# DUNE Software Training

Heidi Schellman, Oregon State University
David DeMuth, Valley City State University

October 19, 2024

# Outline

- Description of our basic training and development of a second module for the hardware database. Credit to Dave Demuth for these slides

- Some feedback from our user base

- Note: Neutrino physicists tend to come from smaller experiments with less formal structure.

# DUNE

~1500 collaborators
~500 active FNAL accounts
~1281 in Slack
~37 countries
> 200 institutions
> 40 compute sites
~ 10 data sites

~ 10 PB/year → 30 PB/year

Oregon State University

Valley City State University

DUNE

# DUNE computing model

- Based on services not tiers

  - ~40 CPU sites, OSG/WLCG >= 2000 MB memory/core

  - HPC

  - ~10 Rucio controlled disk sites: FNAL, CERN, UK sites, NL … (need good network speeds)

  - Local caches

  - Tape store at FNAL, CERN, UK sites, IN2P3

  - Analysis facilities under construction

- CPU anywhere is useful for simulation, reconstruction over xrootd

- Our model tries to keep useful reconstructed samples on disk so dedicated disk resources near CPU are exceptionally valuable for fast analysis.

Oregon State University

Valley City State University

DUNE

# Communication channels

DUNE wiki (collab only)

DUNE docdb (collab only)

Slack (collab only)

Github pages (public)

> https://dune.github.io/computing-basics

> https://github.com/orgs/DUNE/projects/19

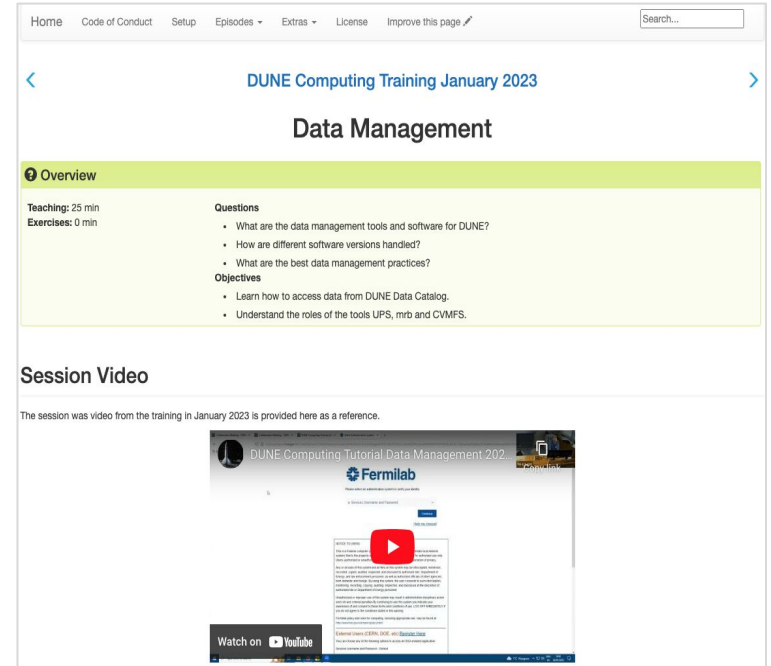Having docs behind a firewall makes early training difficult

We do have some issues getting google to index the public information

# Documentation DUNE Wiki

DUNE Computing hosts a one-stop for computing tutorials in several formats (pdf, wiki, html)

https://wiki.dunescience.org/wiki/Computing_tutorials

- Computing Basics
- FAQ
- The Justin Workflow System
- LArSoft training materials
- General HEP Training Courses
- Metacat documentation
- DUNE Hardware Database Training
- 2x2 and Minerva reconstruction tutorial



An overview of the Computing Basics training and the machinery to produce the lessons follows.

# DUNE Computing Tutorials

- DUNE Computing has coordinated 11 training events since 2016, with 400+ participants.
- Tutorials focused on three topics:
  - Data storage and management,
  - art & LArSoft,
  - Job submission and monitoring.
- The goal is to verify access to DUNE resources, understand the basics of logging in, storage areas, running applications, code modifications, submitting and monitoring batch jobs.
- Simultaneous face2face and Zoom modes, and post-production to allow asynchronous study.

  https://dune.github.io/computing-basics/

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

**Organizers**

Claire David
York University / FNAL

David DeMuth
Valley City State University

**DUNE Computing Consortium Lead**

Heidi Schellman
Oregon State University

**Lecturers**

Mike Kirby
FNAL

Steven Timm
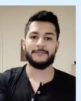FNAL

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. Other training events were two half day, and one day events.

