

Data Analysis Training in CMS

Gabriele Benelli (Brown University)

HSF Training Pre-CHEP Workshop, Oct 19th 2024



Data Analysis Training in CMS



- A quick overview of CMS and the scope of the trainings
- Some details on our main training programs:
 - CMS Data Analysis Schools (DAS)
 - Hands-on Advanced Tutorial Sessions (HATS)
 - Graduate-level academic lectures
 - Topical workshops (Statistical tools, Machine Learning, Trigger)
- Some reflections based on our experience in CMS

- Compact Muon Solenoid
- Large General-Purpose LHC experiment analyzing proton-proton collisions at the highest center-of-mass energies
- Running for over a decade, currently in Run 3 with 13.6 TeV collisions wrapped up for this year and Heavy Ion collisions starting
- Over 1000+ publications

CMS DETECTOR

Total weight : 14,000 tonnes
Overall diameter : 15.0 m
Overall length : 28.7 m
Magnetic field : 3.8 T

STEEL RETURN YOKE
12,500 tonnes

SILICON TRACKERS
Pixel ($100 \times 150 \mu\text{m}^2$) $\sim 1.9 \text{ m}^2$ $\sim 124\text{M}$ channels
Microstrips ($80\text{--}180 \mu\text{m}$) $\sim 200 \text{ m}^2$ $\sim 9.6\text{M}$ channels

SUPERCONDUCTING SOLENOID
Niobium titanium coil carrying $\sim 18,000 \text{ A}$

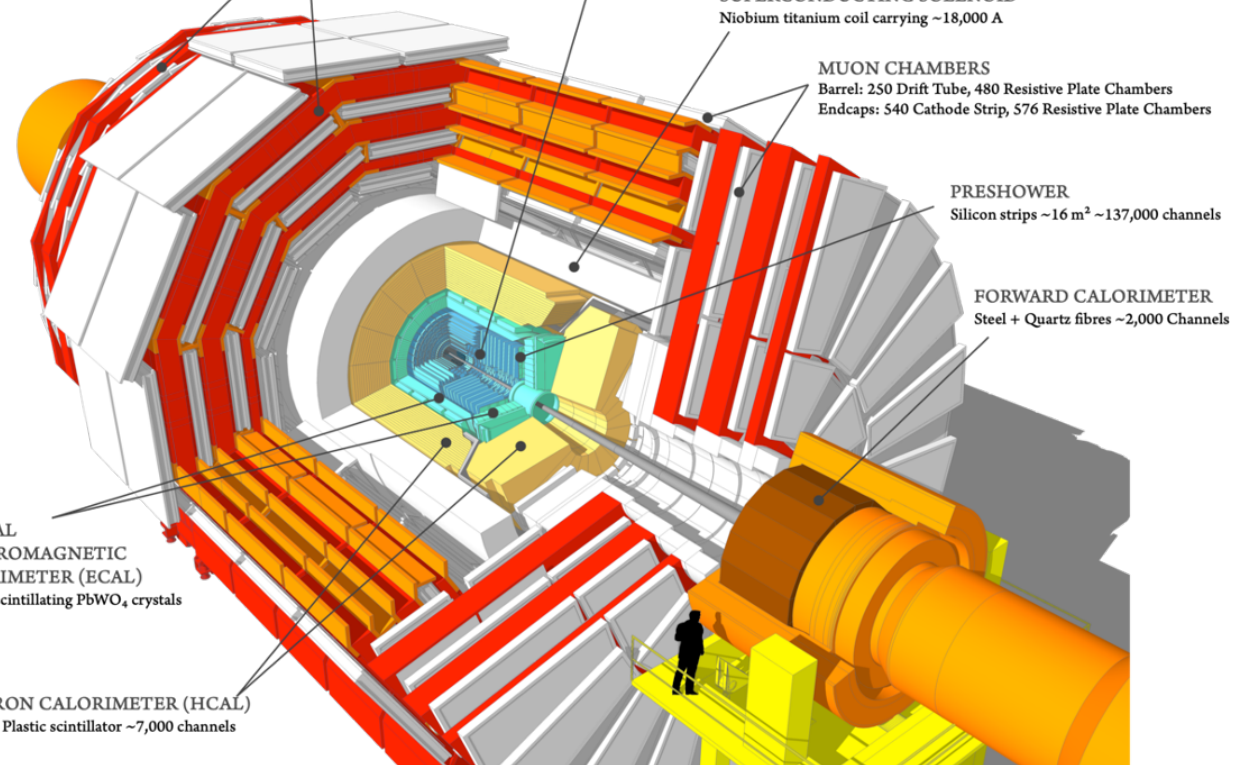
MUON CHAMBERS
Barrel: 250 Drift Tube, 480 Resistive Plate Chambers
Endcaps: 540 Cathode Strip, 576 Resistive Plate Chambers

PRESHOWER
Silicon strips $\sim 16 \text{ m}^2$ $\sim 137,000$ channels

FORWARD CALORIMETER
Steel + Quartz fibres $\sim 2,000$ Channels

CRYSTAL
ELECTROMAGNETIC
CALORIMETER (ECAL)
 $\sim 76,000$ scintillating PbWO_4 crystals

HADRON CALORIMETER (HCAL)
Brass + Plastic scintillator $\sim 7,000$ channels



The CMS Collaboration



3394

PHYSICISTS
(1228 STUDENTS)

