# Physics analysis workflows and pipelines for the HL-LHC

**Alexander Held[1]**, Oksana Shadura[2]

[1] University of Wisconsin–Madison

[2] University of Nebraska–Lincoln

*ISGC 2023*
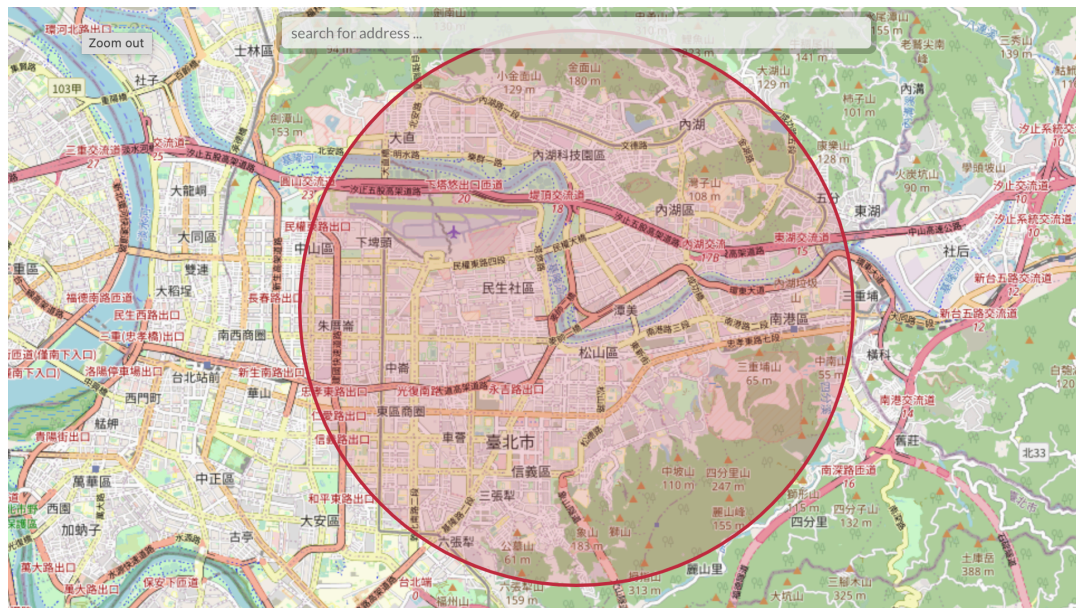
https://indico4.twgrid.org/event/25/

March 23, 2023

# The Large Hadron Collider

- The **Large Hadron Collider (LHC)** collides protons with 13.6 TeV in a 27 km tunnel, 100 m underground

- **Collisions recorded** by multiple detectors and (after a lot of processing) available to physicists as **columnar data**

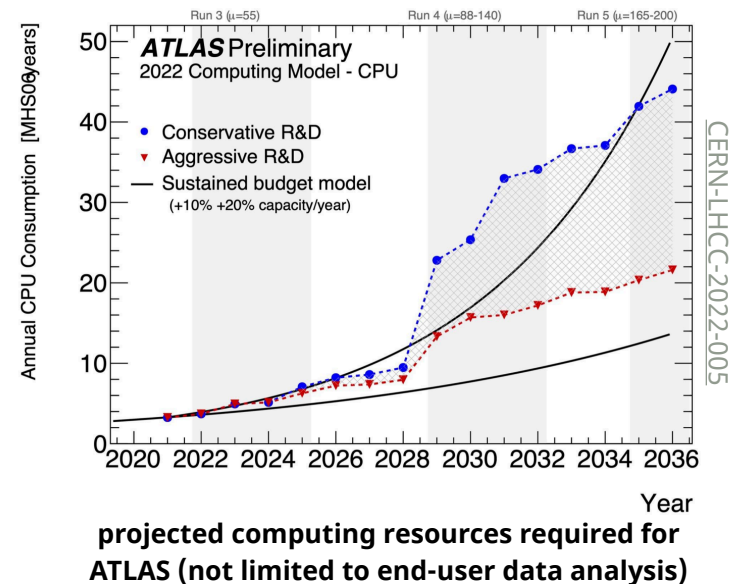  - 1 row per collision event, filled with nested data characterizing collision

