



FAS Research Computing Division of Science https://rc.fas.harvard.edu



# Python and Multiprocessing on the FASRC clusters





#### Learning objectives

- Python using CLI (Command Line Interface)
  - Interactive
  - Sbatch
- Python Package installation
- Python using OOD (Open On Demand)
- Jupyter Notebook
  - Create conda environment (i.e., jupyter kernel)
- Accelerate Python
  - Multiprocessing
  - Other Tools





### Python Programming Language

- High-level, general-purpose, and object-oriented programming language with emphasis on code readability and use of significant indentation.
- Ideal for scripting and rapid application development given its dynamic typing, elegant syntax, and automatic memory management (garbage collection).
- Has a comprehensive standard library. Also known as "batteries-included" language.
- Python's implementation is mostly in C.
  - Python's core interpreter, CPython, written in C.
- Interpreted language, hence slower than compiled languages, like C and Fortran.
  - Compiled generates executable
  - Interpreted executes instructions directly on the fly without compiling a program into machine language





#### Python using CLI - Interactive

- Login to Cannon: ssh <username>@login.rc.fas.harvard.edu
- Change to a desired location if don't want *\$HOME*: pwd
- Go to a compute node on the test partition:

salloc -p test --nodes=1 --cpus-per-task=2 --mem=12GB --time=00:30:00

• Check Python modules available on Cannon:

module spider python

• Get detailed information on specific module, e.g.:

module spider python/3.10.13-fasrc01

Load the latest (usually also the default) Python module:

module load python





#### Python using CLI - Interactive

# Login to Cannon
ssh <username>@login.rc.fas.harvard.edu

# Change to a desired location if don't want \$HOME: cd <location>

# Go to a compute node on the test partition: salloc -p test --nodes=1 --cpus-per-task=2 --mem=12GB --time=00:30:00

# Check Python modules available on Cannon: module spider python

# Get detailed information on specific module, e.g.: module spider python/3.10.13-fasrc01

# Load the latest (usually also the default) Python module: module load python





#### Python using CLI - Interactive

- Check Python version: python --version
- Invoke Python interpreter: python
- Execute Python programming interactively:

```
def square(x):
    """square a number"""
    return x ** 2
for N in range(1, 4):
    print(N, "squared is", square(N))
```

- o Exit Python: exit()
- Or run a python script interactively:

python myscript.py

*myscript.py* available in User\_Codes





#### Python using CLI - sbatch

• FASRC Examples repository - <u>GitHub - fasrc/User\_Codes</u>

ssh <username>@login.rc.fas.harvard.edu

O pwd (cd to a desired location if not *\$HOME*)

git clone <u>git@github.com</u>:fasrc/User\_Codes.git

```
mkdir python-training
```

cd python-training

```
cp -r ../User_Codes/Languages/Python .
```

cp -r ../User\_Codes/Parallel\_Computing/Python/Python-Multiprocessing-Tutorial

#### CLI Training:

 $\bigcirc$ 

https://docs.rc.fas.harvard.edu/wp-content/uploads/2013/10/Getting-started-on-FASRC-clusters-with-CLI.pptx.pdf





#### Example 1 - Pi using Monte-Carlo

#### https://github.com/fasrc/User\_Codes/tree/master/Languages/Python/Example1

cd Python/Example1 sbatch run.sbatch

- run.sbatch: Batch-job submission script for queuing the job
- mc\_pi.py: Source code for calculating Pi using Monte-Carlo method

#### #!/bin/bash

#SBATCH -J mc\_pi #SBATCH -o mc\_pi.out #SBATCH -e mc\_pi.err #SBATCH -N 1 #SBATCH -C 1 #SBATCH -c 1 #SBATCH -p serial\_requeue #SBATCH -t 0-00:30 #SBATCH --mem=4000

# Load required modules
module load python

# Run program
srun -c 1 python mc\_pi.py

# job name # standard output file # standard error file # number of nodes # number of cores # partition # time in D-HH:MM # memory in MB





#### Example 2 - Figures with LaTex Font

#### https://github.com/fasrc/User\_Codes/tree/master/Languages/Python/Example2

cd Python/Example2 sbatch run.sbatch

 tex\_demo.py: source code for generating figures using LaTex fonts

```
#!/bin/bash
#SBATCH –J tex demo
#SBATCH -o tex_demo.out
#SBATCH -e tex_demo.err
#SBATCH -c 1
#SBATCH -p serial_requeue
#SBATCH -t 0-00:30
#SBATCH --mem=4000
# Load required modules
module load python/3.10.12-fasrc01
module load texlive/2018.06.15-fasrc01
# Run program
srun -c 1python tex_demo.py
```





#### Python Package Installation

#### • Go to a compute node on the test partition:

salloc -p test --nodes=1 --cpus-per-task=2 --mem=12GB --time=01:00:00

• Create a vanilla mamba/conda environment (for multiprocessing exercise):

```
module load python
conda create --prefix=/n/holylabs/LABS/<desired-folder>/multiproc_env
python=3.11 -y
```