1102

ENGINEERS

282

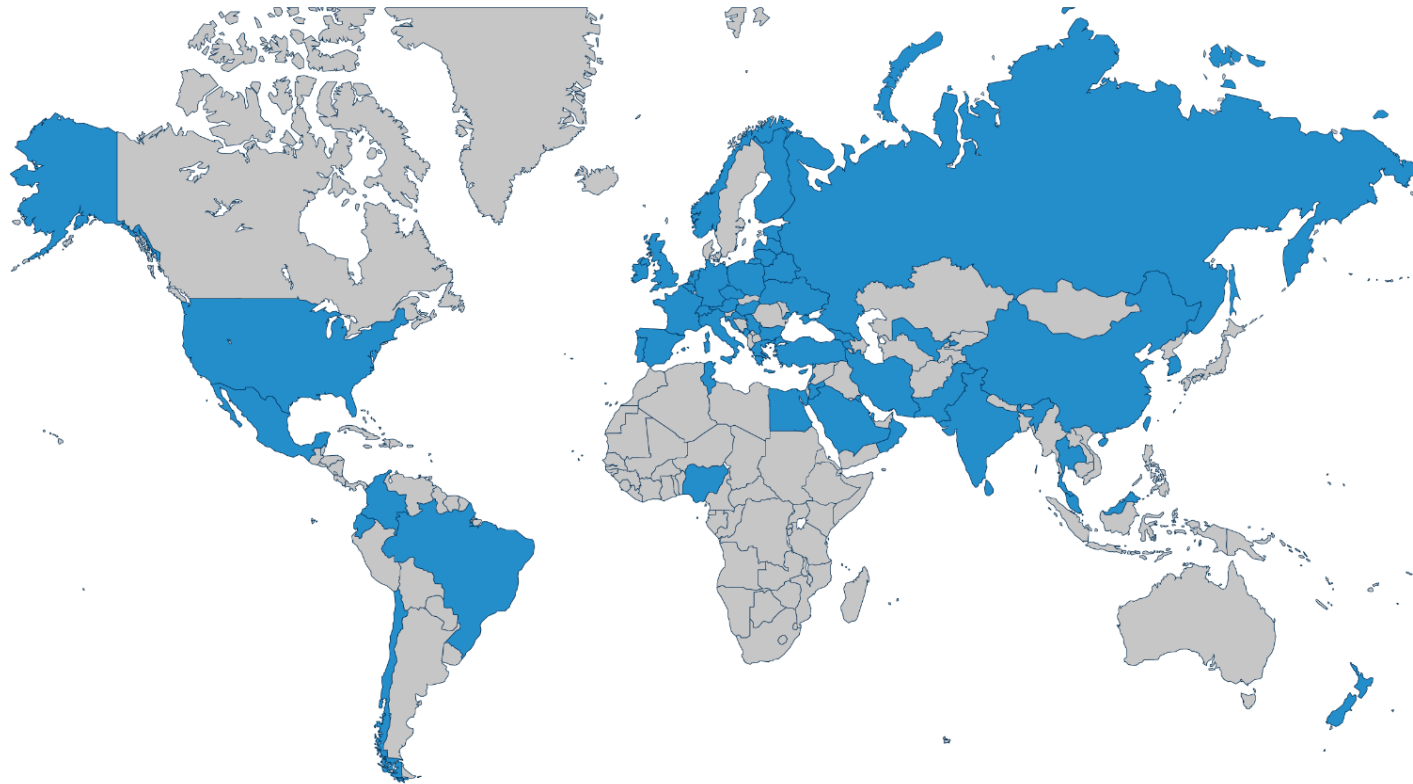
TECHNICIANS

247

INSTITUTES

57

COUNTRIES &
REGIONS



- Large span of timezones (19hrs!)
- Large number of newcomers every year

Data Analysis Training in CMS

- CMS Induction Courses (at CERN, once or twice a year since 2014, so far 16, next one Feb 5-7, 2025)
 - Newcomers introduction to the experiment and the collaboration
 - Talks/discussions (from Physics/Data Analysis to Tracking/Calorimetry/Trigger/Machine Learning/Offline and Computing, from Data Taking/Run Coordination to the various committees/institutions within CMS, Diversity, Collaboration Board, Secretariat, Communication, CMS visits etc.)
 - Split in two days, currently hybrid (in-person+Zoom), fully recorded
- **CMS Data Analysis Schools [DAS]** (2 or 3 a year since 2011, so far 30!)
 - The prototype of Data Analysis Training in CMS
 - Taking newcomers from zero to a full physics analysis
 - From basics computing skills, through all ingredients for analysis, to actually participate in a physics analysis exercise, presenting results
- **Hands-on Advanced Tutorial Sessions [HATS]** (15/20 a year since 2013)
 - Focused tutorials on specific physics objects or tools for data analysis
 - Offered at Fermilab LPC (LHC Physics Center) in Spring/Summer

Data Analysis Training in CMS

- CMS Physics Object Schools [CMSPOS] (only three times)
 - Emphasis on training new contributors to Physics Object Groups (POGs, B-tagging, Tracking, EGamma, Muon etc) and Detector Performance Groups (DPGs, Tracker, ECAL, HCAL, CSC, DT, RPC, etc)
- CMS Upgrade Schools [CUPS] (only once)
 - Emphasis on training new contributors to Phase II/ HL LHC upgrades (test-beam, sensor characterization, thermal/mechanical lab measurements, tuning of operational parameters, test DAQ systems, design of tracking/muon systems)
- Graduate-level advanced courses (1 or 2 a year)
 - Hybrid (in person+remote) lectures with homework and exams
 - Wide range of analysis related topics: Computational Physics, Statistics, Detectors, Machine Learning, etc.
- Dedicated workshops/hackathons
 - Driven by DPG or POG or other groups in preparation to major developments or at critical times (before data-taking, after release of new tools etc)
 - Statistical tools, GPU/heterogeneous computing, Trigger, Data Quality Monitoring, Machine Learning

CMS Data Analysis Schools (DAS)



- Concept started before actual data taking started, but first official CMSDAS was in January 2011
- Emphasis on teamwork and hands-on learning
- Full coverage of all data analysis aspects
- Evolved through the years and adapted to the Covid-19 era
- **Offered yearly in January at Fermilab LPC**, usually a couple more per year in other locations (Pisa, Taipei, DESY, Kolkata, Bari, Daegu, Beijing, CERN)
- Typical attendance 50 to 70 students (from undergraduates to faculty) and around 50 facilitators
- Major logistical and organizational effort



XXIX CMS Data Analysis School
Jan 8-12, 2024 at the Fermilab LPC

 “13 TeV Physics” 

A school designed to teach CMS members how to perform data analyses with the CMS analysis software.

Hands-on exercises will cover all physics objects, trigger, visualization, statistics and participants will engage in full-fledged physics analyses with CMS 13 TeV collision data, ranging from standard model measurements to searches for physics beyond the standard model.

<https://indico.cern.ch/e/cmsdas2024>

Organizing Committee:
Gabriele Benelli (Brown)
Kevin Black (Wisconsin)
Bo Jayatilaka (Fermilab)
Marguerite Tonjes (UIC)
David Yu (Nebraska)

CMS Schools Committee:
Lothar Bauerdick (FNAL)
Gabriele Benelli (co-chair, Brown)
Nitish Dhingra (Chandigarh, India)
Alexander Grohsjean (DESY)
Mohsen Khakzad (IPM, Tehran)
Santeri Laurila (CERN)
Qiang Li (Peking University)
Sudhir Malik (co-chair, UPRM)
Kajari Mazumdar (TIFR, Mumbai)
Aruna Nayak (NISER)
Andre Sznajder (UERJ, Brazil)
Michael Tytgat (Vrije Universiteit Brussel)
Jian Wang (FSU)

Administrative Support:
Carrle Farver, Terry Grozis, Frankie Kelly, Terry Read



CMS Data Analysis Schools (DAS)



- Collaborative and active learning



- From our Welcome to CMS DAS presentation:
- “A high-intensity, engaging camp, that will charge you up, change your mindset, and give you a reason to engage in cutting edge basic research and continue with it happily ever after (we hope)!”

