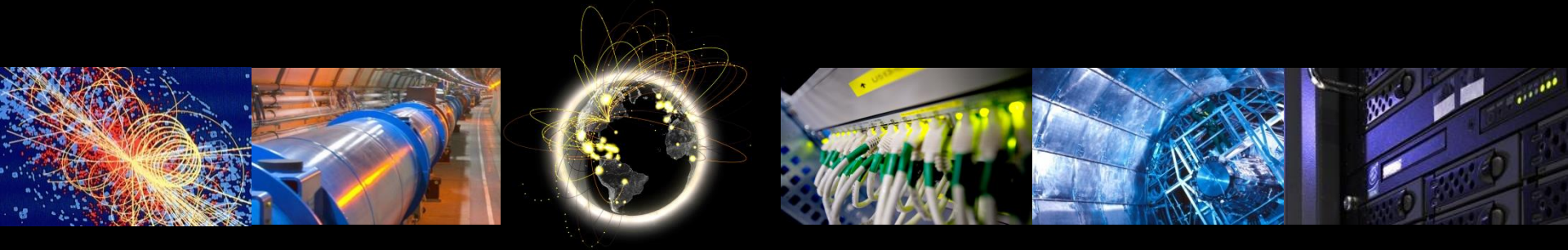# WLCG-HSF Workshop, Adelaide
## Analysis Systems: From Future Facilities to Final Plots

Ian Collier

ian.collier@stfc.ac.uk

UKRI-STFC Rutherford Appleton Laboratory

GDB, December 11th 2019

WLCG
Worldwide LHC Computing Grid

CERN

# Overview

- Workshops are the place where more of the collaboration gathers

- Workshops part of a process
  - Evolve as our infrastructure and computing challenge evolves
  - Joining WLCG & HSF firmly established now

- A first to focus workshop on just one area
  - Brainwriting format on day 2 led to much more active discussion and drew out ideas from many more attendees than usual format

# Workshop Motivation

- Short time available, focus on one topic: **Analysis**

- We have a situation where event rates are climbing precipitously for LHC experiments
- This has a direct impact on analysis
- In a flat budget model we need new ideas!
- Analysis a rich and layered topic:
  – Data storage and access at sites
  – Analysis model workflows as an integrated part of the production workflow
- Including realtime workflows
  – Good languages and analysis ergonomics are vital to productivity
- Resource scaling and responsiveness beyond the laptop…
  – Machine learning becomes more and more important

# Format

- Detailed agenda at:
  - https://indico.cern.ch/event/805983/
- 164 people registered
- First day
  - Plenary talks on key topics
- Second day
  - Group discussions looking at three broad challenge areas – 'brainwriting' format
    - Facilities
    - Machine Learning
    - Analysis Model

# Facilities ideas

- Flexibility
  - Tension between Interactive vs batch
  - Hardware flexibility & heterogeneity
  - Agile services
    - Facilities organised in layers of services
    - E.g. K8s for agile deployment
- Facilities themselves
  - Simultaneous convergence and divergence
    - Trend towards specialised facilities on one hand
    - While boundaries between Tiers blurring
      - focus on capabilities
  - Communities are heterogeneous & diverse
    - ML may be niche for us but more significant for other groups
    - Similarly for HPC & MPI jobs
- Funding
  - Evolution of pledging and funding models
  - How to get credit for e.g. GPUs, super I/O
    - Cost/ contribution should not just be measured in PB or CPU cycles
    - Credit for R&D & testbeds
- Training
  - Diversity of resources will mean
    - different support / training needs
    - Opportunities for site operators to develop skills knowledge

# Machine Learning

- Access to resources (particularly GPUs)
  - Two viewpoints, it's easy if you have access, impossible if you don't
  - A clear need for many of the ML enthusiasts to simply be able to work
  - Resources are there, we need to make them accessible
    - See Facilities Funding credit

