# Training Activities and Experiences in the HSF.

**Delivering sustainable software training**

Michel Hernandez Villanueva (DESY),
and many more! (see slide 12)

On behalf of the HSF Training WG

**FH Sustainability Forum meeting**
June 5, 2023

# Definition

## Or "Why am I here today?"

sustainability

*noun*

the ability to be maintained at a certain rate or level.
"the sustainability of economic growth"

*What we mean by "sustainable training" in other talks*

# Definition
## Or "Why am I here today?"

🔊 sustainability

*noun*

the ability to be maintained at a certain rate or level.
"the sustainability of economic growth"

- <u>avoidance</u> of the <u>depletion</u> of natural resources in order to maintain an <u>ecological</u> balance.
"the pursuit of global environmental sustainability"

*What we will also discuss today*

# Software Development in HEP
## As a key for a successful scientific program

- Scientific collaborations are **big and growing**.
  - O(1000) collaborators in hundreds of institutes around the world.

# Software Development in HEP
## As a key for a successful scientific program

- Scientific collaborations are **big and growing**

    - O(1000) collaborators in hundreds of institutes around the world

- High Energy Physics (HEP) and Nuclear Physics (NP) are **computationally intensive** and **data driven** fields

    - A full physics potential requires investment into the software used to collect, process, and analyse data

- **Developers with strong foundation** are critical resources in the success of the current and future experiments

    - The researchers must be brought up to date with new software technologies, concurrent programming, and artificial intelligence

    - They must maintain, improve, and sustain the software

# Training and Onboarding Initiatives in HEP

## How do experiments teach software?

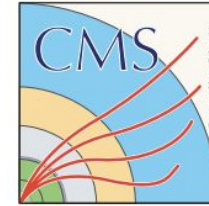Virtual                  Hybrid                In person
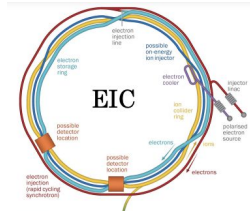


Online book



Starter Kit



Data Analysis Schools

Online tutorials

Software tutorials

Synchronous tutorials "Carpentries-style"

**"Software is different, but challenges are common"**

Reinsvold Hall, Allison (US Naval Academy), CHEP 2023
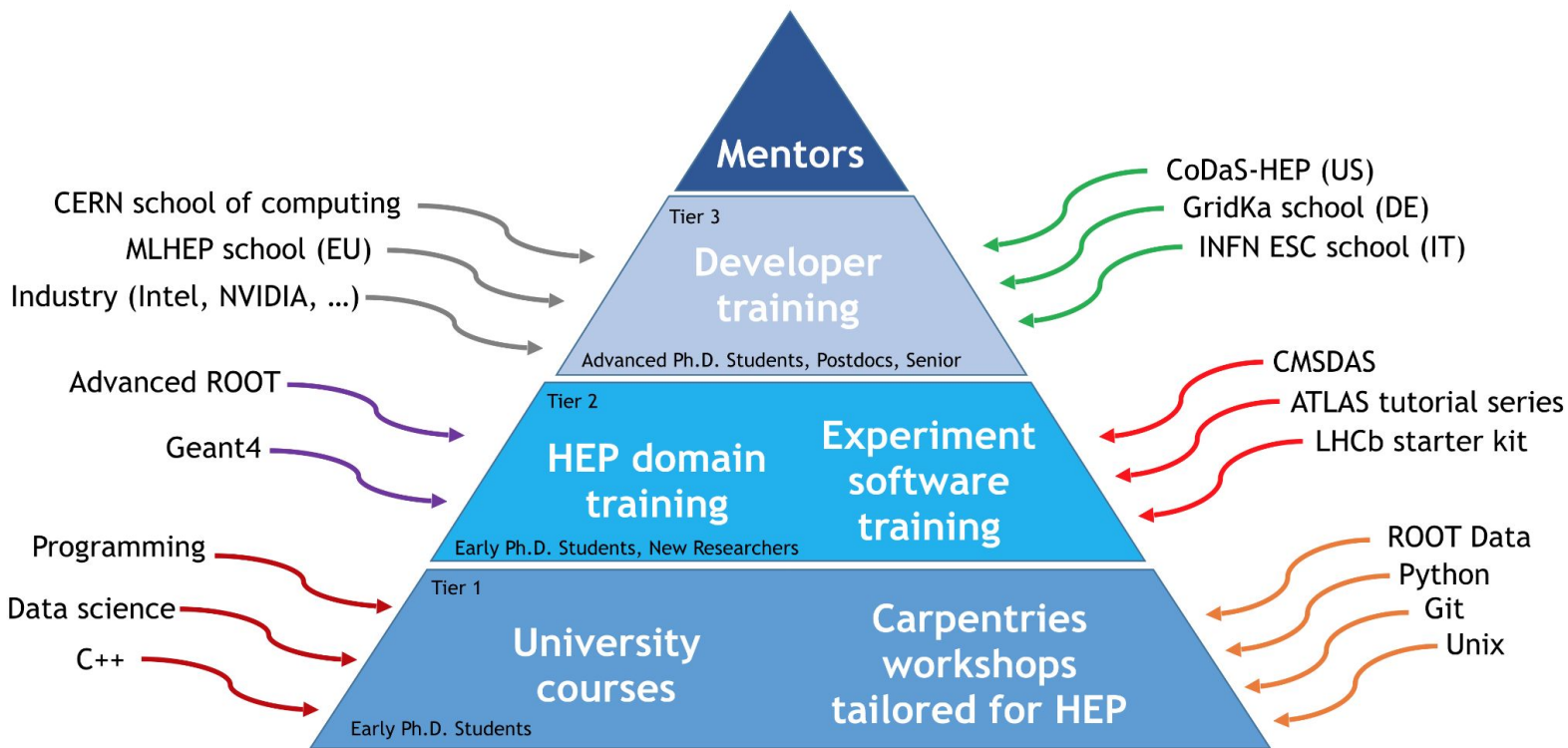
# HEP Software Training

