# Scalable Data Science with Dask

COILED

# Hi!
# I'm Pavithra :)

Evangelist at Coiled

FOSS contributor - Bokeh & Wikimedia

CSE Student, India

@pavithraes on Twitter

hi@pavithraes.me

# What's your experience with Scalable Data Science?

*On a scale of 1 - 5*

*1: What is scalable compute?*

*5: I understand the challenges and have a solution that works for me*

# We'll talk about

- Big data - What is it?

- Parallel and distributed computing

- Dask for scaling data science

- Coiled - Dask on the cloud

**Slides and notebook at:**
*bit.ly/pyladies-berlin-dask*

# Big data

# What is big data?

- Doesn't fit on your local machine

- Traditional tools and methods fail

COILED

# What is big data?

- Doesn't fit on your local machine

- Traditional tools and methods fail

## Characteristics:

- Volume

- Velocity

- Variety

- Veracity

COILED

# Scalable compute

# Parallel computing

- Working in parallel

- Use all CPU cores
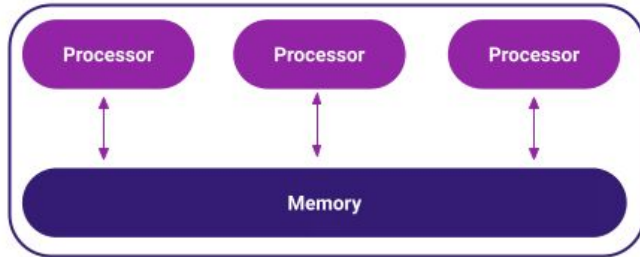
- Multiple processes and shared memory
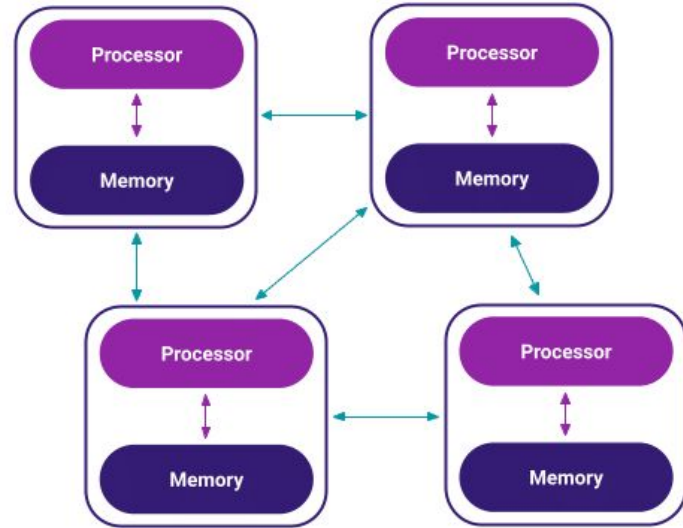


7 am

2 pm

COILED

# Distributed computing

- Using groups of machines

- Each machine has processors and
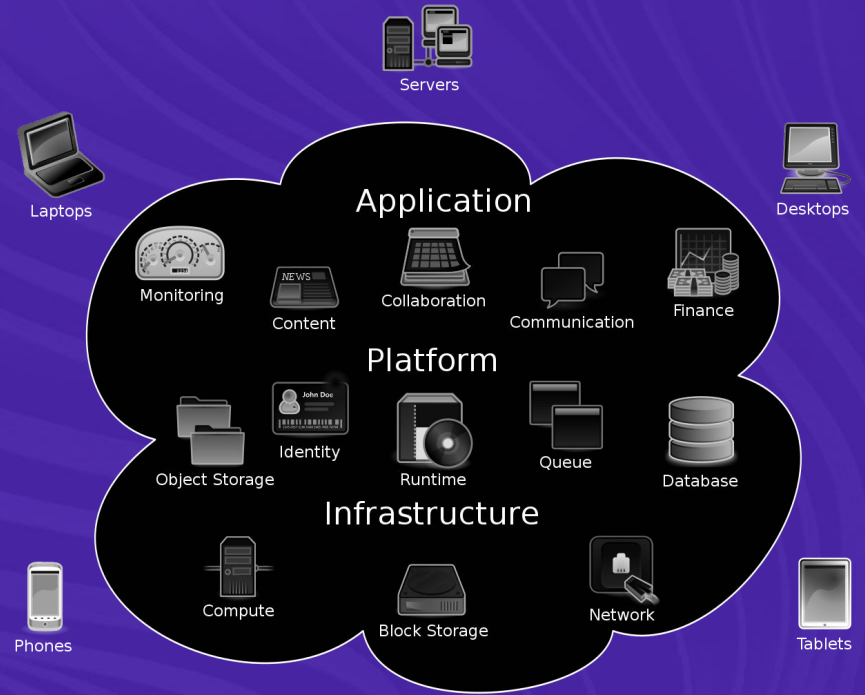  memory

