator completely ignores claimed startds when matching
- › Makes matching faster
- › Startds can still evict jobs, then be rematched

# Concurrency Limits

- › Manage pool-wide resources
  - E.g. software licenses, DB connections
- › In central manager config
  - › `FOO_LIMIT = 10`
  - › `BAR_LIMIT = 15`
- › In submit file
  - › `concurrency_limits = foo,bar:2`



# Summary

- › Many ways to schedule