**Why not leave everything to the experiments?**

- O(10k) HEP people worldwide need to be trained in software engineering & computing

- Common challenges faced:
  - Most people developing code have non-permanent positions with contracts of 2 - 4 years
  - Training activities are undervalued in making career steps and by funding agencies
  - Material for training is a moving target as technology evolves (e.g., ML, GPUs, FPGAs, …)

- **This should be a community effort!**

# HEP Software Training
## We can cover more ground together!



CERN school of computing
MLHEP school (EU)
Industry (Intel, NVIDIA, …)

Advanced ROOT
Geant4

Programming
Data science
C++

**Mentors**

Tier 3
**Developer training**
Advanced Ph.D. Students, Postdocs, Senior

Tier 2
**HEP domain training**
**Experiment software training**
Early Ph.D. Students, New Researchers

Tier 1
**University courses**
**Carpentries workshops tailored for HEP**
Early Ph.D. Students

CoDaS-HEP (US)
GridKa school (DE)
INFN ESC school (IT)

CMSDAS
ATLAS tutorial series
LHCb starter kit

ROOT Data
Python
Git
Unix

https://iris-hep.org/ssc.html

# HEP Software Foundation
## (HSF)

- The role of the <u>HEP Software Foundation</u>, started in 2015, is to facilitate coordination and common efforts in software and computing across HEP in general

- The goal is to describe a global vision for software and computing for the current and future experiments

  - Working groups cover Training, Analysis, Generators, Simulation, Reconstruction and Software Triggers, etc.

- The HSF's role is one of an information conduit and meeting point

  - Report on interesting and common work being done

  - Forum for technical comments and discussion

  - Encourage cooperation across experiments and regions

  - Motivate the publication of summary documents or papers for reference

HSF

HEP Software Foundation

# HSF Software Training
## Organization

- Established 4 years ago

- Develops material for an introductory software curriculum
  - And teaches this curriculum to scientists

- Focuses on common software material across HEP
  - From basic core software skills
  - To advanced training required in software and computing

- Remote weekly public meetings (via Zoom) to plan and assess progress
  - Led by four co-conveners
  - Proposals are discussed and events are planned

- Engages with different experimental collaborations and initiatives
  - IRIS-HEP, FIRST-HEP, and The Carpentries

**Join an event!**
Discover new topics together with mentors and peers!

**Self study!**
Learn at your own pace. No matter if you want to get a quick overview or dive in the details, this is for you!