CMS Data Analysis Schools (DAS)



- Over the years the structure has evolved slightly, but this is the basic structure:
 - Pre-exercises
 - Mandatory with deadline before start of the school
 - Plenary Lectures
 - Introduction to CMS Physics, LHC, CMS Detector, Software/ Analysis Tools, Diversity and Inclusion, Communications/ Outreach
 - Short Exercises
 - Covering objects and basic analysis ingredients
 - Writers PUB
 - Going over the publication process of a paper in CMS
 - Long Exercises
 - A set of complete physics analysis exercises
 - Mini-Symposium
 - Presentation of the results from all teams



CMS Data Analysis Schools (DAS)



- Networking is a key ingredient of CMS DAS:
 - All participants are organized in small teams, based on the long exercise they are assigned to

January 4					January 10	January 10-14	
Short Exercises					Long Exercise Assignment		
Period 1	Period 2	Period 3	Period 4	Period 5	Writers PUB		
Tuesday 10:45->11:30	Tuesday 11:35->12:20	Tuesday 13:20->14:05	Tuesday 14:10->14:55	Tuesday 15:10->15:55	Monday 19:00->20:45	Mon 11:00->Fri 12:00	7 short exercises
Statistics	Visualization	Tagging	Generators	Machine Learning	Sarah Eno	Double Higgs to 4b final state	Jets
Statistics	Visualization	Tagging	Generators	Machine Learning	Sarah Eno	Double Higgs to 4b final state	Statistics
Statistics	Visualization	Tagging	Generators	Machine Learning	Sarah Eno	Double Higgs to 4b final state	Tagging
Statistics	Visualization	Tagging	Generators	Machine Learning	Sarah Eno	Double Higgs to 4b final state	Machine Learning
Statistics	Visualization	Tagging	Generators	Machine Learning	Sarah Eno	Double Higgs to 4b final state	Generators
Statistics	Visualization	Tagging	Generators	Machine Learning	Sarah Eno	Double Higgs to 4b final state	Visualization
Statistics	PU/MET	Tagging	Jets	Machine Learning	Sarah Eno	Double Higgs to 4b final state	PU/MET
Statistics	PU/MET	Tagging	Jets	Machine Learning	Sarah Eno	Double Higgs to 4b final state	
Statistics	PU/MET	Tagging	Jets	Machine Learning	Sarah Eno	Double Higgs to 4b final state	
Statistics	PU/MET	Tagging	Jets	Machine Learning	Sarah Eno	Double Higgs to 4b final state	
Statistics	PU/MET	Tagging	Jets	Machine Learning	Sarah Eno	Double Higgs to 4b final state	

- Their commitment and responsibility to the team (and not to facilitators or organizers) is the key to each team (and team member) success
- While all team members participate to the same long exercise, team members are assigned different short exercises so that the team as a whole can have full coverage of the needed tools for the long exercise analysis
- All team members participate both in the slides preparation and in the final presentation at the Mini-Symposium at the end



CMS DAS during the pandemic

- Evolution during Covid-19:
 - Anticipated the Pre-Exercises release (essential to make sure people can hit the ground running)
 - Mattermost support
 - Factored out some short exercises as Offline short exercises (ROOT, NanoAOD, PPD, Luminosity), fully remote some with video recordings, all with Mattermost support
 - Stretched the length of DAS to two full weeks
 - First week kick-off and asynchronous short exercises
 - Introductions, team building, plenaries
 - Short Exercises kick-off (live)
 - Asynchronous work through material with Mattermost support (some with extra live office hours)
 - Short Exercises wrap-up (live)
 - Second week more plenaries and live long exercises
 - Recordings of all live sessions
 - Heavy use of Mattermost understanding about support expectations from facilitators
 - Release logistical constraints
 - Social events replacement (Scribble I/O, Rubik's cube competition etc)

CMS DAS during the pandemic

- In 2022, for example, we offered 13 short exercises:

January 4					January 12
Short Exercises					
Period 1	Period 2	Period 3	Period 4	Period 5	Writers PUB
Tuesday 10:45->11:30	Tuesday 11:35->12:20	Tuesday 13:20->14:05	Tuesday 14:10->14:55	Tuesday 15:10->15:55	Tuesday 19:00->20:00
Statistics	Forward Protons	Tagging	Jets	Machine Learning	Sarah Eno
Tracking/Vertexing	PU/MET	Muon	Generators	Electron/photon	Jacobo Konigsberg
Triggers	Visualization	Tau			

- And 8 long exercises (which resulted in as many teams):

Double Higgs to 4b final state
B2G Search for $b^* \rightarrow tW$ all hadronic
TTGamma cross-section
Exclusive production of lepton pairs
Top quark mass at 13 TeV
SUSY hadronic with top tagging
Contact Interactions
Z \rightarrow tau tau cross-section at 13TeV

Post-pandemic CMS DAS

- Last year, we offered 12 short exercises

January 8		January 9			
Short Exercises					
Period 1	Period 2	Period 3	Period 4	Period 5	Writers PUB
Monday 1:30->3:30pm	Monday 4:00->6:00pm	Tuesday 11:00am->1:00pm	Tuesday 2:00->4:00pm	Tuesday 4:30->6:30pm	Tuesday 7:45->9:00pm
Statistics	Statistics	Statistics	Statistics	Tagging	Nadja Strobbe
Jets	Jets	Common Analysis Tools	Machine Learning	Machine Learning	Keith Ulmer
Tracking/Vertexing	Electron/photon	PU/MET	Tau	Electron/photon	Joel Butler
Triggers	Generators	Muon	Muon	Common Analysis Tools	Sridhara Dasu

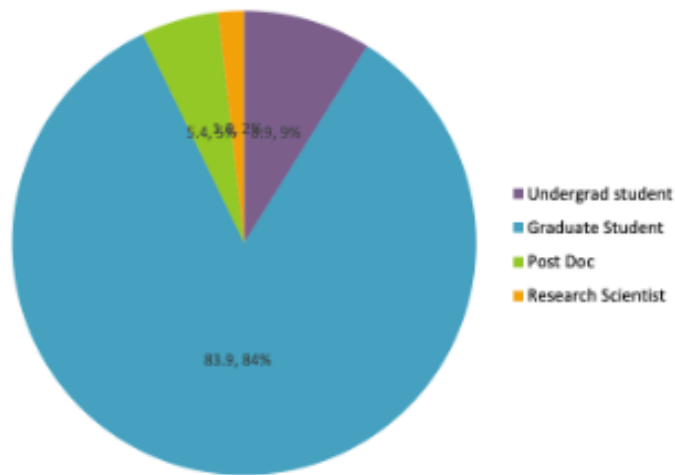
- And 5 long exercises:

Search for structures near the J/psiJ/psi mass threshold
B2G Search for $b^* \rightarrow tW$ all hadronic
TTGamma cross-section
Search for long-lived particles with muon detector shower
Z \rightarrow tau tau cross-section at 13TeV

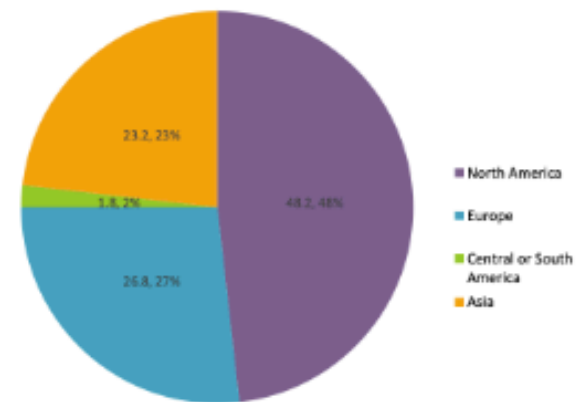
CMS DAS during the pandemic

- Feedback

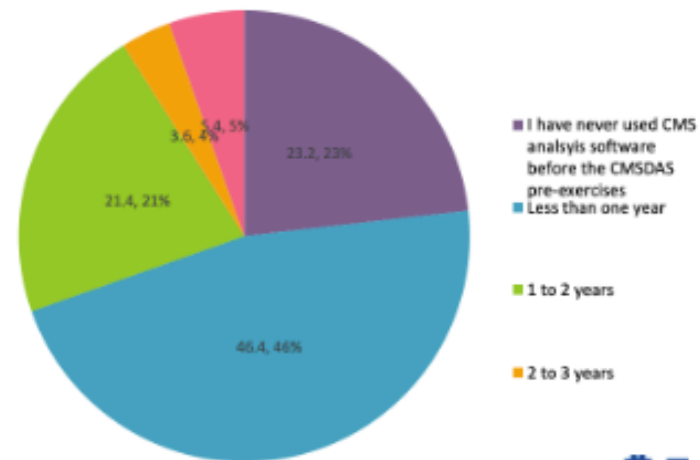
CMS DAS participation level



In which region is your institute located?