Oregon State University

VALLEY CITY STATE UNIVERSITY

DUNE

# Training Logistics

- Event registration and communications are managed by Indico.
- Participants must verify their ability to use the [Unix Shell](#).
- And verify access the DUNE general purpose virtual machines at FNAL or CERN
- Much of the tutorial can be done using CERN resources.
- 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:
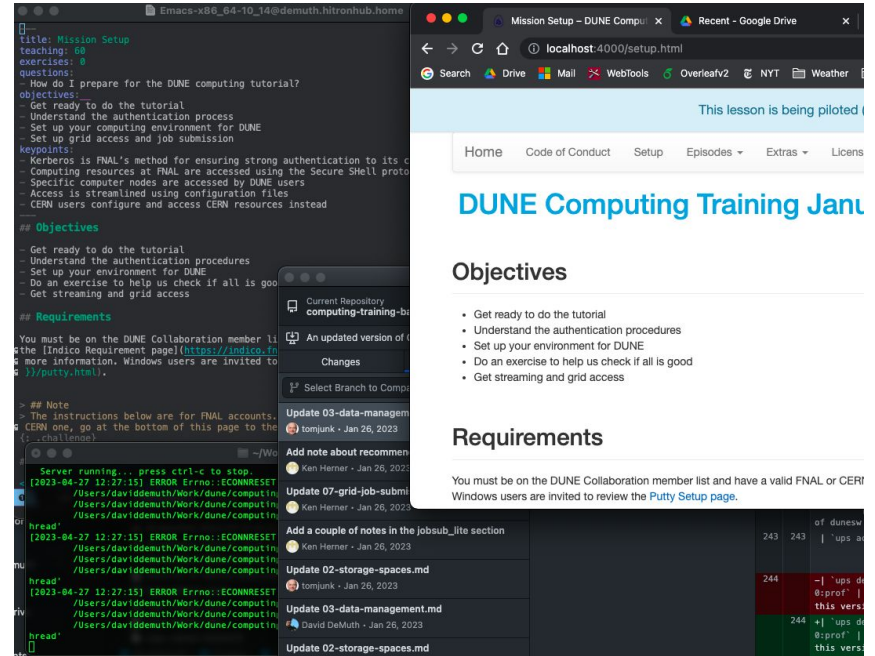
**Bash**
```
kinit -f username@FNAL.GOV
```

In advance, students demonstrate an understanding of using the Unix shell to access secure VMs.

Oregon State University   VALLEY CITY STATE UNIVERSITY   DUNE

# Lesson Development

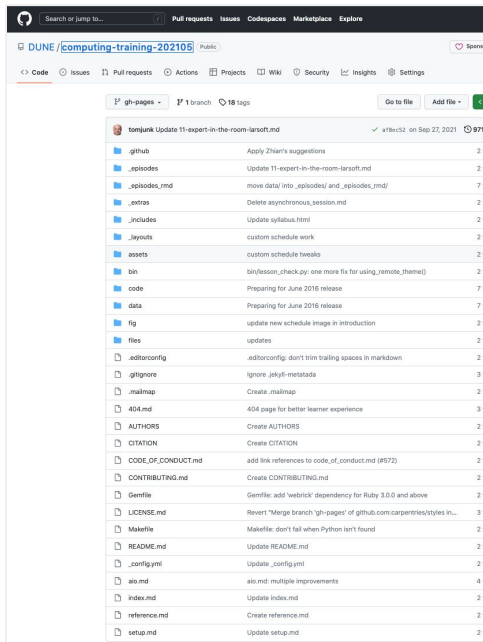The infrastructure to develop lessons is provided in the **Software Carpentries** framework:

- 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.
- Viewing edits in a localhost browser uses a Ruby/Jekyll rendering engine.
- 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.

Oregon State University · VALLEY CITY STATE UNIVERSITY · DUNE

# Lesson Deployment

**End users see:**



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

# Support

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.
  - As questions are volleyed, experts interact with individuals to find solutions.
  - Livedocs can be studied by all, solutions can be contrasted, and becomes a resource.
- **Slack channels** are set up for each training event.
- **GitHub Issues** is also used to support learning.



