# Parsl: Enabling Scalable Interactive Computing in Python

Kyle Chard ([chard@uchicago.edu](chard@uchicago.edu))

Yadu Babuji, Anna Woodard, Ben Clifford, Lukasz Lacinski, Mike Wilde, Dan Katz, Justin Wozniak, Ian Foster

http://parsl-project.org

THE UNIVERSITY OF CHICAGO

Argonne NATIONAL LABORATORY

ILLINOIS UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN

# An increasingly common story…

- I'm developing an application and I need to link together external tools + functions
  - (where each tool is dependent on data from the previous tool)
- I have a notebook that does *X* and I need to run it on a cloud, cluster, supercomputer
- I need to run my analysis using a range of local and distributed datasets
- …
- And I want to do this in an interactive environment

# Parsl: Interactive parallel scripting in Python

Annotate functions to make Parsl *apps*
- Python apps call Python functions
- Bash apps call external applications

Apps return "futures": a proxy for a result that might not yet be available

Apps run concurrently respecting data dependencies.
Natural parallel programming!

Parsl scripts are independent of where they run. Write once run anywhere!

```
pip install parsl
```

```python
@python_app
def hello ():
    return 'Hello World!'

print(hello().result())
```
Hello World!

```python
@bash_app
def echo_hello(stdout='echo-hello.stdout'):
    return 'echo "Hello World!"'

echo_hello().result()

with open('echo-hello.stdout', 'r') as f:
    print(f.read())
```
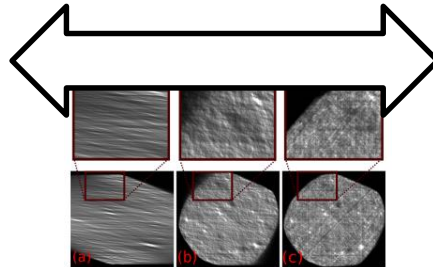Hello World!

parsl
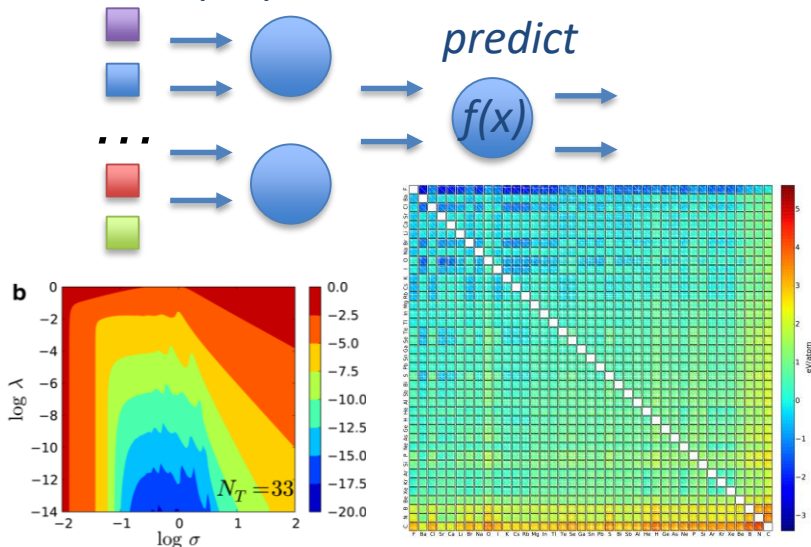
# Scientific workflows are not just for MTC/HPC/HTC
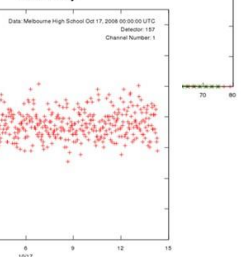
Online computing