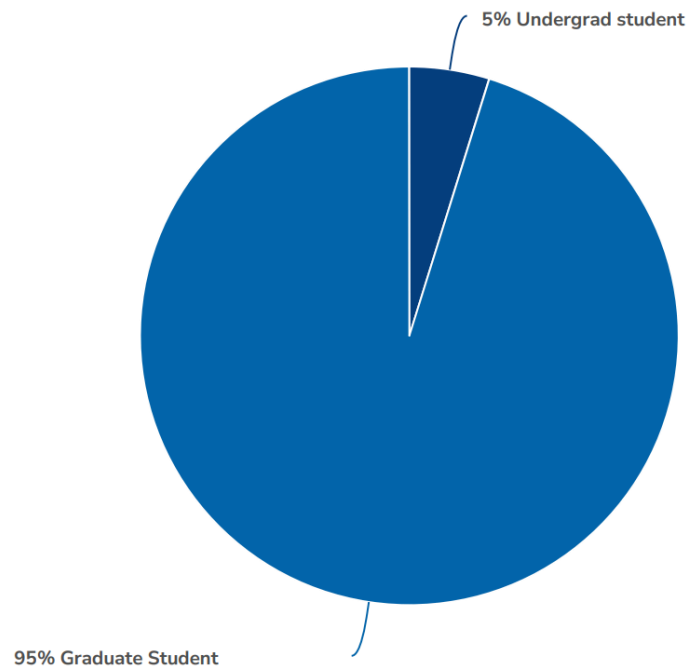
For how many years have you been using CMS software?



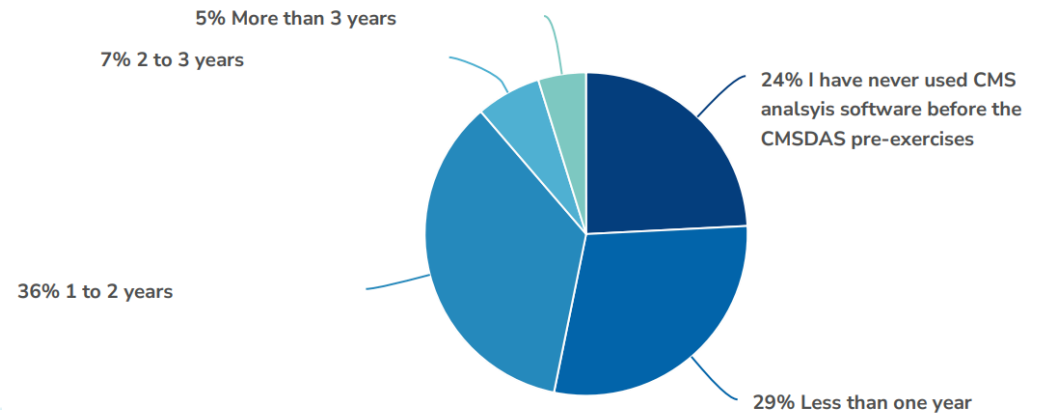
CMS DAS @LPC 2024

- Feedback

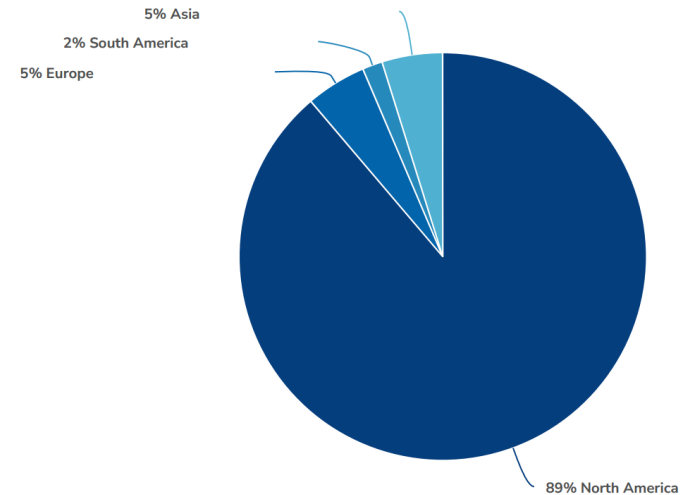
CMS DAS participation level



For how many years have you been using CMS software?



In which region is your institute located?

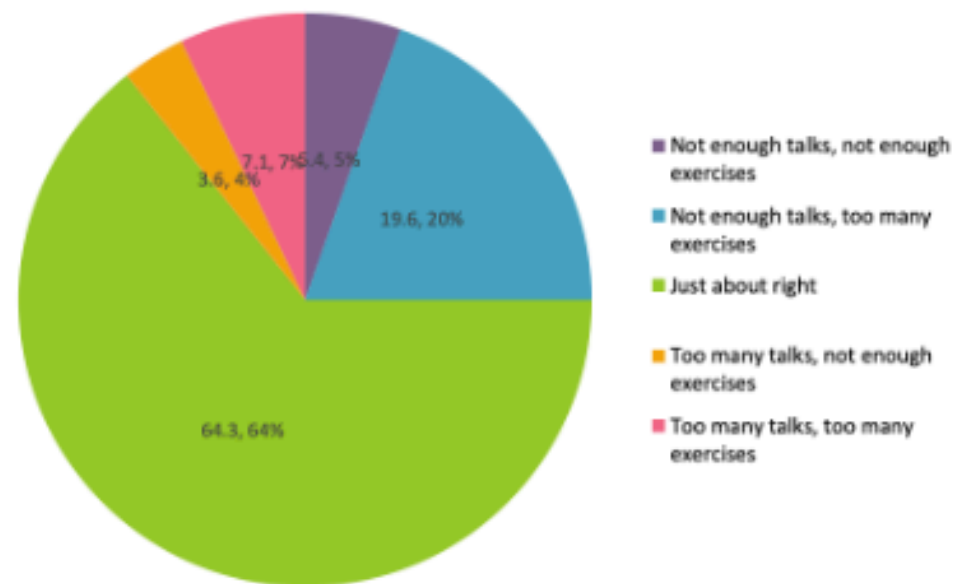
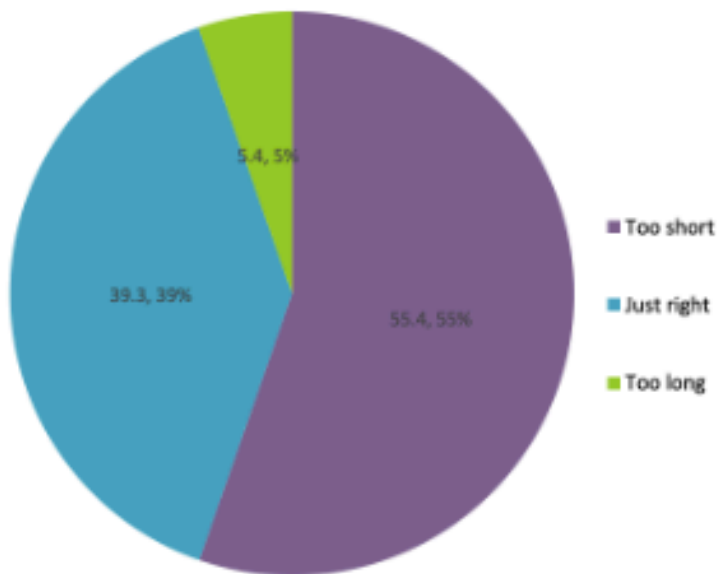


