# DUNE Computing Tutorials

David DeMuth, Valley City State University
CHEP 2023, Norfolk, Virginia, USA
For the DUNE Collaboration
May 8, 2023
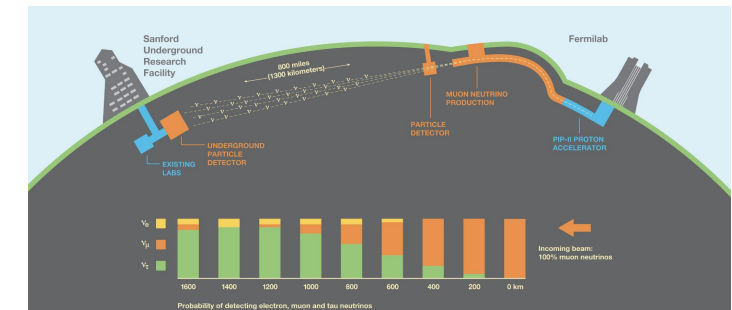
https://bit.ly/chep23-dune-training

# DUNE Computing

- The physics realms of the Deep Underground Neutrino Experiment presents a unique challenge for data analysis and data processing software frameworks:
    - Precision neutrino oscillation measurements,
    - Searches for nucleon decay,
    - Sensitivity to nearby Supernova explosions, and more…
- Particle interactions will span timescales of nanoseconds to 100's of seconds, where event energies and associated processing timelines require a **wide arsenal of software tools**.
- Moreover, **data storage is cloud-based and worldwide**, and require failsafe high bandwidth networks with strong security and monitoring protocols.
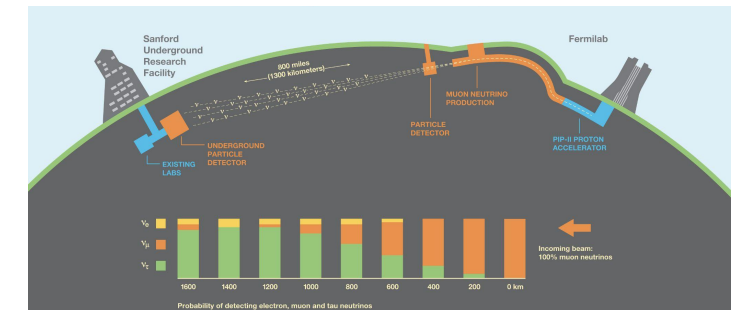


Worldwide and sophisticated, DUNE Computing Tutorials are used regularly to train new colleagues.

DUNE Computing CDR: https://arxiv.org/abs/2210.15665

**VALLEY CITY STATE UNIVERSITY**

# DUNE Computing

- With the first of four detector designs prototyped, blasting of the deep underground caverns underway at Lead, SD, detector installation anticipated for 2027, and science starting in 2029, computing system development has been vigorous.
- ArgoNeuT, and MicroBooNE are LArTPc experiments that have produced software that DUNE Computing has adopted to ensure efficient liquid argon detector designs, those **modelling techniques are highly numerical**, and sophisticated to use.
- Graduate students, post-docs, young and senior researchers are among the hundreds of collaborators added each year, with most **requiring general and specific software training**.
- While selected training can be outsourced to an evolving curricula provided by the HEP Software Foundation (HSF),  some training remains specific to DUNE.



Worldwide and sophisticated, DUNE Computing Tutorials are used regularly to train new colleagues.

