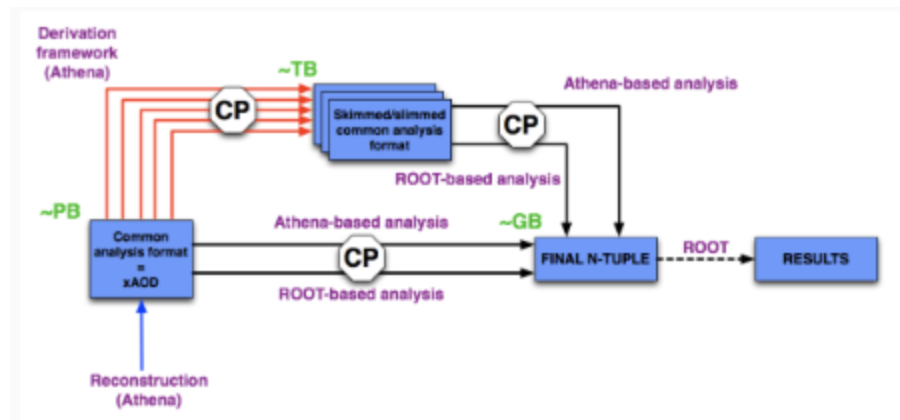


- Need to prepare questions/tracks for both sets of speakers, e.g.
 - **Questions for the DOMA talks:**
 - How things are changing (e.g. datalake)? What does it mean for users?
 - What are the variables that matter for distributed vs local storage?
 - How do you see latencies evolving?
 - What do you think is lowest latency achievable today (are there test systems in WLCG that can provide better than average performance)
 - What do you think is an affordable setup (e.g. comparing to google cloud 2 gbit/core)?
 - “Per institute” analysis machine vs “central HPC like setup”
 - What are the technologies we can expect to have in the future? E.g. HPC vs Grid vs local ... how this changes amount of storage, cpu, latencies, IO bandwidth
 - Which questions do you have for analysis experts on needs for HL-LHC?
 - **Questions to the “analysis experts”:**
 - Usage patterns: e.g. how often do you read xx events, how quick you need to get the result back, how much cpu time per event, which fraction of events is sufficient 90%, 99.9%?
 - What are the different data formats used in analysis?
 - How large are they per event and how often will they re-created, do different analysis teams use overlapping samples and if so how many?
 - How often will a given version be read?
 - What fractions of files will be read
 - Subset of events?

- Subset of data?
 - Lifetime of versions?
 - For the different data formats, what are the average computing times per event?
 - The above questions are very detailed, how certain are answers to these questions
 - Or in other words, how much flexibility in the fabric is needed?
 - Is it possible to (roughly) classify different analysis classes?
 - Data intensive / Compute bound?
- ATLAS has described in the past their analysis model with abstract diagrams and some additional information like the following:

(https://indico.cern.ch/event/773049/contributions/3476057/attachments/1935466/3207215/a_051119.pdf)



- More refined diagrams can be found in Johannes CHEP talk
 - Could the experiments, together with the computing people, come up with a common way to **describe** workflows, so that these descriptions can be used to better understand the requirements and constraints?
- Size of the input data that should be pre-staged/hot on

cache/living on SSD

- How much would that data be shared with other analysis
- Planning of resource usage vs “I want it all and I want it now”, avoid cycles wasting (develop+test+deploy also for analysis?)
 - “Test” requires what fraction of the data?
- Do most analysers plan their “big” run well and only need access to “big” resources at well-defined times? How often? If not, could they work like that?
- What should NOT be naively extrapolated from the current analysis model or data access with compact dAODs to the HL-LHC era ?
- What do you plan to do with compact dAOD (nanOAD, DAOD_PHYSLITE) which are made available on a Grid storage ?
 - Copy once to local storage and never access Grid
 - Produce Ntuple from Analysis facility and store locally NTUP format (very light format)
 - Concrete examples today of Grid access for heterogeneous hardware for tests ?
- Do we want to discuss architectures? <--- this is not a question to the speakers :-)
- Should we include questions of code quality (estimate how much optimisation could be been done)?
 - In this context, how long is analysis code in use (relative to the lifetime to reconstruction code??)
- Which questions do you have for computing experts on resources for HLLHC?

Notes on discussion, 6 March 2020

What does data lake & interface mean? How will it change analyses?

- * Ideally not? Adaptation to handle different facility models
- * Grid data access may be more restricted — no local storage per user

How will access patterns change, what will this mean for distributed storage (e.g. thanks to very small analysis formats)

What do we want to answer?

- * How much data needs to be read by analyses? (E.g. PB/week for different formats)
- * Latency sensitivity
- * Throughput on small event formats

Do we need solid numbers on sizes, rates etc?

- * How frequently will we run AOD->NanoAOD
- * How frequently will we run AOD->DAOD (!= NanoAOD)

Are we limited if we optimise too hard for one architecture (~NanoAOD)?

- * Need model that works for both small formats and something resembling AOD
- * Will there be a different/additional format?

DOMA talks —> what capabilities will future facilities have, which might enable analysis to be approached differently?

“Encouragement” from DOMA for how analyses might be able to structure for better efficiency?

Finalise questions Mon 9 Mar

Plan next discussion in May?

Can we come up with some prototypes/proposals for testing even this year?