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

# Lessons Learned

- SWC template is elegant and functional for delivering active learning materials.
- It is also very popular for asynchronous training!
- A common look and feel for training materials has value.
- Markdown lessons allows experts to provide content without worrying (too much) about format
- **Pre-event homework that checks that participants can access FNAL servers is a must.**
- It's easy to be (over) ambitious with the material for events.
- Mentors ensure skill development and understanding.
- Coffee breaks allow students to catch-up on lesson activities, some assisted by mentors.
- Hybrid synchronous delivery of lessons works well.
- Zoom captures on YouTube allow asynchronous access, and a record of the event.
- Indico and Github sites require coordination.
- Providing build support streamlines all aspects of the training event.
- Undergraduates can contribute to lesson designs.

Oregon State University      VALLEY CITY STATE UNIVERSITY      DUNE

# Case Study: HWDB

- Utilizing the same lesson templates and editing environment, a two half-day hardware database training was recently offered by the University of Minnesota (H. Muramatsu, A. Wagner, U. Ekka): Indico 65297

- Event management, lesson development environment setup, GitHub site management, were among the centrally coordinated tasks provided to the lead instructors: https://github.com/DUNE/computing-HWDB

- A four person team benefitted by weekly meetings over a two month development cycle; key was undergraduate U. Ekka's markdown and styling contributions, e.g. image shading and pop-ups.

- The HWDB training event evidences that other topics could be developed similarly, e.g., condoDB.

**Centralizing DUNE Computing Training, A Case Study - HWDB.**

DUNE Computing has established a teaching and learning framework for new collaborators so as to jumpstart their analysis and simulation research.

Using Indico, computing tutorial events are scheduled for the May collaboration meetings for each year. These are delivered simultaneously in face-to-face and online formats, and are adapted with a third asynchronous modality.

Lesson development and delivery toolkits are provided by The Software Carpentry (SWC), and the deployment and use are consistent with the training activities of the HEP Software Foundation (HSF). HSF offers regular opportunities for training that include Unix Shell, GitHub, Python, and more advanced topics such as CI/CD; these we announce collaboration wide as encouragement for participation.

DUNE-specific training to date has included lessons on data storage, data management, LArSoft coding, and job management. Expanding DUNE-specific training HEP-wide by partnering with the HSF is a stretch goal.

Expanding DUNE-specific training is a nearer goal. With a DUNE hardware database training anticipated for May, 2024, we propose the event be centrally coordinated by the DUNE training group as a service, ensuring the following deliverables:

- Event management: working the the HWDB group, establish dates for the training, and send out a Save the Date.
- Schedule: In Zoom meetings, meet with HWDB training stakeholders to sort out the schedule of the event (lesson ordering and lengths), identify instructors, mentors, and target students.
- Lessons: Establish learner outcomes for each lesson. Ask instructors to sketch lessons in ways in which they are accustomed, beginning with the development of Powerpoint slides for example, and sharing those to seed a Markdown version produced for them (as a service).
- Provide support on the use of Markdown and the use of SWC templates, including the localhost rendering of episodes using Ruby/Jekyll technologies.
- Provide training on live coding techniques (side by side).
- Provide technical support as needed.
- Develop invitation for the training event participation.
- Establish pre- and post-event survey instruments.
- Consider near and long term support strategies that include Slack channels, GitHub FAQ, MediaWIki documentation, and direct discourse with lead engineers.
- Solicit participation via Indico event registration.
- Establish key points for each lesson.
- Determine pre-event training requirements, and ensure these are completed by participants, including the pre-event survey.

Oregon State University    VALLEY CITY STATE UNIVERSITY    DUNE

# Results from a survey

- Survey done before our last collaboration meeting in September 2024

- 66 responses from 20 countries

answers from

BR(2), CERN(3), CL(2), CO(2), CZ(1), DE(1), ES(1), FR(1), HU(1), IL(3), IN(1), IT(6), JP(1), MA(1), NL(2), PT(1), RO(1), TR(1), UK(3), US(30)

10/16/24

Oregon State University

Valley City State University

DUNE

# What career status are you?

66 responses



- 🔵 Undergraduate
- 🔴 Graduate Student
- 🟠 Postdoc
- 🟢 Senior Scientist
- 🟣 Engineer/Computer Scientist/Technician
- 🔵 Faculty
- 🔴 Prof.
- 🟢 Assistant Professor