**VALLEY CITY STATE UNIVERSITY**

DUNE

# DUNE Trainings

- DUNE Computing first sponsored training in 2016, and since coordinated 10 events in various formats, with ~400 participants.
- Tutorials have focused on three topics:
  - Data storage and management,
  - LArSoft (software for LArTPCs),
  - Job submission and monitoring.
- The goal is to certify that new colleagues have access to DUNE computing resources and understand the basics of logging in, storage areas, running applications, code modifications, submitting and monitoring batch jobs.

| Sessions | Wednesday, May 12 | Thursday, May 13 | Friday, May 14 |
|---|---|---|---|
| 8:00 - 8:15 | **Welcome + announcements** C. David & D. DeMuth | **Grid job submission + common errors** Lecture + hands-on + exercises *Follow-up: see "Expert in the room" Friday late morning* K. Herner | "Expert in the room" **LArSoft: How to modify a module** T. Junk |
| 8:15 - 9:00 | **Storage spaces** Lecture + hands-on M. Kirby | | |
| 9:00 - 10:00 | **Data management** Lecture + hands-on S. Timm | | **Code-makeover** Switch to POMS K. Herner |
| 10:00 - 10:30 | Coffee break! | Coffee break! | Coffee break! |
| 10:30 - 11:00 | **QUIZ!** Storage spaces data management | **QUIZ!** Grid job submission | **QUIZ!** Best programming practices |
| 11:00 - 12:15 | **Intro to art/LArSoft** ← lecture **Exploring fcl files** ← hands-on *Follow-up: see Friday morning* T. Junk | **Code-makeover** How to improve your code for better efficiency T. Junk | "Expert in the room" **Grid & batch job submission** K. Herner |
| 12:15 - 12:30 | | | **Closing remarks** C. David & D. DeMuth |

**Organizers**

Claire David
York University / FNAL

David DeMuth
Valley City State University

**DUNE Computing Consortium Lead**

Heidi Schellman
University of Oregon

**Lecturers**

Mike Kirby
FNAL

Steven Timm
FNAL

Tom Junk
FNAL

Kenneth Herner
FNAL

**Mentors**

Amit Bashyal
ANL

Carlos Sarasty
U. of Cincinnati

The May 2021 training was offered as a three day event, each of four lecturers and two mentors doing the bulk of the work.

VALLEY CITY STATE UNIVERSITY

DUNE

# Training Logistics

- Event registration and communications are managed using Indico.
- Participants must verify their ability to use the Unix Shell.
- **Pre-event homework is vital** as it includes verification that participants can access the DUNE general purpose virtual machines, and ideally the corresponding CERN VM's (LX+).
- Livecoding, quizzes, expert in the room sessions, and assigned mentors ensure the event as hands-on.
- Each session is delivered and captured via Zoom, then embedded into the Software Carpentries lesson framework (SWC) which is hosted at DUNE Computing's GitHub site for review.

## DUNE Computing Training May 2021 edition: Mission Setup

### Objectives

- Get ready to do the tutorial
- Understand the authentication procedures
- Set up your environment for DUNE
- Do an exercise to help us check if all is good
- Get streaming and grid access

### Requirements

You must be on the DUNE Collaboration member list and have a valid FNAL or CERN account. See the Indico Requirement page for more information.

> ✏️ **Note**
>
> The instructions below are for FNAL accounts. If you do not have a valid FNAL account but a CERN one, go at the bottom of this page to the "Setup on CERN machines".

### 1. Kerberos business

If you already are a kerberos-aficionado, go to the next section. If not, we give you a little tour of it below.

**What is it?** Kerberos is a computer-network authentication protocol that works on the basis of tickets.

**Why does FNAL use Kerberos?** Fermilab uses Kerberos to implement strong authentication, so that no passwords go over the internet (if a hacker steals a ticket, it is only valid for a day).

**How it works?** Kerberos uses tickets to authenticate users. Tickets are made by the kinit command, which asks for your kerberos password (info on kerberos password here). The kinit command reads the password, encrypts it and sends it to the Key Distribution Centre (KDC) at FNAL. The Kerberos configuration file, which lists the KDCs, is stored in a file named krb5.conf. On Linux and Mac, it is located here:

```
Code
/etc/krb5.conf
```

If you do not have it, create it. A FNAL template is available here for each OS (Linux, Mac, Windows). More explanations on this config file are available here if you're curious.

