

Reduction of Workflow Resource Consumption Using a Density-based Clustering Model

Qimin Zhang, Nathaniel Kremer-Herman, Benjamin Tovar, and Douglas Thain

WORKS Workshop at Supercomputing 2018



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.

Makeflow = Make + Workflow



- Provides portability across batch systems.
- Enables parallelism (but not too much!)
- Fault tolerance at multiple scales.
- Data and resource management.
- Transactional semantics for job execution.



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

Classic Make Language

MEMORY=8 CORES=4 DISK=32

output.5.txt: simulate.py /simulate.py –n 5 > output.5.txt

output.6.txt: simulate.py /simulate.py –n 6 > output.6.txt SLURM Submit Files

JX (JSON + Expressions)

```
"command": "./simulate.py -n "+N+" > output."+N+".txt",
 "inputs": [ "simulate.py" ],
 "outputs": [ "output."+N+".txt" ],
 "resources" : {
    "memory" : 8,
    "cores" : 16,
    "disk" : 128
} for N in range(1,1000)
```

```
• • •
```

Makeflow Shapes a Workflow



Example: Species Distribution Modeling



Full Workflow: 12,500 species
x 15 climate scenarios
x 6 experiments
x 500 MB per projection
= 1.1M jobs, 72TB of output



Small Example: 10 species x 10 expts







The Job Sizing Problem

How should one select the resources (cores, memory, disk) required by each job in a workflow?

- A high throughput workflow:
- hundreds of thousands of jobs;
- complex distributions of resource consumption.

Too small: the job may fail due to resource exhaustion.

Too large: the job may succeed, but resources will be wasted, fewer jobs will run, and throughput will be reduced.

The Job Sizing Problem



Job Sizing is a Client's Problem (Not the Facility's Problem.)

6

User/WMS selects allocation: Too big? Wasted resources. Too small? Job fails, retry.

Facility places the allocation: Too big? Get paid. Too small? Still get paid!

Facility handles scheduling (placement) but not the job sizing problem.



Facility



Job





Other Approaches and Their Limitations

Static Resource Allocation:

- An expert user sets the maximum resources used for all tasks.
- Easy to implement.
- Ignores different kinds of tasks.
- Ignores distribution of tasks.
- Expert knowledge is hard to come by, and may be incorrect.

Dynamic Resource Allocation:

- Resources are allocated as historical data becomes available.
- Harder to implement.
- User labels the different kinds of tasks.
- Takes into account distribution of tasks.
- Failures while distributions are learned.

Clustering Based Approach



- Max resources are set according to available historical data.
- Easier to implement.
- Different types of tasks can be learned without user intervention.
- Takes into account distribution of tasks.
- Needs apriori historical data.
- History may not represent the future.

Clustering Based Approach



Our approach can be applied When users need to:

1. Run the workflow multiple times.

2. Run a small batch of workflow before running the whole workflow.

Clustering Based Approach



A *resource monitor* tool is used to generate resource usage reports.

A density-based clustering model is applied to discover hidden clusters.

An improved resource allocation strategy is given by calculating the maximum resources of each cluster.

Run the remainder (or next) with the improved resource allocation strategy.

Resource Allocation

Data Clustering

Density-based Clustering Model

Advantages of Density-based Clustering:

- Clusters found by density can be any shape;
- The number of clusters is not required to be specified a priori.



An Example: Clustering Result of BWA-GATK Workflow

Parameters of the Clustering Model

Features:

cores, memory, disk, and wall-time

Minimal samples:

- the minimal number of sample in a cluster
- set to be 1 to ensure each job have its cluster label

Weights: (w)

- represent the importance of features (or axes)
- weights of cores, memory, disk and wall-time were set to be 0.1, 0.7, 0.1, 0.1 respectively

Parameters of the Clustering Model

Distance:

- represents inferent similarity between two points
- The distance D(X,Y) between two points $X = (x_1, x_2, x_3, x_4)$ and

$$Y = (y_1, y_2, y_3, y_4) \text{ is defined as (p=2):}$$
$$D(X, Y) = (\sum_{i=1}^4 w \times |x_i - y_i|^p)^{(1/p)}$$

eps:

- The maximum distance between two samples for them to be considered as in the same cluster
- The value of *eps* should be modified in different workflows to achieve better performance.