• Alternatively, if default *\$HOME* is desired, then do following instead:

```
module load python
conda create --name multiproc_env python=3.11 -y
```

#### See <u>Python – FASRC DOCS</u>





#### Python Package Installation

• Activate conda/mamba environment:

source activate /n/holylabs/LABS/<desired-folder>/multiproc\_env

- O Or if \$HOME used, then: source activate multiproc\_env
- Install relevant python packages (Mamba recommended):

conda install numpy pandas matplotlib -y pip install jupyterlab swifter

- Always pip install inside a conda environment to avoid package conflicts
- Deactivate the environment: Conda deactivate
- See <u>https://docs.rc.fas.harvard.edu/kb/python/#Mamba</u>





### Python Using Open OnDemand (OOD)

- Open-source web portal to access clusters
- Web-based, no software needs be installed on your local laptop/desktop (except for a modern browser like Google Chrome, Mozilla Firefox)
- Easy to learn and simple to use
- Very similar to desktop applications
- $\circ$  ~ The easiest way to run GUI applications remotely on a cluster
- $_{\odot}$   $\,$  Safari is not recommended for OOD  $\,$
- OOD Training: <u>https://docs.rc.fas.harvard.edu/wp-content/uploads/2013/10/Getting-started-on-FASRC-clusters-with-OOD-May20</u> <u>24.pdf</u>





#### How to access OOD on FASRC Clusters

- Accessing OOD from Cannon
  - Connect to FASRC VPN <u>Virtual Desktop (VDI) through Open OnDemand FASRC</u> <u>DOCS</u>
  - Then go to <u>https://rcood.rc.fas.harvard.edu</u>
- Accessing OOD from FASSE
  - Connect to FASSE VPN FASSE VDI Apps FASRC DOCS
  - Then go to <u>https://fasseood.rc.fas.harvard.edu</u>





#### Filling a form to launch an app

- Request the resources that you need
   (If you don't know for a first trial run, use similar resources as your laptop/desktop)
  - Partition (Name): depends on <u>Cannon</u> (URL) vs <u>FASSE</u> (URL)
  - Memory (RAM): amount of memory in GB
  - Number of cores: recommended at least 2
  - Number of GPUs: if >= 1, make sure you **select** a gpu partition
  - \_Allocated time: time you would like your session to run\_\_\_\_\_\_
  - Email for status notification: to know when job starts, ends
  - Reservation: if you have a special reservation (this requires approval from FASRC)
  - Account: use this if you have more than one PI\_lab affiliation

the minimum and/or maximum values of each field depends on the selected partition





### Jupyter Notebook

- Launch new Jupyter Notebook session (existing session will not work!)
- $\circ$   $\,$  Select newly created conda environment as the kernel
  - a. Open a notebook
  - b. On the top menu, click Kernel -> Select Kernel -> Click on OOD\_env
  - c. Note: kernels is the same as conda, python, mamba environment







### Closing running OOD windows/tabs

- In most OOD apps, you can close the browser tab while the code is running, and the code will continue to run on the background
- Jupyter Notebook will not! The cell that is running will lose the data and output files will not be written
  - Solution: run Remote Desktop app and launch Jupyter Notebook from within Remote Desktop
  - Documentation: <u>Open OnDemand (OOD/VDI) Remote Desktop: How to open</u> <u>software – FASRC DOCS</u>





FAS Research Computing Division of Science https://rc.fas.harvard.edu

### FASSE proxy

Documentation: FASSE Proxy Settings – FASRC DOCS

- $\circ$   $\,$  You may need to set FASSE proxy on
  - Firefox (web browsing)
  - Jupyter Notebook
  - Access Github
  - (Basically, anything outside of FASSE)





# Multiprocessing in Python

- Ability of a system to run multiple processors at one time
- Allows several processes to run simultaneously
- Multiprocessing module allocates tasks to different processors and makes better use of a multi-core machine
- Different from threading, which is subject to Global Interpreter Local (GIL) and one thread can execute only one Python code
- <u>Multiprocessing vs Threading Python Stack Overflow &</u>
- o <u>multiprocessing Process-based parallelism Python 3.12.4 documentation</u>





### Multiprocessing in Python

- On the cluster, difference between number of CPUs allocated to the job vs total number of CPUs available on the node
- Go to a compute node on the test partition requesting 10 cores:

salloc -p test --nodes=1 --cpus-per-task=10 --mem=12GB --time=01:00:00

• See total number of cores available on the node:

scontrol show node <nodename>

Execute cpu-count.py to see which command gives you the number of cores allocated to your job:
 Cd Python-Multiprocessing-Tutorial python cpu-count.py

• See <u>How to find out the number of CPUs using python - Stack Overflow</u>





### Multiprocessing - Process-based Parallelism - Basic

- <u>Multiprocessing in Python MachineLearningMastery.com</u>
- Two functions declared to execute print statements after sleeping for 2 & 3 seconds, resp.
- 3 processes created using *multiprocessing.Process* inside *main()*
- The *Process()* utilizes *target* argument to run target process
- Processes are run using *start()*
- Use *join()* to run & exit a processes
   before the main program process



