# Introduction to Makeflow and Work Queue with Containers

Nick Hazekamp and Kyle Sweeney
University of Notre Dame
nhazekam|ksweene3@nd.edu

# Go to http://ccl.cse.nd.edu and Click on Container Camp Tutorial

## The Cooperative Computing Lab
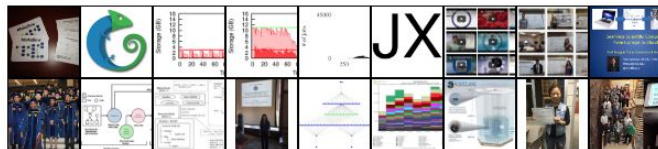
Software | Download | Manuals | Forum | Papers

Go to the CyVerse Container Camp 2018 Tutorial on Makeflow and Work Queue, March 9th!

### About the CCL

We design software that enables our collaborators to easily harness large scale distributed systems such as clusters, clouds, and grids. We perform fundamental computer science research that enables new discoveries through computing in fields such as physics, chemistry, bioinformatics, biometrics, and data mining.

### CCL News and Blog

- Submit Your CCL Highlight (15 Jan 2018)
- CCL on Chameleon Cloud with ACIC (04 Dec 2017)
- TPDS Paper: Storage Management in Makeflow (04 Dec 2017)
- CCL at Supercomputing 2017 (13 Nov 2017)
- TPDS Paper: Job Sizing (26 Oct 2017)
- Makeflow Feature: JX Representation (18 Oct 2017)
- Announcement: CCTools 6.2.0 released (09 Oct 2017)
- 2017 DISC Summer REU Conclusion (30 Aug 2017)
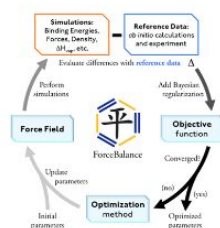- Announcement: CCTools 6.1.6 released (29 Aug 2017)
- (more news)

### Community Highlight

ForceBalance is an open source software tool for creating accurate force fields for molecular mechanics simulation using flexible combinations of reference data from experimental measurements and theoretical calculations. These force fields are used to simulate the dynamics and physical properties of molecules in chemistry and biochemistry.

The Work Queue framework gives ForceBalance the ability to distribute computationally intensive components of a force field optimization calculation in a highly flexible way. For example, each optimization cycle launched by ForceBalance may require running 50 molecular dynamics simulations, each of which may take 10-20 hours on a high end NVIDIA GPU. While GPU computing resources are available, it is rare to find 50 available GPU nodes on any single supercomputer or HPC cluster. With Work Queue, it is possible to distribute the simulations across several HPC clusters, including the Certainty HPC cluster at Stanford, the Keeneland GPU cluster managed by Georgia Tech and Oak Ridge National Laboratories, and the Stampede supercomputer managed by the University of Texas. This makes it possible to run many simulations in parallel and complete the high level optimization in weeks instead of years.

- Lee-Ping Wang, Stanford University

(Submit Your Story Here)

| Research | Software | Community | Operations |
|---|---|---|---|
| • Papers | • Download | • Annual Meeting | • Work Queue Display |
| • Projects | • Manuals | • Workshops | • Condor Display |

# 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://ccl.cse.nd.edu

# Outline

**First Session**

- Thinking Opportunistically
- Overview of the Cooperative Computing Tools
- Makeflow
- Makeflow + Work Queue
- Hands-On Tutorial

**Second Session**

- Containers in Makeflow
- Hands-On Tutorial
- Work Queue API

# Thinking Opportunistically

## Opportunistic Computing

- Much of scientific computing is done in conventional computing centers with a fixed operating environment with professional sysadmins.

- But, there exists a large amount of computing power available to end users that is not prepared or tailored to your specific application:

  - ▷ National HPC facility

  - ▷ Campus-level cluster and batch system.

  - ▷ Volunteer computing systems: Condor, BOINC, etc.

  - ▷ Cloud services.

- Can we effectively use these systems for "long tail" scientific computing?

# Opportunistic Challenges

- When borrowing someone else's machines, you cannot change the OS distribution, update RPMs, patch kernels, run as root…

- This often puts important technology just out of reach of the end user, e.g.:

  ▷ FUSE might be installed, but without setuid binary.

  ▷ Docker might be available, but you aren't a member of the required Unix group.

- The resource management policies of the hosting system may work against you:

  ▷ Preemption due to submission by higher priority users.

  ▷ Limitations on execution time and disk space.

  ▷ Firewalls only allow certain kinds of network connections.

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

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

# Cooperative Computing Tools

## Our Philosophy

- Harness all available resources: 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, FreeBSD, …

- Interoperates with many distributed computing systems.

  - Condor, SGE, Torque, Globus, iRODS, Hadoop…

- Components:

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

  - Makeflow – A portable workflow manager.

  - Work Queue – A lightweight distributed execution system.

  - Parrot – A personal user-level virtual file system.

  - Chirp – A user-level distributed filesystem.

- Provides portability across batch systems.

- Enable parallelism (but not too much!)

- Fault tolerance at multiple scales.

- Data and resource management.

**Makeflow**

| Local | Condor | SGE | Work Queue |

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

# Work Queue API

```
#include "work_queue.h"
while( not done ) {

        while (more work ready) {
        task = work_queue_task_create();
                // add some details to the task
                work_queue_submit(queue, task);
        }

        task = work_queue_wait(queue);
        // process the completed task
}
```

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

# Parrot Virtual File System

Unix Appl

Capture System Calls via ptrace

Custom Namespace

/home = /chirp/server/myhome
/software = /cvmfs/cms.cern.ch/cmssoft

Parrot Virtual File System

File Access Tracing
Sandboxing
User ID Mapping
. . .

