

Harnessing the Power of HPC for Machine Learning

Tools and Techniques

Charles Peterson

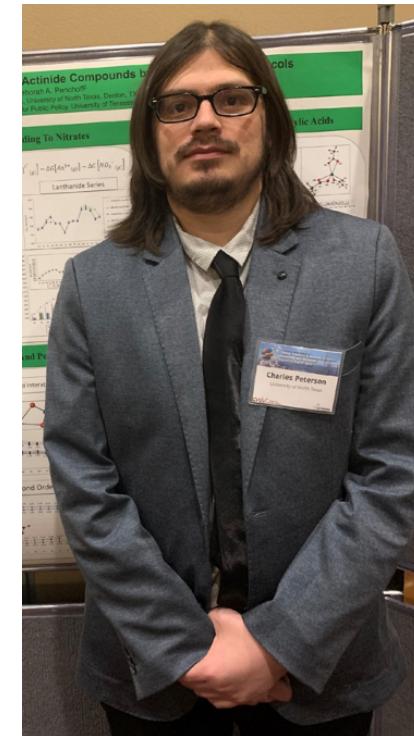
Welcome Everyone!

This workshop provides an overview of topics and practical examples of using Machine Learning tools on HPC resources. 

Key Topics:

- Python/R in HPC (Hoffman2)
- ML Package Installation
- Interactive and Batch Job Submission
- Big Data Insights

For suggestions: cpeterson@oarc.ucla.edu





Access the Workshop Files

This presentation and accompanying materials are available on  [UCLA OARC GitHub Repository](#)

You can view the slides in:

-  PDF format - WS_MLonHPC.pdf
-  HTML format: [Workshop Slides](#)

Each file provides detailed instructions and examples on the various topics covered in this workshop.

Note:  This presentation was built using [Quarto](#) and RStudio.

Machine Learning and HPC



Machine Learning Basics

- What is Machine Learning?
 -  **Machine Learning (ML)** is a subset of artificial intelligence (AI) focused on building systems that learn from and make decisions based on data.
- Key Concepts:
 - **Data:** The foundation of any ML model. It can be labeled (supervised learning) or unlabeled (unsupervised learning).
 - **Algorithms:** Procedures or formulas for solving a problem. Common ML algorithms include linear regression, decision trees, and neural networks.
 - **Training:** The process of teaching a machine learning model to make predictions or decisions based on data.
 - **Inference:** Applying the trained model to new data to make predictions.
- Types of Machine Learning:
 -  **Supervised Learning:** The model learns using labeled data (e.g., spam detection).
 -  **Unsupervised Learning:** The model identifies patterns in data without any labels (e.g., customer segmentation).

? What is HPC and Why Should I Care

-  HPC uses **MANY computers** to solve large problems faster than a normal computer.
-  If your task takes a long time to run on a laptop or a lab's server, HPC can ‘speed up’ your application.
-  You can store **LARGE amounts of data**, too big for your laptop.



General vs High-Performance Computing

General Purpose Computing

- 🚧 Only one person at a time.
- 💻 Calculations run on the machine directly.
- 🗅️ Can only run 1-2 calculations at a time.