*preprocess*

*predict*

$f(x)$

$N_T = 33$
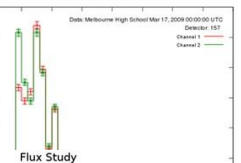
Machine learning

Detector    Detector

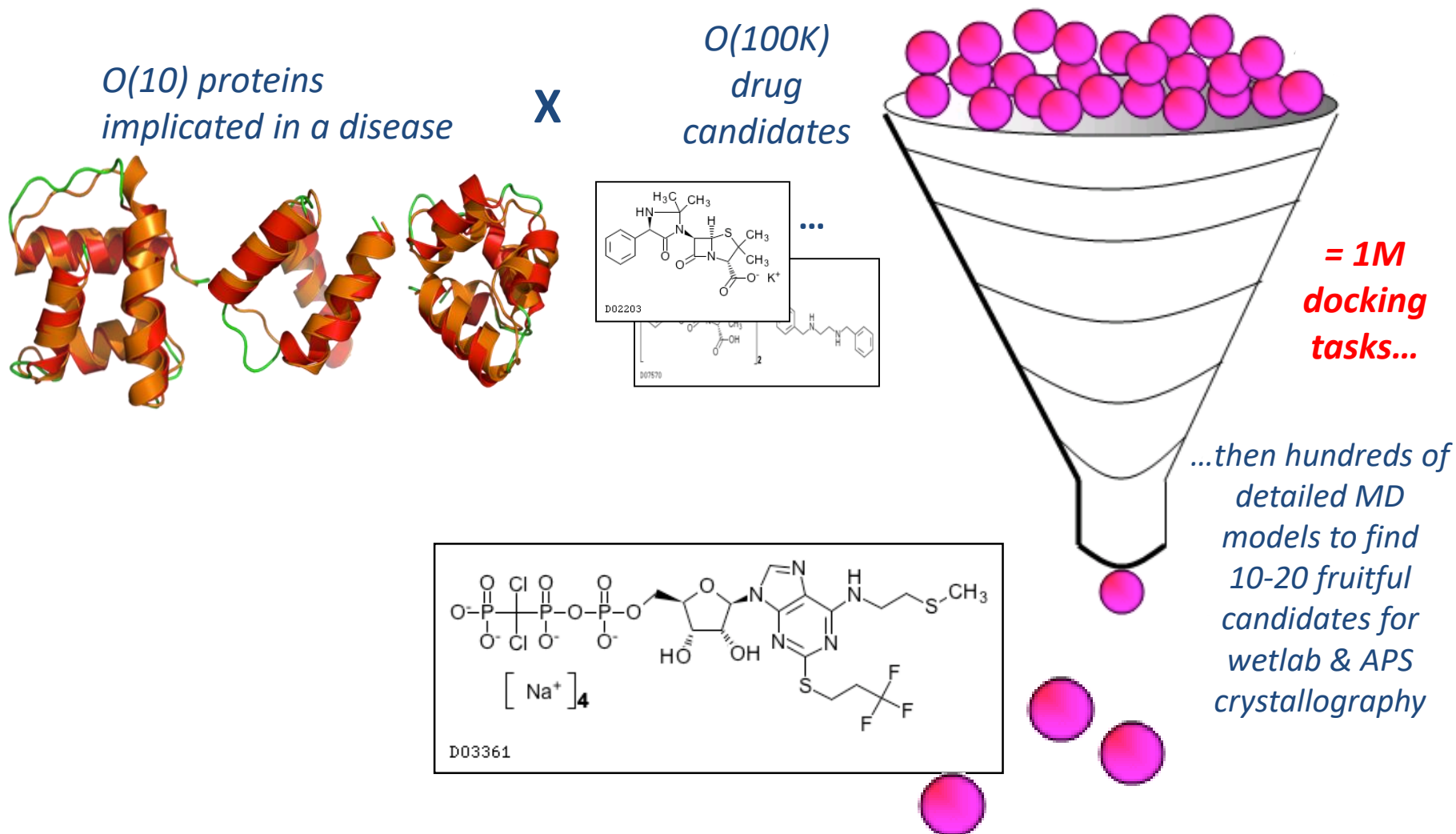Identify shower candidates

Curate

Analyze

Performance Study

Flux Study

Interactive computing

parsl

# When do you need automated workflow?
## Example application: protein-ligand docking for drug screening

*O(10) proteins implicated in a disease*

**X**

*O(100K) drug candidates*

*...*

*= 1M docking tasks...*

*...then hundreds of detailed MD models to find 10-20 fruitful candidates for wetlab & APS crystallography*



parsl

# Expressing a many task workflow in Parsl

*1) Wrap the protein docking code:*

```
@bash_app
def dock(p, c, minRad, maxRad)
  return 'dock.sh {0} {1} {2} {3}'.format(p,
   c ,minRad, maxRad)
```

parsl

# Expressing a many task workflow in Parsl

*2) Execute the protein docking workflow:*

```
for p in proteins:
    for c in ligands:
        structure[p][c] =
          dock(p, c, minRad, maxRad)

scatter_plot = analyze(structure)
```

parsl

# Brief history: the Swift parallel scripting language

- 10+ years of development

- C-like language with implicit parallelism

- Applied in dozens of scientific domains

- Data management, multi-site execution, coasters, etc.

