

# DASK FOR PARALLEL COMPUTING

## CHEAT SHEET

See full Dask documentation at: <http://dask.pydata.org/>

These instructions use conda environment manager. Get yours at <http://bit.ly/getconda>

*TIP: Use `help(object)` to get help about any Python object*

### DASK QUICK INSTALL

|                         |   |
|-------------------------|---|
| Install Dask with conda | <code>conda install dask</code>         |
| Install Dask with pip   | <code>pip install dask[complete]</code> |

### DASK COLLECTIONS

#### DASK ARRAYS

|  |   |
|--|---|
| Import dask.array library                            | <code>import dask.array as da</code>  |
| Create a Dask Array from Numpy-like array            | <code>x = da.from_array(d, chunks=(m, n, ...))</code>   |
| Example. Dask Array from HDF5 file                   | <pre>import h5py  f = h5py.File('datafile.hdf5', 'r') x = f['/group1/dataset1'] d = da.from_array(x, chunks = (1000, 1000))</pre>             |
| Store Dask Array in array-like object                | <code>da.store(x, array)</code>   |
| Example. Store Dask Array into HDF5 file             | <pre>x = da.random.normal(10, .3, size=(5,5), chunks=(5,1)) f = h5py.File('myfile.hdf5') dset = f.create_dataset(...) da.store(x, dset)</pre> |
| Arithmetic element-wise and scalar operations        | <code>*, +, -, **, /, exp, log</code>   |
| Example. Arithmetic element-wise & scalar operations | <code>y = da.sin(x)**2 + da.cos(x)**2</code>  |
| Reduction along axes                                 | <code>sum(), prod(), mean(), std()</code>   |
| Example. Sum reduction along t                       | <code>y = x.mean(axis=t)</code>   |
| Matrix multiplication and dot product                | <code>dot(), tensordot()</code>   |
| Axis reordering                                      | <code>transpose()</code>  |
| Slicing  | <code>x[:5, 20:10:-1]</code>  |
| Fancy indexing                                       | <code>x[[1, 3], :]</code>   |

#### DASK BAGS

|   |   |
|---|---|
| Import dask.bag library   | <code>import dask.bag as db</code>  |
| Create Dask Bag from a sequence   | <code>db.from_sequence(seq, npartitions)</code>   |
| Example.  | <code>b = db.from_sequence([1, 2, 3, 4, 5, 6], npartitions=2)</code>  |
| Create Dask Bag from text files   | <code>b = db.from_filenames('data.*.json')</code>   |
| Map function across all elements in a Dask Bag                                | <code>map()</code>  |
| Example: use <code>from_filenames</code> and <code>json.loads</code> together | <pre>import json b = db.from_filenames('data.*.json.gz').map(json.loads) compute()</pre>                    |
| Trigger computations  |   |
| Example.  | <pre>b = db.from_sequence([2, 3, 5, 7, 11, 13], npartitions=2) c = b.map(lambda x: x + 1) c.compute()</pre> |

## DASK COLLECTIONS (CONTINUED)

### DASK BAGS (CONTINUED)

Some useful functions supported by Dask Bags

max(), min(), mean(), sum(), std(), filter(), fold(), foldby(), frequencies(), groupby(), join(), pluck(), product(), remove(), take(), topk(), var()

Convert to Dask DataFrame

to\_dataframe()

Write Dask Bag to disk

to\_textfiles('path')

### DASK DATAFRAMES

Import dask.dataframe library

import dask.dataframe as dd

Create Dask DataFrame from CSV files

read\_csv('filenames\*.csv')

Example

```
df = dd.read_csv('file-*.csv.gz', compression='gzip')
df.head()
```

Element-wise operations

\*, +, /, -

Row-wise selection

df[df.x > 0]

Selection by label

df.loc['2015-01': '2015-05']

Common aggregations

max(), min(), mean(), std(), sum(), count(), var()

Pandas operations supported by Dask DataFrames

groupby(), value\_counts(), drop\_duplicates(), merge(), set\_index()

Trigger computations

compute()

Example

```
df = dd.read_csv('filenames*.csv')
df.sample(frac=0.1).groupby(df.timestamp.day).value.
mean()
df.compute()
```

## GRAPHS

*TIP: Use single-threaded scheduler for debugging, dask.set\_options(get=dask.async.get\_sync)*

Scheduler backed by thread pool

dask.threaded.get()

Scheduler backed by process pool

dask.multiprocessing.get()

Synchronous scheduler

dask.async.get\_sync()

Example

```
from dask.threaded import get
from operator import add
dsk = {'a': 1,
      'b': 2,
      'c': (add, 'a', 'b')}
get(dsk, 'c')
```

## MORE RESOURCES

Support

<http://bit.ly/anacondasupport>

Training

<http://bit.ly/continuumtraining>

Consulting

<http://bit.ly/continuumconsulting>

Dask gitter chat room

<https://gitter.im/dask/dask>

Report a bug

<https://github.com/dask/dask/issues>

Dask mailing list

<https://groups.google.com/a/continuum.io/forum/#!forum/blaze-dev>

Join the [#AnacondaCrew!](#)

Connect with other talented, like-minded data scientists and developers while contributing to the open source movement. Visit <https://continuum.io/community>