



# **Intermediate Condor**

## **Monday morning, 10:45am**

Alain Roy <roy@cs.wisc.edu>  
OSG Software Coordinator  
University of Wisconsin-Madison

# Before we begin...

---

- Any questions on the lectures or exercises up to this point?





Open Science Grid

---

**How can my jobs access their data?**



# Access to Data in Condor

---

- Use shared filesystem if available
  - Not available for today's exercises
- No shared filesystem?
  - **Condor can transfer files**
    - Can automatically send back changed files
    - Atomic transfer of multiple files
    - Can be encrypted over the wire
    - This is what we'll do in the exercises
  - Remote I/O Socket
  - Standard Universe can use remote system calls (more on this later)



# Condor File Transfer

---

- `ShouldTransferFiles = YES`
  - Always transfer files to execution site
- `ShouldTransferFiles = NO`
  - Rely on a shared filesystem
- `ShouldTransferFiles = IF_NEEDED`
  - Will automatically transfer the files if the submit and execute machine are not in the same `FileSystemDomain`

```
Universe      = vanilla
```

```
Executable   = my_job
```

```
Log           = my_job.log
```

```
ShouldTransferFiles = IF_NEEDED
```

```
Transfer_input_files = dataset$(Process), common.data
```

```
Queue 600
```



**Some of the machines in the Pool  
do not have enough memory or  
scratch disk space to run my job!**



# Specify Requirements

---

- An expression (syntax similar to C or Java)
- Must evaluate to True for a match to be made

```
Universe      = vanilla
Executable    = my_job
Log           = my_job.log
InitialDir    = run_$(Process)
Requirements  = Memory >= 256 && Disk > 10000
Queue        600
```

# Specify Rank

---

- All matches which meet the requirements can be sorted by preference with a Rank expression.
- Higher the Rank, the better the match

```
Universe      = vanilla
Executable    = my_job
Log           = my_job.log
Arguments     = -arg1 -arg2
InitialDir    = run_$(Process)
Requirements  = Memory >= 256 && Disk > 10000
Rank          = (KFLOPS*10000) + Memory
Queue        600
```



# My jobs run for 20 days...

---

- What happens when they get pre-empted?
- How can I add fault tolerance to my jobs?



# Condor's Standard Universe to the rescue!

---

- Condor can support various combinations of features/environments in different “Universes”
- Different Universes provide different functionality for your job:
  - Vanilla: Run any serial job
  - Standard: Support for transparent process checkpoint and restart

# Process Checkpointing

---

- Condor's process checkpointing mechanism saves the entire state of a process into a checkpoint file
  - Memory, CPU, I/O, etc.
- The process can then be restarted from right where it left off
- Typically no changes to your job's source code needed—however, your job must be relinked with Condor's Standard Universe support library



Open Science Grid

# Relinking Your Job for Standard Universe

---

To do this, just place “condor\_compile” in front of the command you normally use to link your job:

```
% condor_compile gcc -o myjob myjob.c
```

- OR -

```
% condor_compile f77 -o myjob filea.f  
fileb.f
```

# Limitations of the Standard Universe

---

- Condor's checkpointing is not at the kernel level. Thus in the Standard Universe the job may not:
  - fork()
  - Use kernel threads
  - Use some forms of IPC, such as pipes and shared memory
- Many typical scientific jobs are OK
- Must be same gcc as Condor was built with

# When will Condor checkpoint your job?

---

- Periodically, if desired (for fault tolerance)
- When your job is preempted by a higher priority job
- When your job is vacated because the execution machine becomes busy
- When you explicitly run:
  - `condor_checkpoint`
  - `condor_vacate`
  - `condor_off`
  - `condor_restart`

