# Getting up to speed with Dask

Aaron Richter
July 2020

https://github.com/rikturr/getting-up-to-speed-with-dask

Saturn Cloud

# Hi!

## Aaron Richter

Senior Data Scientist @ Saturn Cloud

- I work to make data scientists faster and happier

PhD in Machine Learning

aaron@saturncloud.io
rikturr.com
@rikturr

# Saturn Cloud

## Bringing together the fastest hardware + OSS

**DASK**

- Pythonic parallelism
- Rapidly scale PyData

**RAPIDS**

- Multi-GPU computing
- The future of HPC

**PREFECT**

- Workflow orchestration
- Flow insight and mgmt

**kubernetes**

- Fast setup
- Enterprise secure

# Dask

Parallel computing for Python people

# Dask

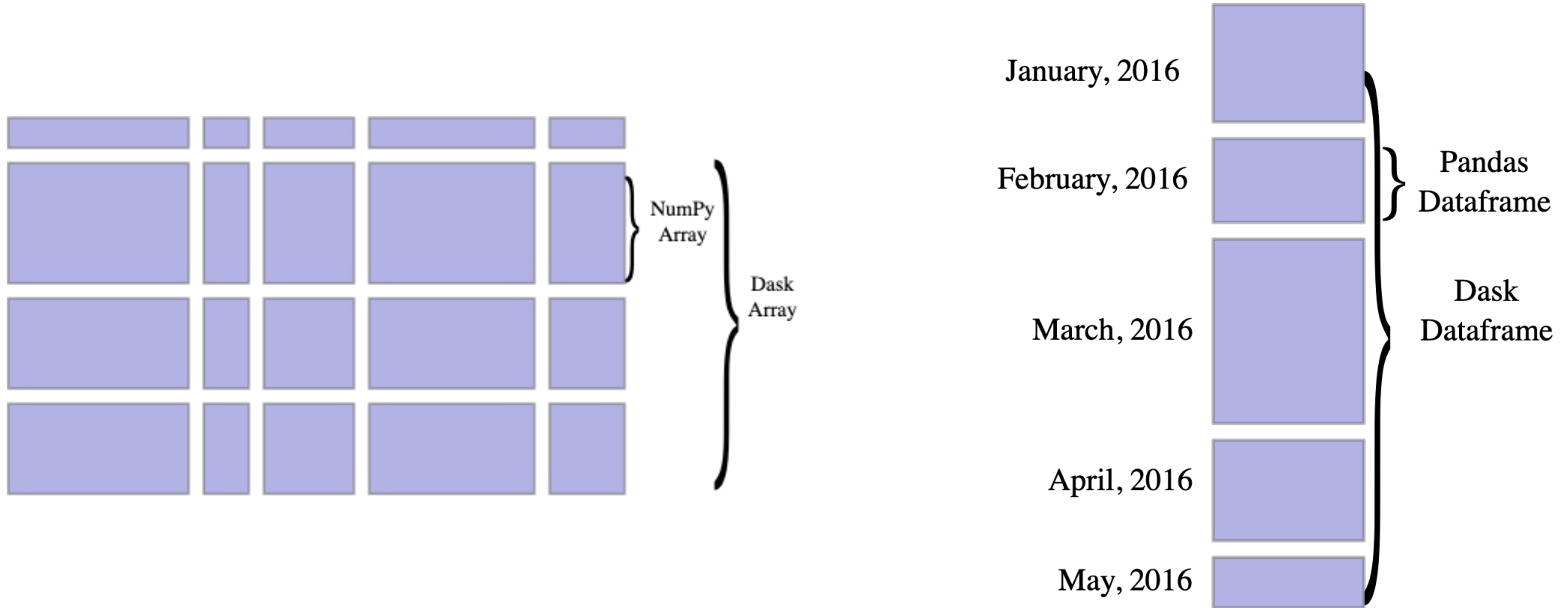## What does it do?