- Ecosystem and Collaboration
  - Software should take as much advantage as possible of industry standard tools (not just e.g. TMVA wrappers)
    - Cannot stay up to date otherwise
    - Challenges for the interfaces to framework but many benefits
      - Cutting edge available
      - Collaboration easier
    - Our data is highly structured in a unique way, needs *thought* to make appropriate translations
  - Format conversions are critical
    - For tooling, ease of collaboration – need to be efficient
  - Training
    - ML has its own language, so does HEP
      - Common tooling reduces the barriers
      - Still HEP-specific use cases (uncertainties) that need to feed back into the tools
    - Using industry standard we benefit from industry docs and support

# Analysis Workflow Ideas

- Virtual datasets
  - Trade CPU for disk, recompute instead of storing intermediate data
  - Add columns from a separate file to avoid duplication
  - Old ideas but maybe time to revisit some, e.g. H1 used a system with four parallel files/trees
- Move analysis to use a workflow paradigm
  - Mantra - analysis will be a workflow, start like that on day 1 !
  - Encourages analysis to always present compute problems as a workflow
    - Not as individual jobs that you then scale up
- Analysis preservation, distinguish short term and long term views
  - Analysis languages (i.e. more declarative-like) help - we should use these !
  - Short term benefits from containers et al
  - Long term should avoid commercial (possibly expiring) tech
    - Good point but not clear there was a real solution proposed here!
  - Metadata challenges either way, in containers, in databases, in both

WLCG
Worldwide LHC Computing Grid

# General Observations

- It turns out analysis means different things to different people – leading to some confusion at time
    - There isn't one definition of analysis! There should be a follow-up meeting in the HSF Data Analysis working group where we try to identify use cases that are distinct from the computing perspective
    - Make sure we map this information back to other parts of the community in the language they understand!
- One interesting discussion was on attitudes toward "Complete"
    - Allowing partial (>90%) completion for most purposes seems reasonable
        - Metadata challenges anyway need to be solved for cases of incomplete statistics
    - Having zero tolerance for "final" run also seems reasonable but was more controversial
        - My (Paul Laycock) two cents - we can forget about analysis preservation if we have such a pessimistic attitude to being able to read data! That final data could come later during the extensive review process. Never being able to read it implies we are really in trouble.

WLCG
Worldwide LHC Computing Grid

# Proposals

- Poll both sides (users, facilities) - what do we think we need to support?
  - A dedicated day at the next WLCG-HSF workshop on these themes
    - Either we sponsor analysts to go to Lund or we hold a different meeting for that, the latter may work better! (Although we should still encourage analysts to attend)
  - GDB / preGCB meeting (to collect the known use cases from sites?)
  - Workshop at CERN would be able to attract may more active analysts

- More focused workshops where users and facilities meet
  - ML is one hot topic - we should better understand the scale of the problem here
    - What do people need? What would larger ML resources be used for?
  - Interactive analysis as another hot topic
    - It would be useful to understand what expectations are
      - For latency
      - For #cores
    - Then test drive pre-empting vs scheduling approaches
  - Test drive existing beta future systems with proponents like IRIS-HEP
    - What breaks and why? If we sponsored analysts to go to Lund, this could potentially be a breakout session for them. "Test Drive Future Analysis!".

# CWP

- Time for a lightweight review of CWP?
  - Check what progress has been made
  - Adjust goals as appropriate
  - Do changes in landscape mean that some adjustment of direction of travel is needed?
  - Are new issues/possibilities emerging?

- Should not be a major exercise
  - Rather check and adjustment
  - Perhaps as part of other planned review activity

WLCG
Worldwide LHC Computing Grid

# Finally

- We aim make use of diverse and most appropriate tools for technical work
- Should do the same for our meetings and workshops
- Brainwriting exercise is just one example
  - Encouraged contributions from all – many who might not speak in a room of 50 or 250
  - More inclusive and ensures more good ideas are heard
- Should consider more in this direction
  - In GDB and pre-GDBs as well as in workshops
- We cannot afford to miss out on good ideas