**Parallel Computing** / **Distributed Computing**

Source: Coiled.io

# Cloud computing

- Using cloud resources

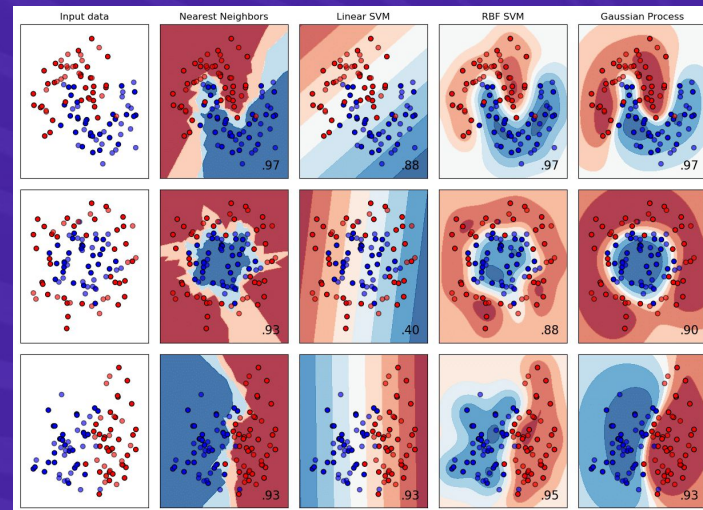- AWS, Azure, GCP

- Lots of storage and computational power



Application

Servers

Laptops

Desktops

Monitoring

Content

Collaboration

Communication

Finance

Platform

Object Storage

Identity

Runtime

Queue

Database

Infrastructure

Phones

Compute

Block Storage

Network

Tablets

*Source: Wikimedia Commons*

COILED

# Dask

# Dask

Library for parallel and distributed

computing in Python

# Dask

Makes it easy to scale-up your workflows to use all cores in your local machine
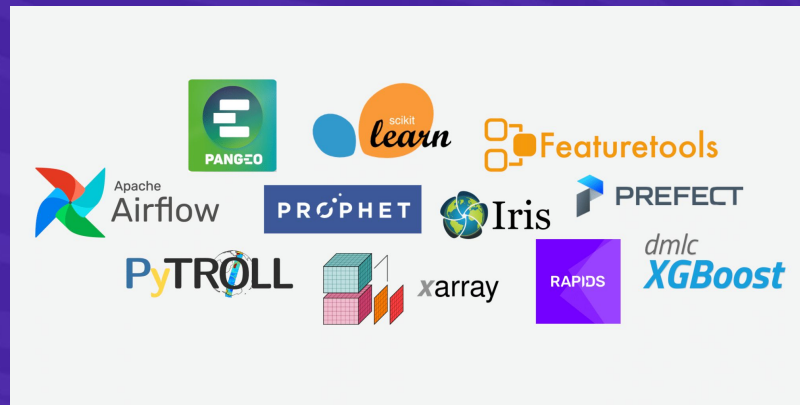
# Dask

Provides a distributed computing framework

Powers tools like RAPIDS, Airflow, PyTorch, and more!