| Local | iRODS | Chirp | HTTP | CVMFS |

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

# Lots of Documentation



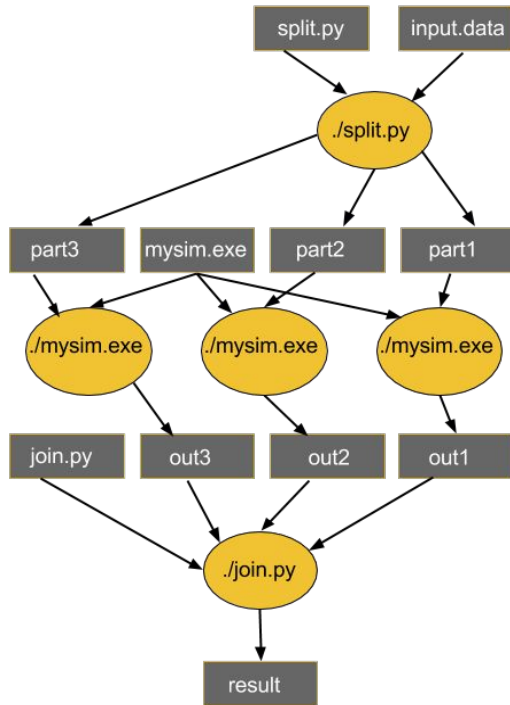http://ccl.cse.nd.edu

# Makeflow

A Portable Workflow System

# 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.
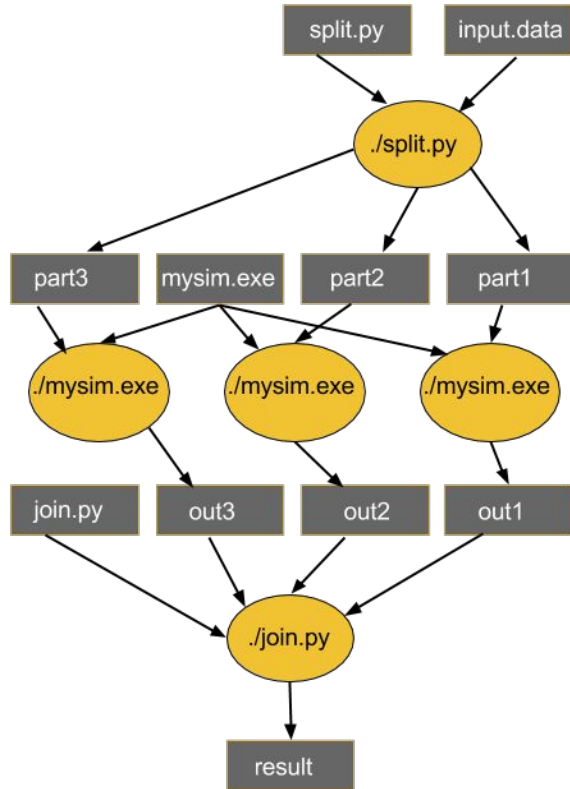
- Data and resource management.

## Makeflow

| Local | Condor | SGE | Work Queue |

# MAKEFLOW (MAKE + WORKFLOW)
# BASED OFF 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
```
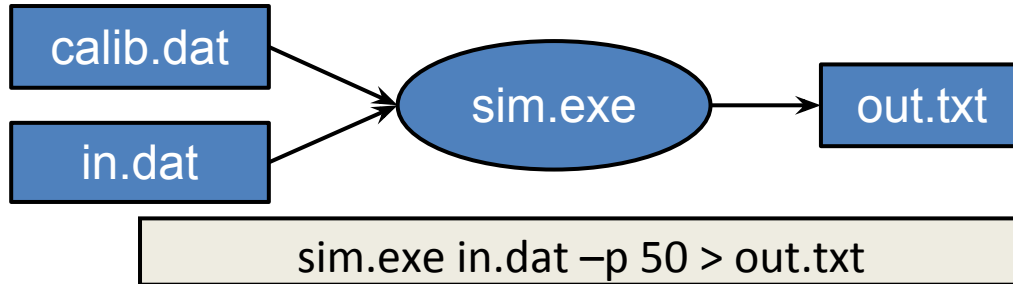
**[output files] : [input files]**

  **[command to run]**

One Rule

calib.dat

in.dat

sim.exe

out.txt

sim.exe in.dat –p 50 > out.txt

**out.txt : in.dat calib.dat sim.exe**

  **sim.exe in.data –p 50 > out.txt**

