# Training on analysis reproducibility in the HSF

Michel Hernandez Villanueva (BNL), and many more!

On behalf of the HSF Training WG

**Workshop on ATLAS Computing and Software Activities at BNL**

Mar 18, 2024

1

# HEP and Nuclear Physics Software
## As a key for a successful scientific program

- O(10k) HEP/NP 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

    - Material for training is a moving target as technology evolves (e.g., ML, GPUs, FPGAs, …)

    - Training activities not in the highest priority when making career steps

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



**Brookhaven**
National Laboratory

# HSF Software Training
## Organization

- The role of the HSF is to **coordinate common efforts** in software and computing across HEP in general

- HSF Training was established in 2018

- Develops material for an introductory software curriculum

- Focuses on **common software material across HEP and NP**
  - From basic core software skills, to advanced training required in software and computing
  - Experiment agnostic, minimal dependencies on having an account at X site

- 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
HEP Software Foundation

Brookhaven
National Laboratory

# The HSF Software Training Platform

## We can cover more ground together!

### Weekly meetings


### Monthly Hackathons


### Platforms

GitHub

### Community pages


### How-to guides


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! | |

- We build a **unified**, **scalable**, and **sustainable** software training framework

**Brookhaven** National Laboratory

# Analysis Reproducibility
## Train to Sustain

- **Data is expensive, experiments are unique**

  - Preserving the knowledge around them is a must

- We have developed modules to the learn the tools & consider analysis reproducibility right from the beginning



- Developed by the HEP community during [Containerization & Analysis Preservation Hackathons](#)

- Using **[CMS Open Data](#)**

- **Analysis preservation is work-in-progress in most of the collaborations**
  - We do not focus in a particular reproducibility scenario
  - Instead, we review the commonly used tools

# Analysis Reproducibility
## Training modules

We take a quick tour learning the basic functionality of tools popular in analysis preservation and reproducibility.

- **Containerization technologies**
  Podman & Docker
  Apptainer (Singularity)
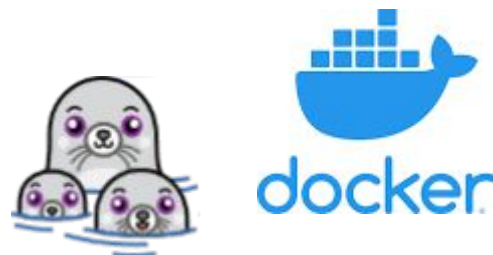
- **Continuous Integration/Deployment (CI/CD)**
  GitLab pipelines
  GitHub actions

- **Analysis platforms and workflows** (in preparation)
  REANA
  Snakemake
  ...

Brookhaven
National Laboratory

# Analysis Reproducibility
## Virtual Events

- Organized in [2023](#) & [2024](#)
  - Monday to Friday
  - **> 100 registrations per event**





Training on Analysis Preservation (Virtual)

16–21 Jan 2023
Virtual

Enter your search term

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

**Note: this tutorial used to be called: "Preservation" → "Pipelines" → "Reproducibility"**

- ATLAS
- CMS
- ALICE
- LHCb
- Belle II
- Theory
- Dune

13.3%
7.4%
7.4%
13.3%
24.5%

1/5

Brookhaven
National Laboratory

## Monday
*Welcome*

## Tuesday to Thursday
*Work on your own, when you want*

## Friday
*Hands-on sessions*

Kickoff/Orientation
[15:00 CET]

Analysis
Preservation@CMS
[15:10 CET]

REANA
[15:40 CET]

Help with Setup
[16:10-17:00 CET]

Watch and work through tutorials:
Indico Agenda

Block 1:
[10-12 CET]

Block 2:
[13-15 CET]

Block 3:
[17-19 CET]

Block 4:
[21-23 CET]

Brookhaven
National Laboratory

# Monday
*Welcome*

# Tuesday to Thursday
*Work on your own, when you want*

# Friday
*Hands-on sessions*

**Kickoff/Orientation**
[15:00 CET]

**Analysis Preservation@CMS**
[15:10 CET]

**REANA**
[15:40 CET]

**Help with Setup**
[16:10-17:00 CET]

Brookhaven
National Laboratory

> *"Your closest collaborator is you six months ago…*

Theory model/ LHE file

CMS GEN-SIM step

geometry/ conditions database

**Your work the of the past N years**

CMS specific

CMS DIGI step

CMS RECO step

ntuplisation/selection

Internal documentation

PUBLISHED

Analysis specific

Statistical analysis

Experiment data

**Reproducible**

**Currently not preserved**

**Preserved**

Clemens Lange — Analysis preservation at CMS

## Monday
*Welcome*

## Tuesday to Thursday
*Work on your own, when you want*

## Friday
*Hands-on sessions*

Kickoff/Orientation
[15:00 CET]

Analysis Preservation@CMS
[15:10 CET]

REANA
[15:40 CET]

Help with Setup
[16:10-17:00 CET]

Watch and work through tutorials:
Indico Agenda

Virtual support on Slack.