COILED

# Dask features

*Familiar API*
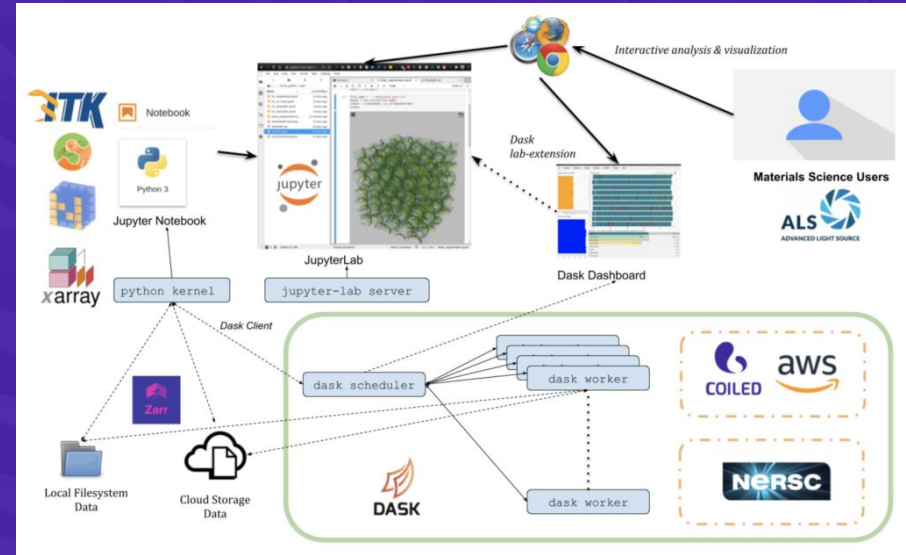
Resembles normal pandas, NumPy,

scikit-learn code

```python
import pandas as pd
df = pd.read_csv("data_taxi/yellow_tripdata_2019-01.csv")
df.groupby("passenger_count").tip_amount.mean()
```

```python
import dask.dataframe as dd
df = dd.read_csv("data_taxi/yellow_tripdata_2019-*.csv")
mean_amount = df.groupby("passenger_count").tip_amount.mean()
mean_amount.compute()
```

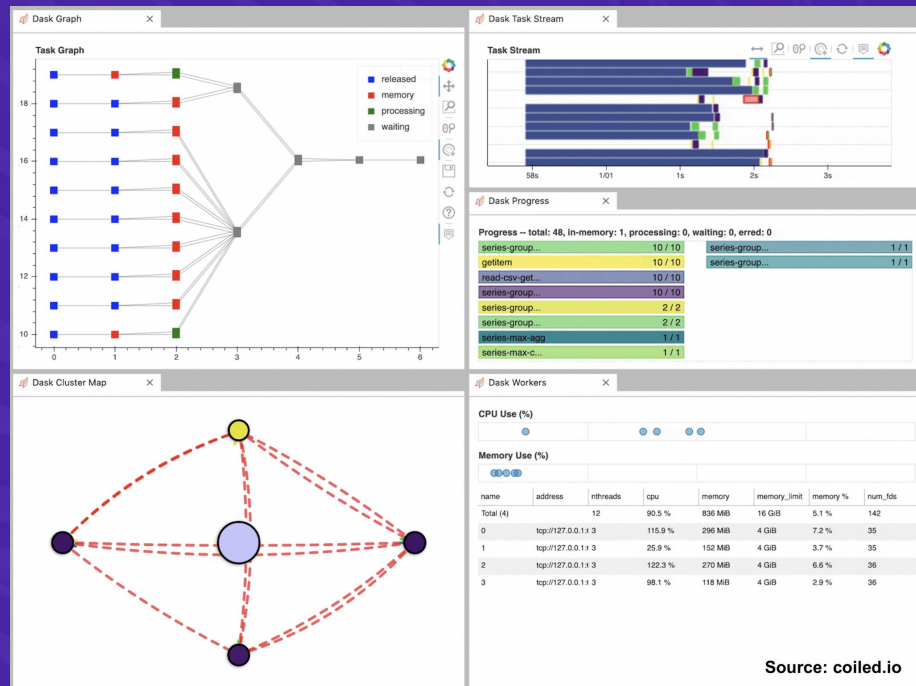COILED

# Dask features

*Flexible*

Local machine, on-prem, on the cloud, anywhere



*3D microstructure interactive image analysis and visualization system architecture.*
*Source: Article presented at Super Computing 2020*