```
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
```

- Run a workflow locally (multicore?)

  `makeflow  -T local sims.mf`

- Clean up the workflow outputs:

  `makeflow –c sims.mf`

- Run the workflow on Torque:

  `makeflow –T torque sims.mf`

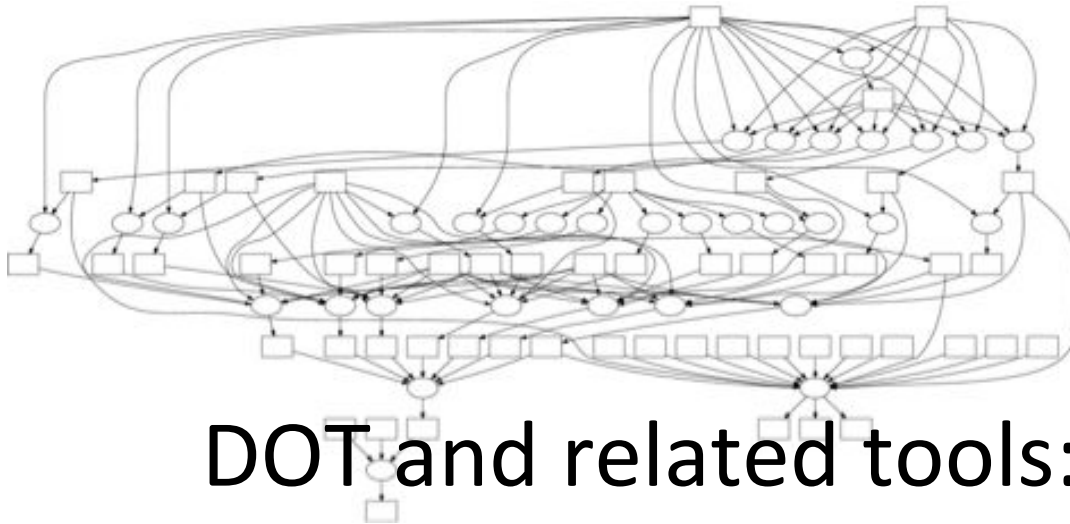- Run the workflow on Condor:

  `makeflow –T condor sims.mf`

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



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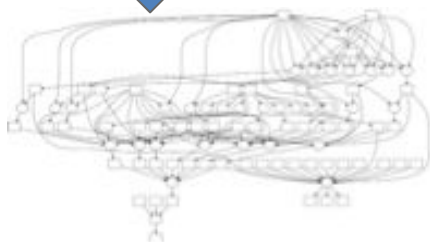http://www.graphviz.org

# Makeflow Shapes a Workflow

# Example: Biocompute Portal



BLAST
SSAHA
SHRIMP
EST
MAKER
…

Progress Bar

Generate Makeflow

Transaction Log

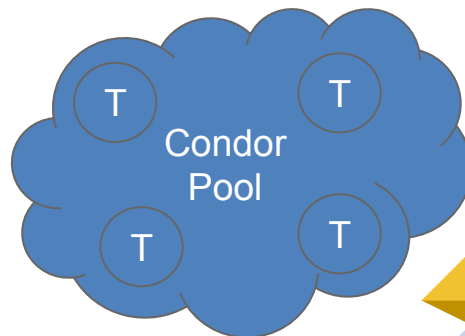Run Makeflow

Update Status

Makeflow

Condor Pool

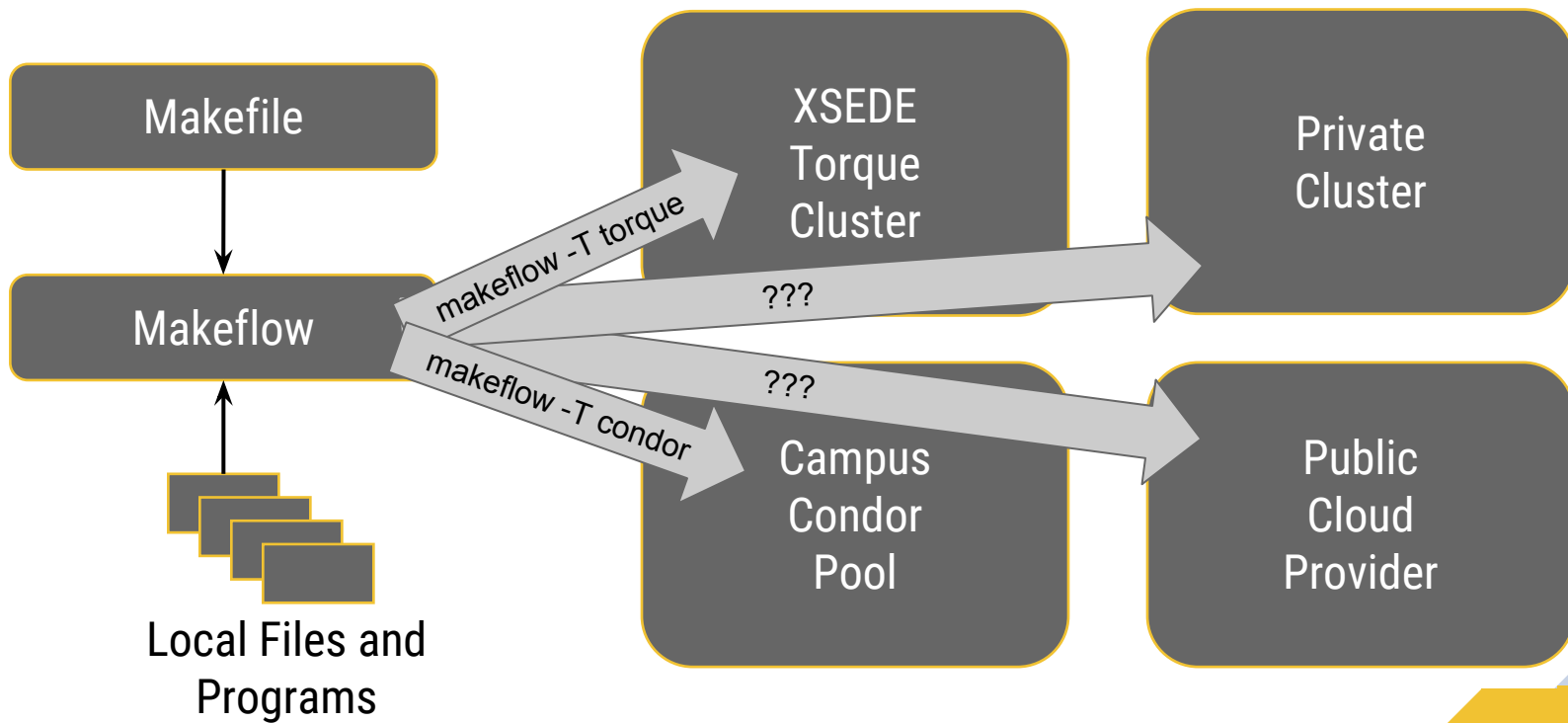T    T    T    T

# Makeflow + Work Queue

A Portable Workflow System

Makefile

Makeflow

Local Files and Programs

makeflow -T torque

makeflow -T condor

XSEDE Torque Cluster

Campus Condor Pool

Private Cluster