- We are leveraging lessons and components to build Parsl

```
type file;

app (file o) simulation (int sim_steps, int sim_range, int sim_values)
{
  simulate "--timesteps" sim_steps "--range" sim_range "--nvalues" sim_values
stdout=filename(o);
}

app (file o) analyze (file s[])
{
  stats filenames(s) stdout=filename(o);
}

int nsim   = toInt(arg("nsim","10"));
int steps  = toInt(arg("steps","1"));
int range  = toInt(arg("range","100"));
int values = toInt(arg("values","5"));

file sims[];

foreach i in [0:nsim-1] {
  file simout <single_file_mapper; file=strcat("output/sim_",i,".out")>;
  simout = simulation(steps,range,values);
  sims[i] = simout;
}

file stats<"output/average.out">;
stats = analyze(sims);
```
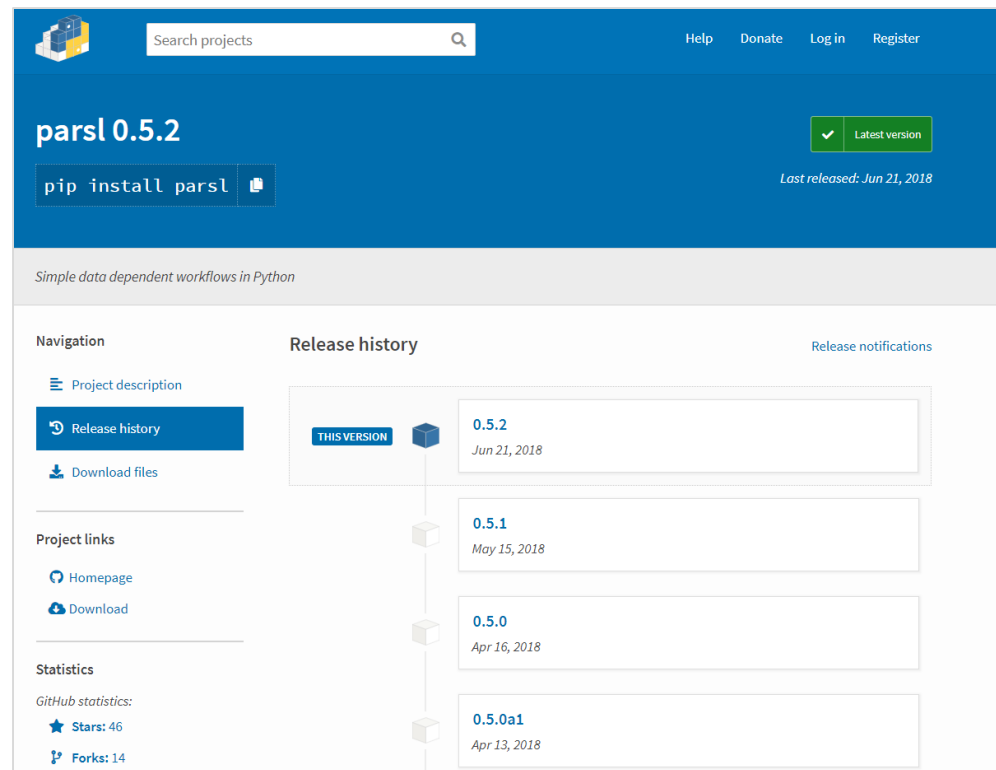
# Parsl in action: dynamic dataflow execution

# Parsl is Python

```
pip3 install parsl
```

- Use Python libraries natively

- Stage Python data transparently

- Integrates with Python ecosystem

# Parsl scripts are execution provider independent

- The same script can be run locally, on grids, clouds, or supercomputers
  - Works directly with the scheduler (no HTC-like setup)
- Containers can be used for per-app execution or repeated invocation of the same app
- Parsl builds on libsubmit
  - https://github.com/Parsl/libsubmit
- Currently supported execution providers:
  - Local, Cloud (AWS, Azure, private), Slurm, Torque, Condor, Cobalt