- Parallel machine learning (scikit)

- Parallel dataframes (pandas)

- Parallel arrays (numpy)

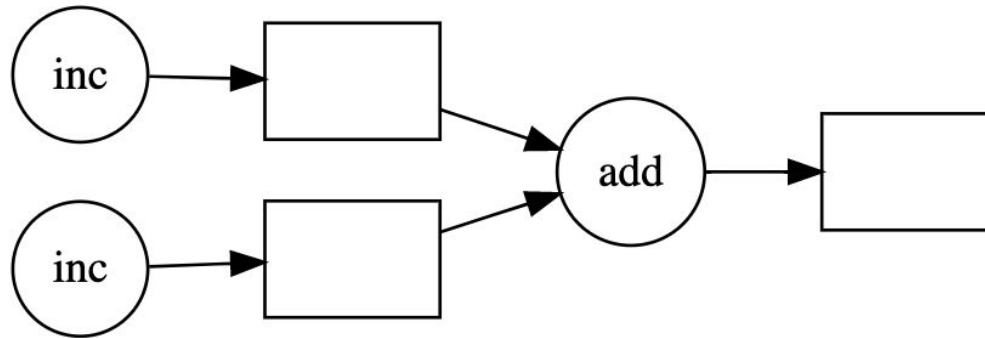- Parallel anything else

# What does it do?

## Arrays and Dataframes
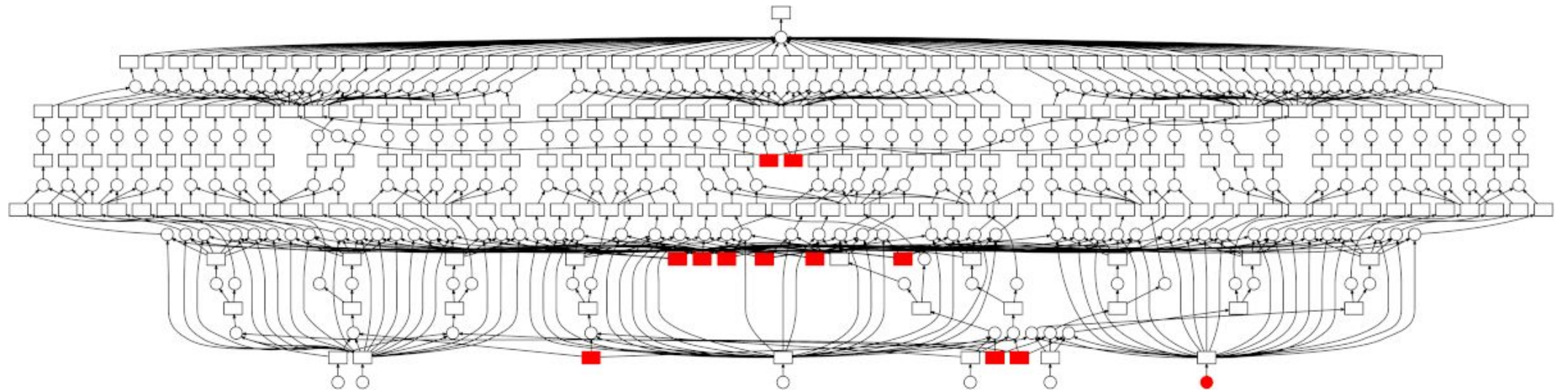
# What does it do?

## Anything else!



```
>>> x = dask.delayed(inc)(1)
>>> y = dask.delayed(inc)(2)
>>> z = dask.delayed(add)(x, y)
>>> z.compute()
5
>>> z.visualize()
```

# What does it do?

## Anything else!

# Dask

## Why should I use it?

- Python native

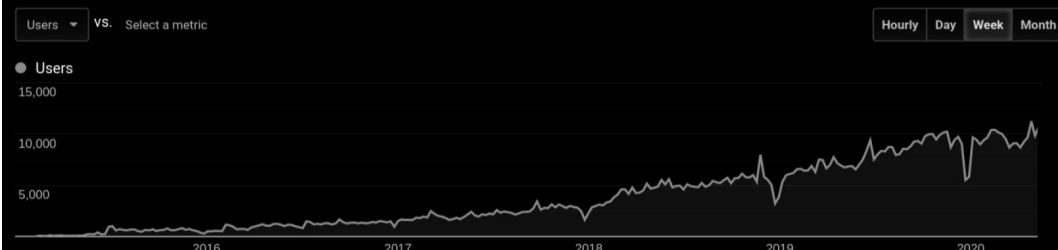- Strong ecosystem (PyData)

- Easily scalable

# Dask

## Why should I use it?

- Make your Python faster with a "pip install"