https://hepsoftwarefoundation.org/training

# HSF Software Training
## Principles

We need a **unified**, **scalable**, and **sustainable** software training framework

**Unified**
- Material and events should be **centrally listed** & **discoverable**
- Concentrate efforts by developing **cross-experiment** content
- A **community** must guide, support, and coordinate

**Scalable**
- Material must be teachable by **multiple instructors**
- **Self-study** must not be an afterthought

**Sustainable**
- Material must be **open source** and **maintained collaboratively**
- **Incentives & recognition** important motivators

# HSF Software Training

## The community

- An active community of members supporting training on voluntary basis

  - Coming from multiple collaborations, adding value to the training from different environments

- Profile of each tutor that contributes is included in the HSF training page

  - Public recognition of their capability, skills and contribution



https://hepsoftwarefoundation.org/training/community.html

# HSF Software Training Platform

## We can cover more ground together!

### Weekly meetings

### Monthly Hackathons

HSF iris hep

✨ Containerization Training Hackathon ✨

🔥 The big goal!

we do a little hacking

### Platforms

GitHub

### Community pages

Our community

Amber Roepe (she/her)
Andrea Valassi
Clemens Lange
Dan Guest

Daniel S. Katz
David Chamont
David Yakobovitch
Giordon Stark

Graeme A Stewart
Henry Schreiner
Jackson Burzynski
Jim Pivarski

### How-to guides

HSF Training Workshop Checklist ✅

Let's streamline our organization and make sure we don't forget anything!

Note: there's also a Hackathon checklist.

Before the workshop

Setting up documents and more

HSF
HEP Software Foundation

iris hep

---

## Software Development and Deployment

| **Version controlling with git** | **Advanced git** | **CI/CD (gitlab)** |
|---|---|---|
| Track code changes, undo mistakes, collaborate. This module is a must. | Learn to work with branches and more with this interactive webpage. | Continuous integration and deployment with gitlab. |
| 📕 Start learning now! | 📕 Start learning now! | 📕 Start learning now! |
| | | 📹 Watch the videos! |
| 🔧 Contribute! | 🔧 Contribute! | 🔧 Contribute! |
| **CI/CD (github)** | **Docker** | **Singularity** |
| Continuous integration and deployment with github actions. | Introduction to the docker container image system. | Introduction to containerization with Singularity/Apptainer. |
| | | ❇ Status: Beta testing |
| 📕 Start learning now! | 📕 Start learning now! | 📕 Start learning now! |
| 📹 Watch the videos! | 📹 Watch the videos! | 📹 Watch the videos! |
| 🔧 Contribute! | 🔧 Contribute! | 🔧 Contribute! |
| **Unit testing** | **Level up your python** | |
| Unit testing in python. | Advanced bits of python (testing, debugging, logging, and more) | |
| ❇ Status: Beta testing | | |
| 📕 Start learning now! | 📕 Start learning now! | |
| 🔧 Contribute! | 🔧 Contribute! | |

https://hepsoftwarefoundation.org/training/
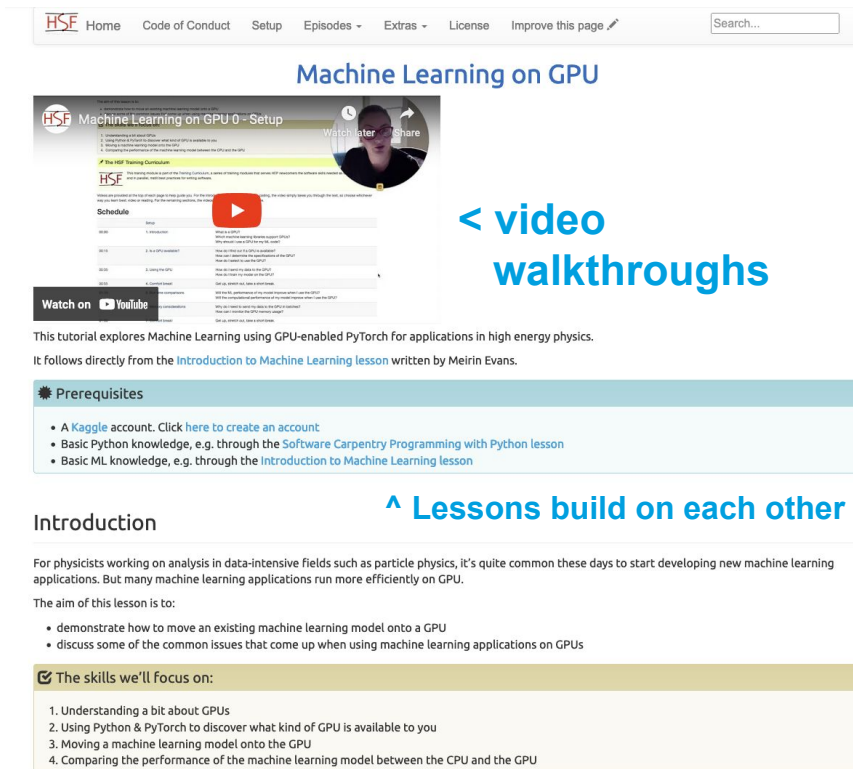
# HSF Software Training Center
## For HEP newcomers