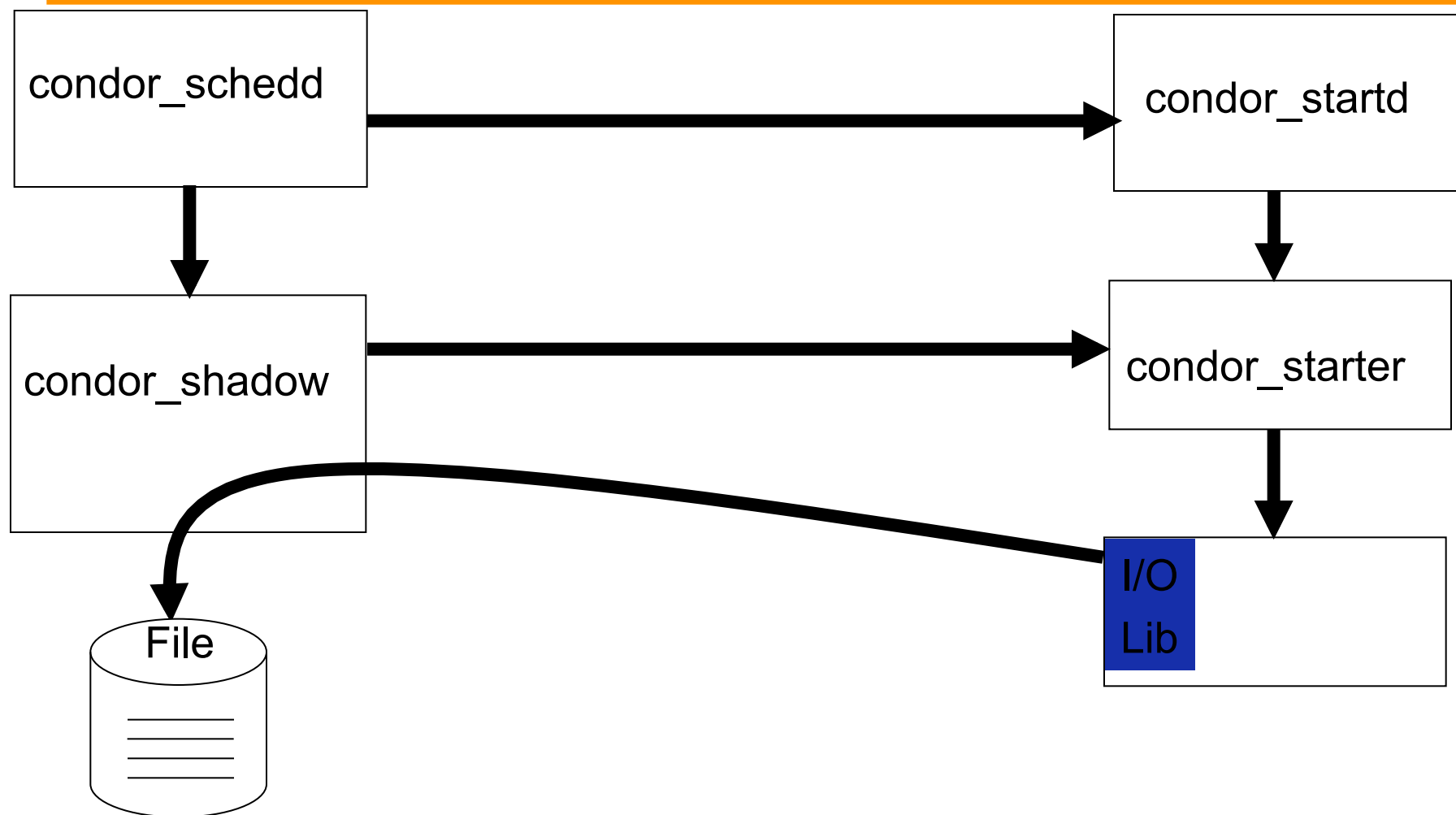


# Remote System Calls

---

- I/O system calls are trapped and sent back to submit machine
- Allows transparent migration across administrative domains
  - Checkpoint on machine A, restart on B
- No source code changes required
- Language independent
- Opportunities for application steering

# Remote I/O







# Clusters and Processes

---

- If your submit file describes multiple jobs, we call this a “cluster”
- Each cluster has a unique “cluster number”
- Each job in a cluster is called a “process”
  - Process numbers always start at zero
- A Condor “Job ID” is the cluster number, a period, and the process number (“20.1”)
- A cluster is allowed to have one or more processes.
  - There is always a cluster for every job



Open Science Grid

# Example Submit Description File for a Cluster

---

```
# Example submit description file that defines a
# cluster of 2 jobs with separate working directories
Universe      = vanilla
Executable    = my_job
log           = my_job.log
Arguments     = -arg1 -arg2
Input         = my_job.stdin
Output        = my_job.stdout
Error         = my_job.stderr
InitialDir    = run_0
Queue         Becomes job 2.0
InitialDir    = run_1
Queue         Becomes job 2.1
```