CMS DAS during the pandemic

- Feedback

Is the two weeks (including the asynchronous time) the right length of time for a virtual CMSDAS?

The number of talks and number of exercises at CMSDAS were

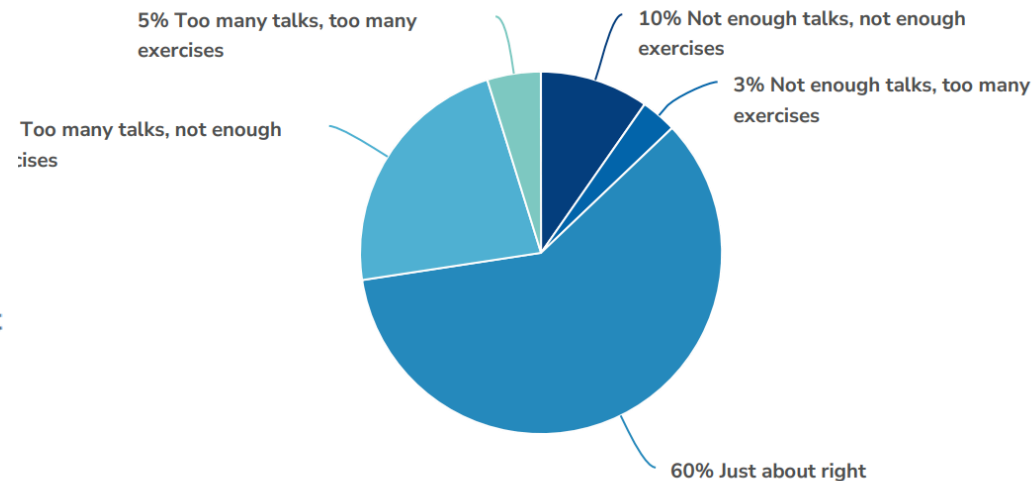
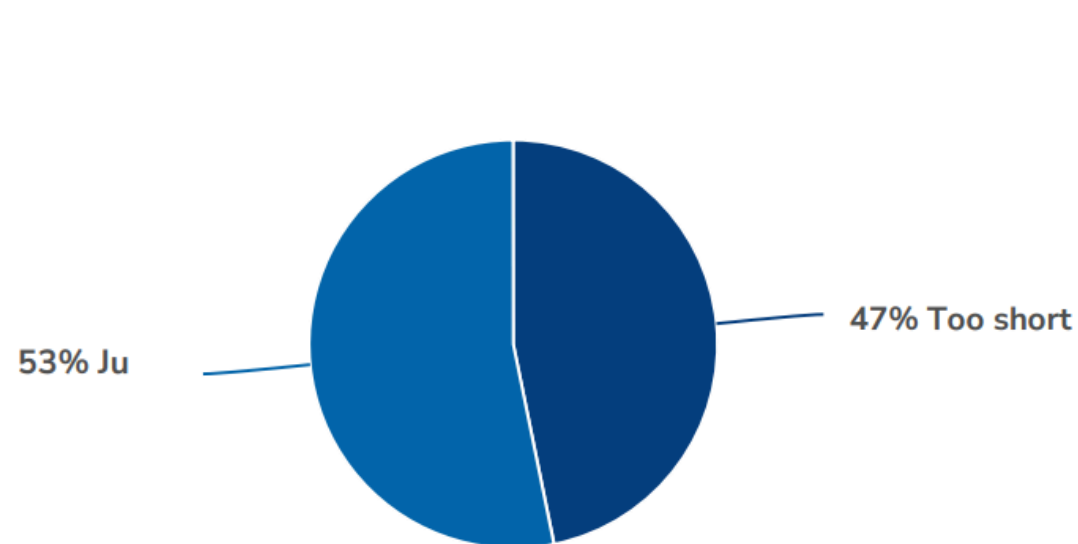


CMS DAS @LPC 2025

- Feedback

Is the two weeks (including the asynchronous time) the right length of time for a virtual CMSDAS?