**LHC ring overlaid over Taipei**



https://natronics.github.io/science-hack-day-2014/lhc-map/

# Data analysis at the LHC and the HL-LHC

- Focusing here on the **final steps of data analysis**

  ‣ physicists turning columnar data into results ready for publication

    - input: nested data structure per collision (~billions of rows)

    - output: results of statistical inference, figures, tables, …

- The upcoming **High Luminosity LHC** poses **computational challenges**

  ‣ significant data volume increases

  ‣ R&D required to scale to the data analysis demands

**projected computing resources required for ATLAS (not limited to end-user data analysis)**
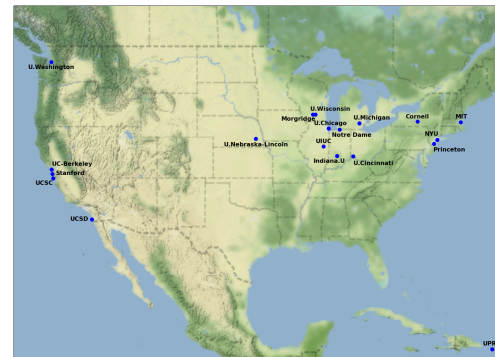
# The Analysis Grand Challenge (AGC)

• The **"Analysis Grand Challenge" (AGC)** aims to help **address the computing challenges** of the HL-LHC

• The AGC has **two components**

    1. define a physics analysis task of realistic scope & scale

    2. develop an analysis pipeline that implements the task

      - find & address performance bottlenecks & usability concerns

# IRIS-HEP and the Analysis Grand Challenge



- **IRIS-HEP**: *"Institute for Research and Innovation in Software for High Energy Physics"*

  ‣ software institute funded by the US National Science Foundation

  ‣ research & development for the HL-LHC

    - innovative algorithms for data reconstruction & triggering

    - analysis systems to reduce time-to-insight and maximize physics potential

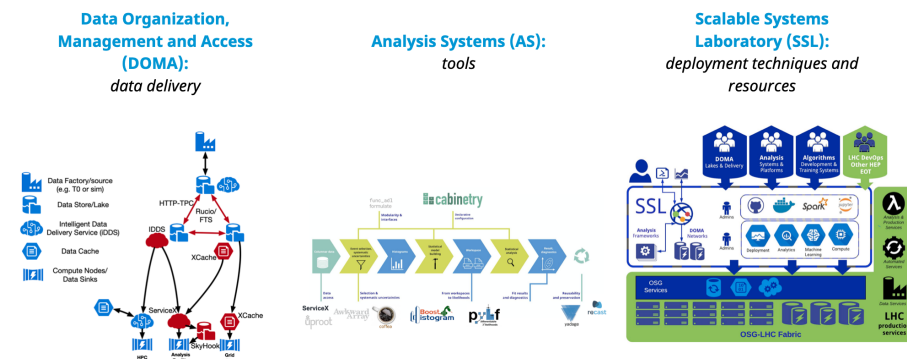    - data organization, management and access systems

  ‣ more information: https://iris-hep.org/



institutes participating in IRIS-HEP

# IRIS-HEP and the Analysis Grand Challenge

- **AGC**: *"Analysis Grand Challenge"*

  ‣ historically, an integration exercise

  - test realistic end-to-end analysis pipelines aimed at HL-LHC use

  - combine technologies being developed in various ares of IRIS-HEP & adjacent ecosystem