Public Cloud Provider

???

???

Application

API

Tasks

submit tasks

Master

XSEDE Torque Cluster

W

W

Private Cluster

W

W

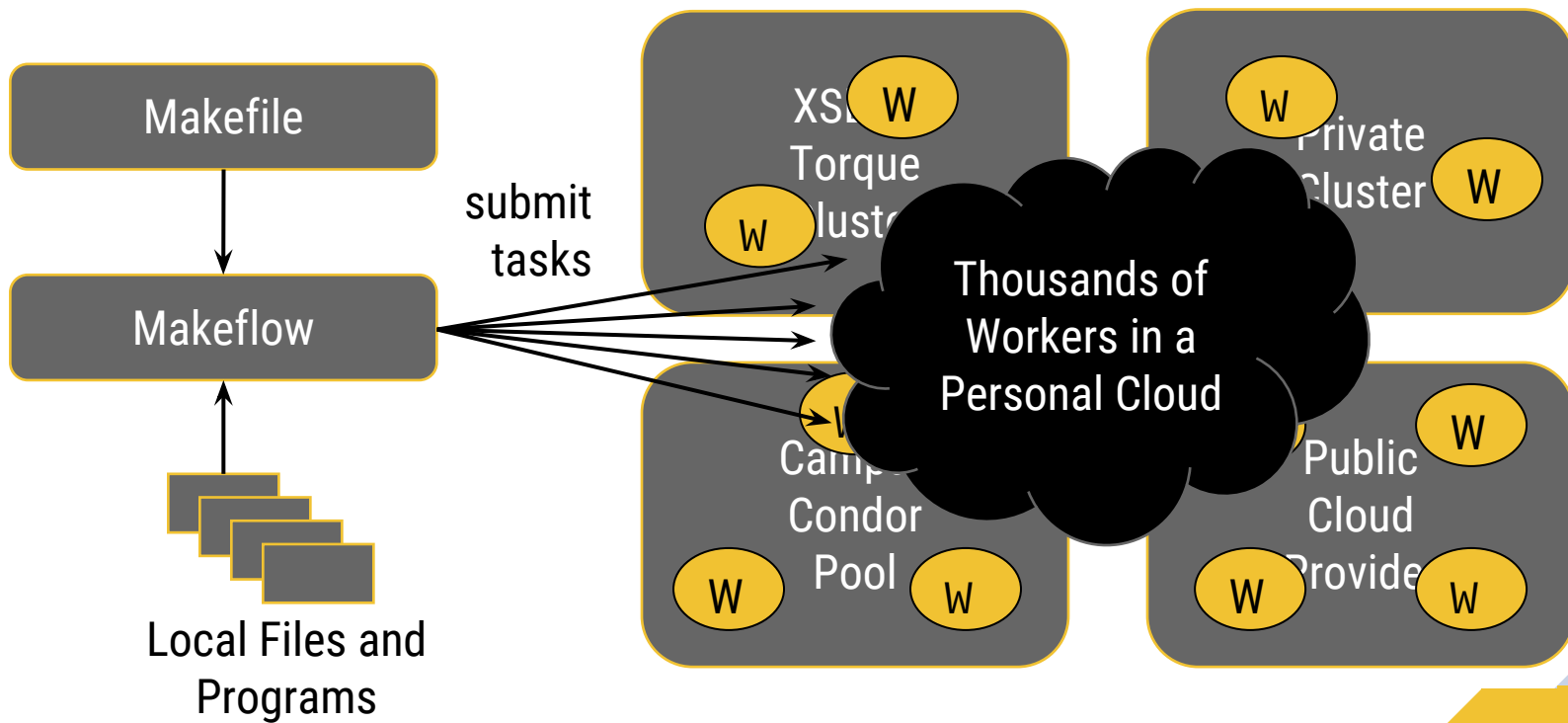Thousands of Workers in a Personal Cloud

Campus Condor Pool

W

W

W

Public Cloud Provider

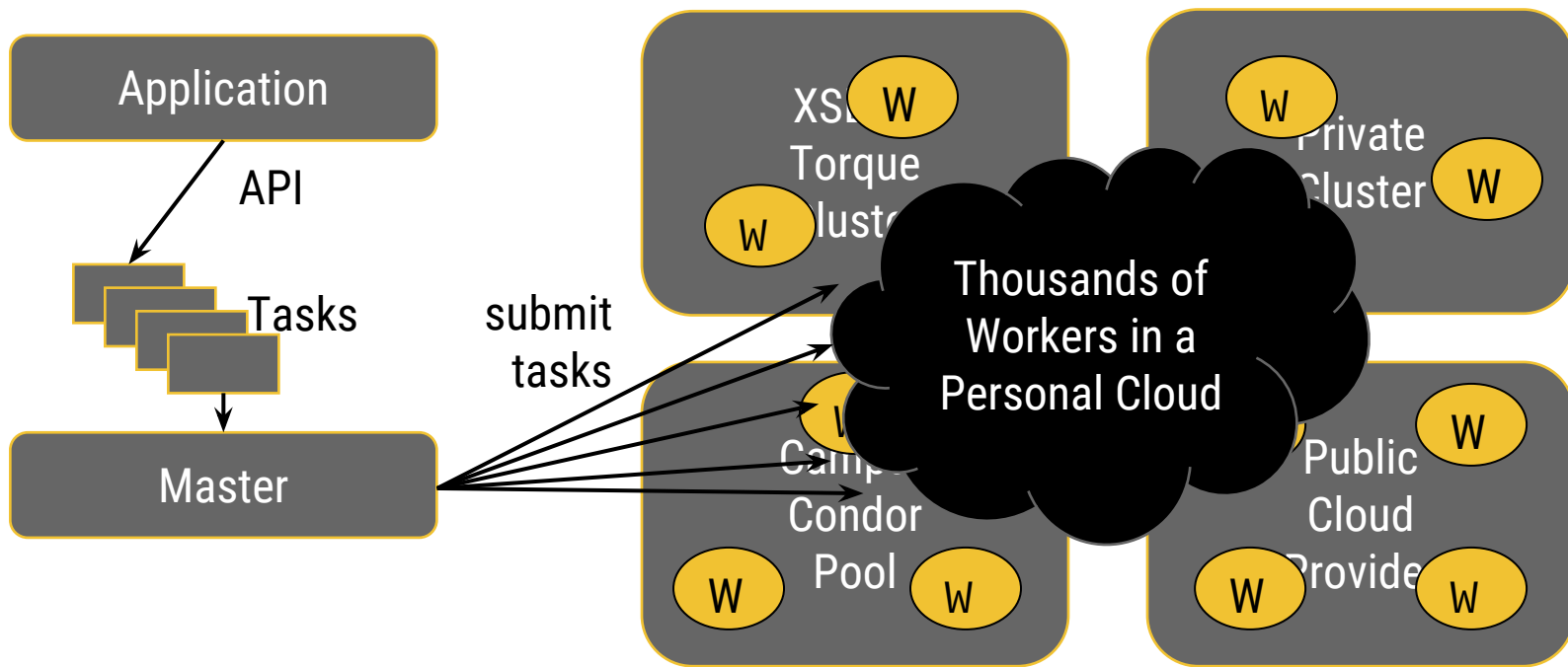W

W

W

# 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.

To start the Makeflow

% makeflow –T wq  sims.mf

Could not create work queue on port 9123.


% makeflow –T wq –p 0 sims.mf

Listening for workers on port 8374…


To start one worker:

% work_queue_worker  master.hostname.org 8374

Work Queue Factory:

work_queue_factory -T slurm -w 5 -W 25

-T : specify the batch system

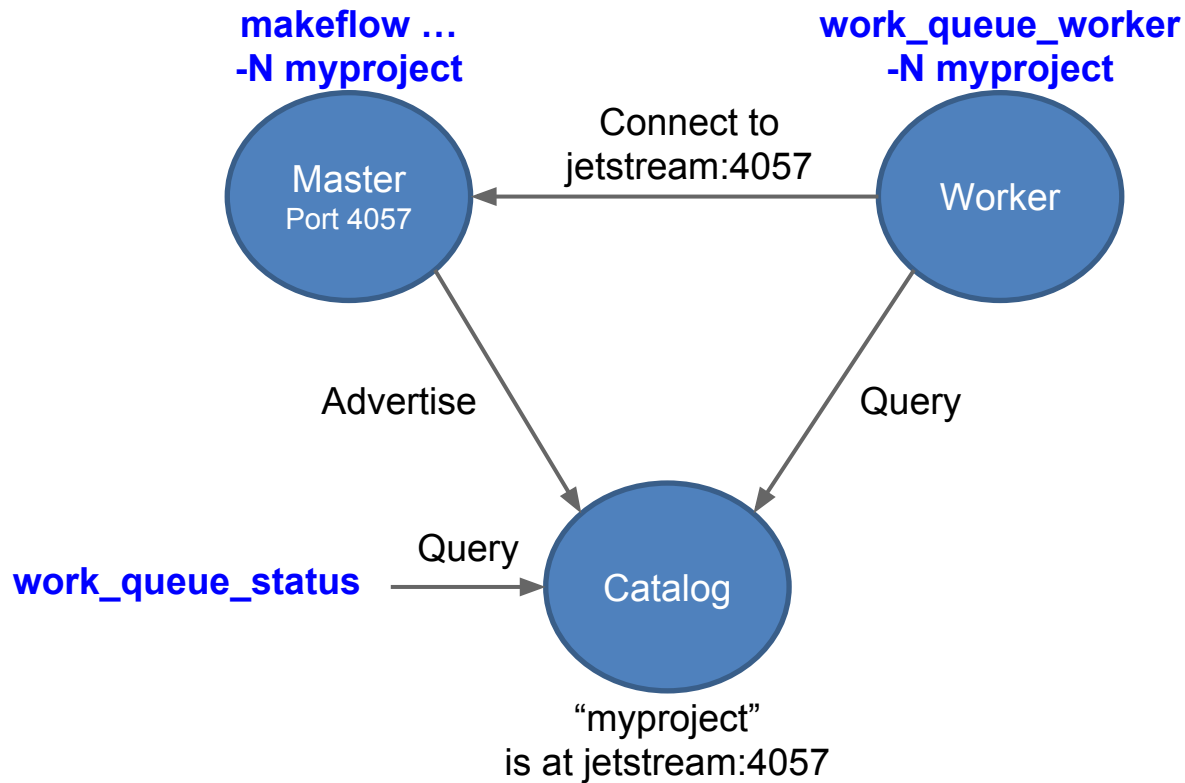-w : Set the lower limit of workers to upkeep

-W : Set the upper limit of workers to submit

# Keeping track of port numbers gets old fast...

# Project Names

## 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:

% work_queue_factory -T slurm –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
%
```