- <u>HSF Training Center</u> currently lists 21 training modules
  - Basics: Bash, Git, Python, Matplotlib
  - HEP basics: ROOT, Uproot
  - Software development, ML, devops, etc

- Goal: Training Center as a focal point for all HEP Training activities
  - Free and experiment-agnostic

- Embrace the framework of The Carpentries
  - Built from markdown files
  - Source at GitHub (**Anyone can contribute!**)
  - Verbose and self-study ready



**< video walkthroughs**

**^ Lessons build on each other**

**^ Enough verbosity for self-study**

# HSF Software Training Center
## For HEP newcomers

- **The big picture: scientists with skills for delivering high-quality code**

- We must aim to train our community with the best practices for sustainable software development
A few examples:
  - Continuous Integration
  - Testing, testing, testing
  - Reproducibility, preservation
  - Project development methodologies
  - Green coding practices: efficient algorithms and data structures, reduce memory consumption and network traffic…

- Large impact at computing centers in the long term!

# HSF Software Training Center
## For HEP newcomers

- **The big picture: scientists with skills for delivering high-quality code**

- We must aim to train our community with the best practices for sustainable software development
  A few examples:

  - Continuous Integration

  - Testing, testing, testing

  - Reproducibility, preservation

  - Project development methodologies

  - Green coding practices: efficient algorithms and data structures, reduce memory consumption and network traffic…

- Large impact at computing centers in the long term!

We are halfway on this list.

**Reaching the bottom needs support from the HEP community
(For example, keeping communication with this forum 😄)**

DESY

# Analysis Preservation
## An example on how-to "train to sustain"

- **"Preservation reduce the resource footprint of our analyses"** [Yves Kemp et al. "Sustainable computing in HEP"]

- Last year, we developed modules to teach how to consider analysis preservation right from the beginning



- Developed by the HEP community during Containerization & Analysis Preservation Hackathons

- Teaching Docker, Singularity/Apptainer, CI/CD with github/gitlab, REANA (soon)

- Using CMS Open Data

- Emphasis on **self study with videos**

- Material can be used in training events

DESY.

# Analysis Preservation
## Virtual workshop

- A week

- 100 participants

- A good example on how the interest of the community can drive towards training events related to environmental sustainability



Training on Analysis Preservation (Virtual)

16–21 Jan 2023
Virtual

Learning the tools to make your analysis last to infinity and beyond!

https://indico.cern.ch/event/1219810/

# Training Events

## In person

- HSF Training software tutorials through 2020:

  - In-person participation only

  - Approximately 35 participants per workshop

- Impact on ecological and social sustainability:

  - Travel limits the accessibility to research groups with access to sufficient funding

  - Typical carbon footprint ~0.5 t $CO_2$e / person:

    - Intra-continental travel: 0.4 t $CO_2$e per person

    - Hotel stays: ~25 kg $CO_2$e per person per night

  - Compare with estimated average EU (US) annual carbon footprint of 7 (16) t $CO_2$e per person

  - **A workshop can increase one's footprint by 5% to 10%**

# Training Events
## Virtual

- Holding Virtual events since 2020.