- 2 to 50 times faster than Spark*
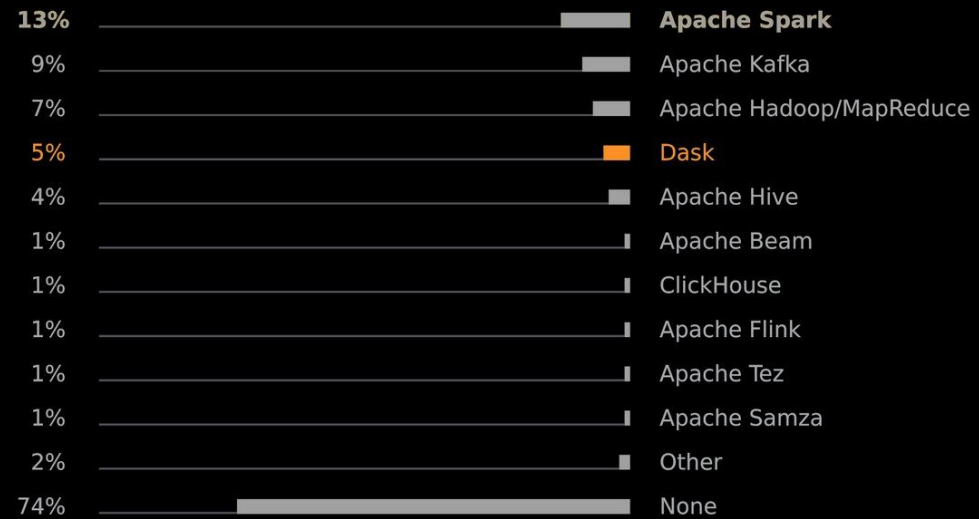
- Will bring you fame and fortune**

*stay tuned          **YMMV

# Dask

## Why should I use it?

10,000 Documentation Visitors

Unique visitors on a weekly basis

5% of Python developers
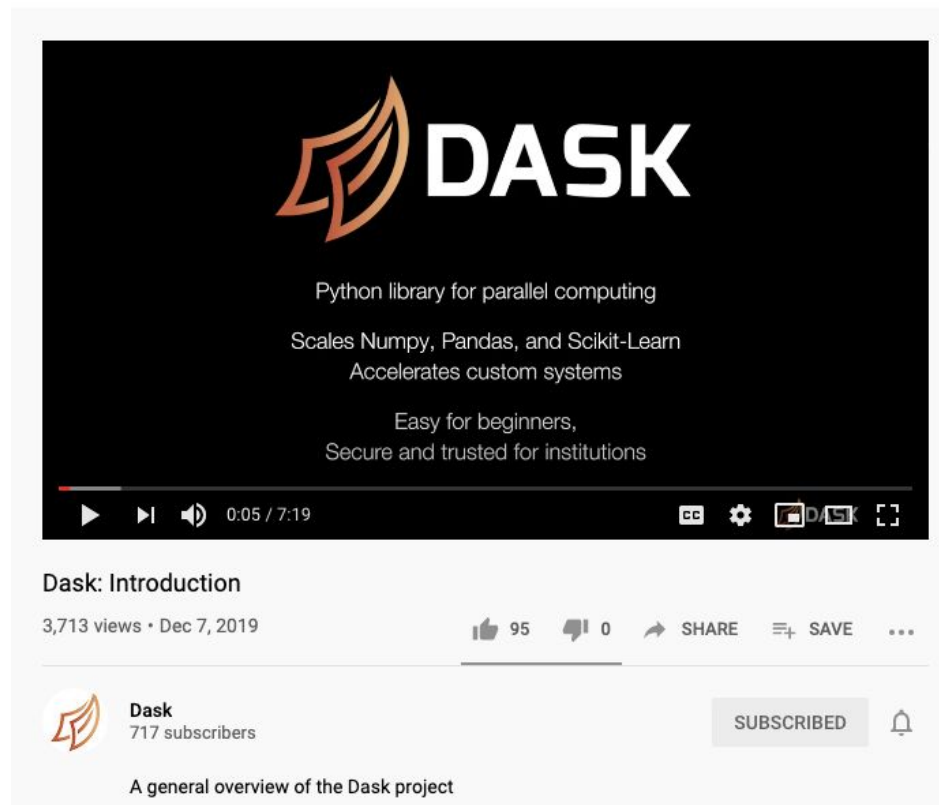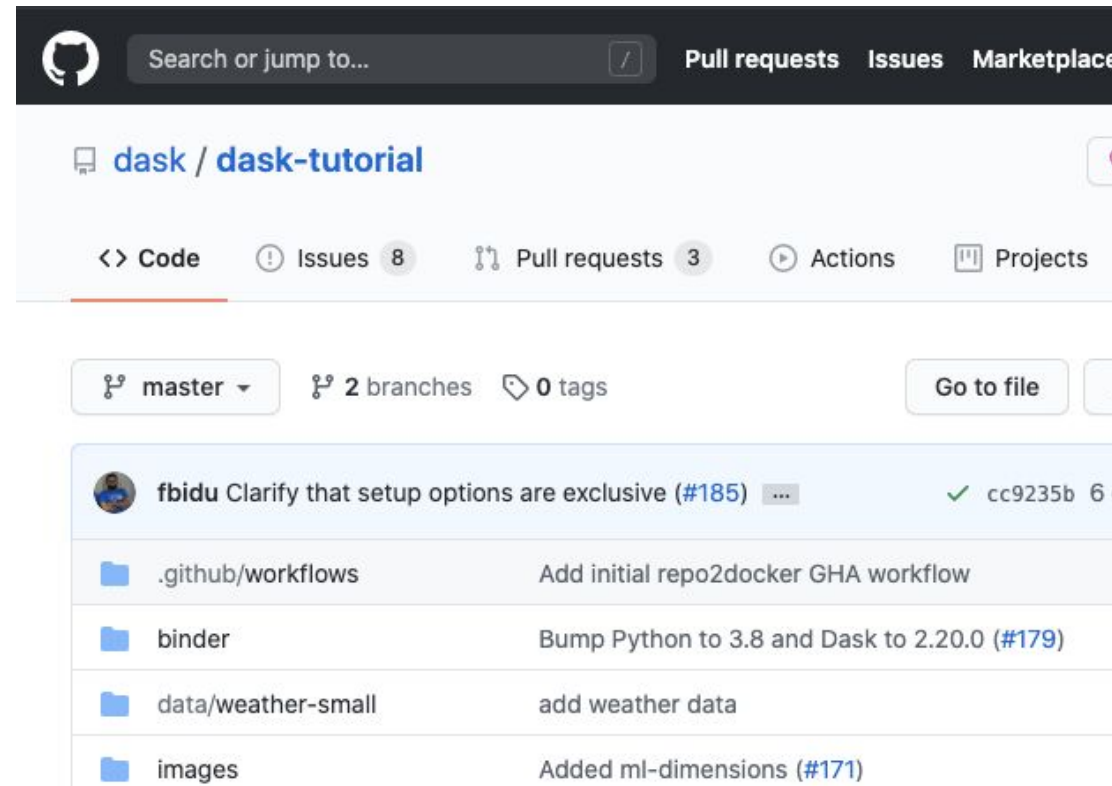(among those who take the Python survey)