# Advantages of Work Queue

- 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 detects failure and restarts the task elsewhere.

  - Workers can be added and removed at any time during workflow execution.

  - 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.

# Alternative Makeflow Formats
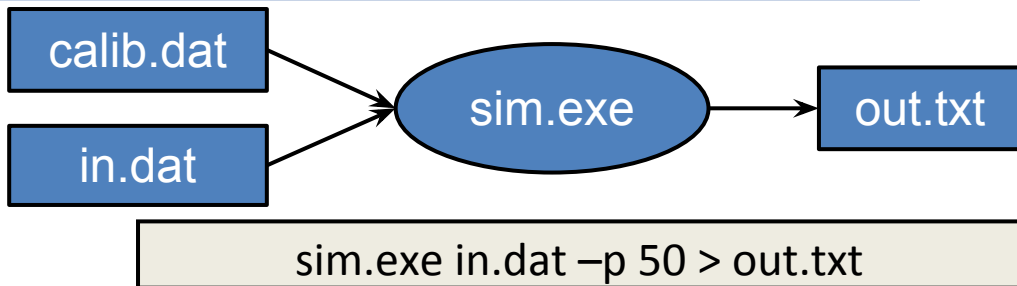
Utilizing JSON/JX for easier scripting

# Makeflow JSON Syntax

- Verbose flexible structure

- Familiar structure

- Consists of four items:

  - "categories": Object<Category>

  - "default_category": String
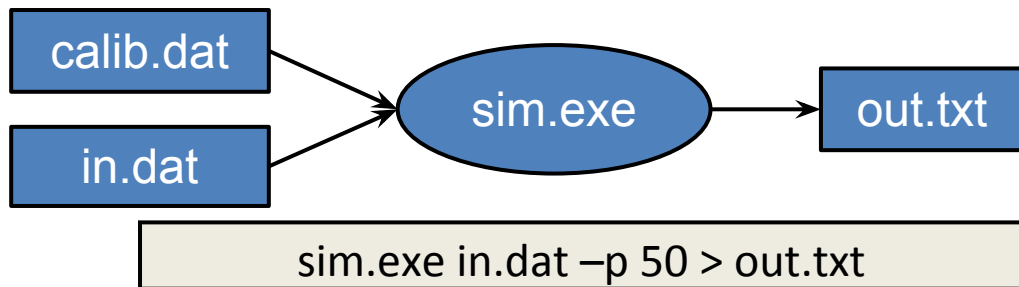
  - "environment": Object<String>

  - "rules": Array<Rule>

```
{
    "outputs": [{"dag_name": "out.txt"}],
    "inputs": [ {"dag_name": "in.dat"},  {"dag_name": "calib.dat"},
                    {"dag_name": "sim.exe"}]
    "command": "sim.exe –p 50 in.data > out.txt",

}
```