To log in to a machine, you need to have a valid kerberos ticket. You don't need to do this every time you login, only when your ticket is expired. Kerberos tickets last for 26 hours. To create your ticket:

```
Bash
kinit -f username@FNAL.GOV
```
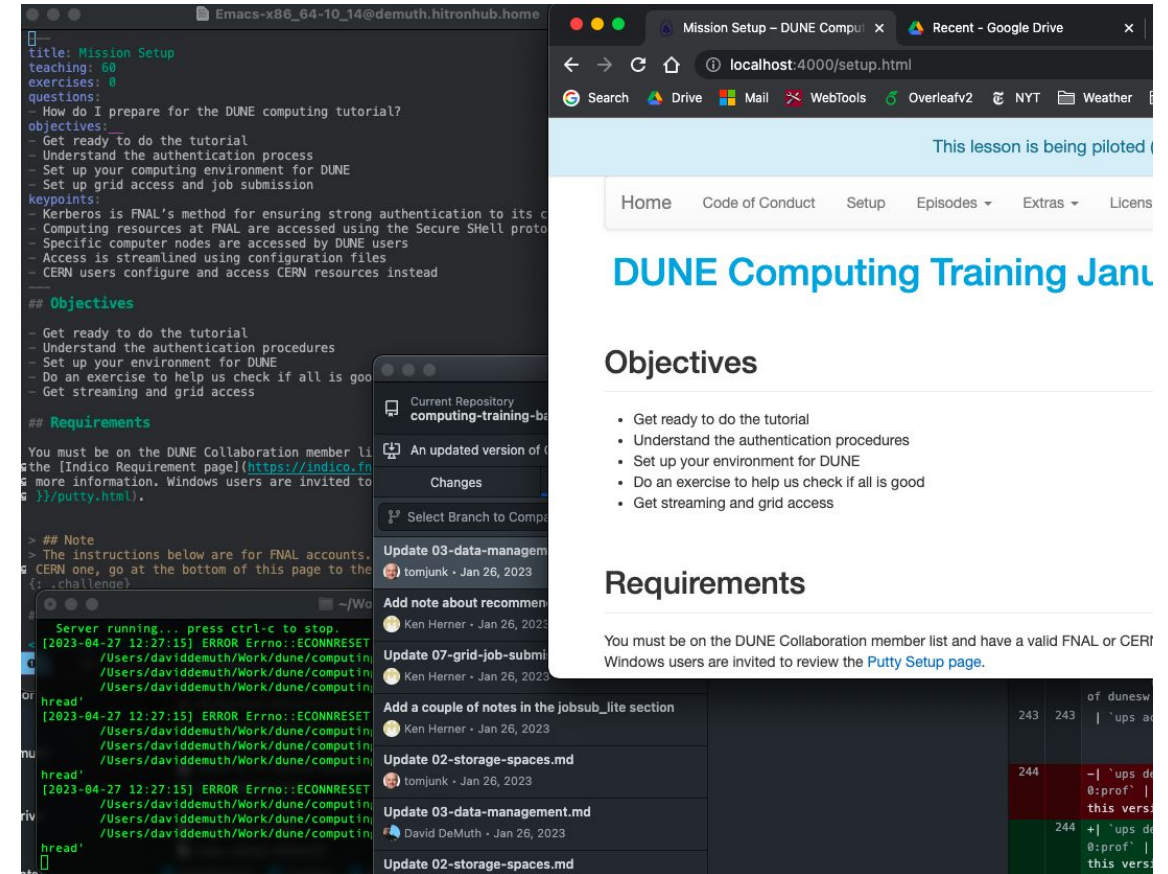
In advance of the opening salvo, students must demonstrate an understanding of using the Unix shell to access secure VMs.