High-Performance Computing

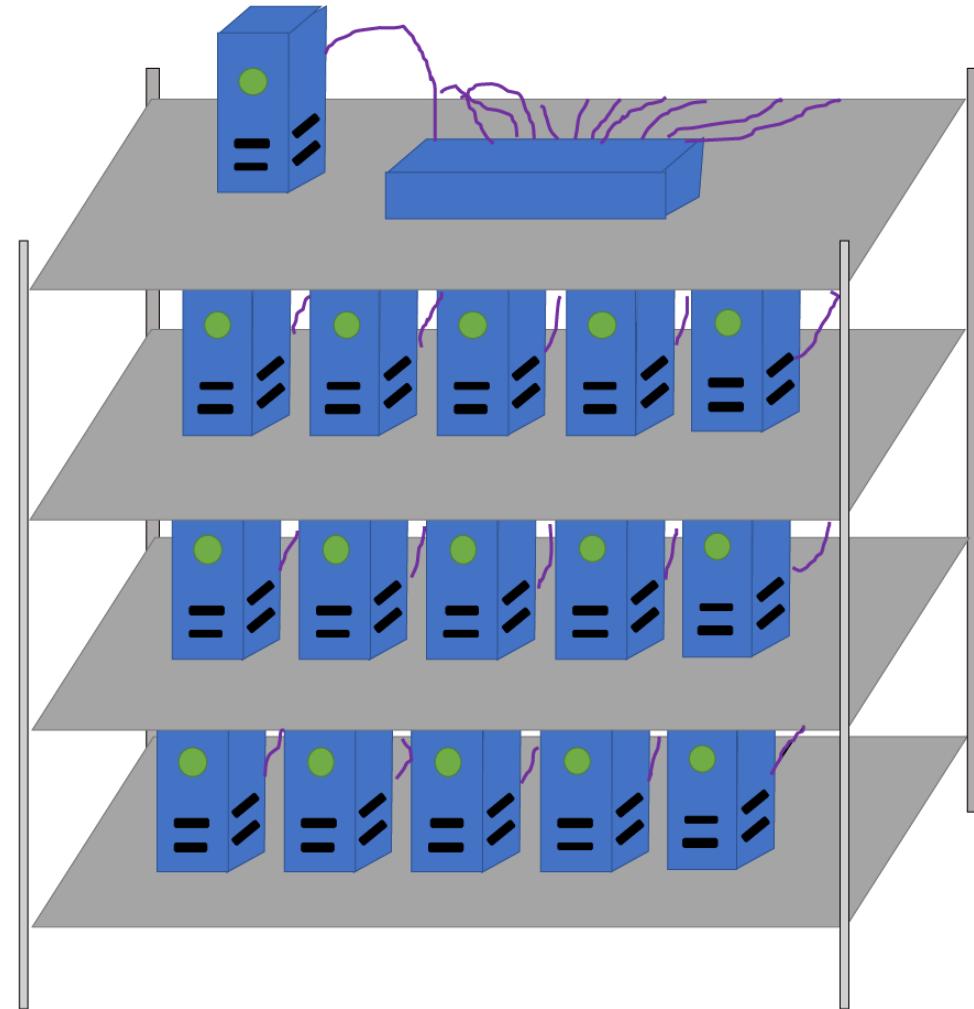
- 👤 Multiple people can log in at one time.
- ⏳ Calculations are ‘scheduled’ to run on a different machine.
- 💻 Can run hundreds of calculations at one time.

🤝 HPC provides an excellent platform for collaborations and faster results.

HPC Overview

Beowulf style cluster

Multiple computers
Single computing resource



The Power of HPC

Single Computer Limitations:

-  Only one CPU.
-  Large problems cannot fit.
-  Long processing times.
-  Limited memory and disk space.

HPC Solutions:

-  More CPUs for faster processing.
-  More memory to handle bigger tasks.
-  More disk space for extensive data.
-  Access to GPUs for advanced computations.



The Power of HPC: Parallelization

Single CPU Program

- Task A ➔ Task B ➔ Task C
- ⏳ Total Time: 3 hours

Parallel Tasks

- Task A on CPU 1 ⏳ 1 hour
- Task B on CPU 2 ⏳ 1 hour
- Task C on CPU 3 ⏳ 1 hour
- ⚡ Total Time: 1 hour

🤖 Most Machine Learning packages can utilize multiple CPUs and GPUs to run your models in parallel!

