

Finding, downloading, and applying software on CURC

August 14, 2024

Brandon Reyes

Slides & Exercises

https://github.com/ResearchComputing/hpc_fundamentals_micro_credential

- In “installing_software” directory



Learning Objectives

- Learn about different methods to install and use software on Alpine

Session Overview

- The Module System (Lmod)
- Building Software on CURC Systems
- Simplifying Source installations with Spack
- Virtual Environments with Anaconda
- Containerization with Apptainer
- Requesting Software Installations



The Module System

In most cases, a supercomputer has far more software installed than the average user will ever use.

- Users may need different versions of the same software, which in general cannot be installed nor used in parallel on the same system.
- The requirements for one package may adversely affect another package or even be mutually exclusive.

The Module System



HPC centers manage this complexity with **environment module systems**.

CURC uses the Lmod system.



The Module System

Setting up for today's session.

1. Log in to CURC HPC system
2. get on an Alpine compute node

```
$ module avail  
$ acompile --help  
$ acompile --time=2:00:00  
$ module avail
```

Login nodes do not have
the full software stack!

The Module System



Live Demo: Loading and unloading modules will dynamically change the software environment on the cluster.

```
$ module purge  
$ module load intel/2022.1.2  
$ module avail  
$ module load impi  
$ module avail
```

```
$ module purge  
$ module load gcc  
$ module avail  
$ module load openmpi  
$ module avail
```


The Module System



Live Demo: Loading and unloading modules will set (and reset) important environment variables for you.

```
$ module purge  
$ module load intel  
$ module load hdf5  
$ module display hdf5  
$ env | grep HDF5
```

```
$ module purge  
$ module load gcc  
$ module load hdf5  
$ module display hdf5  
$ env | grep HDF5
```



The Module System

Live Demo: Useful Lmod commands

<code>module spider</code>	<code># list all available modules</code>
<code>module avail</code>	<code># list modules available to you</code>
<code>module load <package/version></code>	<code># load a module into your env</code>
<code>module purge</code>	<code># unload all modules</code>
<code>module list</code>	<code># list currently loaded modules</code>
<code>module display <package></code>	<code># display module info/help</code>
<code>module spider <package></code>	<code># view info for all version</code>
<code>module spider <package/version></code>	<code># view info for specific version</code>



The Module System



Points to note about CURC-managed modules:

- CURC does not update system modules; we do fresh installs of new versions and change the default when that is appropriate
- Sometimes when a module is outdated or problematic we will remove it from the software stack

Take home: pay attention to what modules you are loading, as this may be important for reproducibility!

Building Software on CURC Systems

- Definitions
 - **Building**- a generic term describing the overall installation process that includes compiling
 - **Compiling**- the process of converting source code to an executable
 - **Linking**- the process of combining pieces of code and data into a single file that can be loaded into memory and executed
 - **Installing**- any process that results in executables

Building Software on CURC Systems

- There are numerous ways to install software on CURC systems
 - Grab pre-compiled binaries
 - Within virtual environments (using Conda, Miniconda, or Mamba)
 - Using containers (Apptainer)
 - From source
 - Using a package manager for HPC systems (Spack)

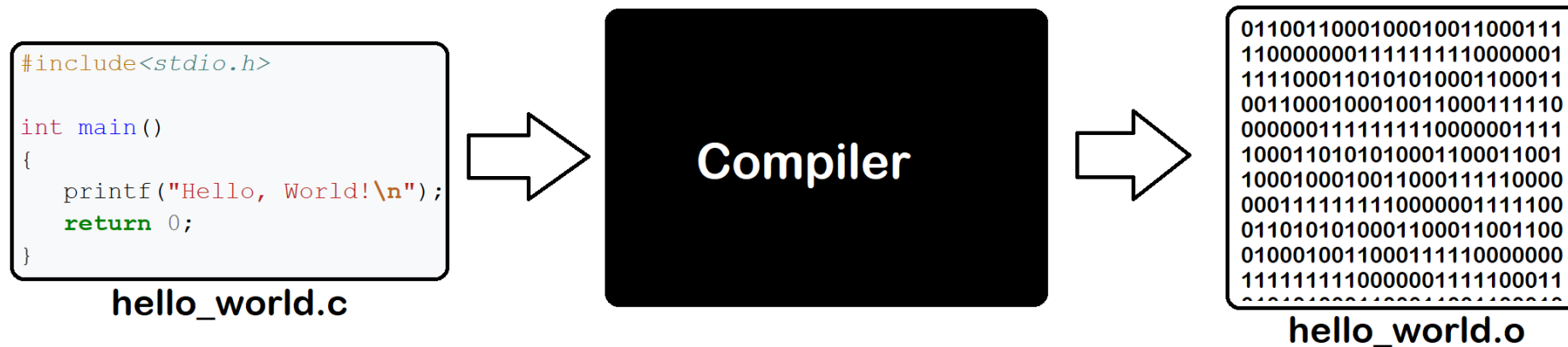
Building Software on CURC Systems

Why compile a research application manually from source?