**VALLEY CITY STATE UNIVERSITY**  DUNE

# Lesson Development

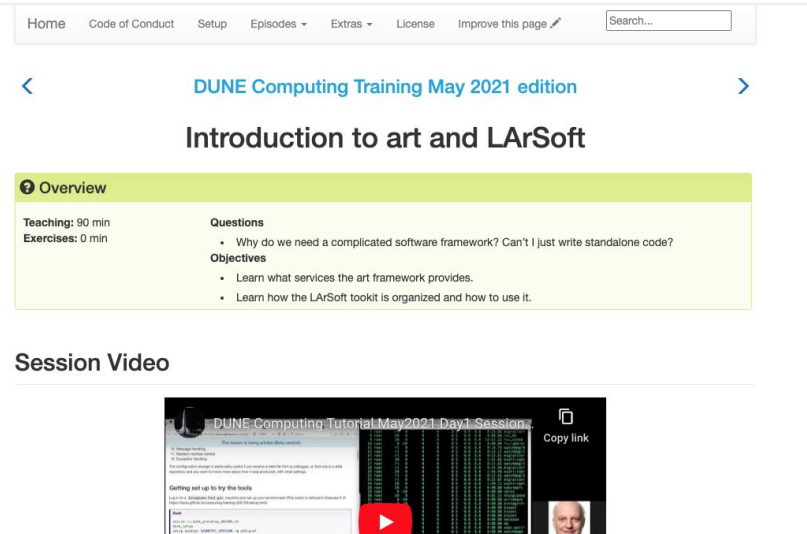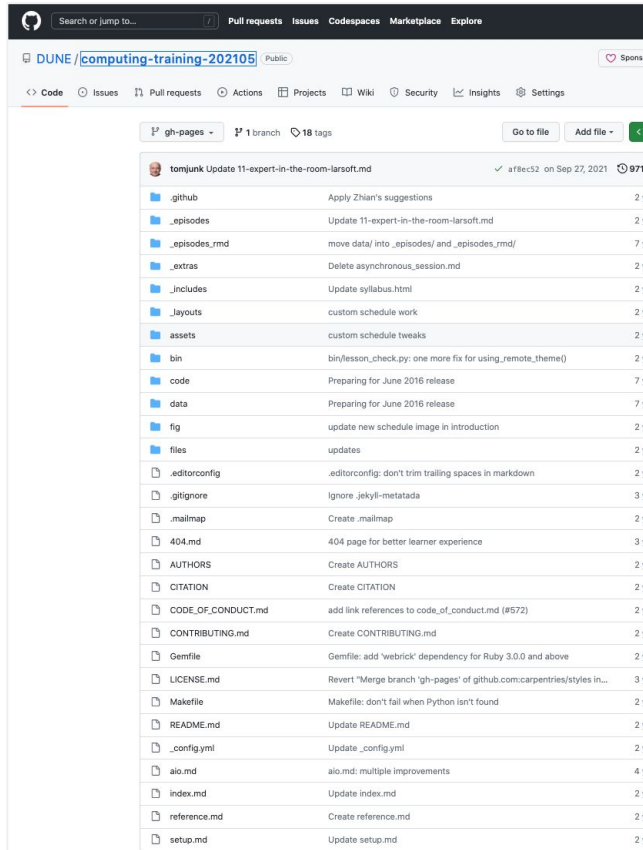The infrastructure to develop lessons is provided in the Software Carpentries framework:

- A  lesson template is imported as a new DUNE GitHub repo, configuration via a _config.yml file, and main lesson content as markdown files (.md) located in _episodes/.
- GitHub Desktop is used to manage the repository locally.
- Raw editing of files on github.com is frowned on, editing .md files locally is encouraged.
- Viewing edits in a localhost browser uses a Ruby/Jekyll rendering engine.
- As edits are verified, lessons are pushed to the site's main branch for access by multiple authors, and the public.
- Lessons are rendered elegantly on the web via GitHub Pages as a free service.