Multi-tasking with Machine Learning Jobs

-  With HPC resources, you can have **multiple machine learning jobs running concurrently**.
-  This parallel processing greatly increases efficiency and productivity.
-  Ideal for complex computations and large-scale data analysis.

```
[ccrp0051@l3-login1 ~]$ squeue -u ccrp0051
   JOBID PARTITION  NAME     USER ST      TIME  NODELIST(REASON)
14225  public yac3_h3-  ccrp0051 PD      0:00  1 (QOSMaxNodePerUserLimit)
14226  public yac3_h3-  ccrp0051 PD      0:00  1 (QOSMaxNodePerUserLimit)
14227  public yac3_h3-  ccrp0051 PD      0:00  1 (QOSMaxNodePerUserLimit)
14228  public yac3_h3-  ccrp0051 PD      0:00  1 (QOSMaxNodePerUserLimit)
14229  public yac3_h3-  ccrp0051 PD      0:00  1 (QOSMaxNodePerUserLimit)
14230  public yac3_h3-  ccrp0051 PD      0:00  1 (QOSMaxNodePerUserLimit)
14231  public yac3_h3-  ccrp0051 PD      0:00  1 (QOSMaxNodePerUserLimit)
14232  public yac3_h3-  ccrp0051 PD      0:00  1 (QOSMaxNodePerUserLimit)
14233  public yac3_h3-  ccrp0051 PD      0:00  1 (QOSMaxNodePerUserLimit)
14234  public yac3_h3-  ccrp0051 PD      0:00  1 (QOSMaxNodePerUserLimit)
14235  public yac3_h3-  ccrp0051 PD      0:00  1 (QOSMaxNodePerUserLimit)
14236  public yac3_h3-  ccrp0051 PD      0:00  1 (QOSMaxNodePerUserLimit)
14237  public yac3_h3-  ccrp0051 PD      0:00  1 (QOSMaxNodePerUserLimit)
14238  public yac3_h3-  ccrp0051 PD      0:00  1 (QOSMaxNodePerUserLimit)
14239  public yac3_h3-  ccrp0051 PD      0:00  1 (QOSMaxNodePerUserLimit)
14240  public yac3_h3-  ccrp0051 PD      0:00  1 (QOSMaxNodePerUserLimit)
14241  public yac3_h3-  ccrp0051 PD      0:00  1 (QOSMaxNodePerUserLimit)
14242  public yac3_h3-  ccrp0051 PD      0:00  1 (QOSMaxNodePerUserLimit)
14243  public yac3_h3-  ccrp0051 PD      0:00  1 (QOSMaxNodePerUserLimit)
14244  public yac3_h3-  ccrp0051 PD      0:00  1 (QOSMaxNodePerUserLimit)
14245  public yac3_h3-  ccrp0051 PD      0:00  1 (QOSMaxNodePerUserLimit)
14224  public yac3_h3-  ccrp0051 R  3:13:16  1 compute-1_22-3
14223  public yac3_h3-  ccrp0051 R  3:54:19  1 compute-2_24-1
14222  public yac3_h3-  ccrp0051 R  11:57:35  1 compute-1_38-3
12458  public lac3_h3-  ccrp0051 R  7:01:09:46  1 compute-1_16-4
12445  public lac3_h3-  ccrp0051 R  8:16:31:81  1 compute-1_26-1
14221  public yac3_h3-  ccrp0051 R  2:14:54:45  1 compute-11_14-3
14228  public yac3_h3-  ccrp0051 R  2:15:58:37  1 compute-11_20-3
14219  public yac3_h3-  ccrp0051 R  2:15:58:59  1 compute-1_32-1
14218  public yac3_h3-  ccrp0051 R  2:16:01:05  1 compute-3_16-2
14217  public yac3_h3-  ccrp0051 R  2:16:04:31  1 compute-11_14-1
14216  public yac3_h3-  ccrp0051 R  2:16:06:36  1 compute-3_24_1
14215  public yac3_h3-  ccrp0051 R  2:16:09:14  1 compute-2_16-3
14214  public yac3_h3-  ccrp0051 R  2:16:16:37  1 compute-2_18-1
14213  public yac3_h3-  ccrp0051 R  2:22:02:23  1 compute-3_26-1
14212  public yac3_h3-  ccrp0051 R  3:21:37:34  1 compute-3_18-4
14211  public yac3_h3-  ccrp0051 R  4:00:49:58  1 compute-2_30_3
14209  public yac3_h3-  ccrp0051 R  4:02:52:15  1 compute-3_12-4
14206  public yac3_h3-  ccrp0051 R  4:03:27:42  1 compute-3_30_1
14205  public yac3_h3-  ccrp0051 R  4:03:30:36  1 compute-11_32-1
14204  public yac3_h3-  ccrp0051 R  4:03:31:57  1 compute-3_14-1
14203  public yac3_h3-  ccrp0051 R  4:05:52:32  1 compute-2_18-3
14202  public yac3_h3-  ccrp0051 R  4:08:19:03  1 compute-3_14-2
[ccrp0051@l3_login1 ~]$ _
```



Python

Hoffman2 supports running 🐍 Python applications.

Hoffman2 supports 🐍 Python applications, and it is **HIGHLY** recommended to use Python versions built and tested by Hoffman2 staff.

🚫 Avoid using system python builds (e.g., `/usr/bin/python`). Instead, use **module load** commands to access optimized versions.

- To see all Python versions installed on Hoffman2:

```
1 modules_lookup -m python
```

- Load a Python module

```
1 module load python/3.7.3
2 which python3
```

- This example shown:
 - Python version 3.7.3
 - Location of python
 - `/u/local/apps/python/3.7.3/gcc-4.8.5/bin/python3`
 - (Location of the Hoffman2 installed python)



Python Packages in Machine Learning



Scikit-learn:

- Versatile tools for machine learning, including classification, regression, and clustering.



TensorFlow:

- Google's library for deep learning and neural networks.



Keras:

- Python interface for neural networks, primarily an interface for TensorFlow.



PyTorch:

- Flexible deep learning library by Facebook's AI Research lab.

 **XGBoost:**

- Efficient gradient boosting library, ideal for structured data.

 **Pandas:**

- Essential for data manipulation and analysis, a cornerstone in machine learning.

 **NumPy:**

- Fundamental for scientific computing, supports large arrays and matrices.

 **SciPy:**

- For scientific and technical computing, extends NumPy's capabilities.



Python Installation on Hoffman2

Basic Builds on Hoffman2: The Python builds on Hoffman2 include only the basic compiler/interpreter and a few essential packages.

User-Installed Packages:

- 📦 Most Machine Learning Python applications will require additional packages, installed by the user.
- 🔍 Hoffman2 staff do **not** install extra packages in the supported Python builds to avoid conflicts.