The number of talks and number of exercises at CMSDAS were



CMS DAS



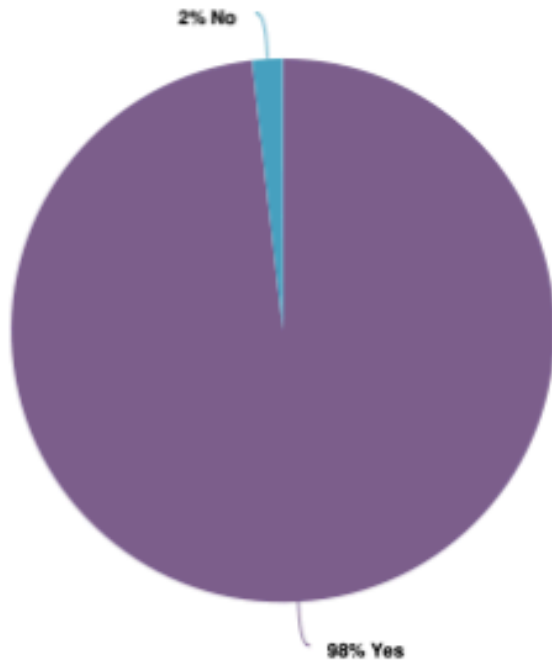
- Feedback

	Timing	Content/Support		★★★★☆ Count: 53 Not Applicable: 2	★★★★★ Count: 52 Not Applicable: 3
Kick-off day (Tue Jan 4)	★★★★☆ Count: 48 Not Applicable: 6	★★★★★ Count: 47 Not Applicable: 5	Short Exercise Wrap-up (Fri Jan 7)		
Plenaries (Tue Jan 4, Fri Jan 7, Mon Jan 10, Tue Jan 11)	★★★★☆ Count: 51 Not Applicable: 4	★★★★★ Count: 54 Not Applicable: 1	Writers PUB (Mon Jan 10)	★★★★☆ Count: 40 Not Applicable: 13	★★★★★ Count: 36 Not Applicable: 18
Long Exercise Meet up (Tue Jan 4)	★★★★☆ Count: 52 Not Applicable: 3	★★★★★ Count: 52 Not Applicable: 3	Communication and Outreach (Tue Jan 12)	★★★★☆ Count: 48 Not Applicable: 5	★★★★★ Count: 50 Not Applicable: 5
Short Exercise Kick-off (Tue Jan 4)	★★★★☆ Count: 48 Not Applicable: 4	★★★★☆ Count: 50 Not Applicable: 4	Optional Social events (Wed Jan 12, Thurs Jan 13)	★★★★☆ Count: 39 Not Applicable: 14	★★★★☆ Count: 41 Not Applicable: 14
Asynchronous short exercises (Tue - Wed Jan 5-6)	★★★★☆ Count: 53 Not Applicable: 1	★★★★☆ Count: 55 Not Applicable: 1	Long Exercises (Mon-Thurs Jan 10-13)	★★★★☆ Count: 54 Not Applicable: 0	★★★★☆ Count: 55 Not Applicable: 0
Offline Short Exercises (released Mon Dec 13)	★★★★☆ Count: 43 Not Applicable: 11	★★★★☆ Count: 44 Not Applicable: 11			

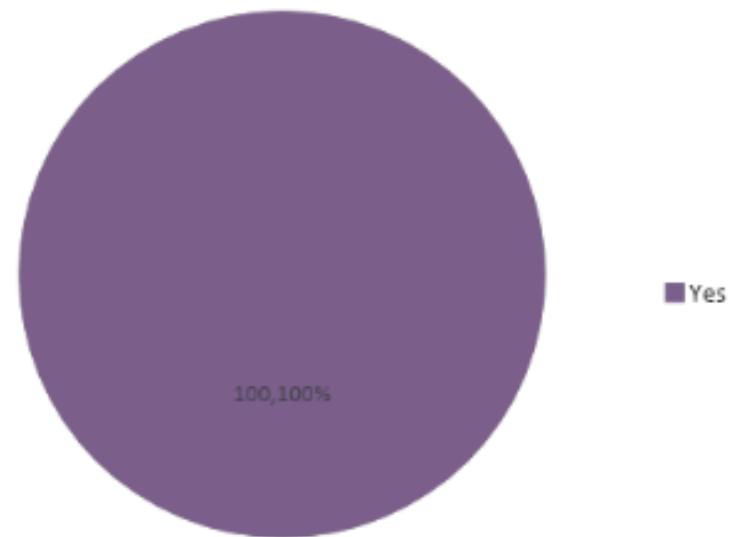
CMS DAS during the pandemic

- Feedback

Has CMS DAS been a valuable experience for you?



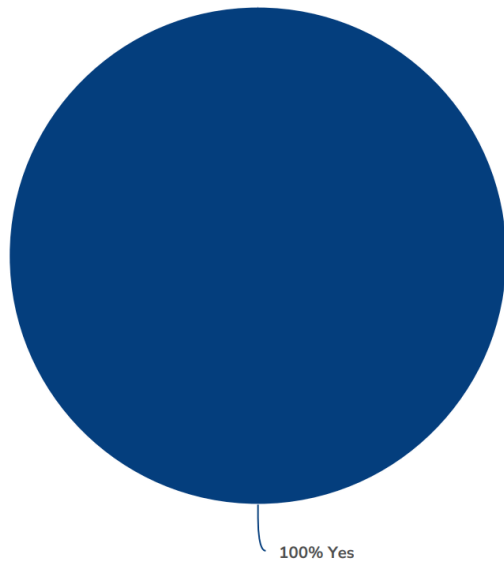
Would you recommend CMSDAS to your colleagues?



CMS DAS @LPC 2024

- Feedback

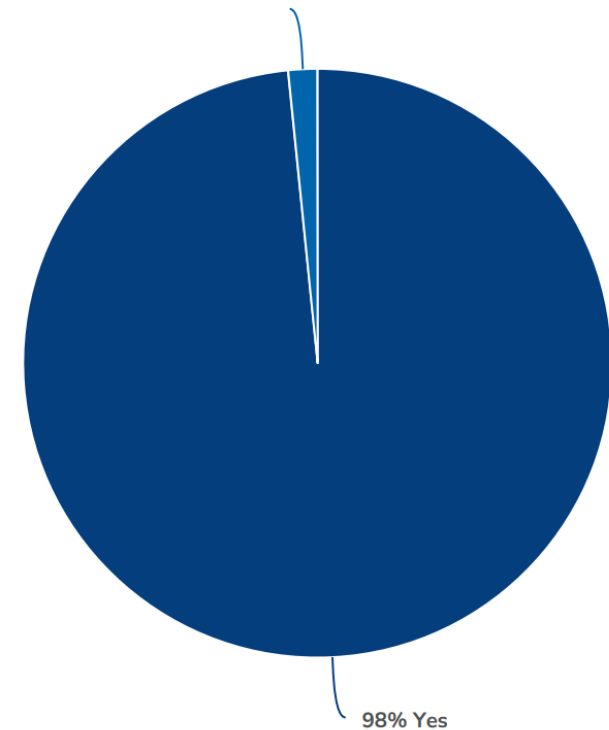
Has CMS DAS been a valuable experience for you?



100% Yes

Would you recommend CMSDAS to your colleagues?

2% No



98% Yes

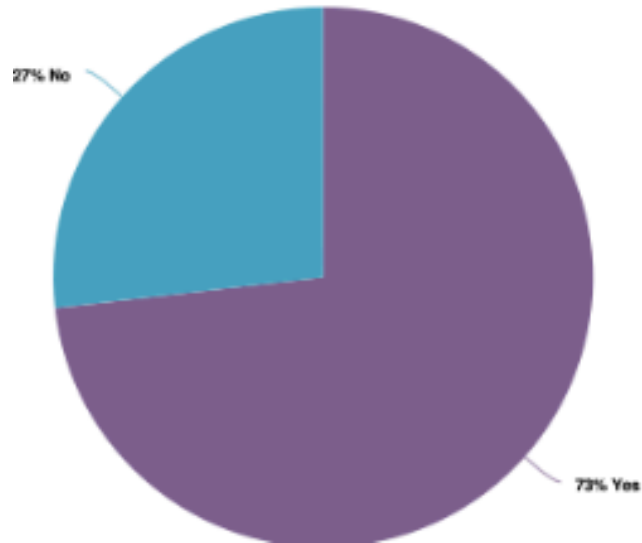
Value	Percent	Responses
Yes	100.0%	62

Totals: 62

CMS DAS during the pandemic

- Feedback

Has CMSDAS enabled you to make new connections to CMS members or groups that you envisage will be helpful in your future analysis and other work for CMS?