Once installed on a curriculum designers local machine, the lesson production environment is slick.

VALLEY CITY
STATE UNIVERSITY

DUNE

# Lesson Deployment

**End users see:**



Markdown language features:

- Simplified editing structure for web browsers,
- Bash script blocks easily copied and pasted in terminal windows,
- Drop down Quiz blocks that can be opened on demand.

GitHub's gp-pages are rendered seamlessly, and for free.

VALLEY CITY STATE UNIVERSITY

DUNE

# Getting Help

As expert instructors work through a lesson plan, often live coding in an adjacent window, multiple methods are used to ensure near- and long-term support:

- Instructors encourage participants to pose questions via a shared and world editable Google Doc (**Livedoc**) which is monitored by mentors and other experts.
  - Each participant selects a unique color to improve clarity.
  - As questions are volleyed, experts will form a dialogue with individuals to reason solutions; often a second expert will contribute.
  - Livedocs can be studied by others, solutions can be contrasted, and the Livedoc becomes a permanent educational resource.
- **Slack channels** are set up for each training event, where questions can be posed asynchronously, when questions receive attention from a cadre of experts.
- **GitHub Issues** has proven to be another method to field questions, provide mentorship, and ensure the novice become proficient in DUNE computing.

Livedocs, Slack, and GitHub Issues are used to assist DUNE colleagues with computing questions.

VALLEY CITY STATE UNIVERSITY

# Student Snapshot

## Where are you career wise?
DUNE Computing Tutorials May 2022 (N=35)

- graduate student — 51.4%
- postdoc — 28.6%
- senior scientist — 8.6%
- computer specialist
- undergraduate student — 8.6%

## How familiar are you with running large batch jobs
35 responses

- Expert
- Sorta — 34.3%
- Maybe — 28.6%
- No — 31.4%

**I am familiar with**
35 responses

- UNIX shells like bash — 28 (80%)
- GitHub — 24 (68.6%)
- Python — 28 (80%)
- C++ — 31 (88.6%)
- Databases (including SQL) — 7 (20%)
- Jupyter notebooks — 19 (54.3%)
- ROOT — 29 (82.9%)
- GEANT simulations — 9 (25.7%)
- neutrino event generators — 8 (22.9%)
- The Art framework — 8 (22.9%)

**What physics/algorithm group are you working with, so far?**
31 responses

- ProtoDUNE analysis — 5 (16.1%)
- Far detector non-oscillation — 3 (9.7%)
- Oscillation physics — 4 (12.9%)
- Cross sections — 3 (9.7%)
- Near detector sim+reco — 4 (12.9%)
- BSM — 2 (6.5%)
- Calibrations — 1 (3.2%)
- Not started yet — 16 (51.6%)
- Coldbox analysis — 1 (3.2%)
- Good ol' blasters. — 1 (3.2%)
- Beam working group — 1 (3.2%)