calib.dat → sim.exe → out.txt
in.dat → sim.exe

sim.exe in.dat –p 50 > out.txt

```
{
    "outputs": ["out.txt"],
    "inputs": [ "in.dat", "calib.dat", "sim.exe"]
    "command": "sim.exe –p 50 in.data > out.txt",
}
```

```
{
    "outputs": [{"dag_name": "out_10.txt"}],
    "inputs": [ {"dag_name": "in.dat"},  {"dag_name": "calib.dat"},
                {"dag_name": "sim.exe"}]
    "command": "sim.exe –p 10 in.data > out_10.txt",
},
{
    "outputs": [{"path": "out_20.txt"}],
    "inputs": [ {"dag_name": "in.dat"},  {"dag_name": "calib.dat"},
                {"dag_name": "sim.exe"}]
    "command": "sim.exe –p 20 in.data > out_20.txt",
},...
```
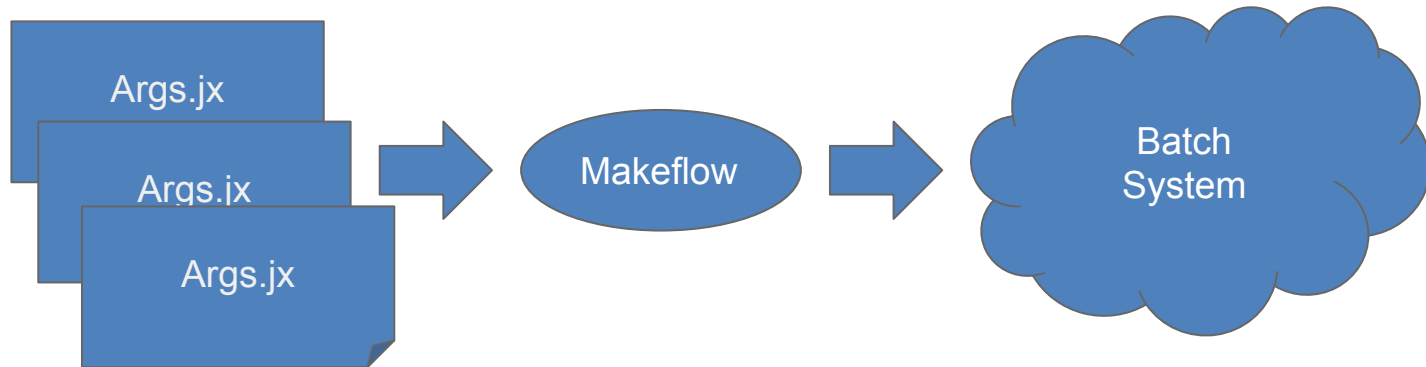
# Makeflow JSON Rule

- "inputs": Array<File>
- "outputs": Array<File>
- "command": String
- "local_job": Boolean
- "category": String
- "resources": Resources
- "allocation": String
- "environment": Object<String>

# Makeflow JX Syntax

- Allows for more compact makeflows.
  - ▷ Provides functions for expanding tasks: range, variables, etc…
- Can be used as templates in conjunction with an arguments file.
- Useful for consistently structure data and different data.

Args.jx

Args.jx

Args.jx

Makeflow

Batch System

# Makeflow JX Syntax

```
{
     "outputs": [{"dag_name": format("out_%d.txt", i)}],
     "inputs": [ {"dag_name": "in.dat"},  {"dag_name": "calib.dat"},
                      {"dag_name": "sim.exe"}]
     "command": format("sim.exe –p %d in.data > out_%d.txt", i),
} for i in range(10, 30, 10),
```

# How to run a Makeflow

- Run a workflow from json
  - makeflow  --json sims.json
- Clean up the workflow outputs:
  - makeflow –c --json sims.json
- Run the workflow from jx:
  - makeflow --jx sims.jx
- Run the workflow with jx and args:
  - makeflow --jx sims.jx --jx-args args.jx
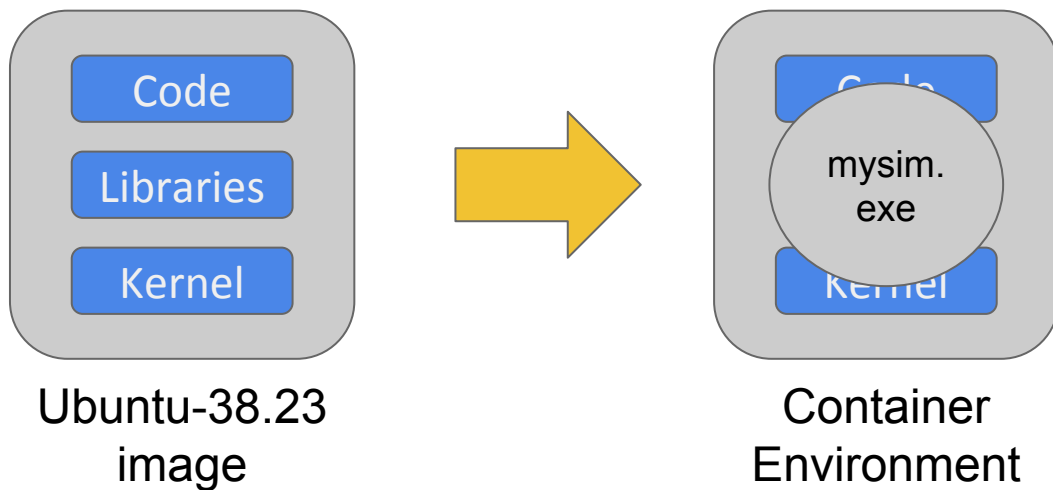
# Hands-on Tutorial

