



# 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>
<b>Example</b> Dask Array from HDF5 file	<code>import h5py</code>  <code>f = h5py.File('datafile.hdf5', 'r')</code> <code>x = f['/group1/dataset1']</code> <code>d = da.from_array(x, chunks = (1000, 1000))</code>  <code>da.store(x, array)</code>
Store Dask Array in array-like object	
<b>Example</b> Store Dask Array into HDF5 file	<code>x = da.random.normal(10, .3, size=(5,5), chunks=(5,1))</code> <code>f = h5py.File('myfile.hdf5')</code> <code>dset = f.create_dataset(...)</code> <code>da.store(x, dset)</code>
Arithmetic element-wise & scalar operations	<code>*, +, -, **, /, exp, log</code>
<b>Example</b> 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>
<b>Example</b> 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>
<b>Example</b>	<code>b = db.from_sequence([1, 2, 3, 4, 5, 6], npartitions=2)</code>
Create Dask Bag from text files	<code>b = db.read_text('data.*.json')</code>
Map function across all elements in a Dask Bag	<code>map()</code>
<b>Example</b> use <code>read_text</code> and <code>json.loads</code> together	<code>import json</code> <code>b = db.read_text('data.*.json.gz').map(json.loads)</code> <code>compute()</code>
Trigger computations	
<b>Example</b>	<code>b = db.from_sequence([2, 3, 5, 7, 11, 13], npartitions=2)</code> <code>c = b.map(lambda x: x + 1)</code> <code>c.compute()</code>

DASK COLLECTIONS (CONTINUED)

DASK BAGS (CONTINUED)

Some useful functions supported by Dask Bags	<code>max()</code> , <code>min()</code> , <code>mean()</code> , <code>sum()</code> , <code>std()</code> , <code>filter()</code> , <code>fold()</code> , <code>foldby()</code> , <code>frequencies()</code> , <code>groupby()</code> , <code>join()</code> , <code>pluck()</code> , <code>product()</code> , <code>remove()</code> , <code>take()</code> , <code>topk()</code> , <code>var()</code>
Convert to Dask DataFrame	<code>to_dataframe()</code>
Write Dask Bag to disk	<code>to_textfiles('path')</code>

DASK DATAFRAMES

Import <code>dask.dataframe</code> library	<code>import dask.dataframe as dd</code>
Create Dask DataFrame from CSV files	<code>df = dd.read_csv('filenames*.csv')</code>
Element-wise operations	<code>*</code> , <code>+</code> , <code>/</code> , <code>-</code>
Row-wise selection	<code>df[df.x &gt; 0]</code>
Selection by label	<code>df.loc['2015-01': '2015-05']</code>
Common aggregations	<code>max()</code> , <code>min()</code> , <code>mean()</code> , <code>std()</code> , <code>sum()</code> , <code>count()</code> , <code>var()</code>
pandas operations supported by Dask DataFrames	<code>groupby()</code> , <code>value_counts()</code> , <code>drop_duplicates()</code> , <code>merge()</code> , <code>set_index()</code>
Trigger computations	<code>compute()</code>
Example	<pre>df = dd.read_csv('filenames*.csv') df.sample(frac=0.1, replace=True)     .groupby(df.timestamp.day)     .value.mean().compute()</pre>

GRAPHS

<i>TIP: Use single-threaded scheduler for debugging, <code>dask.set_options(get=dask.async.get_sync)</code></i>	
Scheduler backed by thread pool	<code>dask.threaded.get()</code>
Scheduler backed by process pool	<code>dask.multiprocessing.get()</code>
Synchronous scheduler	<code>dask.async.get_sync()</code>
Example	<pre>from dask.threaded import get from operator import add dsk = {'a': 1,       'b': 2,       'c': (add, 'a', 'b')} get(dsk, 'c')</pre>

MORE RESOURCES

Support	<a href="https://www.continuum.io/support-plan">https://www.continuum.io/support-plan</a>
Training	<a href="http://bit.ly/continuumtraining">http://bit.ly/continuumtraining</a>
Consulting	<a href="http://bit.ly/continuumconsulting">http://bit.ly/continuumconsulting</a>
Dask gitter chat room	<a href="https://gitter.im/dask/dask">https://gitter.im/dask/dask</a>
Report a bug	<a href="https://github.com/dask/dask/issues">https://github.com/dask/dask/issues</a>
Dask mailing list	<a href="https://groups.google.com/a/continuum.io/forum/#!forum/blaze-dev">https://groups.google.com/a/continuum.io/forum/#!forum/blaze-dev</a>

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>