# Resource Management for Dynamic Function Distribution with Parsl and Work Queue

Douglas Thain, Ben Tovar,
Thanh Son Phung, and Barry Sly-Delgado
Parsl / FuncX Workshop, 27 October 2021

CCTools

# Parsl + Work Queue for Scalable Apps

**CCTools**

## http://parsl-project.org



**Powerful Pythonic Workflow Programming Model**

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

**Work Queue: A Scalable Master/Worker Framework**

Work Queue is a framework for building large master-worker applications that span thousands of machines drawn from clusters, clouds, and grids. Work Queue applications are written in C, Perl, or Python using a simple API that allows users to define tasks, submit them to the queue, and wait for completion. Tasks are executed by a standard worker process that can run on any available machine. Each worker calls home to the master process, arranges for data transfer, and executes the tasks. The system handles a wide variety of failures, allowing for dynamically scalable and robust applications.
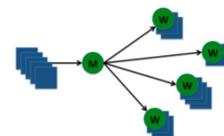
Work Queue has been used to write applications that scale from a handful of workstations up to tens of thousands of cores running on supercomputers. Examples include Lobster, NanoReactors, ForceBalance, Accelerated Weighted Ensemble, the SAND genome assembler, the Makeflow workflow engine, and the All-Pairs and Wavefront abstractions. The framework is easy to use, and has been used to teach courses in parallel computing, cloud computing, distributed computing, and cyberinfrastructure at the University of Notre Dame, the University of Arizona, and the University of Wisconsin - Eau Claire.
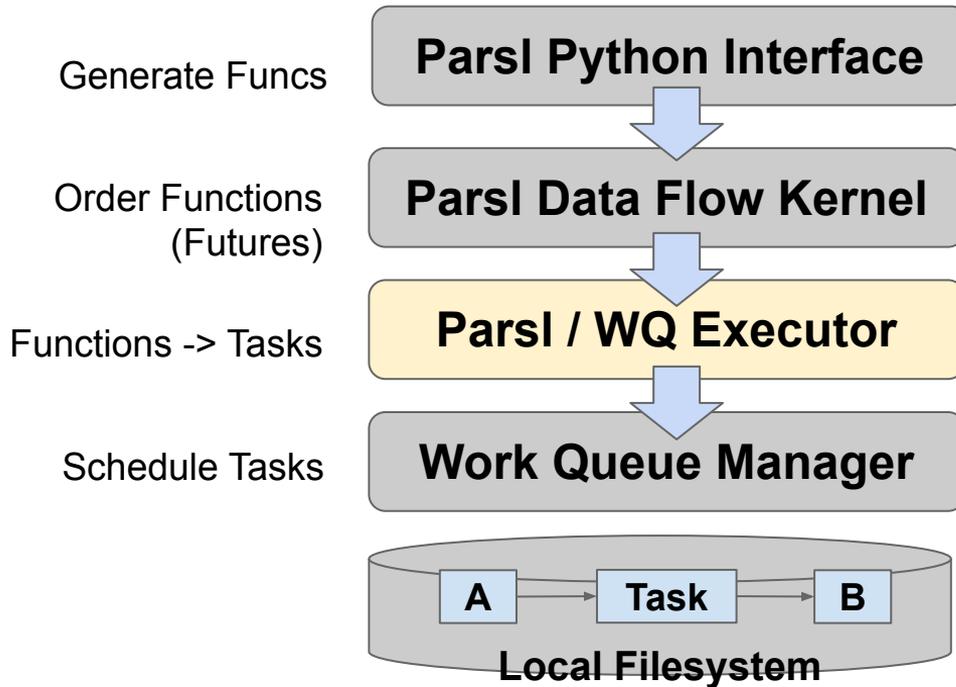
**For More Information**

- Work Queue User's Manual
- Work Queue API (C | Perl | Python)
- Work Queue Example Program (C | Perl | Python)
- Work Queue Status Display
- Download Work Queue
- Getting Help with Work Queue

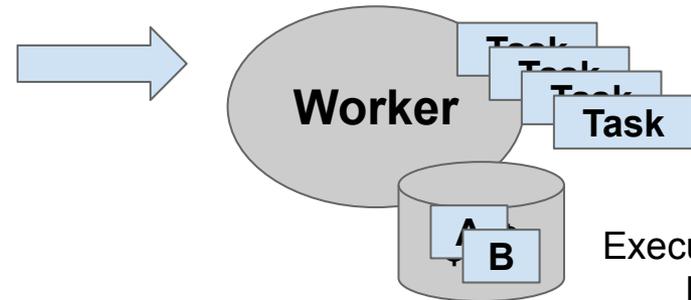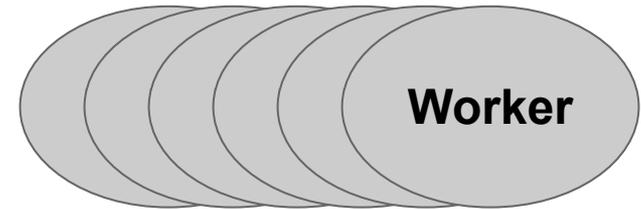**Scalable, Portable, Robust Distributed Execution**