  - identify & address performance bottlenecks and usability issues

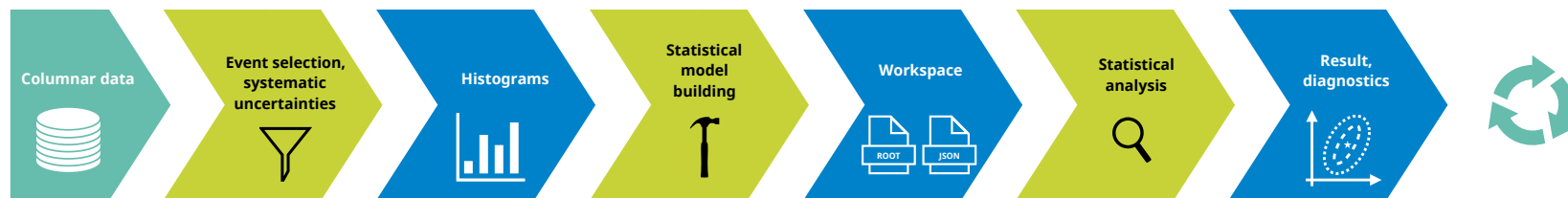  ‣ organized jointly with the US ATLAS & US CMS operations programs
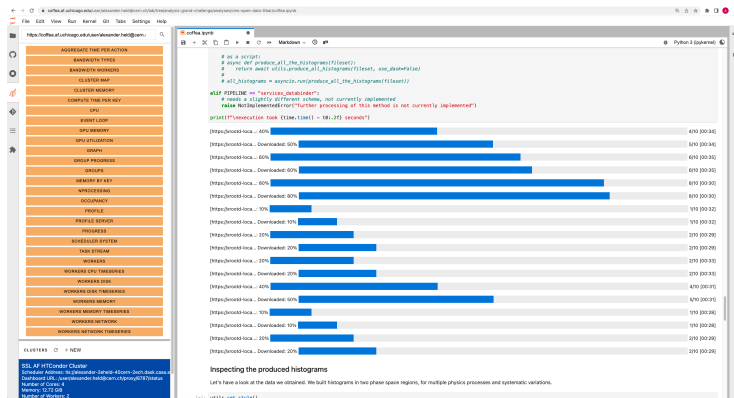


AGC combining IRIS-HEP focus areas

# "Analysis" in the AGC context

- In view of the HL-LHC: "analysis" **starts** from centrally produced **common data samples** (= big tables of information)

- Includes all **subsequent steps** to produce results needed for publication

  ‣ extract relevant data

  ‣ (re-) calibrate objects (groups of columns) & calculate systematic variations (new columns)

  ‣ filter events (rows) & calculate observables (new columns)

  ‣ histogramming (for binned analyses)

  ‣ construct statistical model + perform statistical inference

  ‣ visualize results & provide all relevant information to study analysis details

- Do all these steps in a **reproducible** way

# Moving beyond an integration exercise

- Investigating the possibility of **"interactive analysis"**: turnaround time of minutes or less

  ‣ made possible by highly parallel execution in short bursts, low latency & heavy use of caching

- We hope that the AGC can be **useful to the broader community**!

  ‣ testbed for software library development

  ‣ environment to prototype analysis workflows

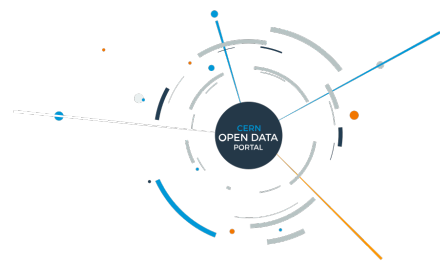  ‣ functionality & integration test for analysis facility development

interactive analysis in a notebook

AGC tools 2022 workshop

# The AGC analysis setup

- Main AGC analysis task: **ttbar cross-section measurement** in single lepton channel

  ‣ includes simple top quark reconstruction

  ‣ setup chosen as it captures relevant workflow aspects and can easily be extended

      - e.g. conversion into a beyond-the-Standard-Model search

  ‣ analysis task prominently features handling of systematic uncertainties

- Analysis is based on **Run-2 CMS Open Data** (~400 TB of data in MiniAOD format)