**Where do you back up your code?**
34 responses

- Local backup drive — 16 (47.1%)
- Personal GitHub or GitLab — 18 (52.9%)
- DUNE GitHub or GitLab — 10 (29.4%)
- Redmine — 0 (0%)
- Dropbox/Box or other cloud ser... — 5 (14.7%)
- Backup? Guess I should be doi... — 6 (17.6%)
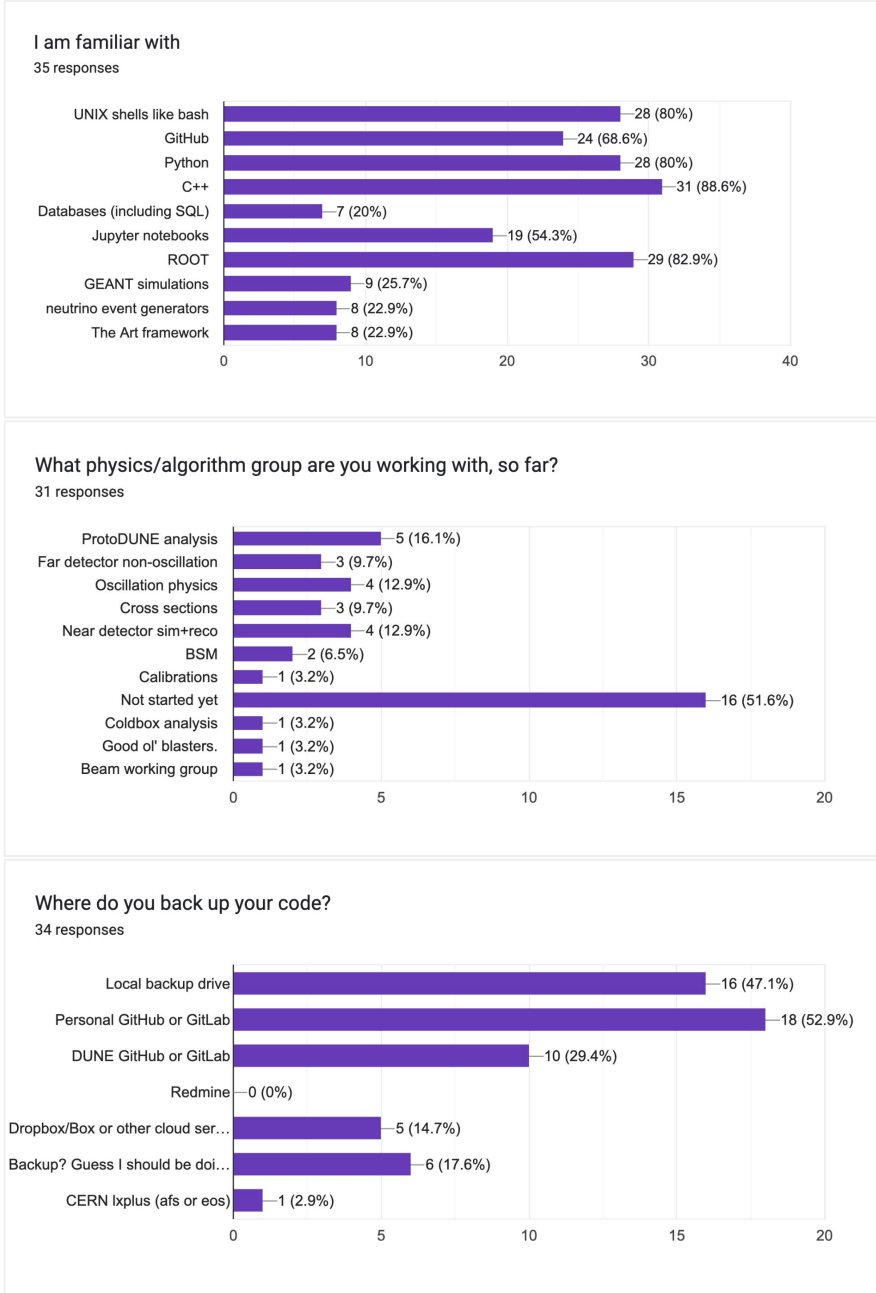- CERN lxplus (afs or eos) — 1 (2.9%)

Graduate students and postdocs make up the majority of our training audience. It is exciting to see undergraduates are participating.

Few are familiar with large batch jobs.

A strong majority know bash, use GitHub, Python, C++, and Root!

Over half have yet to start their physics analyses, and these were who we focused on in the event registration communications.

A healthy number of participants use GitHub for code backups.

**Below, one student provides feedback after participating in a training.**

The FNAL grid job submission has always been a pain for us. It seems very cumbersome to run anything on grid. I found this fact from the oscillation analysis I am working on (some inherited job submission scripts) as well as from the past dune computing tutorials. We need set up so many softwares for LArSoft, dune tpc, mrb as well as things like xrootd, IFDH, SAM input, etc. This maybe very necessary and useful for experienced users, but it's very hard for undergraduate/early graduate to learn. Part of these is due to me not wrapping my head around all the softwares so far used at DUNE (and no one is teaching us).