Installing Machine Learning Packages:

- ⭐ When using Python (or R), you'll need to install the ML packages yourself.
- 📚 We have a workshop covering this topic in detail:
 - 🔗 [Python/R Installation Workshop](#)



User-Installed Packages on Hoffman2

- Users cannot install packages in the main Python build directories.
 - This is to avoid version conflicts and dependencies issues that could break Python.
- Users can install packages in their own directories:
 - `$HOME`, `$SCRATCH`, or any project directories.

Installation Methods:

- **Using pip package manager:** Ideal for standard Python package installations.
- **Using Python Virtual Environments:** Creates isolated environments for specific projects.
- **Using Anaconda:** Suitable for managing complex package dependencies and environments.



Using pip Package Manager

Installing scikit-learn with pip:

- To install the **scikit-learn** package via pip (PyPI) package manager:

```
1 module load python/3.7.3
2 pip3 install scikit-learn --user
```

Understanding the **--user** Flag:

- The **--user** flag ensures the package installs in your \$HOME directory.
- By default, pip tries to install in the main Python build directory, where users lack write access.
- Using **--user**, packages install in \$HOME/.local, avoiding permission errors.



R on Hoffman2

Finding Available Versions of R:

- ★ Hoffman2 supports various versions of R.
- To view all available versions of R on Hoffman2:

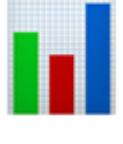
```
1 modules_lookup -m R
```

Loading a Specific Version of R: Example to load R version 4.2.2 with GCC version 10.2.0:

```
1 module load gcc/10.2.0
2 module load R/4.2.2
```

Ensuring Correct Module Loads:

- 🔧 Load the gcc or intel modules first, as indicated by `modules_lookup`. This step ensures that the correct versions of gcc and intel libraries are loaded for R.



R Packages



Caret:

- Framework for building machine learning models. Offers tools for data pre-processing, feature selection, and model tuning.



RandomForest:

- Implements the random forest algorithm. Known for performance in classification and regression.



e1071:

- Contains functions for SVMs, naive Bayes classifier, and more.



nnet:

- For training single-hidden-layer neural networks and multinomial log-linear models.



rpart:

- Recursive partitioning for decision tree models.



xgboost:

- Efficient gradient boosting, effective for large datasets.



glmnet:

- Fitting generalized linear and Cox models via penalized likelihood.



tm:

- Text mining framework, managing and mining text data.



ggplot2:

- Powerful data visualization tool based on the Grammar of Graphics.



dplyr:

- Essential for data manipulation, providing a set of tools for dataset management.



R Package Installation

Standard Installation Command:

- Use the following command to install R packages:

```
1 install.packages('PKG_name')
```

- 🚫 On Hoffman2 (and most other HPC resources), you cannot modify the main R global directory.
- Example Installation:

```
1 install.packages("dplyr")
```

- 🏡 R will suggest a new path in your \$HOME directory, determined by \$R_LIBS_USER.
- Each R module on Hoffman2 has a unique \$R_LIBS_USER to prevent conflicts between different R versions.



Anaconda

Anaconda is a popular Python and R distribution, ideal for simplifying package management and pipelines.

Hoffman2 has Anaconda installed, allowing users to create their own conda environments.

```
1 module load anaconda3
```



Warning



No Need for Other Python/R Modules:

- Your Anaconda environment includes a build of Python and/or R. Loading other modules may cause conflicts.



Note

For more information, we had done a workshop on using [Anaconda on Hoffman2](#) that you can review.



Containers

Containers, like Apptainer and Docker, are excellent for running Machine Learning applications on Hoffman2.

Advantages of Containers:

- 🚧 **Isolated Environments:** Comes with all necessary Machine Learning software pre-installed.
- 🚚 **Portability:** Use the same container on different computers, ensuring version control and reproducibility.

Apptainer on Hoffman2:

- 🔧 Hoffman2 uses Apptainer for running containers.
- 🔎 For more information, refer to our previous workshop:
 - 🔗 [Containers Workshop](#)

Example: Fashion MNIST



Fashion MNIST

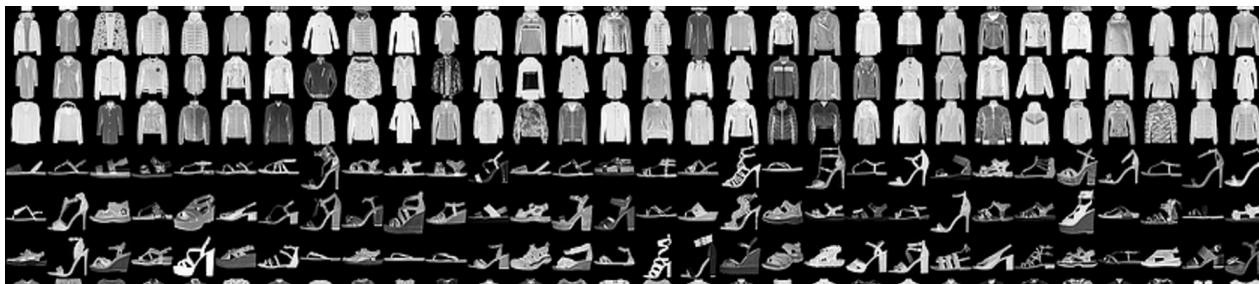
This example focuses on the “Fashion MNIST” dataset, a collection used frequently in machine learning for image recognition tasks.