Block 1:
[10-12 CET]

Block 2:
[13-15 CET]

Block 3:
[17-19 CET]

Block 4:
[21-23 CET]

Brookhaven
National Laboratory

# Monday
*Welcome*

# Tuesday to Thursday
*Work on your own, when you want*

# Friday
*Hands-on sessions*
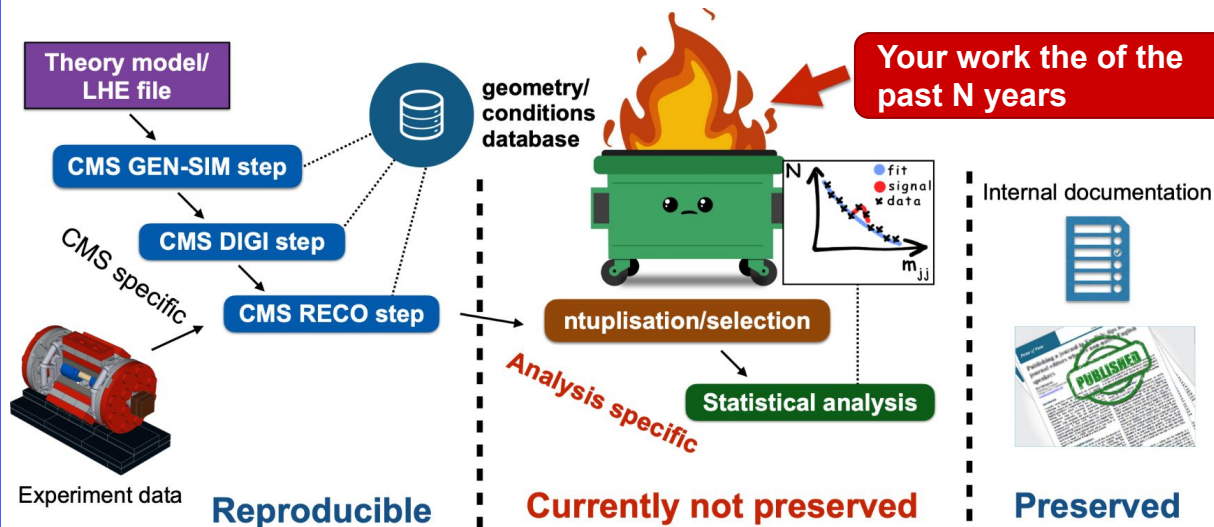
Kickoff/Orientation
[15:00 CET]

Analysis
Preservation@CMS
[15:10 CET]

REANA
[15:40 CET]

Help with Setup
[16:10-17:00 CET]

rials:

## Apptainer/Singularity

- **Easy:** Repeat one of the docker exercises on your cluster with Apptainer/Singularity using an interactive session

- **Medium:** Do the same with a definition file.
  - Option: Use the `%runscript` directive to print out the `uproot` version when the container runs

- **Hard:** Perform the CMS example analysis in a single execution using a definition file. Save the plots in the execution directory.
  - Hint: keep in mind where to store intermediate files.

exercises.pdf

Block 1:
[10-12 CET]

Block 2:
[13-15 CET]

Block 3:
[17-19 CET]

Block 4:
[21-23 CET]

**Brookhaven**
National Laboratory

# Surveys

## Measuring how effective is our training

How comfortable are you with the following docker commands/tasks

### Pre-workshop



How difficult were the lessons/exercises?



How comfortable are you with the following docker/podman commands/tasks

### Post-workshop



Would you like to join the next events as a teacher or mentor? We will be happy to provide all the support necessary for your success.
19 responses



Brookhaven
National Laboratory

12

# Analysis Reproducibility on ATLAS

## Opportunities for collaboration

- The HEP/NP community would benefit greatly from learning about analysis reproducibility in ATLAS

