# Introduction to Makeflow and Work Queue



Prof. Douglas Thain, University of Notre Dame

http://www.nd.edu/~dthain
dthain@nd.edu
@ProfThain

# Lots of information here:
# http://ccl.cse.nd.edu

# The Cooperative Computing Lab

- We *collaborate with people* who have large scale computing problems in science, engineering, and other fields.
- We *operate computer systems* on the O(10,000) cores: clusters, clouds, grids.
- We *conduct computer science* research in the context of real people and problems.
- We *develop open source software* for large scale distributed computing.

## http://www.nd.edu/~ccl

# Science Depends on Computing!



AGTCCGTACGATGCTATTAGCGAGCGTGA...

# Condor Cycles at Notre Dame

# Users of Condor Cycles

# Superclusters at Amazon

I can get as many machines
on the cloud/grid as I want!

How do I organize my application
to run on those machines?

# The Cooperative Computing Tools

# Our Philosophy:

- Harness all the resources that are available: desktops, clusters, clouds, and grids.
- Make it easy to scale up from one desktop to national scale infrastructure.
- Provide familiar interfaces that make it easy to connect existing apps together.
- Allow portability across operating systems, storage systems, middleware…
- Make simple things easy, and complex things possible.
- ***No special privileges required.***

# A Quick Tour of the CCTools

- Open source, GNU General Public License.
- Compiles in 1-2 minutes, installs in $HOME.
- Runs on Linux, Solaris, MacOS, Cygwin, FreeBSD, …
- Interoperates with many distributed computing systems.
  - Condor, SGE, Torque, Globus, iRODS, Hadoop…
- Components:
  - Makeflow – A portable workflow manager.
  - Work Queue – A lightweight distributed execution system.
  - All-Pairs / Wavefront / SAND – Specialized execution engines.
  - Parrot – A personal user-level virtual file system.
  - Chirp – A user-level distributed filesystem.

http://ccl.cse.nd.edu/software

# Makeflow:
# A Portable Workflow System

# An Old Idea: Makefiles



part1 part2 part3: input.data split.py
    ./split.py input.data

out1: part1 mysim.exe
    ./mysim.exe part1 >out1

out2: part2 mysim.exe
    ./mysim.exe part2 >out2

out3: part3 mysim.exe
    ./mysim.exe part3 >out3

result: out1 out2 out3 join.py
    ./join.py out1 out2 out3 > result

# Makeflow = Make + Workflow



- Provides portability across batch systems.
- Enable parallelism (but not too much!)
- Trickle out work to batch system.
- Fault tolerance at multiple scales.
- Data and resource management.

## Makeflow

| Local | Condor | Torque | Work Queue |

http://ccl.cse.nd.edu/software/makeflow

# Makeflow Syntax

[output files] : [input files]
[command to run]

One Rule



out.txt : calib.dat in.dat sim.exe
./sim.exe –p 50 in.data > out.txt

You must state
all the files
needed by the command.

# sims.mf

```
out.10 : in.dat calib.dat sim.exe
        ./sim.exe –p 10 in.data > out.10


out.20 : in.dat calib.dat sim.exe
        ./sim.exe –p 20 in.data > out.20


out.30 : in.dat calib.dat sim.exe
        ./sim.exe –p 30 in.data > out.30
```

# How to run a Makeflow

- Run a workflow locally, using multiple cores:
  - makeflow -T local sims.mf
- Run the workflow on Torque:
  - makeflow –T torque sims.mf
- Run the workflow on Condor:
  - makeflow –T condor sims.mf
- Run the workflow on SLURM:
  - makeflow –T slurm sims.mf

# You should see this:

```
% makeflow -T local sims.mf
parsing sims.mf...
checking sims.mf for consistency...
sims.mf has 3 rules.
starting workflow....
submitting job: ./sim.exe -p 30 in.data > out.30
submitted job 2035
submitting job: ./sim.exe -p 20 in.data > out.20
submitted job 2036
submitting job: ./sim.exe -p 10 in.data > out.10
submitted job 2037
job 2035 completed
job 2036 completed
job 2037 completed
nothing left to do.
```

# If you do the same thing twice:

```
% makeflow -T local sims.mf
parsing sims.mf...
checking sims.mf for consistency...
sims.mf has 3 rules.
recovering from log file sims.mf.makeflowlog...
starting workflow....
nothing left to do.
```

Makeflow keeps a log of operations, so it knows which jobs have been sent to the batch system, and which files have already been created.