🔼 1/2 🔽

Oregon State University

Valley City State University

DUNE

Where are you mainly located

66 responses

Oregon State University

Valley City State University

DUNE

# Where do you do your computing

66 responses

10/16/24

Oregon State University · Valley City State University · DUNE

# What Languages are you familiar with

65 responses



- C++ — 62 (95.4%)
- python — 54 (83.1%)
- bash — 40 (61.5%)
- perl — 10 (15.4%)
- Julia — 3 (4.6%)
- R — 5 (7.7%)
- FORTRAN — 19 (29.2%)
- (unlabeled) — 2 (3.1%)
- Swift — 1 (1.5%)
- embedded C — 1 (1.5%)

Oregon State University    Valley City State University    DUNE

# What software packages are you familiar with?

65 responses

| Package | Responses |
|---|---|
| ROOT | 58 (89.2%) |
| GEANT4 | 27 (41.5%) |
| LArSoft | 20 (30.8%) |
| edep-sim | 10 (15.4%) |
| numpy / PANDAS / R / ... | 44 (67.7%) |
| TensorFlow/CUDA | 12 (18.5%) |
| git, cmake, singularity/apptainer | 1 (1.5%) |
| Self written numerical analysis… | 1 (1.5%) |
|  | 1 (1.5%) |
| Fluka | 1 (1.5%) |
| Eigen, GenFit | 1 (1.5%) |
| Pandora | 1 (1.5%) |

Oregon State University

Valley City State University

DUNE

# How do you do batch processing?

62 responses

10/16/24

Oregon State University    Valley City State University    DUNE

Would you be interested in learning more about:

65 responses

| Topic | Count |
|---|---|
| Unix basics for DUNE | 23 (35.4%) |
| Databases for DUNE | 39 (60%) |
| Batch processing basics | 38 (58.5%) |
| Reconstruction frameworks | 40 (61.5%) |
| Analysis frameworks | 47 (72.3%) |
| Documentation techniques | 22 (33.8%) |
| Code management with git | 26 (40%) |
| Configuration management (… | 25 (38.5%) |
| Continuous integration | 21 (32.3%) |
| Build systems | 21 (32.3%) |
| Parallel processing techniques | 26 (40%) |
| How to do training | 15 (23.1%) |
| DUNE-specific AI techniques | 33 (50.8%) |
| Beamline/Accelerator codes | 19 (29.2%) |
| Containers | 16 (24.6%) |
| LarSoft | 1 (1.5%) |

10/16/24

Oregon State University

Valley City State University

DUNE

How do you get computing help?

65 responses

| Category | Count (Percentage) |
|---|---|
| Ask local expert | 27 (41.5%) |
| DUNE Slack channels | 36 (55.4%) |
| Servicenow tickets | 10 (15.4%) |
| DUNE tutorials | 29 (44.6%) |
| DUNE wiki | 35 (53.8%) |
| Stack Exchange | 26 (40%) |
| Online forums such as root.cer… | 33 (50.8%) |
| Mr. Google/Chat-GPT | 32 (49.2%) |
| Books | 1 (1.5%) |
| Call fermilab help service line (I… | 1 (1.5%) |

Oregon State University

Valley City State University

DUNE

# What format do you like for tutorials?

64 responses

Oregon State University    Valley City State University

Are you interested in helping with a tutorial?

23 responses

# what would you like to help develop?

- When DUNE software is ready and available with Spack I would be happy to help prepare or give some tutorials

- Python

- Databases, anything else I know well enoug to be useful.

- LArSoft Basics

- Software development

- Reco/Sim framework, Analysis framework

- job submission

- Data management

- I'm interested in helping data analysis and simulation/reconstruction.

- GEANT4 and ROOT related topics mainly

- HWDB, but willing to help other DB-related topics

- python, analysis tools, git

- Unix, code management, configuration though mostly wherever my skill set best serves the groups needs.

- C, Python

Oregon State University

Valley City State University

DUNE

# When you joined DUNE what resources were you pointed to?

- DUNE wiki

- DUNE docb

- LArSoft workshop

- dune.github.io/computing-basics (the online tutorial)

- Google (which doesn't work well as most items are now hidden)

- Nothing….

- No mention of the DUNE slack among respondents.

10/16/24

Oregon State University    Valley City State University

# Lessons learned

- Onboarding needs to be more uniform

- People like the tutorial format (Carpentries) but we need to extend to more topics

- DUNE software environment is changing rapidly

  - SAM -> Metacat/Rucio for data ✓

  - POMS☐ Justin for batch 🔋

  - ups ☐ spack for configuration management ☄

- Any documentation we write now will be pretty up-to-date as so much is new.

- Basic training on batch systems (what they even are?) would be useful

- But many things are not yet documented

10/16/24

Oregon State University     Valley City State University     DUNE