Experiments on Scientific Workflows



https://github.com/cooperative-computing-lab/makeflow-examples

Lifemapper-Workflow



Estimated number of clusters: 3



SHRIMP-Workflow



memory

BWA-Workflow



BWA-GATK-Workflow



Cores

Evaluation Results (Lifemapper)

Lifemapper	Job	Resource							Memory	Diek	
		cores	Memory (MB)	Disk (MB)	Time (s)	Memory (GB.min)	Disk (GB.min)	Time saving	saving (GB.min)	saving (GB.min)	Cores saving
1 st	153	8	8000	50000	1394.23	539.23	3370.22	compare compare with 1 st with 1 st 84 83% 72 34%			
2 nd	153	4	8000	50000	542.49	406.41	2540.03		compare	e compare with 1 st	1224
3 rd	153	2	8000	50000	245.87	293.22	1832.64		with 1 st 72 34%		Before clustering After clustering After clustering
After clustering	153	Cluster 1					04.0070	12.0470	00.0070	Cores sawing	
		1	100	3000	211.50	149.13	122.05	compare with 2 nd 61.01%	compare with 2 nd 63.31%	compare with 2 nd 95.59%	612
			Cluster 2								357 255 255 255 51
		2	5000	30000				compare with3 rd 13.98%	compare with 3 rd 49.14%	compare with 3 rd 93.89%	
			Cluster 3								
		2	1500	3000							63.50% 58.33% 16.67%

Evaluation Results

	Resource Saving							
Workflow	Cores	Memory (GB.min)	Disk (GB.min)					
SHRIMP	50.00%	91.14%	92.10%					
BWA	50.00%	95.44%	51.82%					
BWA-GATK	59.10%	94.98%	94.81%					

- SHRIMP, BWA, and BWA-GATK were run on an HTCondor batch system.
- Compared to naive approach, the least resource saving of cores, memory, and disk are 50.00%, 91.14%, and 51.82%, respectively.

Second Order Effect: Smaller Jobs are Easier to Schedule!

Larger resource requesting:

- harder to schedule within the cluster
- longer to wait in work queue



Smaller resource requesting:

- faster scheduling times
- faster to fullfill



Takeaways and Caveats

- Clustering-based resource allocation strategy significantly reduces resource consumption of scientific workflows.
- The density-based clustering model performs well in discovering hidden clusters or relationships between jobs.
- Caveat: User cannot avoid setting some knobs: Relative weights of resources. EPS parameter controlling max distance
- Caveat: Clustering requires some advance resource consumption. (How much pays off?)

Clustering: The Chicken or the Egg?



Unlabeled Workflow

Measure to Find Resources

Cluster Resources to find Groupings

Execute Workflow

Qimin Zhang, Nathaniel Kremer-Herman, Benjamin Tovar, and Douglas Thain, **Reduction of Workflow Resource Consumption Using a Density-based Clustering Model**, WORKS Workshop at Supercomputing, November 2018.

30

Clustering: The Chicken or the Egg?



Benjamin Tovar, Rafael Ferreira da Silva, Gideon Juve, Ewa Deelman, William Allcock, Douglas Thain, and Miron Livny, <u>A Job Sizing Strategy for High-Throughput Scientific Workflows</u>, *IEEE Transactions on Parallel and Distributed Systems*, **29**(2), pages 240-253, February, 2018. DOI: 10.1109/TPDS.2017.2762310

But manual clusters might not be compact!



MB memory (RAM) per task

How to pick the first allocation?