parsl

# Separation of code and execution

```python
@python_app(executors=['midway'])
def midway():
    return 'I am run on midway!'


@bash_app(executors=['local'])
def local():
    return 'I am run locally!'
```
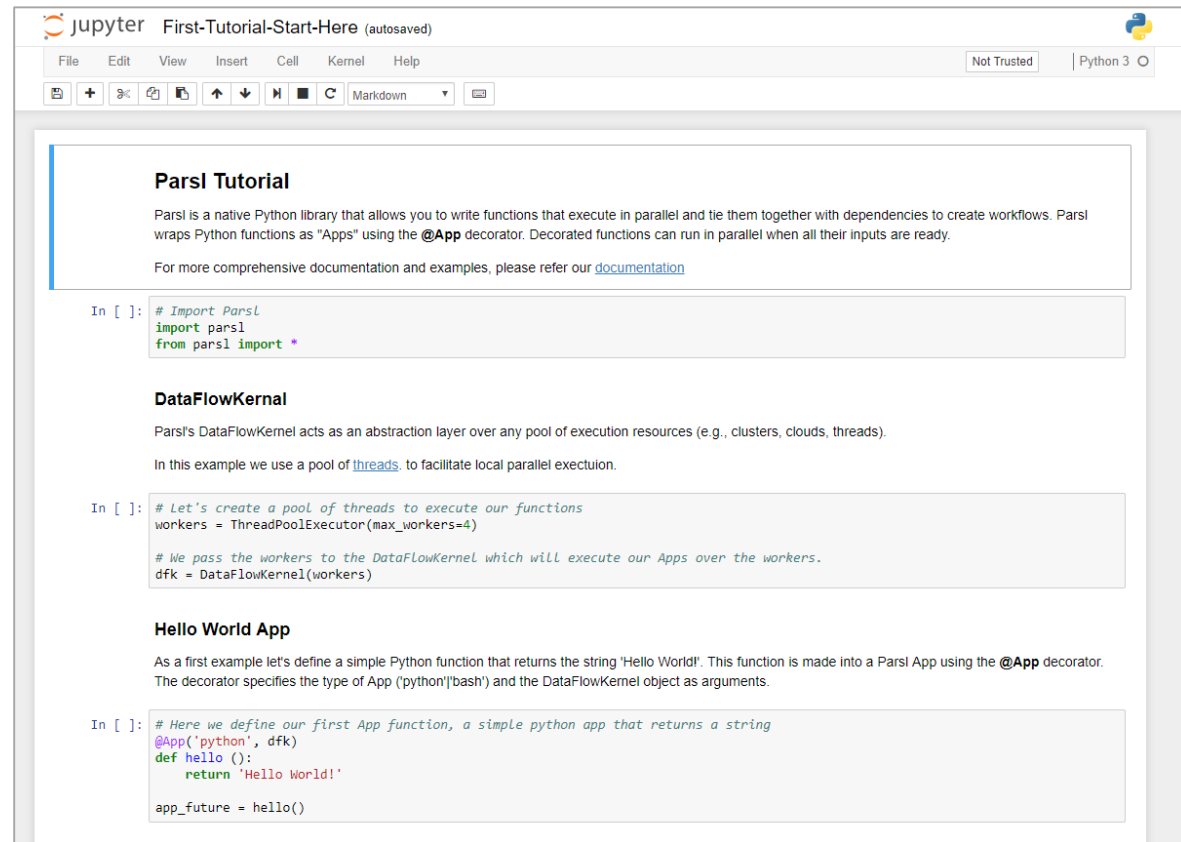
```python
from libsubmit.channels import SSHChannel
from libsubmit.providers import SlurmProvider

import parsl
from parsl.config import Config
from parsl.executors.ipp import IPyParallelExecutor
from parsl.executors.threads import ThreadPoolExecutor

config = Config(
    executors=[
        IPyParallelExecutor(
            label='midway',
            provider=SlurmProvider(
                'westmere',
                channel=SSHChannel(
                    hostname='swift.rcc.uchicago.edu',
                    username='annawoodard'
                ),
                max_blocks=1000,
                nodes_per_block=1,
                tasks_per_node=6,
                overrides='module load singularity; module load Anaconda3/5.1.0; source activate parsl_py36'
            ),
        ),
        ThreadPoolExecutor(label='local', max_threads=2)
    ],
)

parsl.load(config)
```

