

Dask for Parallel Computing Cheat Sheet

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

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

DASK QUICK INSTALL	
Install Dask with conda	<code>conda install dask</code>
Install Dask with pip	<code>pip install dask[complete]</code>

DASK COLLECTIONS	<i>Easy-to-use big data collections</i>
Dask Dataframes	<i>Parallel Pandas dataframes for large data</i>
Import	<code>import dask.dataframe as dd</code>
Read CSV data	<code>df = dd.read_csv('my-data.*.csv')</code>
Read Parquet data	<code>df = dd.read_parquet('my-data.parquet')</code>
Filter and manipulate data with Pandas syntax	<code>df['z'] = df.x + df.y</code>
Standard groupby aggregations, joins, etc.	<code>result = df.groupby(df.z).y.mean()</code>
Compute result as a Pandas dataframe	<code>out = result.compute()</code>
Or store to CSV, Parquet, or other formats	<code>result.to_parquet('my-output.parquet')</code>
EXAMPLE	<pre>df = dd.read_csv('filenames.*.csv') df.groupby(df.timestamp.day)\ .value.mean().compute()</pre>
Dask Arrays	<i>Parallel NumPy arrays for large data</i>
Import	<code>import dask.array as da</code>
Create from any array-like object. This includes HDF5, NetCDF, or other on-disk formats. Alternatively generate an array from a random distribution.	<pre>import h5py dataset = h5py.File('my-data.hdf5')['/group/dataset'] x = da.from_array(dataset, chunks=(1000, 1000)) da.random.uniform(shape=(1e4, 1e4), chunks=(100, 100))</pre>
Perform operations with NumPy syntax	<code>y = x.dot(x.T - 1) - x.mean(axis=0)</code>
Compute result as a NumPy array	<code>result = y.compute()</code>
Or store to HDF5, NetCDF or other non-disk format	<pre>out = f.create_dataset(...) x.store(out)</pre>
EXAMPLE	<pre>with h5py.File('my-data.hdf5') as f: x = da.from_array(f['/path'], chunks=(1000, 1000)) x -= x.mean(axis=0) out = f.create_dataset(...) x.store(out)</pre>

DASK COLLECTIONS (CONT.)	
Dask Bags	Parallel lists for unstructured data
Import	<code>import dask.bag as db</code>
Create Dask Bag from a sequence	<code>b = db.from_sequence(seq, npartitions)</code>
Or read from text formats	<code>b = db.read_text('my-data.*.json')</code>
Map and filter results	<pre>import json records = b.map(json.loads) .filter(lambda d: d["name"] == "Alice")</pre>
Compute aggregations like mean, count, sum	<code>records.pluck('key-name').mean().compute()</code>
Or store results back to text formats	<code>records.to_textfiles('output.*.json')</code>
EXAMPLE	<pre>db.read_text('s3://bucket/my-data.*.json') .map(json.loads) .filter(lambda d: d["name"] == "Alice") .to_textfiles('s3://bucket/output.*.json')</pre>
Advanced	
Read from distributed file systems or cloud storage Prepend prefixes like <code>hdfs://</code> , <code>s3://</code> , or <code>gcs://</code> to paths	<pre>df = dd.read_parquet('s3://bucket/myfile.parquet') b = db.read_text('hdfs:///path/to/my-data.*.json')</pre>
Persist lazy computations in memory	<code>df = df.persist()</code>
Compute multiple outputs at once	<code>dask.compute(x.min(), x.max())</code>
CUSTOM COMPUTATIONS	
Dask Delayed	Lazy parallelism for custom code
Import	<code>import dask</code>
Wrap custom functions with the <code>@dask.delayed</code> annotation Delayed functions operate lazily, producing a task graph rather than executing immediately Passing delayed results to other delayed functions creates dependencies between tasks	<pre>@dask.delayed def load(filename): ... @dask.delayed def process(data): ...</pre>
Call functions in normal code	<pre>data = [load(fn) for fn in filenames] results = [process(d) for d in data]</pre>
Compute results to execute in parallel	<code>dask.compute(results)</code>
Concurrent futures	
Import	<code>from dask.distributed import Client</code>
Start local Dask Client	<code>client = Client()</code>
Submit individual task asynchronously	<code>future = client.submit(func, *args, **kwargs)</code>
Block and gather individual result	<code>result = future.result()</code>
Process results as they arrive	<pre>for future in as_completed(futures): ...</pre>
EXAMPLE	<pre>L = [client.submit(read, fn) for fn in filenames] L = [client.submit(process, future) for future in L] future = client.submit(sum, L) result = future.result()</pre>

SET UP CLUSTER

<https://docs.dask.org/en/latest/setup.html>

Manually

Start scheduler on one machine	<pre>\$ dask-scheduler Scheduler started at SCHEDULER_ADDRESS:8786</pre>
Start workers on other machines Provide address of the running scheduler	<pre>host1\$ dask-worker SCHEDULER_ADDRESS:8786 host2\$ dask-worker SCHEDULER_ADDRESS:8786</pre>
Start Client from Python process	<pre>from dask.distributed import Client client = Client('SCHEDULER_ADDRESS:8786')</pre>

On a single machine

Call Client() with no arguments for easy setup on a single host	<pre>client = Client()</pre>
---	------------------------------

Cloud Deployment

Install Dask on Kubernetes	<pre>pip install dask-kubernetes</pre>
Install Dask on AWS / GCP / AzureML	<pre>pip install dask-cloudprovider</pre>
Multi-user server for managing Dask Clusters	<pre>pip install dask-gateway</pre>

More resources

User Documentation	https://docs.dask.org
Technical Documentation for distributed scheduler	https://distributed.dask.org
Ask for help on StackOverflow	https://stackoverflow.com/questions/tagged/dask
Report a bug	https://github.com/dask/dask/issues

With more than 20 million users, Anaconda is the world's most popular data science platform and the foundation of modern machine learning. We pioneered the use of Python for data science, champion its vibrant community, and continue to steward open-source projects that make tomorrow's innovations possible. Our enterprise-grade solutions enable corporate, research, and academic institutions around the world to harness the power of open-source for competitive advantage, groundbreaking research, and a better world.

Visit anaconda.com to learn more.