Approach:

- We will use a Random Forest algorithm to train a model for predicting fashion categories.

Dataset Overview:

- **Images:** 28x28 grayscale images of fashion products.
- **Categories:** 10, with 7,000 images per category.
- **Total Images:** 70,000.





ML Packages for Python and R

Using Scikit-learn with Python:

- Ideal for algorithms like classification and clustering.
- Useful for preprocessing, model building, and evaluation.

Package Installation:

Python:

- Install Python and Scikit-learn:

```
1 module load python/3.9.6
2 pip3 install sklearn --user
```

R:

- Install R and necessary packages:

```
1 module load gcc/10.2.0
2 module load R/4.2.2
3 # Needed for OpenML package
4 module load libxml2
5 R -e 'install.packages(c("randomForest", "OpenML", "dplyr", "ggplot2", "caret", "farff"), repos = "https://cran.r-project.org/")'
```

Python Example Run

Getting Started with Interactive Compute Node

- Start by requesting an interactive compute node:

```
1 qrsh -l h_data=10G
```

Cloning and Navigating to the Code Repository

- Clone the repository and navigate to the mnist-ex directory:

```
1 cd $SCRATCH
2 git clone https://github.com/ucla-oarc-hpc/WS_MLonHPC
3 cd WS_MLonHPC/mnist-ex
```

Lets look at the code, `mnist.py`

Running the Python Script:

- Load Python module and run the `mnist.py` script:

```
1 module load python/3.9.6
2 python3 mnist.py
```



Parallel Processing with Python

The initial training took about 1 minute over 1 CPU core.

Speeding Up with Parallel Processing:

- Request 10 cores for parallel processing:

```
1 qrsh -l h_data=10G -pe shared 10
```

Note: Use the shared parallel environment as sci-kit learn doesn't support multi-node parallelism.

Code Adjustment for Parallelism:

- Lets look at the code, [mnist-par.py](#)
- The main change is n_jobs option in the Classifier

```
1 clf = RandomForestClassifier(random_state=42, n_jobs=10)
```

Run the code!

```
1 module load python/3.9.6
2 python3 mnist-par.py
```



Batch Submission on Hoffman2

Submitting Non-Interactive Jobs:

- For tasks that don't require interactive sessions, you can submit jobs to be processed in the background.

Command to Submit a Job:

- Use the `qsub` command to submit your job script to the queue:

```
1 qsub mnist-py.job
```

Advantages:

- Efficient for longer or resource-intensive tasks.
- Allows you to free up your session while the job runs in the background.



Running R

Executing Code with a Single CPU:

- Start with requesting an interactive compute node:

```
1 qrsh -l h_data=10G
2 module load gcc/10.2.0
3 module load R/4.2.2
4 Rscript mnist.R
```

Running Code with Parallel Processing (10 CPUs):

- Request multiple cores for parallel execution:

```
1 qrsh -l h_data=10G -pe shared 10
2 module load gcc/10.2.0
3 module load R/4.2.2
4 Rscript mnist-par.R
```

Submitting as a Batch Job:

- For non-interactive execution, submit the job script:

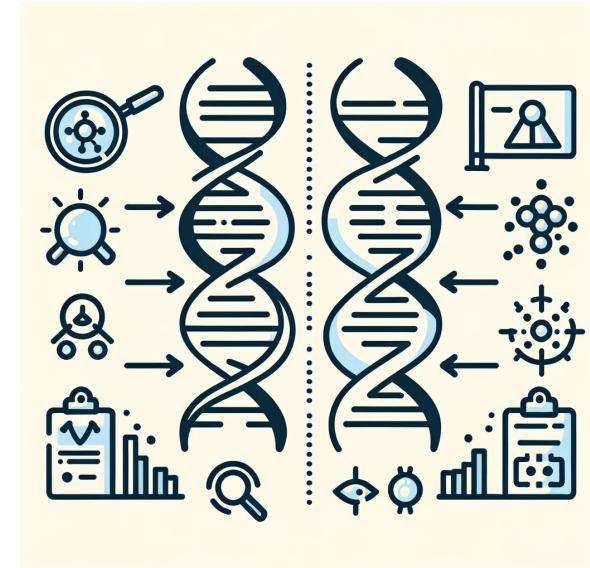
```
1 qsub mnist-R.job
```

Example: DNA Sequence

DNA Sequence classification

DNA Sequence Classification with PyTorch

-  **Objective:** Create a model to classify DNA sequences into 'gene' or 'non-gene' regions.
- **Gene Regions:** Segments of DNA containing codes for protein production.
- **Dataset Creation:** Generate random DNA sequences labeled as 'gene' or 'non-gene'.
-  **Model Development:** Use PyTorch to build a model predicting the presence of 'gene' regions.
-  **Leveraging GPUs:** Utilize the parallel processing power of GPUs for efficient training.