COILED

# Dask features

*Dashboards!*

Real-time visualizations



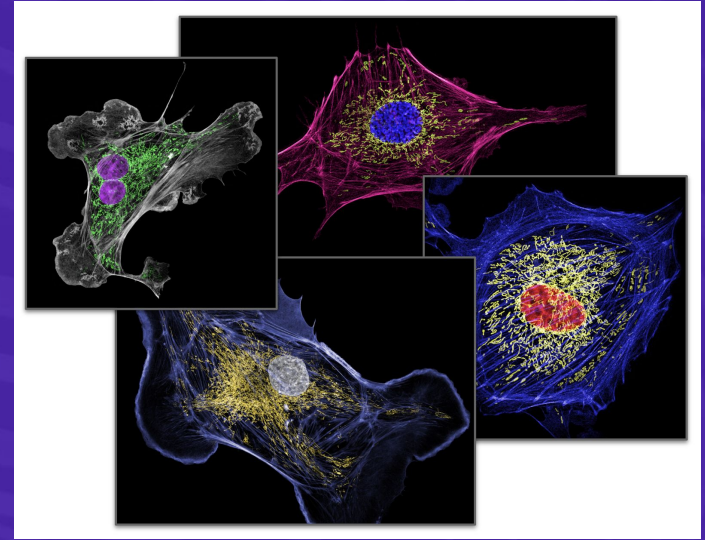Source: coiled.io

COILED

# Dask users in retail

**Walmart**

- demand forecasting

- 500M+ store combinations

- 100x speedup from RAPIDS and Dask



COILED

# Dask users in life science

**High resolution, 4-dimensional, cellular imagery**

- Harvard Medical School
- Howard Hughes Medical Institute
- Chan Zuckerberg Initiative
- UC Berkeley Advanced Bioimaging Center

COILED

# Dask users in finance

**Capital One**

- ETL and ML pipeline speedup

**Barclays**

- Financial system modelling



COILED

# Dask users in Geo

*Farallon Institute , Los Alamos National Labs*

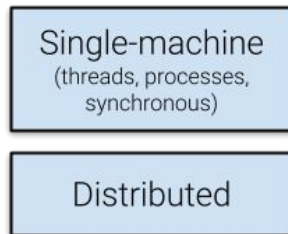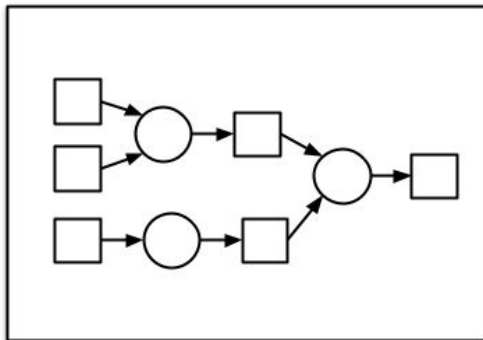- Climate Science
- Energy
- Hydrology
- Meteorology
- Satellite Imaging
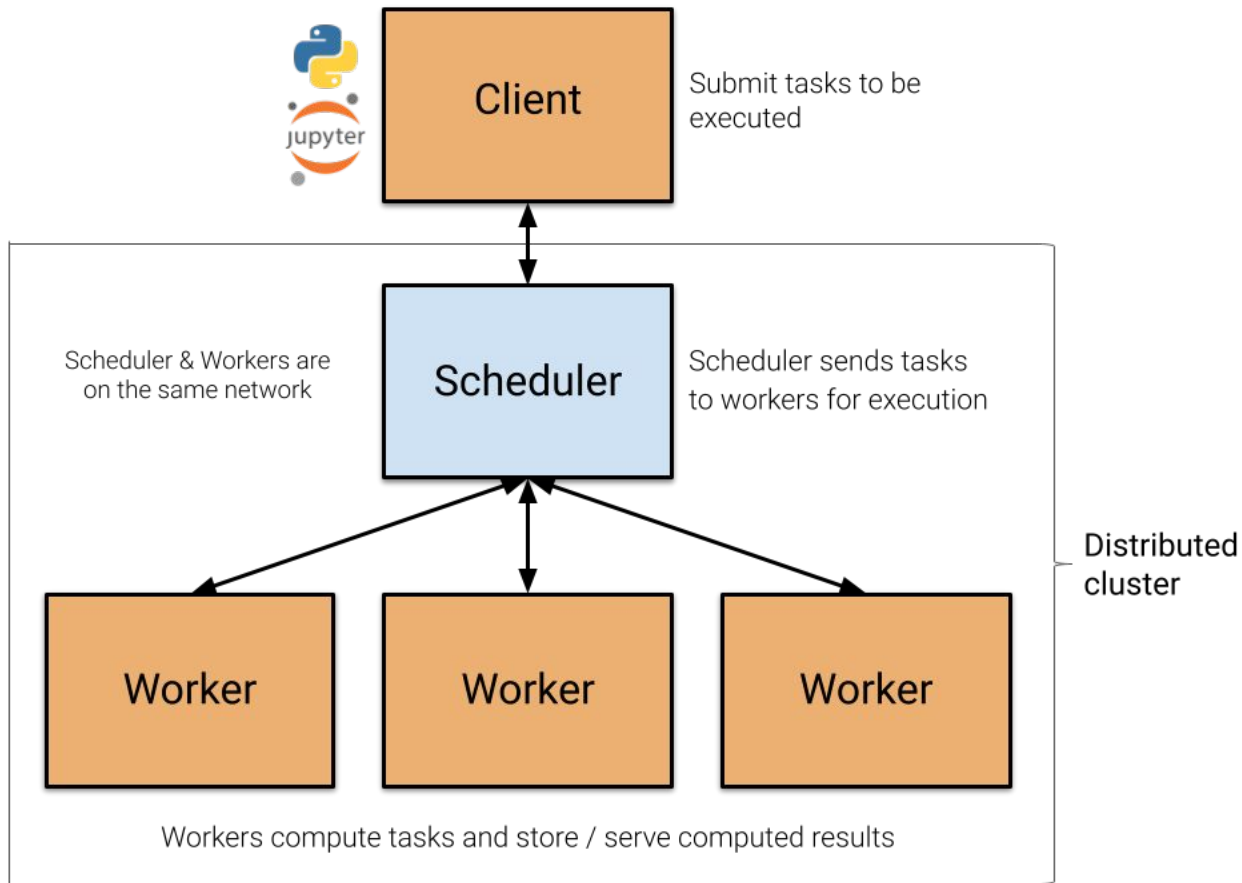