VALLEY CITY STATE UNIVERSITY

# Lessons Learned

- SWC lesson template is elegant, functional, and practical for delivering hands-on learning materials.
- Pre-event homework that includes checking that students can access FNAL servers is a must.
- It's easy to be ambitious with the material for one half day introductions, two day, or three day events.
- Coffee breaks during trainings are important for assimilation.
- Mentors are essential to ensure skill development and understanding.
- Hybrid synchronous delivery is ideal.
- Zoom captures on YouTube allow asynchronous access, and a record of the event.
- Other approaches to training, such as Hackathons, could be used to take lessons to a next step.

**VALLEY CITY**
**STATE UNIVERSITY**

DUNE

# HEP Software Skill Development

- The list of skills that experimental HEP physicists use on a regular basis is lengthy and evolving.
- Mastering data analysis techniques is not specific to one experiment.
- Using HSF training materials such as Unix Shell is effective and efficient.
- DUNE Computing encourages participation in HSF workshops.





This excellent instructional model by the HSF introduces concepts of continuous integration and development using GitHub motivated DUNE Computing's use of the SWC lesson templates.

VALLEY CITY STATE UNIVERSITY

# HEP Software



**Data Management**: XrootD, dCache, HDCondorCE, StashCache. CVMFS, Rucio CLI, SAM, MetaCat, Data Dispatcher

**Reconstruction**: Art, LArSoft, Pandora, Geant4, ROOT, WCT, PyTorch, UPS

**Simulations**: MARS, g4lbnf, Geant4, GDML, DUNENDGGD, GENIE, NuWRO, GIBUU, NEUT, CORSIKA, MUSUN/MUSIC, BXDECAY0, FLUKA, Geant4Reweight, dk2nu, PPFX, larnd-sim, edep-sim, CRY, garg4, ECAL, STT, GRAIN, OverlayGenie, Garfield, SPICE

**Visualization**: Bee, WebEVD, TEve

**Analysis**: CAF, CAFAna, CAFFEA, HighLAnd, Nuisance

**Documentation**: Wiki, GitHub, Read the Docs, Sphynx

Due to the vast suite of software used by the DUNE collaboration for the various analysis steps and sub-detectors, training is important, with some on an ad-hoc basis by physics and hardware groups, e.g. software systems which monitor QA/QC for installation of the far detector.

**VALLEY CITY STATE UNIVERSITY**

DUNE

# Future Work

- HEP experiments host a large number of collaborators, many have limited computer science training.
- The movement is for software to be reused by multiple experiments; common tools are anticipated.
- DUNE Computing's training work thus far has been focused on a few essential topics.
- Use of the Carpentries templates provides a proof of concept for other topics which have wider appeal.
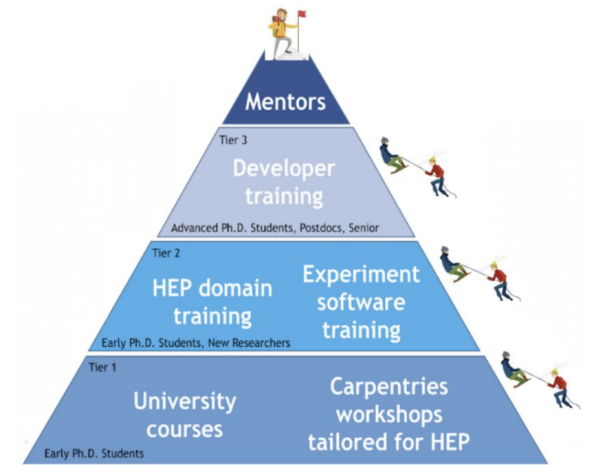


**Figure 1.** Evolution of HEP Education and Training

A rough estimate is that there are ~10K people in the HEP community to train every year, ranging from undergraduates to scientists.

