# High performance Jupyter: Faster workloads with Dask and RAPIDS

JupyterCon 2020



https://github.com/rikturr/high-performance-jupyter

# Hi!

## Aaron Richter, PhD



Senior Data Scientist @ Saturn Cloud
Organizer @ PyData Miami

> *I work to make data scientists faster and happier*

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

**Saturn Cloud**

**PREFECT**
- Workflow orchestration
- Flow insight and mgmt

**kubernetes**
- Fast setup
- Enterprise secure

# Data science with Jupyter

High performance data science 🚀

# Dask

- *Parallel computing for Python people*

- Anaconda, ~2015

- Built in Python; Python API

- Mature, scientific computing communities

- Low-level task library

- High-level libraries for DataFrames, arrays, ML

- Integrates with PyData ecosystem

- Runs on laptop, scales to clusters

https://dask.org/

# RAPIDS

- *GPU accelerated data science*

- NVIDIA, ~2018

- Built in C++(CUDA), Python; Python API

- Large dev team, support from NVIDIA

- Native DataFrames, arrays, ML, graph, streaming, spatial

- Integrates with PyData ecosystem

- ***Scales to clusters with Dask integration***

https://rapids.ai/

**RAPIDS**

# RAPIDS

Single GPU
Data <16GB

100x faster than
CPU

Numpy -> cuPy
Pandas -> cuDF
Scikit-learn -> cuML

# RAPIDS
## DASK

Dask Array [cuPy]
Dask DataFrame [cuDF]
Dask + cuML

1000s GPUs
Data TB++

100x faster +
"Big data"

**Better
hardware**

**Why not
both??**

## PyData

Single machine
Data <200GB

Slow
computation

NumPy    pandas

scikit
learn

**More
hardware**

## DASK

Dask Array [numpy]
Dask DataFrame [pandas]
Dask ML [scikit-learn]

1000s machines
Data TB++

Horizontal
scaling

# Code time!



https://github.com/rikturr/high-performance-jupyter