2

# System Architecture

**Thousands of Workers on National Cyberinfrastructure**

Generate Funcs

**Parsl Python Interface**

Order Functions (Futures)

**Parsl Data Flow Kernel**

Functions -> Tasks

**Parsl / WQ Executor**

Schedule Tasks

**Work Queue Manager**

| A | → | Task | → | B |

**Local Filesystem**

HTCondor, PBS, SLURM, Amazon, Blue Waters, OSG, XSEDE...

**Worker**
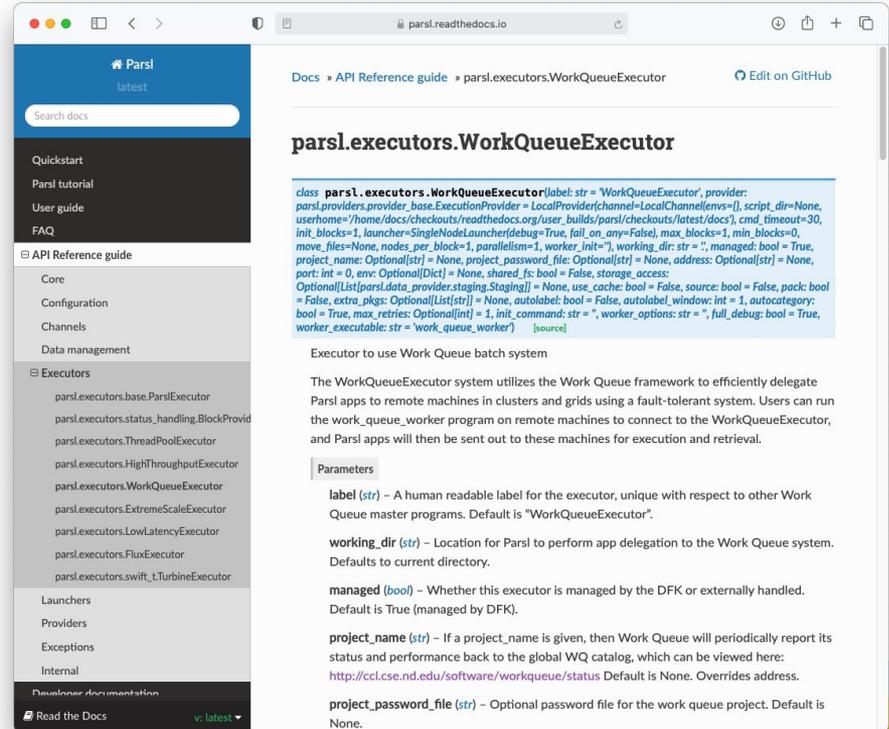
**Worker**

Task
Task
Task
**Task**

A
**B**

Execute Tasks Remotely on Local Disk

CCTools

# Configuring Parsl + WQ



```python
import parsl

from parsl.executors import WorkQueueExecutor


config = parsl.config.Config(

    executors=[

        WorkQueueExecutor(

            label="wq-parsl-app",

            port=9123,

            project_name="wq-parsl-app",

            shared_fs=False,

            full_debug = True,

) ] )
```

CCTools

Two common problems of scaling up:

● What resources should be assigned to a function call?
● What software dependencies does this function need?

How can we solve these problems automatically at runtime, without requiring the user to make advance declarations?

# Packing Functions Into Manycore Nodes

**CCTools**

Allocate 2GB per Function A?

Python App    F_A   F_B

Parsl DFK

Work Queue Manager

F_A    F_B

Work Queue Worker

F_A    F_A    F_A

F_A    F_A    F_A

12 cores and 12 GB RAM

# Packing Functions Into Manycore Nodes



Allocate 4GB per Function A?