- Many possibilities to share the ATLAS experience
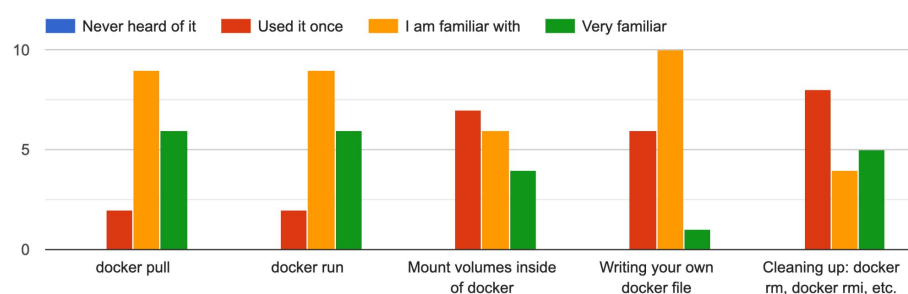
  - An **introductory talk** on analysis reproducibility in the ATLAS collaboration

  - Include our material in ATLAS training events and **improve it based on the feedback**
    - Source code in [GitHub](#), modules written in Markdown
    - A good example is the [feedback from the Analysis Preservation BootCamp @ Valencia](#)

  - Usage of **ATLAS open data** & analysis examples for HSF modules

  - New training modules with tools used at ATLAS

# Summary

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

- We have included a training event teaching Analysis Preservation: containerization & CI/CD with open data
  - **Extend to more related topics 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>.**



**Brookhaven**
National Laboratory

# Join us!

 **@hsf-training**       **hepsoftwarefoundation.org**

# Backup

# Training and Onboarding Initiatives in HEP/NP

How do experiments teach software?

Virtual        Hybrid        In person

 Online book

 Starter Kit
(As a consequence
of covid19)

 Data Analysis
Schools

 Online tutorials

 Software
tutorials

 Synchronous tutorials
"Carpentries-style"

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

Brookhaven
National Laboratory

# Software Training

## The training pyramid



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

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

Advanced ROOT
Geant4

Programming
Data science
C++

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

CMSDAS
ATLAS tutorial series
LHCb starter kit

ROOT Data
Python
Git
Unix

Brookhaven
National Laboratory

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

18

# 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

**Brookhaven**
National Laboratory

# HSF Training Center

- **The big picture: scientists with skills for delivering high-quality and reproducible analysis**

- 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…

We are halfway on this list.

**Reaching the bottom needs support from the community**

- Large impact at computing centers ($$) in the long term!

**Brookhaven**
National Laboratory

# Training Events

## In Person

- HSF Training software tutorials through 2020:

  - In-person participation only

  - Approximately 35 participants per workshop

- Ecological and social impact:

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

  - Typical carbon footprint ~0.5 t $CO_2e$ / person:

    - Intra-continental travel: 0.4 t $CO_2e$ per person

    - Hotel stays: ~25 kg $CO_2e$ per person per night

  - Compare with estimated average EU (US) annual
    carbon footprint of 7 (16) t $CO_2e$ per person

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


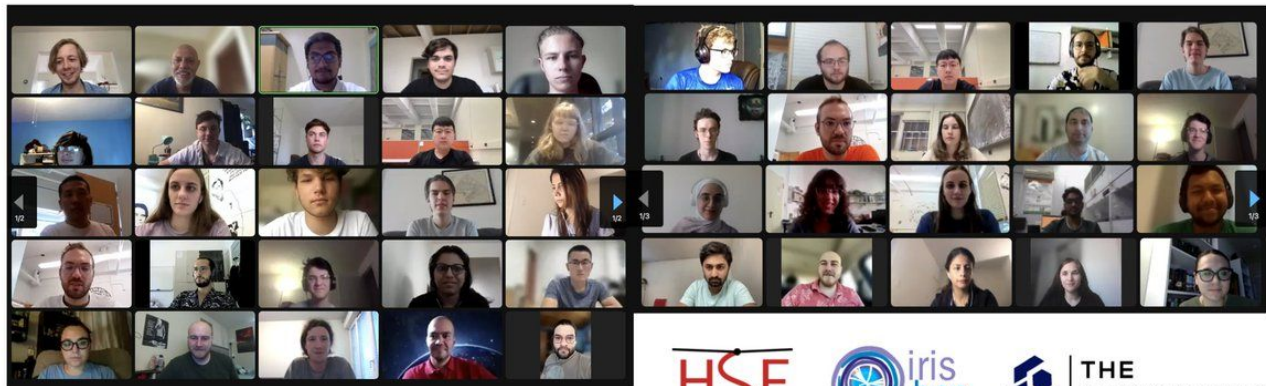


**Brookhaven**
National Laboratory

# 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



**Brookhaven** National Laboratory

# The Training Challenge

## Scaling up

- Proposal to expand the effort in the long-term

  - **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?

Cumulative statistics (w/o carpentry contributions to commits)

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

Normalized in arb. units

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

Jan 2020 — Jul — Jan 2021 — Jul — Jan 2022 — Jul

**Brookhaven**
National Laboratory