  - COVID-motivated, but this training modality is here to stay

- 18 online software trainings, **1300+ participants trained**

  - Logistics are easier. Recordings available

  - Minimum environmental impact

- But also disadvantages:

  - Lower engagement, distractions

  - Meaningful interactions harder

Past Events

- 18 May - 19 May 2023 - HSF/IRIS-HEP Software Basics Training (Virtual) HSF
- 6 Mar - 10 Mar 2023 - 6th HEP C++ Course and Hands-on Training - The Essentials HSF
- 16 Jan - 20 Jan 2023 - Analysis Preservation Workshop HSF
- 11 Oct - 13 Oct 2022 - 5th HEP C++ Course and Hands-on Training - Advanced C++ HSF
- 3 Oct - 8 Oct 2022 - ESC22 EFFICIENT SCIENTIFIC COMPUTING
- 28 Sep - 30 Sep 2022 - HSF/IRIS-HEP Software Basics Training HSF

# Training Events
## Back to in person?

- Discussion on how to quantify how effective **online vs in-person events** are

    - 3 eigenvectors: knowledge exchange, co-creation, community building

    - We need standardise metrics to use throughout the courses

- How to get the advantages of in-person interaction without the environmental impact?
  Several ideas, discussion in progress

    - Self-organized training events (we provide all the material)

    - Events allocated with major conferences

    - Regional training events

- Or the other way around: **improve interaction in online events**

    - Breaking-Ice sessions at the beginning of the event may bring back the joy of the events
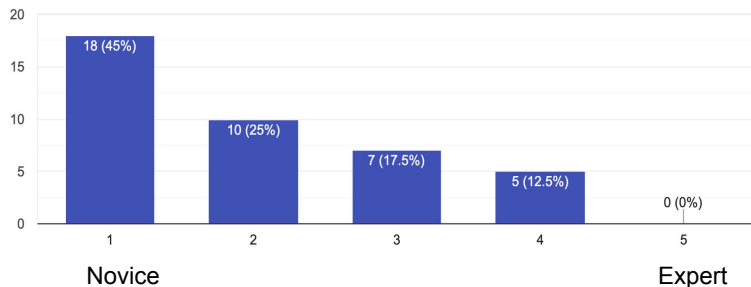      when they are organized online

# Training Events
## Measuring how effective are online events

- Each event, anonymized pre and post-surveys are circulated with the students
  - Pre-survey: Demographics, How much do you know?
  - Post-survey: How much do you know **now**? What can we do better next time?

- We also do our best collecting information of people dropping the event
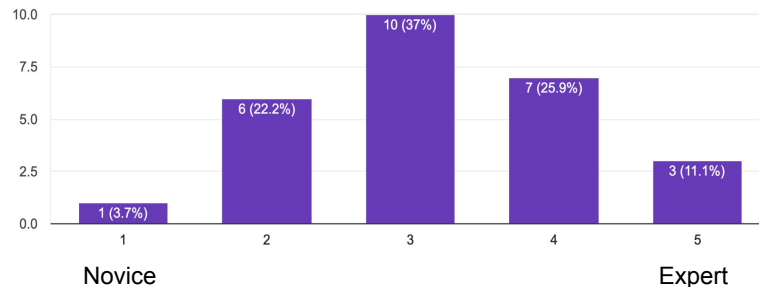
How confident are you in your knowledge and abilities when using Git?
40 responses



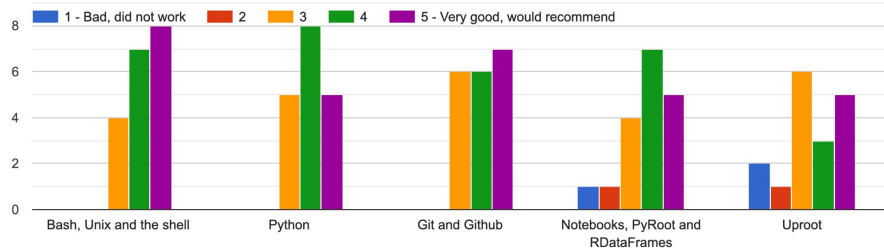How confident are you in your knowledge and abilities when using Git after the workshop?
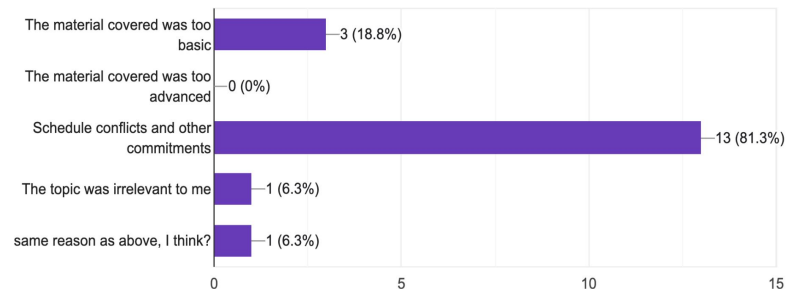27 responses