**Software Training in High Energy Physics,** Michel H. Villanueva, Sudhir Malik, Meirin Oan Evans

# Summary

Students and postdocs who join HEP experiments who have no formal software training but become proficient with structured training, mentorship, and peer support.

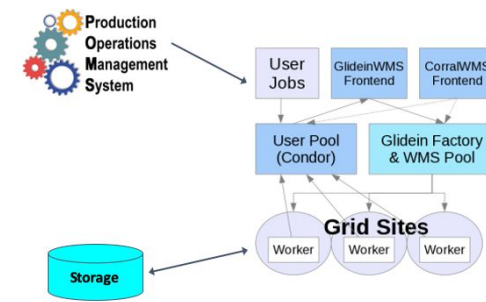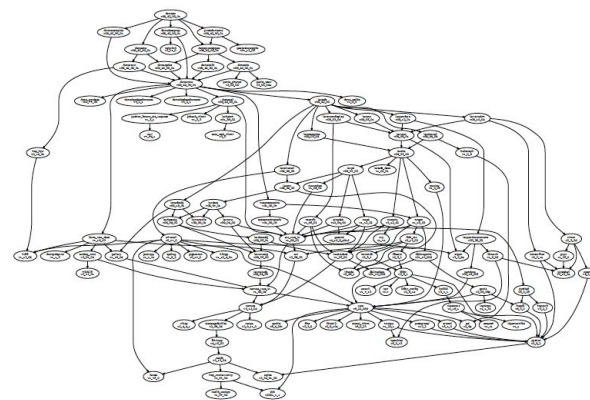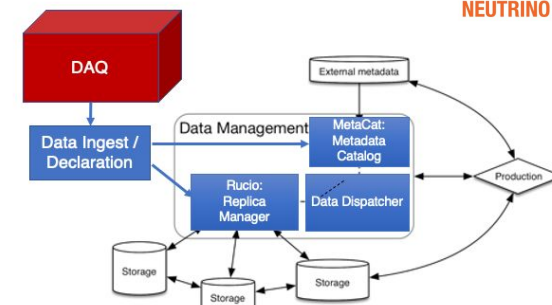**DUNE Computing has developed multi-day training events** to focus on the basics:

- Understanding data storage and best practices for operating in a cloud infrastructure,
- Learn event reconstruction and analysis techniques in liquid argon experiments,
- Master batch submission practices in a HPC grid environment.

Our **aim is to jump start individual's simulation and reconstruction work**.

Expert instructors practice clear communication, produce materials that are straightforward, are open and patient with students, offer access for questions though livedocs, and long term support via Slack channels, GitHub, and email.

Rich with multicultural heritages themselves, DUNE collaborators are conscientious when developing curriculum to ensure the innate diversity becomes a feature for learning, not a bug.

**DUNE Computing is committed to software training**, recognizes that the materials developed for DUNE has pertinence for other experiments, and is working to contribute to the larger open science community.

# CHEP 2023 Abstract

Providing computing training to the next generation of physicists is the principal driver for a biannual multi-day workshop hosted by the DUNE Computing Consortium. Materials are cast in the Software Carpentries templates, and to date topics have included storage space, data management, LArSoft, grid job submission and monitoring. Moreover, experts provide extended breakout sessions to demonstrate the intricacies of the unique software used in HEP analysis. Each workshop session uses live documents for real time correspondence, and are captured on Zoom; afterwards, videos are embedded on the corresponding webpages for review.  As a GitHub repository, shared editing of the learning modules is straightforward, and provides a trusted framework to extend to other training topics in the future. An overview of the machinery will be provided, post workshop statistics will be discussed, with lessons learned will be the focus of this presentation.

**VALLEY CITY**
**STATE UNIVERSITY**

DUNE