#### HARVARD UNIVERSITY



# Multiprocessing - Pooling

- Run 1000 processes together may not be possible
- Create a process pool to limit number of processes that can be run at a time
- Function declared to return the cube
- The *multiprocessing.Process* doesn't work with p.start() & p.join(), would need an output queue as well. But faster than *Pool()*
- The *multiprocessing.Pool* module easier to use, returns ordered result using *pool.map()*, & causes less overhead

```
import multiprocessing
import time
import os
def cube(x):
    return x**3
if name == ' main ':
  # The Process class
  processes = [multiprocessing.Process(target=cube, args=(x,)) for x in
range(1,len(os.sched getaffinity(0)))]
  [p.start() for p in processes]
  result process = [p.join() for p in processes]
  # The Pool class
  pool = multiprocessing.Pool(processes=len(os.sched getaffinity(0)))
  result pool = pool.map(cube, range(1,len(os.sched_getaffinity(0)))]
```

• See Python multiprocessing: How to know to use Pool or Process? - Stack Overflow 21





### Multiprocessing + Numpy

- Using Multiprocessing along with Numpy to accelerate python program
- Go to OOD (Cannon or FASSE) & launch JupyterLab notebook on *test* with
  - 52 CPUs
  - *gcc/12.2.0-fasrc01* loaded as a module
  - multiproc\_env loaded as a kernel
  - In python-training/Python-Multiprocessing-Tutorial
- Problem Statement:
  - A sample data file has 4 columns and 1000 entries. Columns correspond to the time a job was submitted, when it started, when it ended, and number of CPUs allocated.
  - Calculate the total number of CPUs in use by currently running jobs for every submitted job





### Multiprocessing + Numpy

- Convert numerical columns to Numpy arrays.
- Declare a function to calculate CPUs utilized: calculate\_cpus\_utilized()
- $\circ$   $\,$  Multiple methods utilized for the calculation:
  - Use the function over each submitted-job entry
    - Pandas apply()
    - swifter.apply()
  - Using Numpy arrays & for-loop
  - Using Multiprocessing with a pool of processes = #CPUs requested for OOD job
- Run the notebook to see which method gives the fastest result
- Fastest: Combination of Numpy and Multiprocessing





#### Accelerate Python - Other Tools

- o Numba
  - <u>https://numba.pydata.org/</u>
- o Swifter
  - <u>Speed up your Pandas Processing with Swifter | by Cornellius Yudha</u> <u>Wijaya | Towards Data Science</u>
  - <u>GitHub jmcarpenter2/swifter: A package which efficiently applies any</u> <u>function to a pandas dataframe or series in the fastest available manner</u>
- o Dask
  - <u>https://www.dask.org/</u>





# Other Training(s) on Python

Training portal: <a href="https://trainingportal.harvard.edu/Saba/Web\_spf/NA1PRD0068/app/dashboard">https://trainingportal.harvard.edu/Saba/Web\_spf/NA1PRD0068/app/dashboard</a>

#### Data Handling in Python Workshop:

https://github.com/HarvardRC/Python\_Data\_Handling

- Requirement: working FASRC account with cluster access
- Audience
  - Users familiar with command-line interface
  - New to Cannon and FASSE, but familiar with Python
- o **Content** 
  - Python Data Types: Focus on Collections
  - File I/O Operations: Emphasis on profiling & performance
  - Parallel Processing





# FASRC Upcoming Trainings

Training calendar: <a href="https://www.rc.fas.harvard.edu/upcoming-training/">https://www.rc.fas.harvard.edu/upcoming-training/</a>

#### Parallel Job Workflows

This training would focus on best practices for running parallel workflows on FASRC clusters. The module would provide the basic knowledge to execute OpenMP and MPI applications efficiently on the cluster.

Objectives:

- Best Practices
- Brief Introduction to Parallel Computing
- Embarrassingly Parallel Jobs/Workflows
- OpenMP Jobs/Workflows
- MPI Jobs/Workflows
- Hybrid (MPI+OpenMP) Jobs/Workflows





### FASRC documentation

- FASRC docs: <u>https://docs.rc.fas.harvard.edu/</u>
- GitHub User\_codes: <u>https://github.com/fasrc/User\_Codes/</u>
- Getting help
  - Office hours: <u>https://www.rc.fas.harvard.edu/training/office-hours/</u>
  - Ticket
    - Portal: <u>http://portal.rc.fas.harvard.edu/rcrt/submit\_ticket</u> (requires login)
    - Email: <u>rchelp@rc.fas.harvard.edu</u>





#### Survey

Please, fill out our course survey. Your feedback is essential for us to improve our trainings!!

http://tinyurl.com/FASRCsurvey





FAS Research Computing Division of Science https://rc.fas.harvard.edu



#### **Thank you :)** FAS Research Computing