COILED

# Dask users

*Many many more!*

NASA

Software Libraries



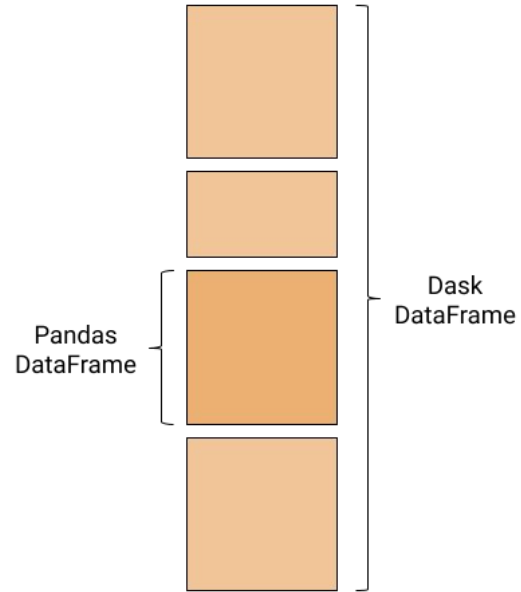COILED

# Demo

**Collections**
(create task graphs)

**Task Graph**

**Schedulers**
(execute task graphs)

Dask Array

Dask DataFrame

Dask Bag

Dask Delayed

Futures

Single-machine
(threads, processes, synchronous)

Distributed

*Source: dask.org*

Source: dask.org

*Shout out to Matthew Rocklin, and the entire Dask community for this material!*

# Coiled

COILED

Built by Dask maintainers, contributors, and enthusiasts.

Open source culture is at the heart of Coiled.

# Cloud computing
## *has some challenges*

- Security concerns

- Managing software environments

- Cost optimization

Coiled tackles these challenges for you.

[welcome.coiled.io](welcome.coiled.io)

# Thank you!

Slides and notebook at:
*bit.ly/pyladies-berlin-dask*

COILED