# Surveys
## Some examples

Please rate how successfully each topic was covered (1 - 5)



If you only attended FEW sessions, why did you skip the other sessions?
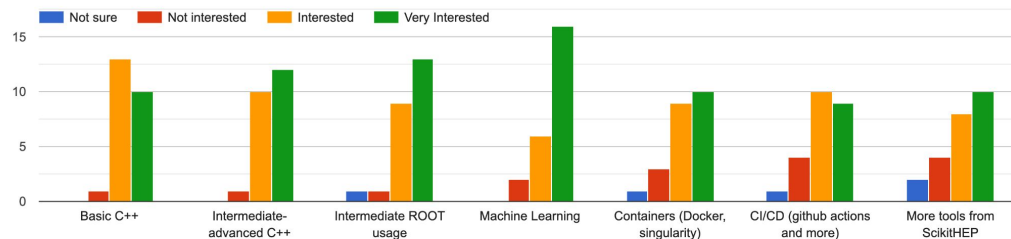
16 responses



How can we improve our bash/shell training?

3 responses

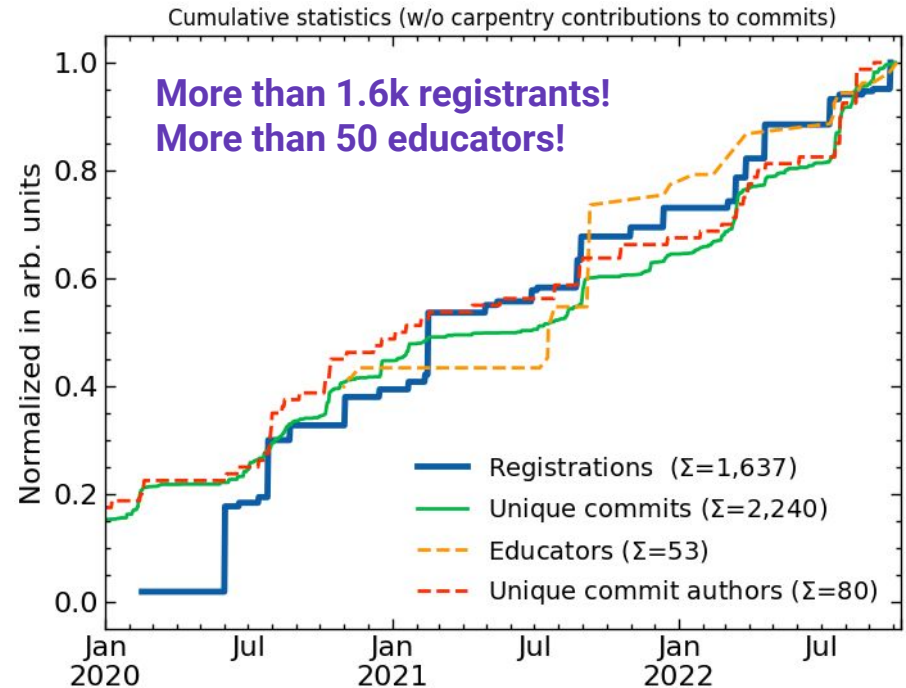How can we improve our python training?

1 response

Please rate your interest in attending future tutorials on the following topics:

# Training Challenge
## Scaling up

- Proposal to expand the effort in the long-term, defining a clear target in form of a Training Challenge.

  - **Scalability**: What is the number of students to reach? How many events does imply?

  - **Sustainability**: How to incentivize new trainers to continually join?

  - … **and Sustainability:** How to minimize the environmental impact, delivering effective training?

  - **Diversity and inclusion**: Everyone feels welcome to participate? How to standardize metrics?