Creating a Conda Environment

Setting Up for GPU-Enabled PyTorch:

- Begin by loading the Anaconda module:
- Create a new Conda environment named biotest with Python, scikit-learn, and scipy:
- Activate the newly created environment:
- Install PyTorch with GPU support using pip:

```
1 module load anaconda3
```

```
1 conda create -n biotest python=3.11 scikit-learn
```

```
1 conda activate biotest
```

```
1 pip3 install torch
```



Running PyTorch on Hoffman2

Code Location and Versions:

- The code for this task is located in the `dna-ex` directory.
- There are two versions:
 - `dna-cpu.py` for the CPU version.
 - `dna-gpu.py` for the GPU version.

Executing the Examples:

- Request a node with GPU resources:
- Run CPU version
- Run GPU version

```
1 qrsh -l h_data=10G,gpu,A100
```

```
1 conda activate biotest
2 python3 dna-cpu.py
```

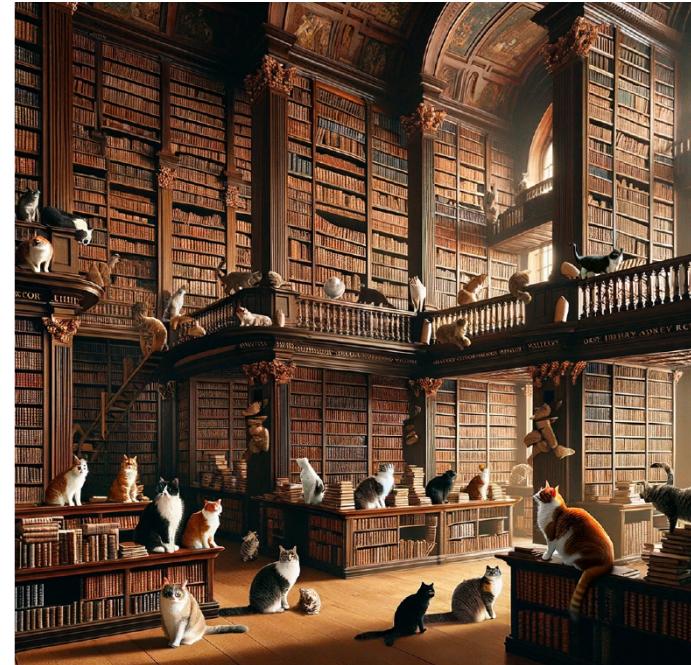
```
1 python3 dna-gpu.py
```

Understanding Big Data

Big Data

The term **Big Data** refers to datasets and data science tasks that become too large and complex for traditional techniques.

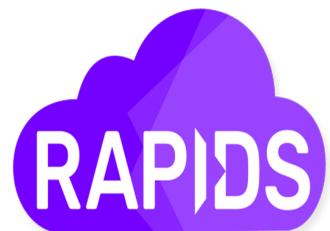
- For comprehensive information, check our workshop on [Big Data on HPC](#)