Pilot jobs on
a cluster

Local threads

\* Config format for Parsl 0.6

parsl

# Interactive supercomputing in Jupyter notebooks

- Parsl can be used within a Jupyter notebook with no modifications necessary

- Tunneling and OAuth-based flows supports remote execution from the notebook

- Visualization of Parsl graph in notebook

# Parsl supports a variety of execution models

- Threads
  - Local execution
- Ipython.parallel
  - Pilot job model
- Swift/T
  - Extreme scale execution

- New execution models can be added

# Authentication and authorization

- A&A is hard today
  - 2FA, X509, etc.
- Integration with Globus Auth to support native app integration for accessing Globus (and other) services
- Using scoped access tokens, refresh tokens, delegation support

# Parsl provides transparent (wide area) data management

- Implicit data movement to/from repositories, laptops, supercomputers, ...

- Globus for third-party, high performance and reliable data transfer
  - Support for site-specific DTNs

- HTTP/FTP direct data download/upload
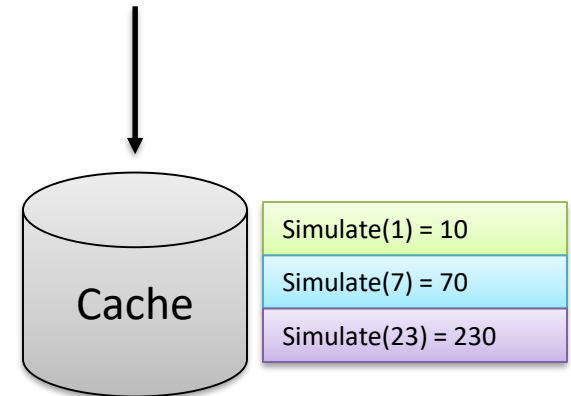
- Compliments node-specific staging and caching models
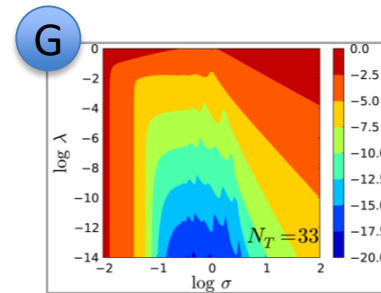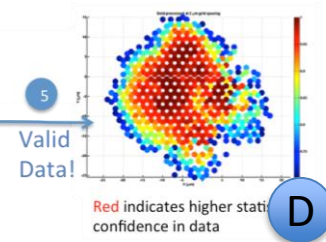
```
parsl_file =
    File(globus://EP/path/file)
```



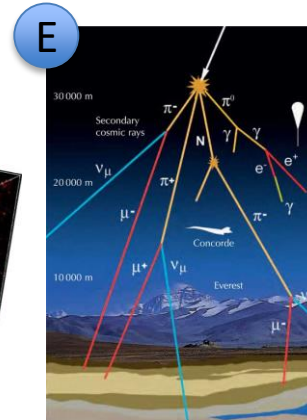www.globus.org

# App caching (memoization)

- Parsl apps are often expensive to recompute
- In many development modes results need not be recomputed
  - During development or interactive workflow
- Memoization optimizes execution by caching app results when called with the same inputs
- Parsl relies on user control to annotate deterministic functions

```python
@python_app(cache=True)
def simulate(input_variable):
    return input_variable * 10
```



Cache

| | |
|---|---|
| Simulate(1) = 10 | |
| Simulate(7) = 70 | |
| Simulate(23) = 230 | |

# Scientific applications using Parsl

A  Machine learning to predict stopping power in materials

B  Protein and biomolecule structure and interaction

C  Information extraction to discovery facts in publications

D  Materials science at the Advanced Photon Source

E  Cosmic ray showers as part of QuarkNet

F  Weak lensing using sky surveys

G  Machine learning and data analytics (DLHub)

# Summary

- Parsl's implicit dataflow model in Python allows for simple expression of complex dependencies
  - Expressed directly in Python
  - Can be used to implement a range of workflow models
- Parsl integrates with the scientific ecosystem
  - Development and execution of scalable applications in Jupyter
  - Use of common SciPy libraries
  - Integration with Globus
- In Parsl, code is separate from the specification of computing resources and data location: this makes Parsl scripts portable and scalable
- Parsl has a number of other important features:
  - app caching, checkpointing, elasticity, container support, data transfer, and more

# Questions?

**http://parsl-project.org**

Try Parsl: http://try.parsl-project.org
https://mybinder.org/v2/gh/Parsl/parsl-tutorial/master