| | |
|---|---|
| 13% | **Apache Spark** |
| 9% | Apache Kafka |
| 7% | Apache Hadoop/MapReduce |
| 5% | Dask |
| 4% | Apache Hive |
| 1% | Apache Beam |
| 1% | ClickHouse |
| 1% | Apache Flink |
| 1% | Apache Tez |
| 1% | Apache Samza |
| 2% | Other |
| 74% | None |

https://www.jetbrains.com/lp/python-developers-survey-2019/

# How do I get started?

## Videos

# How do I get started?

## Tutorial

# How do I get started?

## Docs

# How do I get started?

## Key concepts

- Task graph

- Lazy execution

- Parallel objects are "normal" objects under the hood

# Code time!

https://github.com/rikturr/getting-up-to-speed-with-dask

# Dask

## Should I use it now?

- Use pandas (numpy) until you can't

- Then use Dask on your laptop

- Then try a big machine in the cloud

- *Then* go for clusters (in the cloud)

# Dask

## Running on a cluster

- Dask runs on most cluster/HPC platforms

    - Hadoop/YARN, Kubernetes, SLURM, etc.

- Rent your machines! (AWS, Azure, GCP)

- Managed solutions like Saturn Cloud, Coiled Computing

# What's coming next?

## Exciting stuff!

- High level graph optimization

- Scheduler performance

- Chan Zuckerberg Initiative life science grant



https://threadreaderapp.com/thread/1280885850914553856.html

# Get involved!

Thriving community of open source contributors

- https://docs.dask.org/en/latest/develop.html

- Do you use Dask now?

  - Take the survey: https://dask.org/survey

# Thank you!

aaron@saturncloud.io

https://rikturr.com

@rikturr

Saturn Cloud