



Data Processing Update

Data Production meeting
23/01/2020

Stefano Lacaprara, Marco Milesi
INFN Padova



Proc10 status



- Local (KEKCC) All events **DONE 20/1**
 - Initial ETA was ~4/1
 - On 29/12, b2_prod cores drop from 1500->400
 - New ETA moved to 13/1
 - We found that ~23k jobs failed due to cvmfs issues at KEKCC
 - Not caught at the time by our monitor script
 - Resubmitted at the end (with 400 cores)
- Path:
 - /group/belle2/dataproduct/Data/OfficialReco/proc10/e0007/4S/GoodRuns
 - /group/belle2/dataproduct/Data/OfficialReco/proc10/e0008/4S/GoodRuns
 - /group/belle2/dataproduct/Data/OfficialReco/proc10/e0008/Continuum/GoodRuns
 - /group/belle2/dataproduct/Data/OfficialReco/proc10/e0008/Scan/GoodRuns
- <https://confluence.desy.de/display/BI/Processing+2019a-b#Processing2019a-b-Processing10details>

Issues with Proc10



- Some issues reported for some files (few) by Ami (thanks!)
 - Other failures (glitch) not caught by monitor script
 - Jobs resubmitted and files are now ok
- Fixed our monitor script by using output `json basf2_status` instead log parsing
 - [PR #92](#) (Jake, can you approve it?)
- Good/BadRun list ~finalized
 - Still some discussion with Watanuki and Karim about some early exp7 runs (before fire accident) [BIIDP-2333](#)
- Luminosity computation underway [BIIDP-2338](#)
- **Please report any additional issue you might find.**

Proc10 on the grid

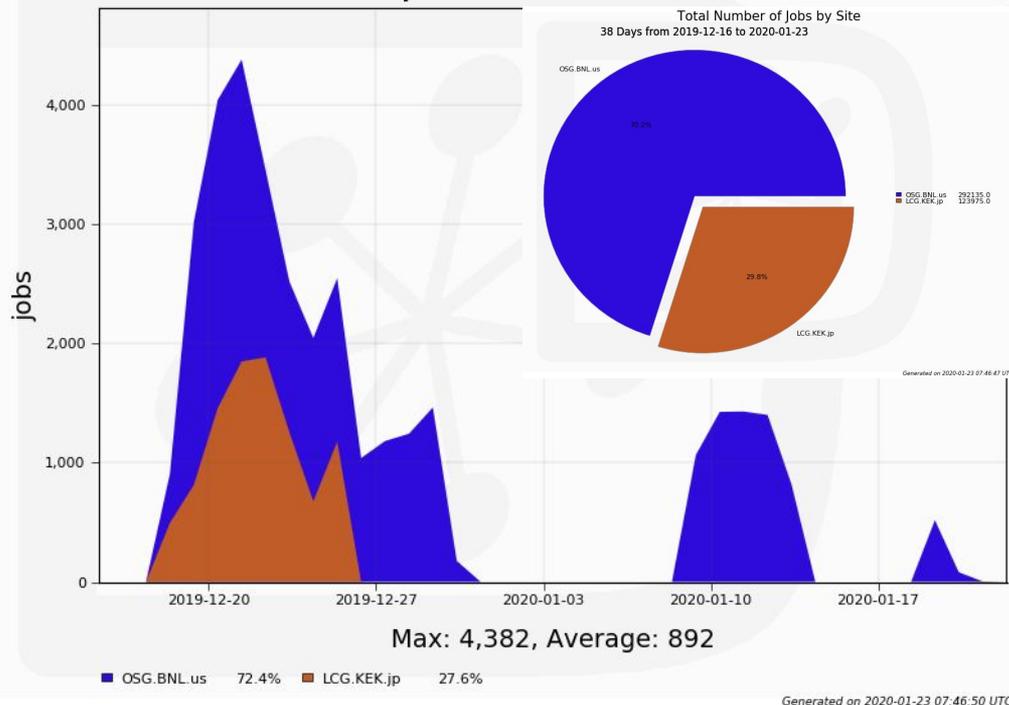


- A long and painful story.
- Multiple prodid submitted (100 runs per Prodid)
 - Additional one due to issue with our script for exp7 run<926
- 4S_offres and 4S_scan runs invalidated and resubmitted with proper metadata and path
- Few RAW files were missing on the grid for exp7 run<925
 - Resubmit after raw has been uploaded
- In total 20 Prodid: 18 valid, 2 cancelled
 - Exp7: 9629 9630 9631 9632 9633 9634 9863 **100% DONE**
 - Exp8 4S: 9635 9636 9637 9638 9639 9640 9641 9642 9643 **100% DONE**
 - Exp8 4S_offres: 9777 **17 merging jobs still waiting**
 - Exp8 4S_scan: 9776 **100% DONE**

Proc10 on the grid (II)