1. It is not distributed as a pre-compiled binary, by any package manager, and is not easily containerized.
2. Compiling from source on the cluster will greatly improve performance.

Building Software on CURC Systems

Compilers are programs that convert code written in high level programming languages (like C/C++ or Fortran) to executable binary files.



Building Software on CURC Systems

Build systems automate the process of compiling and linking.

1. GNU Build System

- your application includes instructions to run `./bootstrap`, `./autogen.sh`, `./configure` or `make` (the latter without a preceding `cmake`)
- `make` is available in `/usr/bin`; Autotools available as a module

```
$ ./configure --prefix=/projects/$USER/software/bin  
$ make  
$ make install
```


Building Software on CURC Systems

Build systems automate the process of compiling and linking.

2. Cmake

- your application includes a cmake step
- module avail cmake

```
$ cmake .. -DCMAKE_INSTALL_PREFIX=$INSTALLDIR \ --  
          DCMCMAKE_CXX_COMPILER=g++ -DREGRESSIONTEST_DOWNLOAD=ON  
$ make -j 8  
$ make install
```

Building Software on CURC Systems

Conventions and best practices

- You will need to adapt installations for **local** or **user installations** (look for these terms in the software's docs)
- Don't install software in /home/\$USER (too small) or scratch (purged every 90 days); /projects/\$USER/software is the way to go!
- Keep your software installations organized by using a consistent file structure and naming convention
- Load the compiler first, MPI implementation second, and third-party libraries last

Building Software on CURC Systems

Conventions and best practices

- Don't install executables to the source directory
 - `cmake -DCMAKE_INSTALL_PREFIX, ./configure --prefix`
- The newest version of a compiler might not be compatible with your application. Read the package documentation and don't be afraid to try different compilers and compiler versions
- Read our 'Compiling and Linking' documentation
<https://curc.readthedocs.io/en/latest/compute/compiling.html>

Building Software on CURC Systems

Conventions and best practices

- Make life easier for yourself by adding executables to PATH and any directories with libraries that your application links to LD_LIBRARY_PATH

```
$ export PATH=/projects/$USER/software/phyloflash/bin:$PATH
$ echo $PATH
$ export
    LD_LIBRARY_PATH=/curc/sw/hdf5/1.10.1/mpi/17.3/intel/17.4/lib
:$LD_LIBRARY_PATH
$ echo $LD_LIBRARY_PATH
```

Want to go the extra mile?

Try our Hands-on exercise #1 provided in EXERCISES.md.
This is not required for the micro-credential.

Objectives:

- 1) Explore CURC compilers and compiler environment variables.
- 2) Perform a simple source installation.

Estimated time to complete: 15 minutes

Simplifying Installations with **Spack**

How can we simplify source installations?

- **Package Managers** – Tools that automate installing, maintaining, and configuring software and any dependencies
- **Environments** – A collection of resources that are available in a self-contained 'bubble'