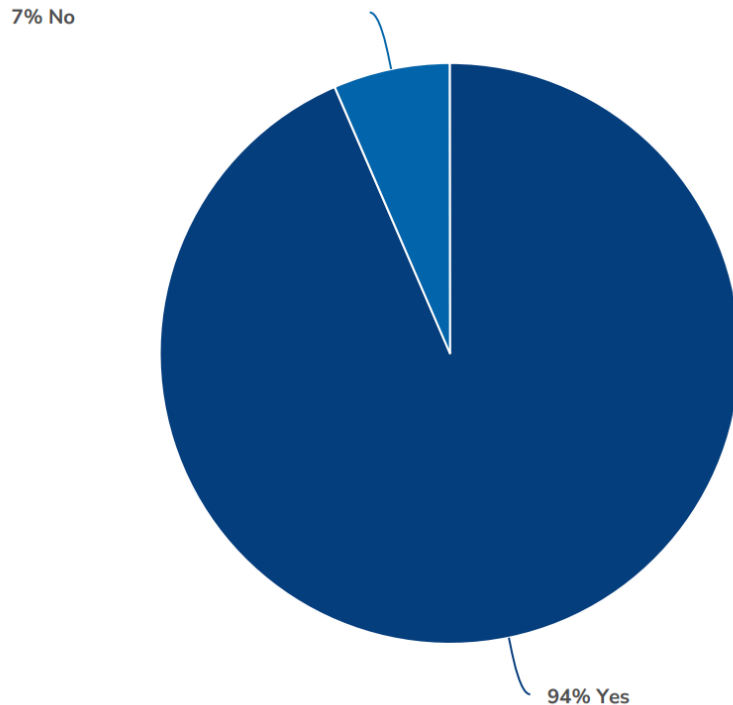
Would you be willing to serve as a facilitator for future schools?



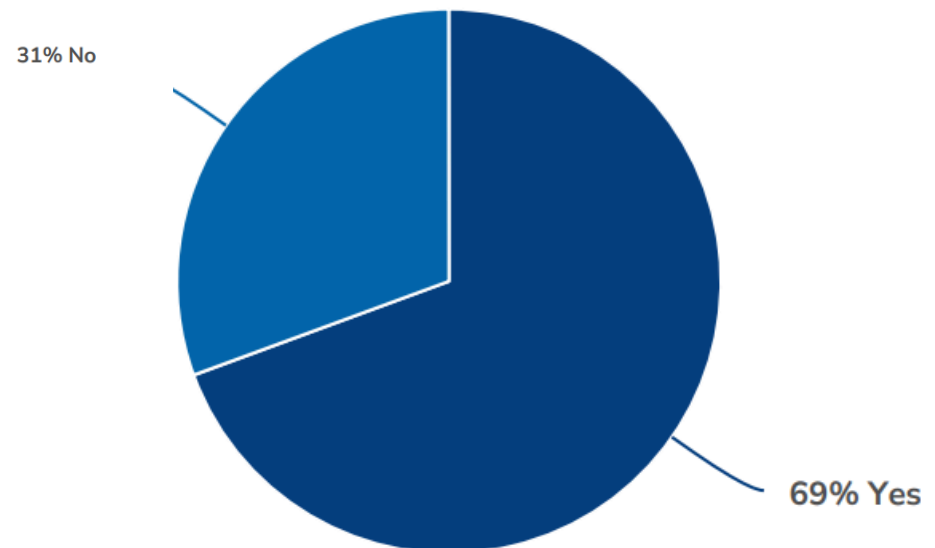
CMS DAS @LPC 2024

- Feedback

Has CMSDAS enabled you to make new connections to CMS members or groups that you envisage will be helpful in your future analysis and other work for CMS?



Would you be willing to serve as a facilitator for future schools?





Hands-on Advanced Tutorial Sessions (HATS)



2024 HATS@LPC (and lectures)

- Particle Flow LEGO HATS (Sep 11-12, 2024) - [Indico Agenda](#)
- Machine Learning HATS (Sep 9, 2024) - [Indico Agenda](#)
- Tagging HATS (Aug 21, 2024) - [Indico Agenda](#)
- Jet Energy Corrections/Resolution and Substructure (Jets II) HATS (Aug 20, 2024) - [Indico Agenda](#)
- Jet Algorithms and Pile-up reweighting/mitigation (Jets I) HATS (Aug 19, 2024) - [Indico Agenda](#)
- Trigger HATS (Aug 16, 2024) - [Indico Agenda](#)
- Electrons and Photons HATS (Aug 13, 2024) - [Indico Agenda](#)
- Muons HATS (Aug 6, 2024) - [Indico Agenda](#)
- CMS Combine HATS (July 23, 2024) - [Indico Agenda](#)
- Matplotlib for HEP HATS@LPC (July 22, 2024) - [Indico Agenda](#)
- Coffea (Columnar Analysis Tools) HATS@LPC (July 19, 2024) - [Indico Agenda](#)
- Dask (Effective Scale Out Techniques) HATS@LPC (July 18, 2024) - [Indico Agenda](#)
- Awkward Array (and uproot) for columnar analysis HATS@LPC (July 17, 2024) - [Indico Agenda](#)
- Containers HATS@LPC (July 16, 2024) - [Indico Agenda](#)
- Git/GitHub HATS@LPC (July 15, 2024) - [Indico Agenda](#)
- "High Energy Physics: Detectors Past, Present, Future" Fall 2024 Course (August 27, 2024 - November 26, 2024, Tuesdays 10:00-11:30am; Zoom) - [Indico Agenda](#) (CMS CERN account required)

2023 HATS@LPC (and lectures)

- Condor/CRAB/CMS Connect/Data Management Tools HATS (September 13, 2023) - [Indico Agenda](#)
- Generators HATS (September 12, 2023) - [Indico Agenda](#)
- Government Outreach HATS (September 7, 2023) - [Indico Agenda](#)
- Trigger HATS (August 18, 2023) - [Indico Agenda](#)
- MET HATS (August 17, 2023) - [Indico Agenda](#)
- Machine Learning HATS (August 16, 2023) - [Indico Agenda](#)
- Jets II HATS (August 15, 2023) - [Indico Agenda](#)
- Jets I HATS (August 14, 2023) - [Indico Agenda](#)
- Tau HATS (August 11, 2023) - [Indico Agenda](#)
- CMS Combine HATS (August 10, 2023) - [Indico Agenda](#)
- Particle Flow Lego HATS@LPC (July 11-12, 2023) (Fermilab in person only)- [Indico Agenda](#)
- Effective Scale Out Techniques HATS@LPC (July 10, 2023) - [Indico Agenda](#)
- Columnar Analysis Tools HATS@LPC (July 7, 2023) - [Indico Agenda](#)
- Uproot and Awkward Array for columnar analysis HATS@LPC (July 6, 2023) - [Indico Agenda](#)
- Git/GitHub HATS@LPC (July 5, 2023) - [Indico Agenda](#)
- Computational Physics 1 Fall 2022 & 2 Spring 2023 (August 29, 2022 - May 21, 2023, Mondays, Wednesdays & Fridays 1pm-1:50pm; Zoom) - [Indico Agenda](#) (CMS CERN account required)



Hands-on Advanced Tutorial Sessions (HATS)