# Automatically clean outputs:

```
% makeflow --clean sims.mf
parsing sims.mf...
checking sims.mf for consistency...
sims.mf has 3 rules.
recovering from log file sims.mf.makeflowlog...
cleaning filesystem...
```
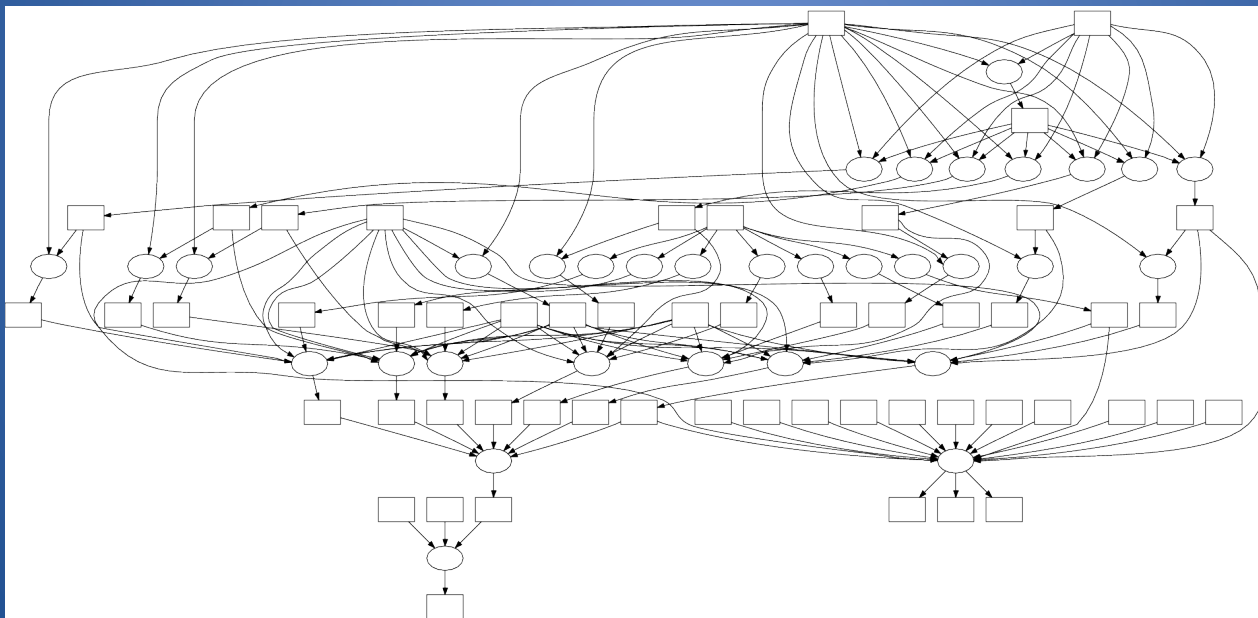
Note that you do *not* have to write a "clean" rule.
Makeflow just figures it out for you.

# Some more handy options:

- Limit the number of jobs running at once:

  --max-local #

  --max-remote #

- Retry jobs that have a tendency to fail:

  --retry-count=5

- Send email when the workflow is done:

  --email user@domain.com

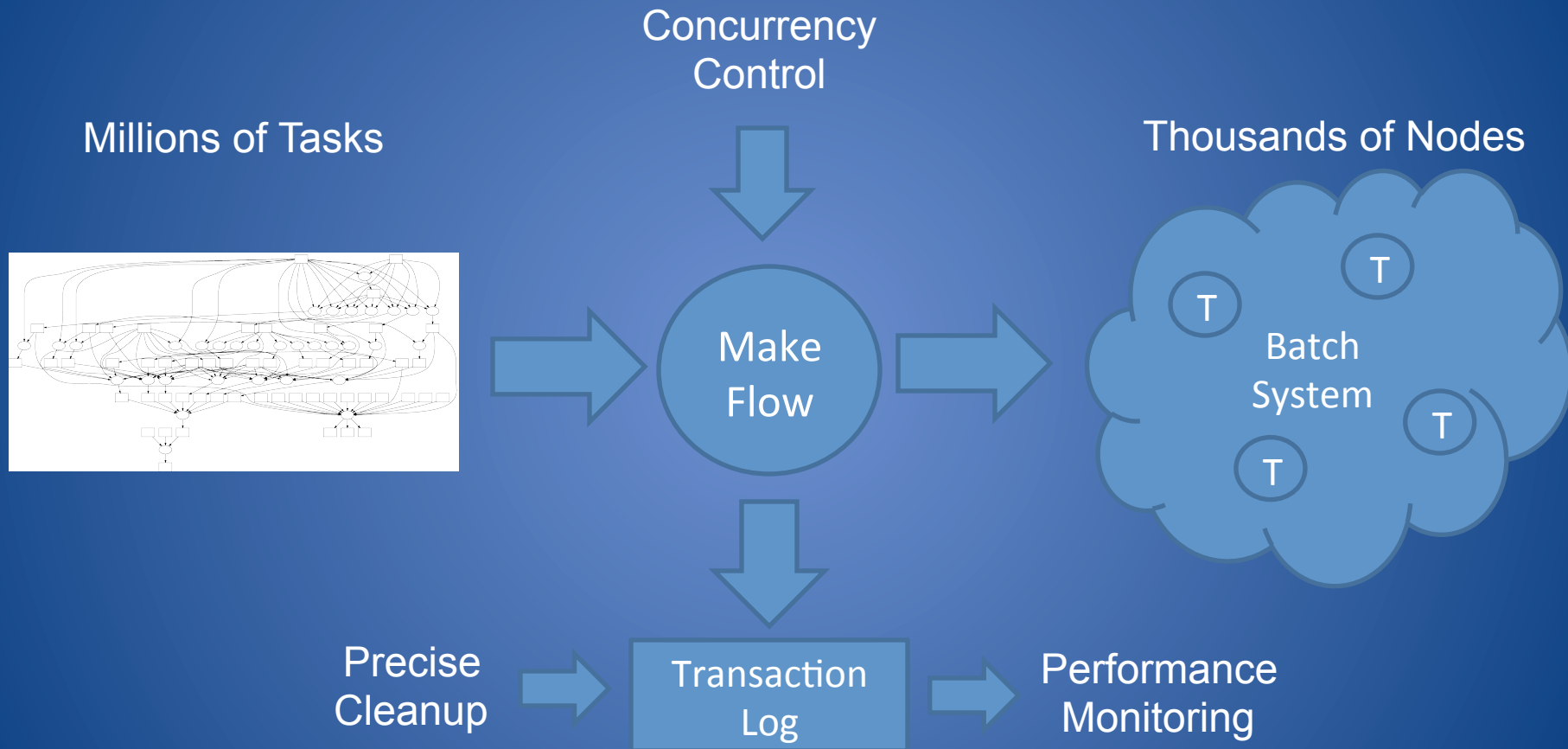- Monitor the resources consumed by each job:

  --monitor <output-dir>

# Visualization with DOT

- makeflow_viz –D example.mf > example.dot
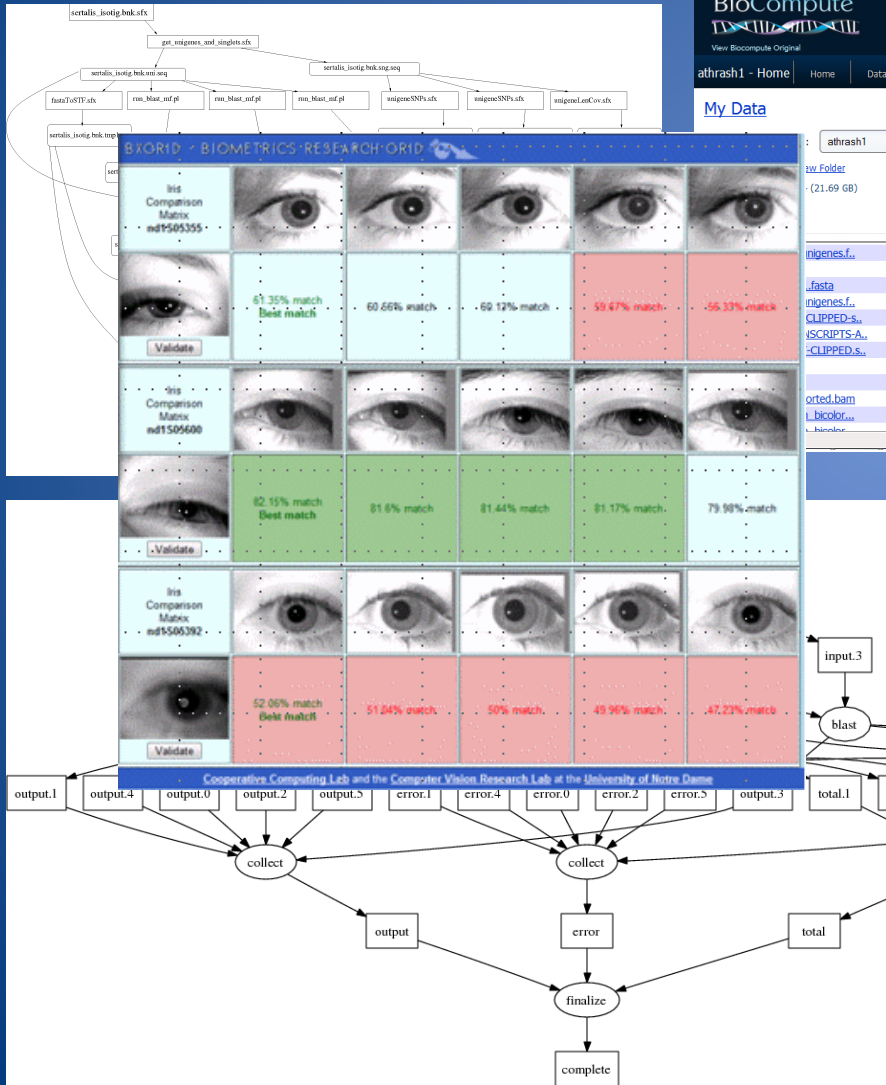- dot –T gif < example.dot > example.gif



DOT and related tools:
http://www.graphviz.org

# Makeflow Shapes a Workflow

Concurrency
Control

Millions of Tasks



Thousands of Nodes

Make
Flow

Batch
System

T
T
T
T

Precise
Cleanup

Transaction
Log

Performance
Monitoring

# Example: Biocompute Portal



BLAST
SSAHA
SHRIMP
EST
MAKER
...

Generate Makefile

Progress Bar

Transaction Log

Update Status

Run Workflow

Make flow

Submit Tasks

Condor Pool

# Makeflow Applications

# Makeflow + Work Queue

Makeflow can send jobs to one batch system at a time.

# Advantages of Work Queue

- Harness multiple resources simultaneously.
- Hold on to cluster nodes to execute multiple tasks rapidly.  (ms/task instead of min/task)
- Scale resources up and down as needed.
- Better management of data, with local caching for data intensive tasks.
- Matching of tasks to nodes with data.

# Makeflow and Work Queue

First, start the Makeflow:

% makeflow –T wq  sims.mf

Could not create work queue on port 9123.

Whoops, try again:

% makeflow –T wq --port 0 sims.mf

Listening for workers on port 8374…

Start one worker and tell it where to find makeflow:

% work_queue_worker  master.hostname.org 8374

# Start 25 Workers in Batch System

Submit workers to Condor:

condor_submit_workers master.hostname.org 8374 25


Submit workers to SGE:

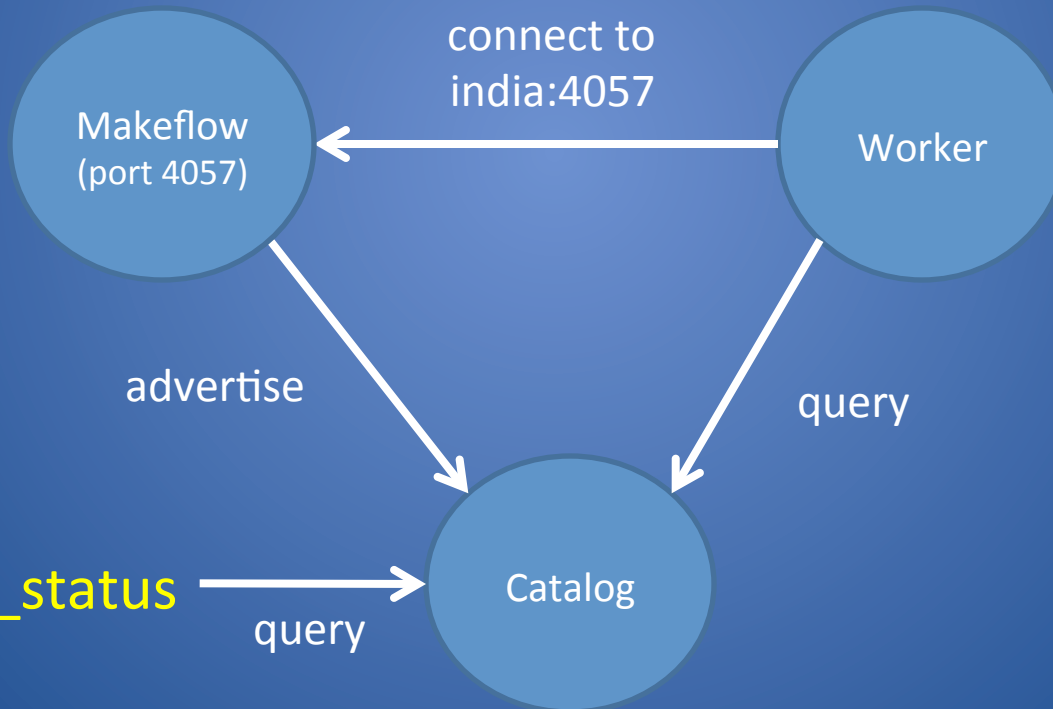sge_submit_workers master.hostname.org 8374 25


Submit workers to Torque:

torque_submit_workers master.hostname.org 8374 25

Keeping track of port numbers gets old fast...

# Project Names

makeflow …
–N myproject

work_queue_worker
–N myproject

connect to
india:4057

Makeflow
(port 4057)

Worker

advertise

query

work_queue_status

query

Catalog

"myproject"
is at india:4057

# Project Names

Start Makeflow with a  project name:

% makeflow –T wq –N myproject  sims.mf

Listening for workers on port XYZ…

Start one worker:

% work_queue_worker -N myproject

Start many workers:

% torque_submit_workers –N myproject  5

# work_queue_status



```
% ./work_queue_status
PROJECT               NAME                   PORT  WAITING  BUSY  COMPLETE  WORKERS
awe-fip35             fahnd04.crc.nd.edu     1024      719  1882   1206967     1882
hfeng-gromacs-10ps    lclsstor01.crc.nd.edu  1024     4980     0   1280240      111
hfeng2-ala5           lclsstor01.crc.nd.edu  1025     2404   140   1234514      140
forcebalance          leeping.Stanford.EDU   5817     1082    26       822       26
forcebalance          leeping.Stanford.EDU   9230        0     3       147        3
fg-tutorial           login1.futuregrid.tacc 1024        3     0         0        0
%
```

# Resilience and Fault Tolerance

- MF +WQ is fault tolerant in many different ways:
  - If Makeflow crashes (or is killed) at any point, it will recover by reading the transaction log and continue where it left off.
  - Makeflow keeps statistics on both network and task performance, so that excessively bad workers are avoided.
  - If a worker crashes, the master will detect the failure and restart the task elsewhere.
  - Workers can be added and removed at any time during the execution of the workflow.
  - Multiple masters with the same project name can be added and removed while the workers remain.
  - If the worker sits idle for too long (default 15m) it will exit, so as not to hold resources idle.

# Lots more information here: http://ccl.cse.nd.edu