  ‣ Open Data is crucial: everyone can participate

  ‣ currently using 4 TB of ntuple inputs (pre-converted, ~1B events before cuts)

- Goal of setup is showing **functionality**, not discovering new physics

  ‣ want to capture *workflow*, but can use made-up tools for evaluating calibrations & systematic uncertainties



reconstructed top mass

# Systematics and other analyzer user experience aspects

- Handling **systematic uncertainties** is a **key challenge** in analysis workflows

  ‣ AGC analysis task includes different types of systematic uncertainties to mirror practical requirements

    - weight-based uncertainties

    - object-based systematic variations affecting kinematics (+ thereby event selection / observables)

    - non-histogram-based uncertainties (e.g. cross-section uncertainties)

- **Metadata** handling

  ‣ capturing various bookkeeping aspects in analysis task

- **Scale-out**: from laptop to analysis facility

  ‣ challenge: write analysis implementation that can run anywhere

**Pain points in analysis user experience, ordered**

1. **Systematics**
   ○ Recurring topic throughout this workshop: this is not solved

2. **Metadata**
   ○ Finding & handling information

3. **Scale-out**
   ○ Prototyping vs scale-out, different implementations / details on different sites
   ○ Need for consistent environments across all resources

Analysis Ecosystem Workwshop II
User experience & Declarative Languages summary

# Tools and services in our implementation

- Employing stack of **Python HEP libraries** for analysis tasks

- **ServiceX** used as data delivery service

- Execution on a **coffea-casa analysis facility**



HEP-specific libraries used for data analysis

data delivery services

optional services

# Analysis Facilities for execution

- **coffea-casa** is a **prototype analysis facility** for the HL-LHC

  ‣ interactive facility for columnar analysis providing analysis tools & scaling to computing resources

  ‣ more information: https://iris-hep.org/projects/coffea-casa.html

**Coffea-Casa**

# Implementation: ttbar analysis in a notebook

- **From data delivery to statistical inference** in a [notebook](notebook)

## multiple supported processing schemes



### reconstructed observables



### nuisance parameter pulls



### coffea processor



### systematic variations



### post-fit distributions

# Benchmarking results

- **Benchmarking AGC implementation performance** at the University of Nebraska–Lincoln CMS Tier-2

  ‣ tested various configurations of hardware, data pipeline and analysis task

  ‣ for more information, see this ACAT 2022 contribution



**good scaling to hundreds of cores**



**efficient resource usage via Dask**

# On-demand columnar data delivery: ServiceX

- **ServiceX** is a **data extraction and delivery** service

  ‣ users provide list of datasets to process + instructions for how to extract data (e.g. declarative)

  ‣ ServiceX can be co-located with input datasets for fast execution

  ‣ columnar data is returned and cached

    - subsequent executions can use the cache!



*input dataset*

*filtered & derived data*

*publication*

ServiceX

further analysis

subsequent analysis iterations use cache for significant speedup

# AGC "version 2" and future work


E. Kauffman at AGC demo day #2

- Development of a **"version 2" of the AGC analysis task** is ongoing

  ‣ expanded task: more complexity and data to process

  ‣ inclusion of machine learning aspects (training & inference)

- **Develop & compare** different **implementations**

  ‣ e.g. implementation using `ROOT RDataFrame`

- **Benchmarking**

  ‣ investigate performance, identify potential additional bottlenecks & implement solutions

- Longer term plan: **differentiable analysis pipeline**

  ‣ investigate end-to-end analysis optimization, evaluate usefulness vs cost of gradient information

# AGC events

- Organizing **yearly workshops**

  ‣ Mix of tutorials, demonstrations, discussions & planning

  ‣ Next workshop: IRIS-HEP AGC workshop on May 3–5

- Recently started bi-monthly **"demo day"** meetings

  ‣ informal, short demos on latest developments

  ‣ broad mix of topics to bring together diverse audience

  ‣ examples: #1, #2 (including recordings)

  ‣ this format works well!