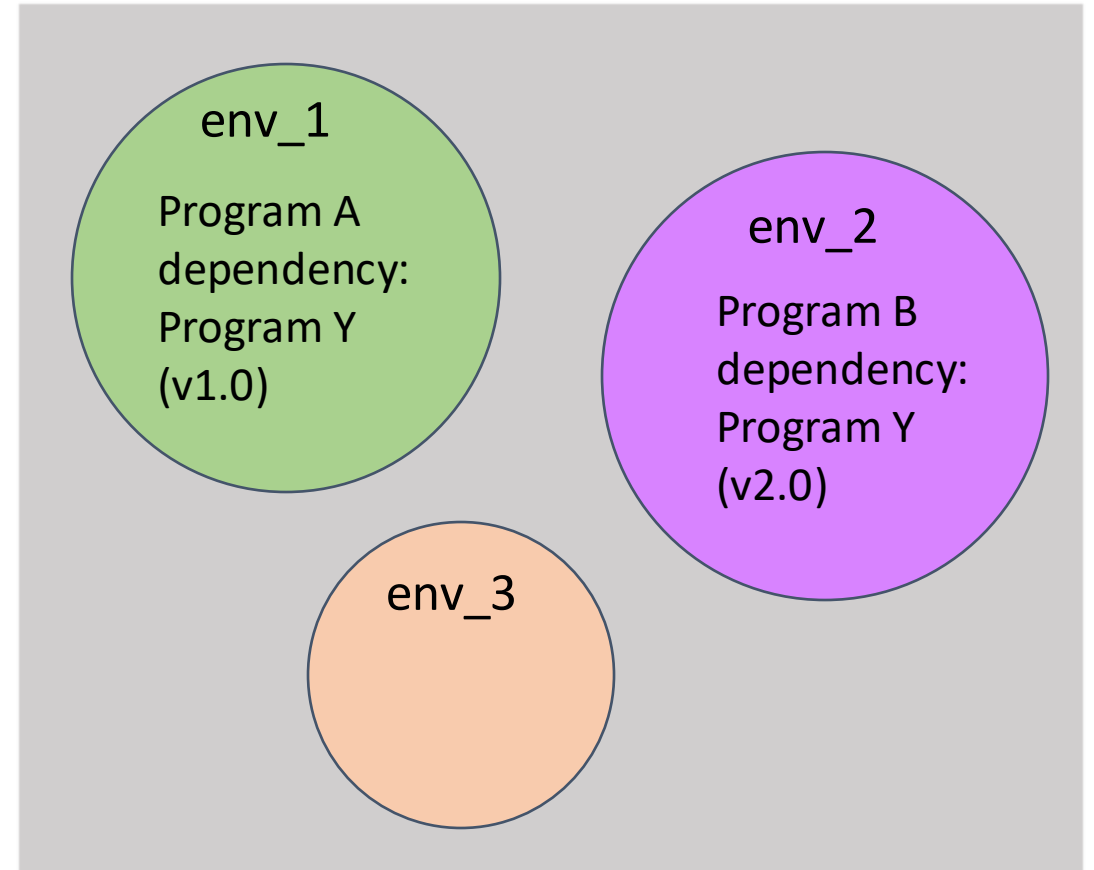
Simplifying Installations with Spack

Think of virtual environments as self-contained bubbles.

env_1 contains all the dependencies of 'Program A'.

env_2 contains all the dependencies of 'Program B'.

The environments do not interact.



Simplifying Installations with Spack

Your workflow requires two programs, 'Program A' and 'Program B'.

- 'Program A' depends on 'Program Y' **v1.0**
- 'Program B' depends on 'Program Y' **v2.0**

What do you do?!

Simplifying Installations with Spack

Environments are created and programs are installed in a few simple steps

```
$ module purge  
$ module load spack/0.20.1  
$ spack env create my_first_env  
$ spack activate my_first_env  
$ spack install --add samtools
```

Don't install packages outside of an environment!

Simplifying Installations with Spack

- Packages are installed within **activated** environments

using spack install

```
$ spack install --add samtools          #install default samtools  
$ spack install --add samtools@1.9     #install specific version
```

Simplifying Installations with Spack

- Spack installations can be slow but will progress more quickly with more cores.
 - Spack builds all packages in parallel. The default parallelism is equal to the number of cores available to the process, up to 16.

Simplifying Installations with Spack

- Useful spack commands

<code>spack env list</code>	<code># list all your environments</code>
<code>spack remove <env></code>	<code># remove an environment</code>
<code>spack uninstall <packagename></code>	<code># remove package</code>
<code>spack env status</code>	<code># check which env you're in</code>
<code>spack info <packagename></code>	<code># prints detailed package info</code>
<code>spack find</code>	<code># show installed packages</code>
<code>despacktivate</code>	<code># deactivate environment</code>
<code>spack spec <packagename></code>	<code># list packages plan</code>

Simplifying Installations with Spack

- Useful spack file paths

```
# root of the spack install tree  
/projects/$USER/software/spack
```

```
# location of package executables - these are symbolically linked to  
the installation tree subdirectory  
/projects/$USER/spack/environments/<env>/ .spack-env/view/bin
```

```
# location of spack config file  
/home/$USER/.spack/config.yaml
```

Want to go the extra mile?

Try our Hands-on exercise #2 provided in EXERCISES.md. This is not required for the micro-credential.

