

# Analysis infrastructure/framework

A collection of  
questions, observations, suggestions  
concerning  
analysis infrastructure and framework

Compiled by Marco van Leeuwen

Collected input from experienced analysers  
Christian KB, Silvia M, Andrea D, Marco v L,  
Constantin L, Pietro A, Laurent A, and friends

Possible **action items** highlighted in red

# Running on the GRID

## (User experiences)

- Individual experiences differ significantly
  - Depends on dataset, format, (person/privileges)
    - Can run over all (ESD) LHC10h data in one day!
    - p+p MC very fragmented lots of small jobs, cumbersome ('recent' option of merging via JDL in plugin helps a lot !)
- Success rate of jobs low (50-70%)
  - Also depends on use-case
  - Resubmission mechanism convenient, but often need >3 resubmissions to reach 90%
  - Not clear whether there is a single main cause for this
    - Large memory usage is an issue (sometimes >2-3 GB)
      - drawback of running trains with many tasks?
    - I/O errors
      - Recent work on R\_\_unzip, xrootd may improve this?
    - Site-issues contribute, but normally transient?

**Low success rate should be a concern – keep monitoring and improving**  
**Ask users for experience/input if needed**

# Stability of central services/grid connection

- Update of alien made aliensh more responsive
- However, still regular hiccups
  - Always quickly fixed (mail on atf list)
  - There may be work ongoing to improve this (e.g. reconfiguring/share of pcapiservers) ?
- It is practically impossible to keep a grid data connection (running merge or a test job on a local machine) open for long (>4 hrs)
  - Not so clear why; site-issues/R\_unzip etc contribute
  - Related to recent observation of expiring certificates after N bytes? (fix ongoing)
  - Does this also affect jobs running on the grid?

No show-stoppers, but definite room for improvement; need to keep an eye on these issues:

- **Keep users informed** and **ask for feedback** on the mailing list (e.g.: we changed ... let us know if you see improvements/problems)
- Suggestion: agree on **scheduled maintenance** times ?

# AOD vs ESD

Many analyses need only a small subset of ESD info

AODs: selected info to reduce resource usage

- AODs actively used in cases where they have distinct advantages:
  - Analysis needs small subset of data (e.g. MUON AODs)
  - Common intermediate analysis objects (jets, Heavy Flavour)
- Other ('minbias') analyses often use ESD:
  - Analysis code development sometimes started on ESD
  - More info available; more cross-checks possible
  - For some analyses, essential info is (not yet) in AOD
- Hybrid situation is 'natural'
  - Different analyses have different needs
  - Putting 'everything' in the AOD would be counter-productive !

**In mean time: keep encouraging people to use AODs**

Requires regular re-production of AODs and more information to the users  
about new sets and features

What is the forum to discuss AOD content and (re)production schedule?

→ Andreas' Analysis and QA meeting?

→ re-establish Analysis discussion during Offline Weekly? Fixed time slot?

# Analysis trains

- Single 'central' analysis train (or 1 per PWG) probably not practical/manageable
  - Too many carts
  - Different analyses run over different inputs
- Currently many analysis trains exist; mostly organised around 'common need'
  - This is probably the natural situation
  - Train operators concentrate knowledge/experience
    - Would like to have a forum to exchange this experience
  - Train operators could/should have special privileges
    - Being realised via 'official' PWG trains?
- Central trains:
  - PWG1/QA train
  - Common AOD production

Suggestion: set up monthly meeting of PWG analysis software coordinators and train operators  
Forum to exchange experiences and report common issues

Need to refresh (and publicize?) list of PWG analysis software coordinators?

# How to keep the information flowing?

## ALICE

A web of interconnected groups/meetings  
+  
central mailing lists, website(s)

- Only way to 'reach everyone' is via mail (or website if people know where to look)
  - If you have/know something of general use, send it to the list; need to encourage everyone !
  - Especially important for **announcements** (from offline or others)
- Alice-analysis-taskforce list works well
  - Only down-side: it is a firehose, with many topics/mails

Suggestion: **separate operational** (day-to-day issues) e-mail **from 'high-level'** questions/discussion  
Set up new mailing list for analysis software discussion

**Next slide: bonus track**

# AOD and ESD consistency, tender and repass

Is there a rule/paradigm that info ESD and AOD should be consistent?  
Can we afford to stick to it, or need rethink?

- Reality: number of improvements in software + calibration since 2010 data were produced.

To make those available for analysis:

- New reconstruction pass? Most straightforward option, but (too?) expensive in terms of CPU and/or storage
- Alternative: use tender (only fixes a subset of issues?)
  - Option I: produce AODs with tender
  - Option II: use ESD+tender in analysis
- Or: use existing ESD+tender to produce new ESD

Need input for further discussion: CPU+storage impact for different options  
Amount of human effort for each option

NB: Hopefully, the current situation is specific to 'first-year' data.  
No guarantees, though (expect sw+calibs to improve further over the years)