Running jobs by Site
38 Days from 2019-12-15 to 2020-01-22



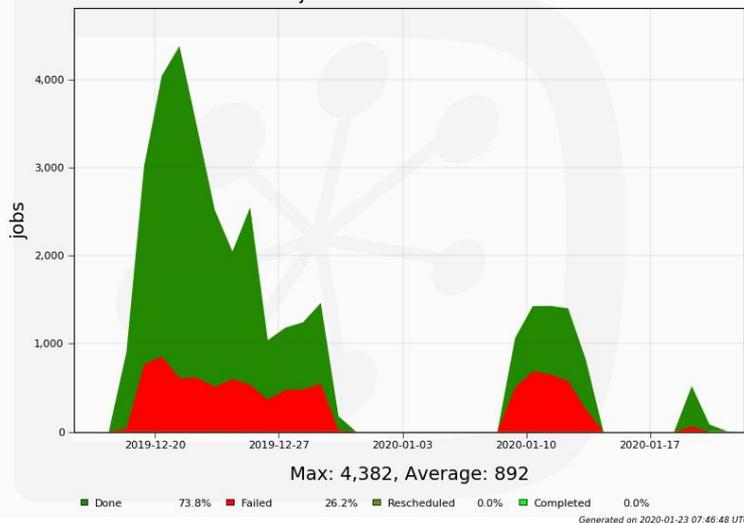
- A nice start (both BNL and KEKCC) 70-30
- Then several issues at BNL
 - Jobs seen as stalled
 - Long ticket BIIDCO-2194
 - Minor priority on Jira (!)
 - Problem not understood
 - Eventually went away
- Prod w/o progress for several days
 - Last peak is the very last processing for run <926

Only RawProcessing jobs shown, don't know how to show merge jobs just for this campaign (and not from MC13 also)

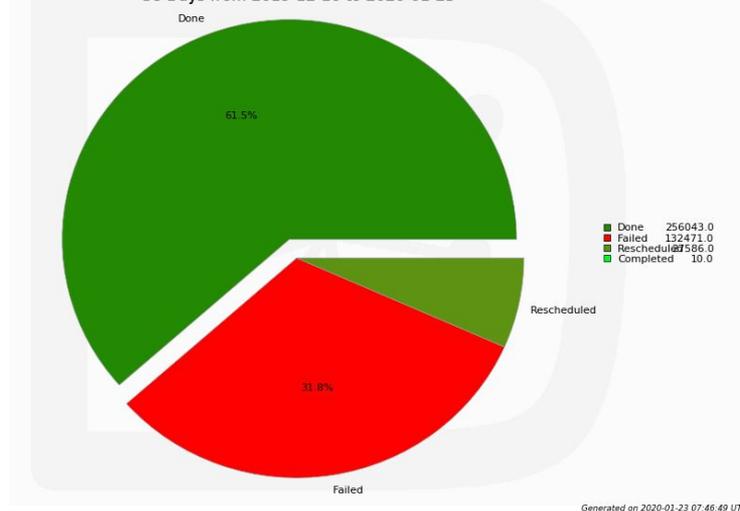
Failure rate: all failures at BNL.



Running jobs by FinalMajorStatus
38 Days from 2019-12-15 to 2020-01-22



Total Number of Jobs by FinalMajorStatus
38 Days from 2019-12-16 to 2020-01-23



- Initially issues with a WN at BNL (cvmsfs) fixed
- Two crashes in basf2: input files removed and tickets opened
- Then tons of stalled jobs killed and resubmitted (automatically by DIRAC)
 - 30% of total jobs

Proc10 on Dataset Searcher



- We (Marco) is uploading all proc10 related files fo Dataset Searcher right now
 - Excluding the 9777 (4S_offres) which is still waiting for 17 merging jobs to run (KMI or Nagoya)
 - `gb2_job_status -p 00105384 --status Waiting -l [...] 17 jobs are selected.`
- Will take some time due to AMGA query, DB uploading, etc

GoodRuns - BadRuns on the grid

- Until integration of DS, RunDB is in place, temporary workaround:
 - Get LPN from DS for all proc10 runs and save to a file
 - Functionality already existing in DS
 - Remove from that list all LPN corresponding to bad runs
 - Publish GoodRun LPN list on confluence
 - User can get that list and pass to `gbasf --input_dslist GoodRunLPNList.txt`
 - Eventually the purging of bad runs can be done automatically by DS querying RunDB
 - It should work (not tested personally)

Bucket 8 processing



- Calibration (including cdst) by AirFlow (Umberto/David)
- Final processing as usual
 - **First HLT_SKIM** (including hlt_hadrons) **at KEKCC**
 - Might consider to run first hlt_hadrons,
 - and then the others (bhabha, gammagamma, mumu2trk) to finish sooner
 - **Then all events at KEKCC and on the grid**
 - Grid processing will start in parallel with hlt_skims
- **Timescale: $L(\text{exp}10)=4 \text{ fb}^{-1}$**
 - Caveat: we still have just 400 cores on **b2_prod**
 - and on I we managed to run only $O(100)$ jobs at once
 - **Hlt_skims** (400 cores) 1.2 fb^{-1} per day => **3.5 days**
 - We will know from cdst processing for calibration
 - **All events** (400 cores): x12 (based on bucket7 experience) => **40 days**
 - **Do we want to do this?**
 - **All events on the grid** (based on bucket7) : **2 weeks**
 - Provided we don't face same issues as for proc10

- As presented in past meeting, we'd like to test hlt_skim processing on the grid
 - Current local workflow: 1 job -> multiple skims. Possible with current production tools?
 - We can test the workflow producing just one skim (e.g. hadrons)
 - ALL_RAW => HLT_hadron_REC
 - And treat the output as it is a "normal" mdst file
 - Need some work on our side to setup and test the workflow
 - In the (long) todo list since some time
 - [Marco's slide at DC meeting 28/11](#)
- Also: discuss how to integrate physics skim in the processing
 - In principle, it would be nice to launch just one procXX which:
 - Reconstruct events for a set of hlt_skims, then all events, then physics skims, ...
 - Publish everything on DS
 - Etc etc
 - And we just have to check the status from time to time and do other stuff
 - (ok, I'm dreaming)

/dataproduct cleanup

DONE

- <https://agira.desy.de/browse/BIIDP-2020>
- If you need some file which has been deleted, well, it is too late now, sorry.