Upcoming AGC workshop in May



First AGC "demo day"

# Summary

- The **Analysis Grand Challenge** is an **integration exercise** to study **HL-LHC analysis workflows**

- Developed **ttbar analysis task & implementation** based on **CMS Open Data**
  - ‣ all data & our implementation are publicly available

- We hope that the **Analysis Grand Challenge** can be **useful to the broader community**
  - ‣ test analysis tools, compare different workflows, test analysis facilities, …

- **Upcoming workshop:** https://indico.cern.ch/e/agc-workshop-2023

- Stay in touch via our **mailing list**
  - ‣ analysis-grand-challenge@iris-hep.org (sign up at this Google group)

# Give it a try!

- You can **run our AGC** analysis pipeline on **Binder**
  - ‣ **Try it out today!**

- All code also available on GitHub

- See also this PyHEP 2022 contribution
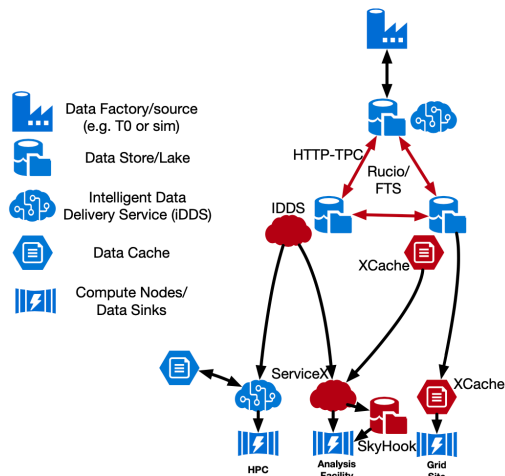  - ‣ includes recording of walkthrough

# Thank you!

- The **AGC is made possible** thanks to the **help of a large number of people** working on many different projects.

- **Thank you** in particular to the teams behind:

  ‣ coffea-casa

  ‣ Scikit-HEP, coffea, IRIS-HEP Analysis Systems

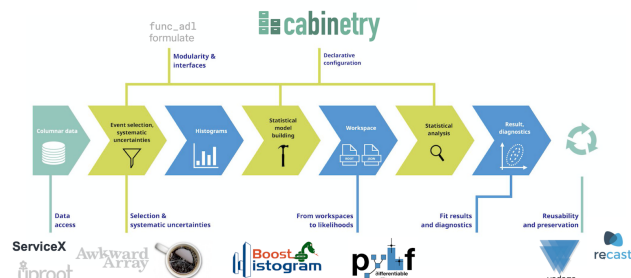  ‣ ServiceX, IRIS-HEP DOMA

  ‣ IRIS-HEP SSL

  ‣ CMS Open Data
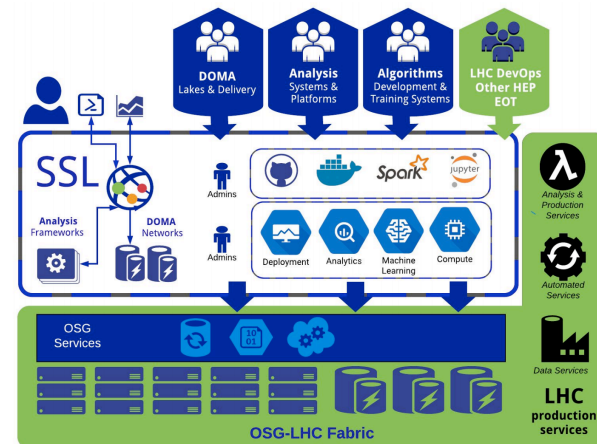
Backup

# Integration: connecting IRIS-HEP focus areas

**Data Organization, Management and Access (DOMA):**
*data delivery*

**Analysis Systems (AS):**
*tools*

**Scalable Systems Laboratory (SSL):**
*deployment techniques and resources*

# Top quark pair production