Big Data Tools

Explore various frameworks, APIs, and libraries for handling Big Data





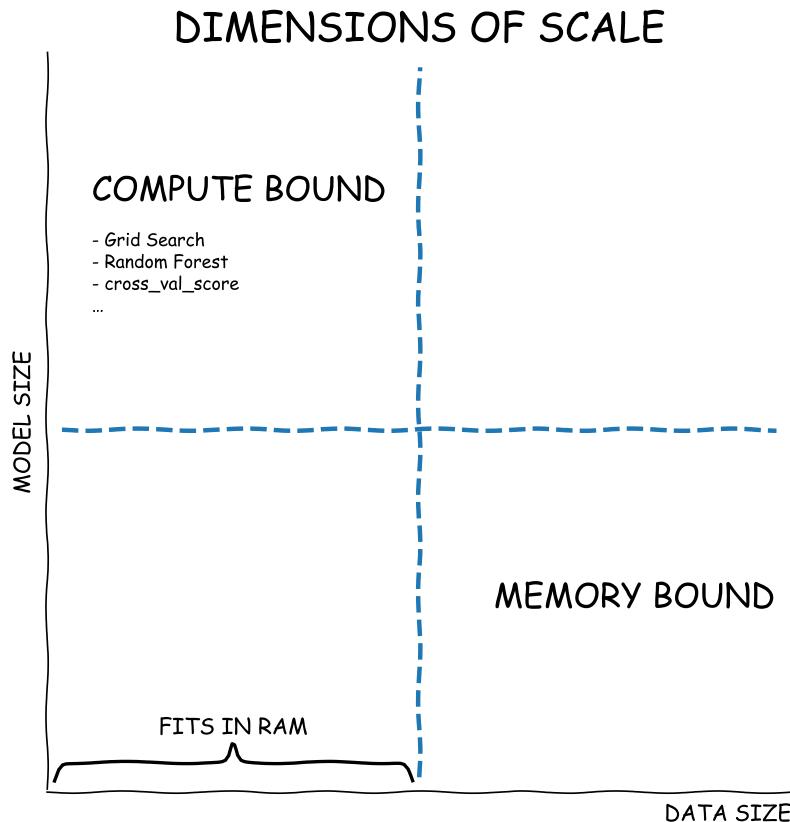
Challenges with LOTS of Data

Dealing with extensive DATA presents unique challenges 😰:

- 🧠 **Insufficient RAM:** Struggling to accommodate large datasets.
- ⏳ **Time-Consuming Processing:**
 - Difficulty in managing large datasets with traditional techniques.
 - Prolonged computation times.
- 🤖 **Complex Machine Learning Models:**
 - Training advanced models requires significant computational power for accuracy.
- 🤖 **Solution: High-Performance Computing (HPC)**
 - HPC resources supercharge solving Big Data challenges with superior computing power 💪
 - Many Big Data tools are designed to run efficiently across multiple compute nodes in HPC systems.



Big Data Challenges



- **Scaling Data Size**
 - Datasets can become so large that they can't fit into RAM
- **Scaling Model/Task Size**
 - Machine Learning or other tasks become so complex that a single CPU core is not adequate

Example: Million Song Dataset

Million Song Example

Using Spark's MLlib for Music Data Analysis:

- This example utilizes Spark's Machine Learning library (MLlib).
- We will analyze data from the [Million Song Subset](#).

Dataset Characteristics:

- 🎵 The subset contains approximately 500,000 songs.
- 📊 Features include:
 - Year of the song.
 - 90 features related to the timbre average and covariance.



Installing Spark and PySpark

Creating and Activating the Conda Environment:

- Load Anaconda and create a new environment named `mypyspark`:
- Installing Spark

```
1 module load anaconda3
2 conda create -n mypyspark openjdk pyspark python \
3                     pyspark=3.3.0 py4j jupyterlab findspark \
4                     h5py pytables pandas matplotlib \
5                     -c conda-forge -c anaconda -y
6 conda activate mypyspark
7 pip install ipykernel
8 ipython kernel install --user --name=mypyspark
```

Environment Features:

- 📚 This Conda environment, `mypyspark`, is configured with Jupyter.
- 🚀 It includes both Spark and PySpark, ready for big data processing tasks.

PySpark: Basic Operations



Let's practice basic PySpark functions with examples.

- Download the workshop content from the GitHub repository
- We'll work with a Jupyter Notebook: `Spark_basics.ipynb`
- Jupyter Notebook: [`MSD.ipynb`](#) from [`MSD_ex`](#)

```
1 cd $SCRATCH
2 git clone https://github.com/ucla-oarc-hpc/WS_MLOnHPC
3 cd WS_MLOnHPC/MSD-ex
```

Downloading the Dataset:

- Retrieve the dataset to your workspace:

```
1 cd $SCRATCH/WS_MLOnHPC/MSD-ex
2 wget https://archive.ics.uci.edu/ml/machine-learning-databases/00203/YearPredictionMSD.txt.zip
3 unzip YearPredictionMSD.txt.zip
```

PySpark: Basic operations: Starting the notebook

We will use the [h2jupynb](#) script to start Jupyter on Hoffman2

You will run this on your LOCAL computer.

```
1 wget https://raw.githubusercontent.com/rdauria/jupyter-notebook/main/h2jupynb
2 chmod +x h2jupynb
3
4 #Replace 'joebruin' with you user name for Hoffman2
5 #You may need to enter your Hoffman2 password twice
6
7 python3 ./h2jupynb -u joebruin -t 5 -m 10 -e 2 -s 1 -a intel-gold\\* \
8           -x yes -d /SCRATCH/PATH/WS_MLonHPC/MSD-ex
```



Note

The `-d` option in the `python3 ./h2jupynb` will need to have the `$SCRATCH/WS_MLonHPC` full PATH directory

This will start a Jupyter session on Hoffman2 with ONE entire intel-gold compute node (36 cores)

More information on the [h2jupynb](#) can be found on the [Hoffman2 website](#)

AutoML with H2O.ai

Introduction to AutoML

AutoML, or Automated Machine Learning, is an innovative approach to automating the process of applying machine learning to real-world problems.