Python App  $F_A$  $F_B$

Parsl DFK

Work Queue Manager

$F_A$  $F_B$

Work Queue Worker

$F_A$  $F_A$  $F_A$

12 cores and 12 GB RAM

# Packing Functions Into Manycore Nodes

**CCTools**

**Mix Function A and Function B?**



Python App    $F_A$    $F_B$

Parsl DFK

Work Queue Manager

$F_A$    $F_B$

Work Queue Worker

$F_A$    $F_A$    $F_B$    $F_B$    $F_B$

12 cores and 12 GB RAM

## LFM - Lightweight Function Monitor

**Ben Tovar**

Python Interpreter

LFM

function

import A    import B

fork (COW)

LFM

function

import A    import B

Resource Usage

Tim Shaffer, Zhuozhao Li, Ben Tovar, Yadu Babuji, TJ Dasso, Zoe Surma, Kyle Chard, Ian Foster, and Douglas Thain, **Lightweight Function Monitors for Fine-Grained Management in Large Scale Python Applications**, IEEE International Parallel & Distributed Processing Symposium, May, 2021. DOI: 10.1109/IPDPS49936.2021.00088

Activate LFMs with an import and the @monitored keyword

```
In [7]:   from resource_monitor import monitored
          from time import sleep
```

```
In [12]:  # declare a function to be monitored with the @monitored() decorator

          @monitored()
          def my_function_1(wait_for):
              sleep(wait_for)
              return 'waitied for {} seconds'.format(wait_for)

          (result, resources) = my_function_1(.1)
          print(result, '{}'.format({'memory': resources['memory'], 'wall_time': resources['wall_time']}))

          waitied for 0.1 seconds {'memory': 49, 'wall_time': 101689}
```

# Example: Colmena-XTB Application



**CCTools**

**Thanh Son Phung**

**Application**: XTB

./xtb-run.sh

**Application Framework**: Colmena

```
@agent
    def producer(self):
...
@agent
    def consumer(self):
```

**Workflow Manager**: Parsl

**Scheduler**: Work Queue

Worker Node 1

Worker Node 2

Tasks:
- Are of two types: inference and simulation
- Display significant differences in resource consumption

Thanh Son Phung, Logan Ward, Kyle Chard, Douglas Thain, **"Not All Tasks are Created Equal: Adaptive Resource Allocation for Hetergeneous Tasks in Dynamic Workflows"**, WORKS Workshop at Supercomputing 2021.

# Memory Consumption of Colmena-XTB's Tasks



Memory consumption over time



Buckets progressing over time

**Problem**

Tasks can consume as low as 2 GBs or as high as 30 GBs of RAM!

**Solution**

Bucket tasks with similar consumption and allocate new tasks accordingly.

# Run Time Dependency Management



**Barry Sly-Delgado**

Manager Environment

How do we ensure that all the tasks get a consistent, minimal environment matching the manager?

# Poncho Toolkit

The Poncho Toolkit allows users to create and deploy self contained Python environments at user level in arbitrary distributed systems via a JSON specification file.

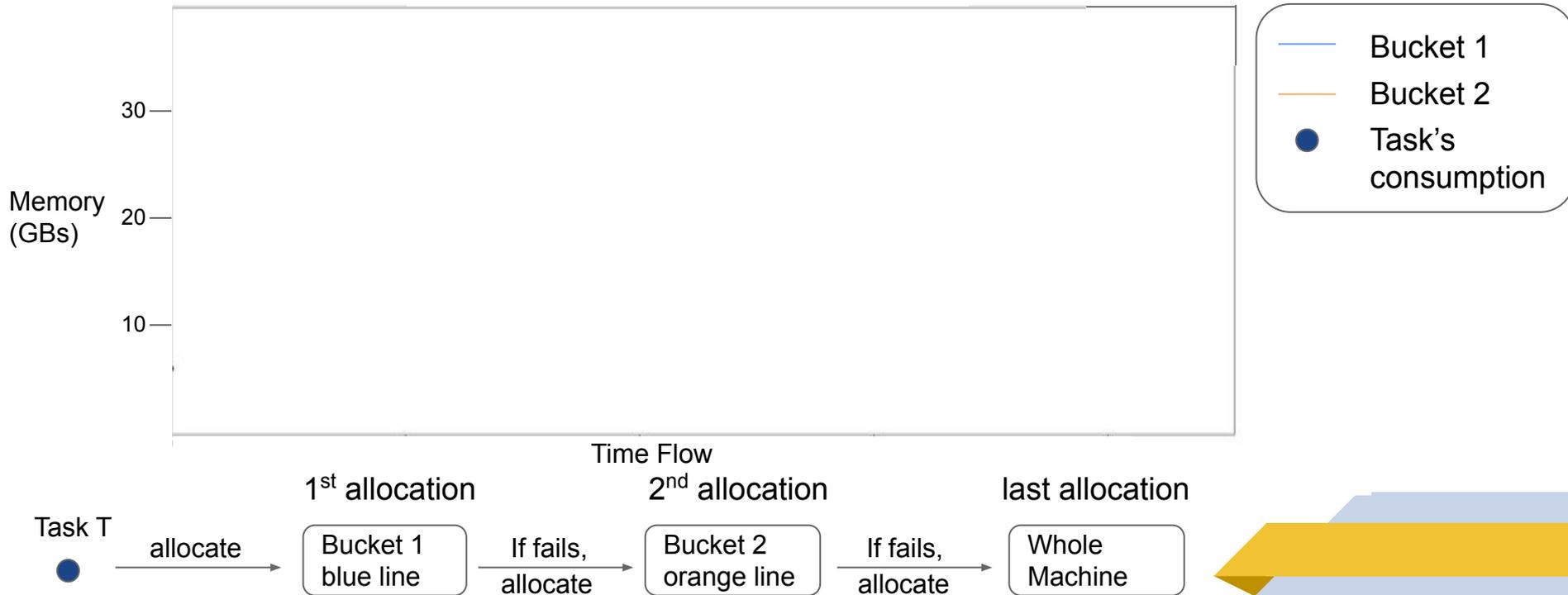- poncho_package_analyze
- poncho_package_create
- poncho_package_run