Short tutorial followed by Lunch

# Container Integration

Providing consistent environments

docker run ubuntu-38.23 mysim.exe



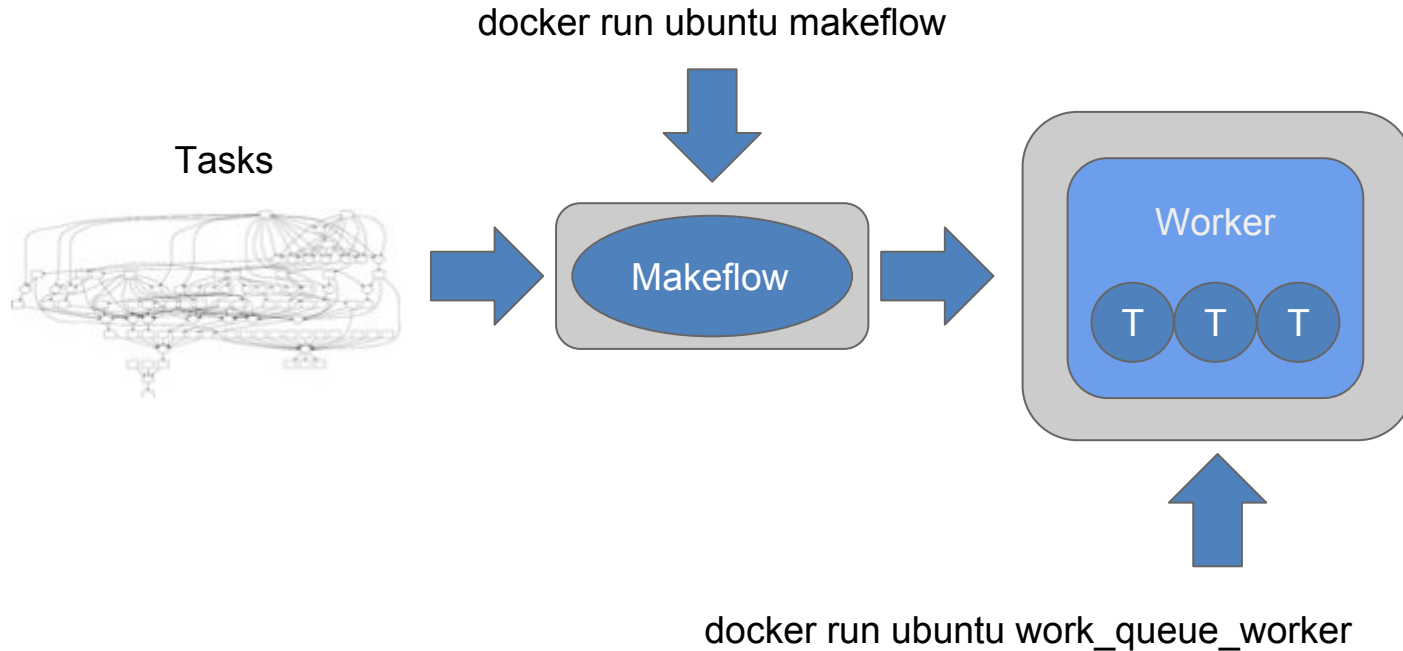Ubuntu-38.23
image

Container
Environment

# Approaches to Containers with Makeflow

- Approach 1:
  - Create containers for starting MF and WQ, then let them run as normal.
  - You are responsible for moving container images responsibly.
- Approach 2:
  - Let MF create containers as needed for each task.
  - Provides more control over moving container images.
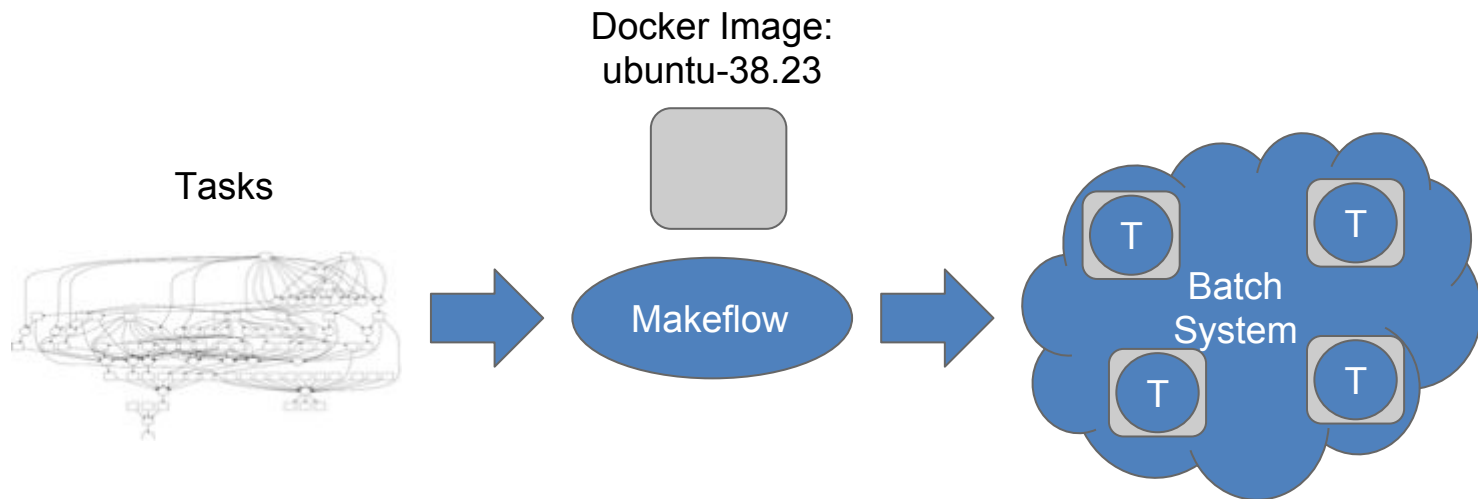  - Sending and storing containers for each task.

docker run ubuntu makeflow

Tasks

Makeflow

Worker

T  T  T

docker run ubuntu work_queue_worker

Docker Image:
ubuntu-38.23

Tasks

Makeflow

Batch
System

T

T

T
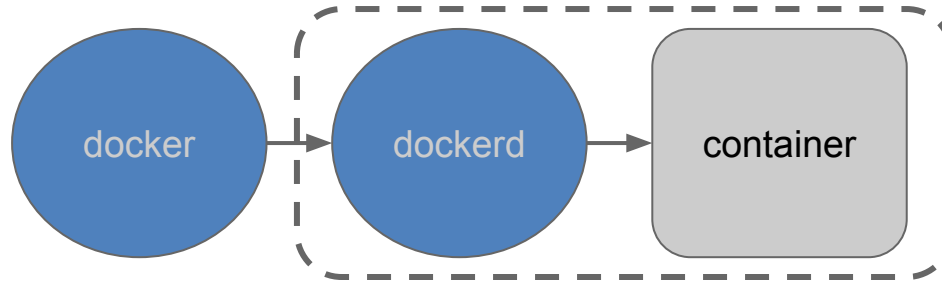
T

makeflow --docker ubuntu-38.23 –T sge . . .

# Container Technology is Evolving



docker.io

docker run ubuntu command

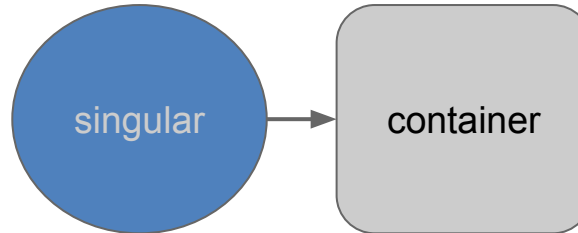docker → dockerd → container

Installed service running as root

singularity exec ubuntu command

singular → container

Container runs directly as a child process
(still needs setuid tool, though)