Key Benefits

- **Efficiency:** Streamlines the model development process.
- **Accessibility:** Makes ML more accessible to non-experts.
- **Optimization:** Automatically selects the best models and parameters.

Components of AutoML

- **Data Preprocessing:** Automatic handling of missing values, encoding, and normalization.
- **Feature Engineering:** Automated feature selection and creation.
- **Model Selection:** Choosing the best model from a range of algorithms.
- **Hyperparameter Tuning:** Optimizing parameters for peak performance.
- **Model Validation:** Ensuring robustness through cross-validation.

AutoML Tools

HPC resources can be used to enhance AutoML since it can be very computationally demanding

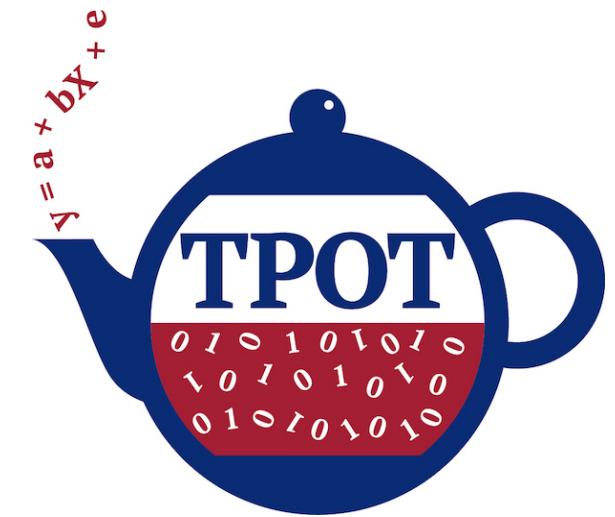
H2O.ai AutoML: An open-source platform that automates the process of training and tuning a large selection of candidate models within H2O, a popular machine learning framework.



Auto-sklearn: An automated machine learning toolkit based on the scikit-learn library, focusing on automating the machine learning pipeline, including preprocessing, feature selection, and model selection.



TPOT (Tree-based Pipeline Optimization Tool): An open-source Python tool that uses genetic algorithms to optimize machine learning pipelines.



MLBox: A powerful Automated Machine Learning python library that provides robust preprocessing, feature selection, and model tuning capabilities.



MLBox,
Machine Learning Box

Auto-Keras: Auto-Keras is a AutoML program built on the Keras platform.





Using H2O.ai for AutoML

Setting Up H2O.ai for Automated Machine Learning:

- Start by loading Anaconda and creating a new environment named `h2oai`:
- Activate the newly created environment:
- Install essential packages including H2O:
- Install IPython kernel and configure it for the `h2oai` environment:

```
1 module load anaconda3  
2 conda create -n h2oai python matplotlib -c conda
```

```
1 conda activate h2oai
```

```
1 pip install requests tabulate future h2o
```

```
1 pip install ipykernel  
2 ipython kernel install --user --name=h2oai
```



Your environment is now set up with H2O.ai, ready for AutoML tasks.



H2O AutoML Example

Exploring AutoML with H2O:

- We will work through an AutoML example from [H2o-tutorials](#).
- The focus is on the [Combined Cycle Power Plant dataset](#).
- **★ Objective:** Predict the energy output of a Power Plant using temperature, pressure, humidity, and exhaust vacuum values.
- This example, we will use the Python API, but H2O.ai has a R API as well



Accessing the Notebook:

- The Jupyter notebook for this example is in the `automl-ex` directory.
- To start Jupyter, execute the following command, adjusting the path as necessary:

```
1 python3 ./h2jupynb -u joebruin -t 5 -m 50 -e 2 -s 1 -a intel-gold\\* \
2                               -x yes -d /SCRATCH/PATH/WS_MLonHPC/automl-ex
```

Wrap-up



Workshop Highlights

- High-Performance Computing (HPC) and Machine Learning:
 -  Introduction to HPC and its benefits for Machine Learning.
 -  Utilizing Python and R on HPC for advanced data processing.
- Key Tools and Frameworks:
 -  Installation and usage of vital Python packages like Scikit-learn, PyTorch.
 -  R package installation and management in HPC environment.
 -  Setting up Anaconda environments for Python and machine learning libraries.
- Big Data and Its Challenges:
 -  Understanding Big Data, its challenges, and tools to handle large datasets.
 -  Introduction to various Big Data frameworks and libraries.
- Conclusion:
 - This workshop offered a comprehensive overview of leveraging HPC

👏 Thanks for Joining! ❤️

Questions? Comments?

- cpeterson@oarc.ucla.edu
- Look at for more [Hoffman2 workshops](#)