- Hands-on learning, typically with some pre-requisite training:
 - Partly lectures
 - Partly exercises
- Format has been evolving to adapt to audience in Covid-19 times, strictly connected to DAS
- All recorded (videos.cern.ch), with transcripts



edit

RN OPEN-VIDEO-2022-200

Git/GitHub HATS@LPC2022 Zoom Recordings

Remote Live Zoom recordings for Git/GitHub HATS@LPC 2022. Please find all the material and links on the indico agenda at: <https://indico.cern.ch/e/githubhats2022>

Videos

Git/GitHub HATS@LPC2022 Recording 1/3

Git/GitHub HATS@LPC2022 Recording 2/3

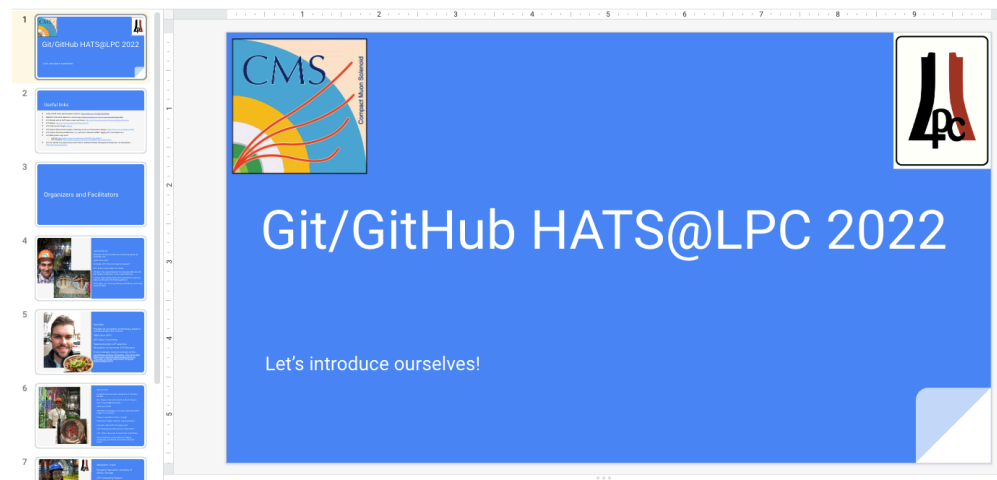
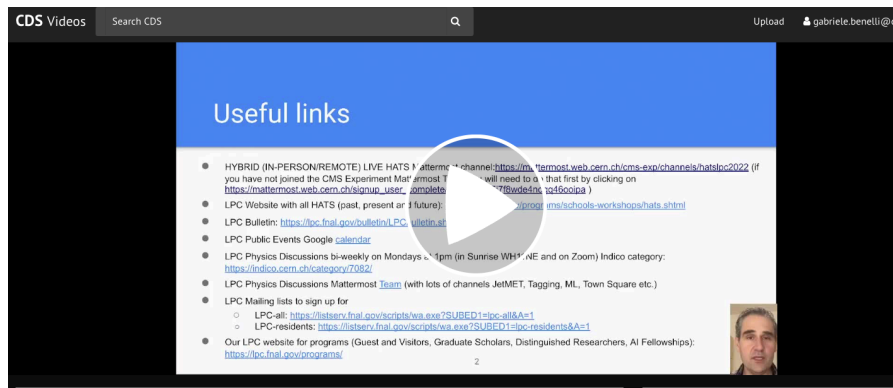
Git/GitHub HATS@LPC2022 Recording 3/3

HATS, HATS@LPC, Git, GitHub

Hands-on Advanced Tutorial Sessions (HATS)



- Covid optimizations:
 - General remarks packaged into a video/set of slides



- Google Slides for introductions
- SliDo/Google Forms/Mattermost polling for interactive polling
- Use of Carpentries format for pre-exercises simplifying maintenance (also for some of the short exercises and HATS)
- Challenge
 - Zoom fatigue, interactions, networking, collaboration spirit and continued interactions

HATS

- Special issues during the pandemic prevented the fully hands-on ones:
 - Particle Flow Lego HATS
 - Use a Lego physical model to do Particle Flow reconstruction by hand and identify an event
 - CMS Upgrade Detectors HATS
 - Experience several hardware hands-on exercises involving oscilloscopes etc (characterization of sensors, production of scintillating fibers, testing of readout electronics)
- Ideas to improve/try/add
 - Luminosity HATS
 - Higgs Combine (statistics tool) HATS
 - Generic Python/Jupyter (synergy with HSF curriculum)
 - Move to Software carpentry model for all to simplify maintenance
 - Automate some of the logistics
- Since 2023 we've been back in-person, allowing a substantial remote participation

Reflections

Brainstorming and Networking

- The value of in-person interaction is made more evident by its lack in the pandemic of years:
- Brainstorming and talking among colleagues in an informal setting where multiple conversations can happen in parallel and there is no meeting goal (having a coffee/tea, or eating lunch) is the seed of many important developments and ideas and it increases dramatically the efficiency of our collaborations
- Networking is a major benefit of in-person interactions, as one of the most important results of participating to DAS and HATS is to know who to ask questions about specific topics, mapping out a network within the collaboration
- Zoom and Mattermost while extremely precious, are not ideal to facilitate those kind of conversations
- The worst damage is to people who do not know what they lost with fully remote collaboration: newcomers who do not have a networking social capital as they navigate these hybrid and remote times.
- The interactions with facilitators, POG/PAG conveners for a remote center like the LPC are fundamental in establishing synergies and collaborations that extend beyond the training, a major benefit (including shifts-taking etc)

Advanced Training for Developers

- A personal opinion on the concept of trainings for developers:
 - While covering the bases with technical knowledge is important (and the previous trainings provide those), ultimately “advanced” becomes very quickly very specific and sitting next to some expert to understand the scope of the problem, draft a solution and interact repeatedly with a few people to get feedback is what is needed to get to the last step and become a developer.
 - Without a specific project the training is not very effective risking to drain more effort from experts than real gain by new experts in training
 - The practical limit to what can be reasonably packaged/worked on as training for developers should be determined by the typical questions/issues that arise most frequently in developers communication channels

More Ideas for improvement

- Interaction with Physics Object Groups and Physics Analysis Groups to be strengthened to get more facilitators, and drive the content relevance
 - Analysis tools evolution
 - Data format evolution
 - Heterogenous resources
- Rely on the lesson learned to be more inclusive (time of day, asynchronous plus synchronous, asynchronous support on Mattermost)
 - But also encourage more in-person participation
 - More regional duplication
- Ideas about collaborating with tools like GatherTown (office hours?)
- Streamline the core curriculum

Reflections

- Challenges
 - Documentation is currently heterogeneous (transitioning to alternatives to traditional Twikis)-> Software Carpentry
 - Training organization is extremely time consuming, lots of communications necessary
 - People's engagement online
 - Getting back to in-person interactions (in a hybrid scenario) after a long time of fully remote operations
- Reasons for hope
 - All of the work described is based on the good will and the collaboration spirit embodied by facilitators and organizers
 - Training, similarly to operations, highlights the collaborative nature of our endeavor, seeding the future of our experiments
 - Attention to diversity and inclusion is growing and definitely training is reflecting it, being much more accessible than in the past, reaching people previously not served.

Reflections



Back-up