singularity.lbl.gov

Singularity mage:
ubuntu.img

Tasks

Makeflow

T   T

Batch
System

T   T

makeflow --singularity ubuntu.img –T sge . . .
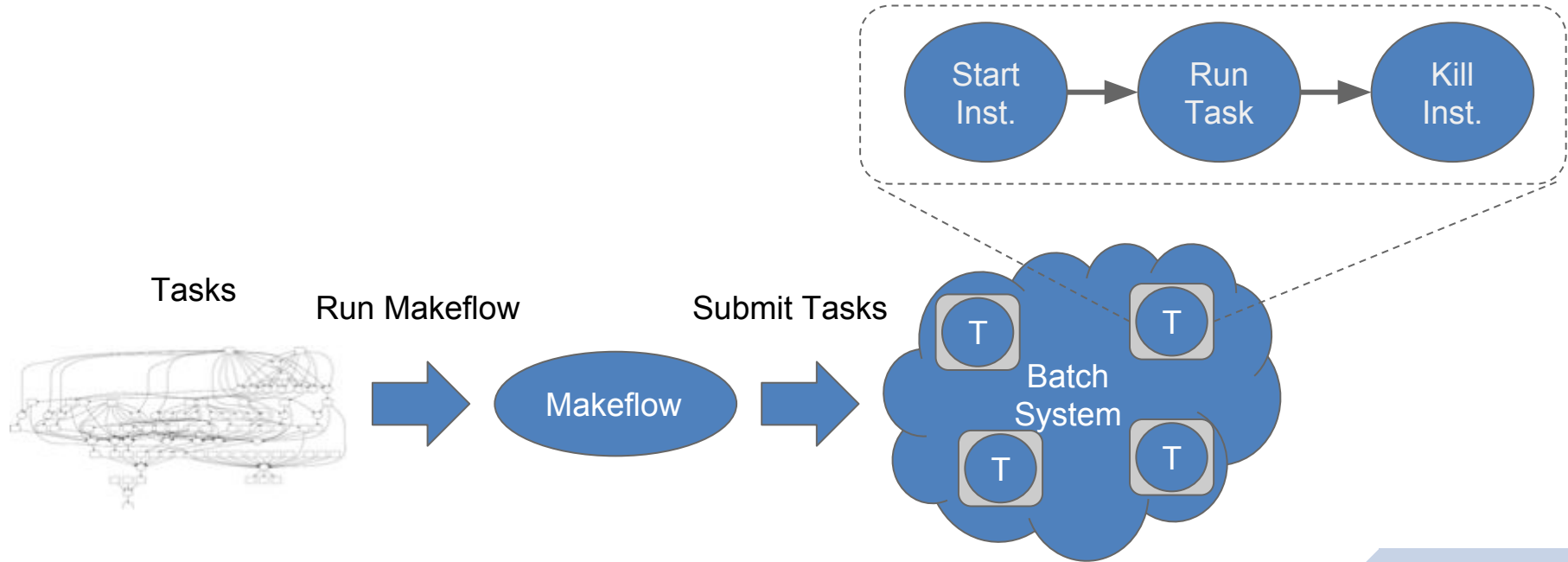
# Cloud Operation

Methods to Deploying

# Approaches to Cloud Provisioning with Makeflow

- Approach 1:

  - MF creates unique instance for each task.

  - Provides complete isolation between tasks.

  - Requires startup and tear-down time of instances.

- Approach 2:

  - Create instances and run WQ Workers on them, submitting to WQ from MF.

  - Relies on WQ for task isolation, but caches shared files.

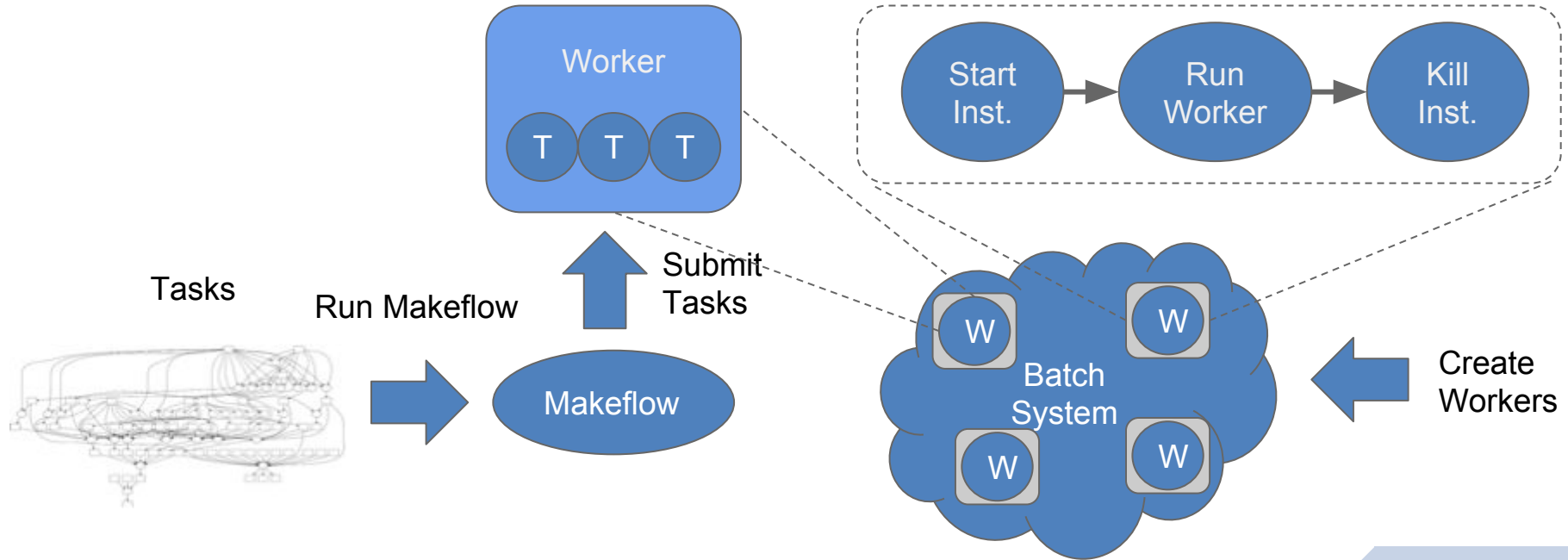  - Instance management relies on the user.

Start Inst. → Run Task → Kill Inst.

Tasks

Run Makeflow

Submit Tasks

Makeflow

Batch System

T T T T

makeflow -T amazon --amazon-config my.config ...

work_queue_factory -T amazon --amazon-config my.config

# Questions?

Nick Hazekamp
Email : nhazekam@nd.edu

Kyle Sweeney
Email:ksweene3@nd.edu

CCL Home : http://ccl.cse.nd.edu
Tutorial Link : http://ccl.cse.nd.edu/software/tutorials/cyversecc2018