https://cctools.readthedocs.io/en/latest/poncho/

```
"conda":{
      "channels":[
            "defaults",
            "conda-forge"
      ],
      "packages":[
            "ndcctools=7.3.0",
            "parsl=1.1.0",
      ]
},
"pip": [
      "topcoffea"
]
"git": {
      "DATA_DIR": {
            "remote": "http://.../repo.git"
      }
```

# Run Time Dependency Management

# Package per Task



JSON Specification

poncho_package_create

Environment Tarball

tgz

Conda environment
Pip packages
Cloned repos
Fetched Files

work_queue_worker

Activate environment for each task.

poncho_package_run

Task

This required modest changes to the Parsl+WQ Executor to wrap each task invocation with additional files and commands.

# Package per Worker

JSON Specification

poncho_package_create

Environment Tarball

tgz

Conda environment
Pip packages
Cloned repos
Fetched Files

poncho_package_run

Activate environment once for whole worker.

work_queue_worker

Task

This requires modest changes to whatever method is used to deploy workers on the batch system, to modify the worker command and input files.

# First Look: Orders of Magnitude

poncho_package_create

| Application | No Versions Specified | All Versions Specified |
|---|---|---|
| TopEFT | 2940s | 170s |
| SHADHO | 257s | 159s |

poncho_package_run

| Application | Size compressed | Size unpacked | Unpack Time |
|---|---|---|---|
| SHADHO | 438MB | 1.4GB | 12s |
| TopEFT | 594MB | 2GB | 21s |
| Colmena-XTB | 1.4GB | 4.8GB | 46s |

# Next Steps...

- Do we really need all this code just to run a function? (maybe)
- Understanding the dependencies **actually used** by a function execution, and how they evolve over time.
- Extending dependency detection to other kinds of resources: databases, executables, file system resources...
- Closing the loop on application configuration: capture discovered resource configurations from multiple runs and use to predict future runs.
- Conveying known application categories from top to bottom through software stacks.

# End to End Integration Testing

## All workflows
Showing runs from all workflows

🔍 Filter workflow runs

341 workflow runs                          Event ▾   Status ▾   Branch ▾   Actor ▾

✅ **CI-Daily**                                        📅 16 hours ago   ...
    CI-Daily #160: Scheduled                          ⏱ 3m 39s

✅ **CI-Daily**                                        📅 2 days ago   ...
    CI-Daily #159: Scheduled                          ⏱ 3m 44s

✅ **drop conda-pack from coffea.sh (comes fr...**   `main`    📅 2 days ago   ...
    CI-Daily #158: Commit 7f26ea3 pushed by btovar           ⏱ 4m 43s

✅ **update python for parsl, drop conda-pack f...**  `main`    📅 2 days ago   ...
    CI-Daily #157: Commit 0b66cc3 pushed by btovar           ⏱ 30m 4s

❌ **CI-Daily**                                        📅 3 days ago   ...
    CI-Daily #156: Scheduled                          ⏱ 30m 25s

❌ **CI-Weekly**                                       📅 4 days ago   ...
    CI-Weekly #25: Scheduled                          ⏱ 30m 0s

❌ **CI-Daily**                                        📅 4 days ago   ...
    CI-Daily #155: Scheduled                          ⏱ 30m 27s

End-to-end daily test that simply installs parsl+workqueue and runs a trivial example out of the manual to see if it gets the right result.

conda-forge dropped support for python 3.6, resulting in attempts to install taking forever while conda tries to solve an unsolvable dependency problem!

# For More Information…

**Quick Start:**
conda install -c conda-forge python=3.9 ndcctools parsl

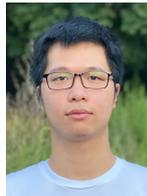**https://cctools.readthedocs.io**
**https://ccl.cse.nd.edu/software/workqueue**

btovar@nd.edu

tphung@nd.edu

bslydelg@nd.edu

dthain@nd.edu