- Active discussion happening right now.

Cumulative statistics (w/o carpentry contributions to commits)

**More than 1.6k registrants!**
**More than 50 educators!**

- Registrations ($\Sigma$=1,637)
- Unique commits ($\Sigma$=2,240)
- Educators ($\Sigma$=53)
- Unique commit authors ($\Sigma$=80)

# Summary
## And how to collaborate

- The HSF Training is a community-driven effort, covering the software training requisites for a sustainable operation of physics experiments

- We have consolidated Software Training events virtually held
  - Discussion on how to sustain & scale up is relevant and happening right now

- We have included a training event teaching Analysis Preservation: containerization & CI/CD with open data
  - **Extend to more topics related to sustainability depends on the motivation of the community**

- Public weekly meetings: Mondays at 4pm CEST
  - https://indico.cern.ch/category/10294/
  - Everyone is welcome to join!

- **Reach us also via the channels shown in <u>our webpage</u>.**

# Join us!

@hepsoftfound       @hsf-training       hepsoftwarefoundation.org

# Backup

# Training Events

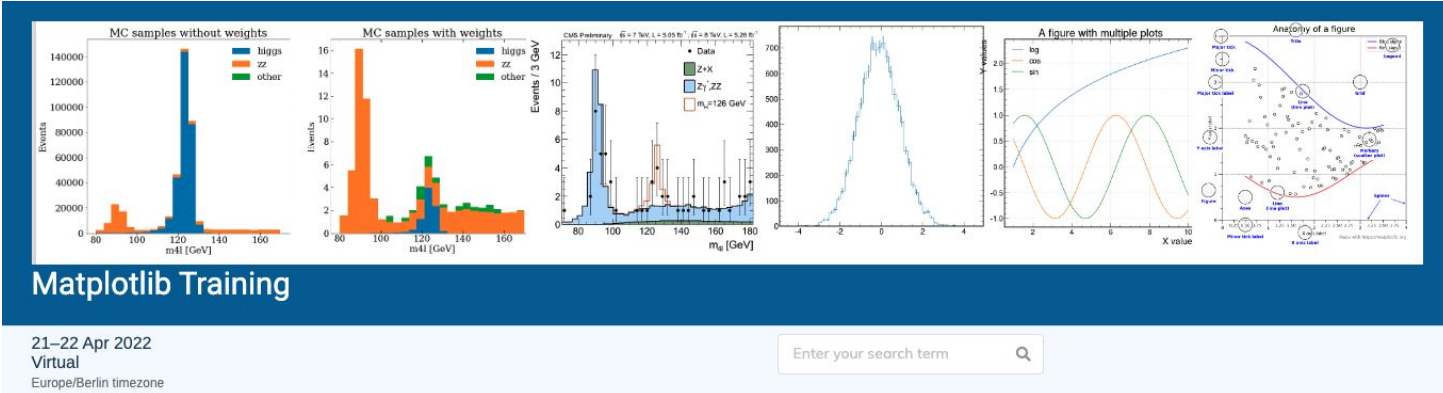## Software Training Basics

- Since Aug 2021, in collaboration with <u>The Carpentries</u>, we have established a training event for newcomers 3 times per year.

  - Agenda of events in 2022: <u>here</u>, <u>here</u> and <u>here.</u>

- Three days event organized as follows:

  - Two days of fast-track to competency with software fundamentals: Bash, Git, Python.

  - One day dedicated to HEP Software:
    <u>ROOT</u> (data analysis framework), <u>Scikit-HEP</u> (data analysis in Python)

- A limit 80 students per event.

  - Instructors for each topic.

  - 5 mentors in average, helping in breakout rooms.

- Material and recordings are preserved on the page of the event.

DESY.

# Intermediate Training
## Development and organization

- Continuously organizing training meetings and hackathons for extending/improving material.

- Established training events in <u>C++</u> and <u>Matplotlib</u>.

- Current topics in development: Docker / Apptainer, CI/CD in GitHub and GitLab.

  - During <u>hackathons</u>, aiming to prepare an Analysis Preservation Training event.

- New ideas are always welcome.



**Matplotlib Training**

21–22 Apr 2022
Virtual
Europe/Berlin timezone

Enter your search term