# Submitting The Job

---

```
% condor_submit my_job.submit-file
```

```
Submitting job(s).
```

```
2 job(s) submitted to cluster 2.
```

```
% condor_q
```

```
-- Submitter: perdita.cs.wisc.edu : <128.105.165.34:1027> :
```

ID	OWNER	SUBMITTED	RUN_TIME	ST	PRI	SIZE	CMD
2.0	frieda	4/15 06:56	0+00:00:00	I	0	0.0	my_job
2.1	frieda	4/15 06:56	0+00:00:00	I	0	0.0	my_job

```
2 jobs; 2 idle, 0 running, 0 held
```



Open Science Grid

# Submit Description File for a BIG Cluster of Jobs

---

- The initial directory for each job can be specified as `run_$(Process)`, and instead of submitting a single job, we use “Queue 600” to submit 600 jobs at once
- The `$(Process)` macro will be expanded to the process number for each job in the cluster (0 - 599), so we’ll have “run\_0”, “run\_1”, ... “run\_599” directories
- All the input/output files will be in different directories!



Open Science Grid

# Submit Description File for a *B/G* Cluster of Jobs

---

```
# Example condor_submit input file that defines  
# a cluster of 600 jobs with different directories
```

```
Universe      = vanilla
```

```
Executable    = my_job
```

```
Log           = my_job.log
```

```
Arguments     = -arg1 -arg2
```

```
Input         = my_job.stdin
```

```
Output        = my_job.stdout
```

```
Error         = my_job.stderr
```

```
InitialDir    = run_$(Process)
```

```
Queue 600
```

·run\_0 ... run\_599

·Creates job 3.0 ... 3.599

# More \$(Process)

- You can use \$(Process) anywhere.

```
Universe      = vanilla
```

```
Executable    = my_job
```

```
Log           = my_job.$(Process).log
```

```
Arguments     = -randomseed $(Process)
```

```
Input         = my_job.stdin
```

```
Output        = my_job.stdout
```

```
Error         = my_job.stderr
```

```
InitialDir    = run_$(Process)
```

```
Queue 600
```

·run\_0 ... run\_599

·Creates job 3.0 ... 3.599



# Sharing a directory

- You don't have to use separate directories.
- `$(Cluster)` will help distinguish runs

```
Universe      = vanilla
Executable    = my_job
Arguments     = -randomseed $(Process)
Input         = my_job.input.$(Process)
Output        = my_job.stdout.$(Cluster).$(Process)
Error         = my_job.stderr.$(Cluster).$(Process)
Log           = my_job.$(Cluster).$(Process).log
Queue 600
```

# Job Priorities

---

- Are some of the jobs in your sweep more interesting than others?
- `condor_prio` lets you set the job priority
  - Priority relative to your jobs, not other peoples
  - Priority can be any integer
- Can be set in submit file:
  - `Priority = 14`





# What if you have **LOTS** of jobs?

- Set system limits to be high:
  - Each job requires a shadow process
  - Each shadow requires file descriptors and sockets
  - Each shadow requires ports/sockets
- Each condor\_schedd limits max number of jobs running
  - Default is 200
  - Configurable: can be quite high (2000+)
- Consider multiple submit hosts
  - You can submit jobs from multiple computers
  - Immediate increase in scalability & complexity
- We constantly strive to improve scalability

# Advanced Trickery

---

- You submit 10 parameter sweeps
- You have five classes of parameters sweeps
  - Call them A, B, C, D, E
- How can you look at the status of jobs that are part of Type B parameter sweeps?

## Advanced Trickery cont.

---

- In your job file:  
+SweepType = "B"
- You can see this in your job ClassAd  
`condor_q -l`
- You can show jobs of a certain type:  
`condor_q -constraint 'SweepType == "B"'`
- Very useful when you have a complex variety of jobs
- Try this during the exercises!
- Be careful with the quoting...

# Time for more exercises!



# Questions?

---

- Questions? Comments?
- Feel free to ask me questions later:  
Alain Roy <roy@cs.wisc.edu>
- Upcoming sessions
  - Now – 12:15
    - Hands-on exercises
  - 12:15 – 1:15
    - Lunch
    - Room 2310