Ben Tovar says: Minimize probability of first attempt succeeding + fallback succeeding, weighted by resources.

$$\begin{split} \mathrm{E}[\mathrm{waste}(\mathbf{r},\tau,a_1)] &= \int_0^\infty \left(\underbrace{\int_0^{a_1} (a_1-\mathbf{r})\tau p(\mathbf{r},\tau) \mathrm{d}\mathbf{r}}_{\mathrm{first-allocation succeds}} \right. \\ &+ \underbrace{\int_{a_1}^{a_m} ((a_m+a_1-\mathbf{r})\tau p(\mathbf{r},\tau) \mathrm{d}\mathbf{r}}_{\mathrm{final allocation succeds}} \right) \mathrm{d}\tau \\ &= a_1 \underbrace{\int_{a_1}^{a_m} \int_0^\infty \tau p(\mathbf{r},\tau) \mathrm{d}\tau \mathrm{d}\mathbf{r}}_{\mathrm{mean wall-time for all tasks}} \\ &+ a_m \int_{a_1}^{a_m} \underbrace{\int_0^\infty \tau p(\tau|\mathbf{r}) \mathrm{d}\tau}_{\mathrm{mean wall-time tasks w. peak r}} p(\mathbf{r}) \mathrm{d}\mathbf{r} \end{split}$$

Big Picture Questions

- What is the minimal amount of useful a priori data?
 - Simple categories provided by the user? Estimates of past resource consumption? Access to runs from "similar" workflows and users? (But what is similar?)
- Can we determine when clustering/predictions fail?
 - Detecting one outlier is easy, but it's not practical to recluster at every step. What's the cost of doing this incrementally?
- Resource description is not entirely separate from scheduling!
 - How much faster can I schedule a 32-core job than a 128

Reduction of Workflow Resource Consumption Using a Density-based Clustering Model

Qimin Zhang, Nathaniel Kremer-Herman, Benjamin Tovar, and Douglas Thain

Support for this work is provided by an iSURE summer fellowship at Notre Dame. We thank Tim Shaffer and Kyle Sweeney, researchers at the University of Notre Dame, for assistance with Makeflow and the HTCondor distributed batch computing system.





••• • • •

ccl.cse.nd.edu

0 1 0

The Cooperative Computing Lab

=

Take the ACIC 2015 Tutorial on Makeflow and Work Queue

About the CCL

Community Highlight

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 in that enables new discoveries through computing in fields such as physics, chemistry, bioinformatics, biometrics, and data mining.

CCL News and Blog

- Global Filesystems Paper in IEEE CiSE (09 Nov 2015)
- Preservation Talk at iPres 2015 (03 Nov 2015)
- CMS Case Study Paper at CHEP (20 Oct 2015)
- OpenMalaria Preservation with Umbrella (19 Oct 2015)
- DAGVz Paper at Visual Performance Analysis Workshop (13 Oct 2015)
- <u>Virtual Wind Tunnel in IEEE CiSE</u> (09 Sep 2015)
- Three Papers at IEEE Cluster in Chicago (07 Sep 2015)
- CCTools 5.2.0 released (19 Aug 2015)
- Recent CCL Grads Take Faculty Positions (18 Aug 2015)
- (more news)





C

Software | Download | Manuals | Papers

software installations. By

Filesystem

(CVMFS), a

using Parrot, CVMFS, and additional components integrated by the Any Data, Anytime, Anywhere project, physicists working in the Compact Muon Solenoid experiment have been able to create a uniform computing environment across the Open Science Grid. Instead of maintaining large software installations at each participating institution, Parrot is used to provide access to a single highly-available CVMFS installation of the software from which files are downloaded as needed and aggressively cached for efficiency. A pilot project at the University of Wisconsin has demonstrated the feasibility of this approach by exporting excess compute jobs to run in the Open Science Grid, opportunistically harnessing 370,000 CPU-hours across 15 sites with seamless access to 400 gigabytes of software in the Wisconsin CVMFS repository.

- Dan Bradley, University of Wisconsin and the Open Science Grid



http://ccl.cse.nd.edu



13 1

9 4

...

New to Twitter?