Objectives:

- 1) Create a Spack environment
- 2) Install fastqc in your Spack environment

Estimated time to complete: 20 minutes

Documentation:

<https://curc.readthedocs.io/en/latest/software/spack.html>

Virtual Environments With

Conda is a package (software) management system

- Installs, runs, and updates packages and their dependencies
- Creates, saves, loads, and switches between virtual environments
- Created for Python programs, but can package and distribute software for any language

Virtual Environments With

For our system, please be sure to add a `.condarc` file in your home directory.

- Prevents package installs from going to your home directory

```
$ cd ~  
$ nano .condarc
```

Add the following content:

```
pkgs_dirs:  
- /projects/$USER/.conda_pkgs  
envs_dirs:  
- /projects/$USER/software/anaconda/envs
```



Virtual Environments With

Environments are created and programs are installed in a few simple steps

```
$ module load anaconda  
$ conda create -n my_first_env python==3.10  
$ conda activate my_first_env  
$ python
```

Don't install packages in your base environment!

Virtual Environments With

Packages are installed within **activated** environments

using conda install (preferred method, when available)

```
$ conda install pandas                #install latest pandas  
$ conda install pandas==0.20.3      #install specific version of pandas
```

Virtual Environments With

Packages are installed within **activated** environments

using pip install (if you must)

```
$ pip install --no-cache-dir pandas          #install latest pandas
```

--no-cache-dir is crucial on CURC systems!

Virtual Environments With

Useful conda commands

```
conda env list           # list all environments
conda list               # list packages in active env
conda env remove -n <envname> # remove an environment
conda config --show channels # view configured channels
conda deactivate         # deactivate environment
conda create --name <clonedenv> / # clone an environment
    --clone <envtoclone>
```

Virtual Environments With

Useful conda file paths

```
# location of python libraries  
/projects/$USER/software/<env>/lib/python3.10/site-packages  
  
# location of package executables  
/projects/$USER/software/<env>/bin  
  
# location of .condarc file  
/home/$USER/.condarc
```

Want to go the extra mile?

Try our Hands-on exercise #3 provided in EXERCISES.md. This is not required for the micro-credential.

Objectives:

- 1) Configure your .condarc file
- 2) Create a conda environment and install samtools
- 3) Activate the environment and run samtools.

Estimated time to complete: 15 minutes

Documentation:

<https://curc.readthedocs.io/en/latest/software/python.html>



Containerization With **APPTAINER**

Containers are portable virtualizations of an operating system, software, libraries, data, and/or workflows

- pros
 - portability- containers can run on any system equipped with its specified container manager
 - reproducibility- because containers are instances of prebuilt isolated software, the software will always execute the same every time
- cons
 - steeper learning curve than conda
 - can be difficult to troubleshoot issues
 - building containers can be tricky for multi-node MPI applications

Containerization With **APPTAINER**

CURC offers Apptainer (formerly Singularity) as container management software

- Apptainer comes pre-installed on all Alpine nodes, so no need to load any specific software

Many common research applications have already been containerized and can be pulled from container repositories (such as Docker Hub).

Containerization With APPTAINER

Useful Apptainer commands

<code>apptainer exec</code>	<code>#Execute a command to your container</code>
<code>apptainer run</code>	<code>#Run your image as an executable</code>
<code>apptainer build</code>	<code>#Build a container</code>
<code>apptainer pull</code>	<code>#pull an image from hub</code>
<code>apptainer inspect</code>	<code>#See labels/environment vars, run scripts</code>
<code>apptainer shell</code>	<code>#Shell into your image</code>

Containerization With APPTAINER

A container has its own file system and so needs help “seeing” files outside the container (on the host system). If not done in the .def file, this can be accomplished at runtime with bind mounting.

```
# bind mount a directory  
apptainer run -B /source/directory:/target/directory sample-image.sif
```

On CURC systems, a running container automatically bind mounts these paths: /home/\$USER, \$PWD. *Note that other locations will need to be manually mounted.*

Want to go the extra mile?

- Try our Hands-on exercise #4 provided in EXERCISES.md. This is not required for the micro-credential.

Objectives:

- 1) Become familiar with basic Apptainer commands.
- 2) Pull an image from a pre-built container, then run the program from the container.

Estimated time to complete: 20 minutes

Documentation:

<https://curc.readthedocs.io/en/latest/software/Containerizationon.html>

Requesting Software Installations

- Is the software already installed on the cluster?
<https://curc.readthedocs.io/en/latest/clusters/alpine/software.html>
- Have you considered its utility and complexity?
 - Are you the only user of this software?
 - How complex or difficult is this software to install?
- Have you tried installing the package on your own?
- Software request form:
<https://www.colorado.edu/rc/userservices/software-request>

Thank you!