# CI pipeline triggering analysis execution on Analysis Facility

**Matteo Bartolini**[1], Mattia Lizzo[1], Lorenzo Viliani[1], Tommaso Tedeschi[2]

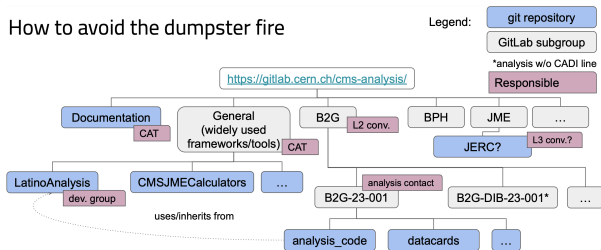[1]Università e INFN, Firenze
[2]Università e INFN, Perugia

Bi-weekly WP2 meeting

5 December 2023

# Continuous integration with CMS dataset



How to avoid the dumpster fire

- The CAT should allow analyzers to setup CI pipelines running on CMS dataset for code checking purposes (e.g. check effect of a commit on cut-flow yields...)
- `https://indico.cern.ch/event/1180058/contributions/5569735/attachments/2718208/4722157/` `CATcmsweekSept2023-2.pdf`
- Issue: difficult, because authentication is often needed in order to access the dataset→ CERN gitlab CI runners typically don't have it
- Solution: leverage EOS tokens, service in place to provide them, can be activated with a couple of lines in .gitlab-ci.yml

# The mkShapeRDF case

- mkShapesRDF, port to RDF of the "latinos" framework, to steer template based analyses with configuration files
- The framework is used, in this case, by the code WpWmJJpolarizations to perform the analysis
- This analysis needs to:
  - access the dataset stored somewhere on EOS
  - submit jobs to condor
  - store the output somewhere
  - Use the output to make plots or run fits
- The condor jobs will run on the Analysis Facility:
  - https://infn-cms-analysisfacility.readthedocs.io/en/latest/introduction/
  - EOS and AFS are not mounted → need xrdfs to access datasets stored on eos

# How do we do it?

These are the steps:

- Create a standalone docker image of mkShapesRDF containing all the libraries needed to run it
- The image is used by the analysis code running on the CI runners everytime a new commit is made.
- The CI runners will submit the condor jobs to workers running on the AF
- The workers on the AF will also run the docker image of the framework and perform all the operations

# The CI tool

- To build a docker image of your framework the CI tool is needed
- This project has the objective to supply an easy to use gitlab CI template to build images



```yaml
1   stages:
2   - build
3
4   include:
5     - project: 'ci-tools/container-image-ci-templates'
6       file:
7         - 'kaniko-image.gitlab-ci.yml'
8
9   variables:
10    CONTEXT_DIR: ""
11    DOCKER_FILE_NAME: "Dockerfile"
12    GIT_SUBMODULE_STRATEGY: recursive
13    PUSH_IMAGE: "true"
14    ACCELERATED_IMAGE: "false"
15    BUILD_ARGS: ""
16    SCAN_IMAGE: "false"
17    REGISTRY_IMAGE_PATH: ${CI_REGISTRY_IMAGE}
18
19  add_to_image:
20    extends: .build_kaniko
21    stage: build
22    tags:
23      - cvmfs
```

# The dockerfile in mkShapesRDF

- The dockerfile used by the CI tool should contain all the commands to build a standalone image of your framework
- In this case our image is built on top of the image of ubuntu 20
- If the building is successfull, the image will be created and pushed to gitlab-registry.cern.ch/lenzip/mkshapesrdf in this case

# The analysis code:WpWmJJpolarizations



```
.gitlab-ci.yml    2.17 KiB

1    default:
2        image:
3            name: gitlab-registry.cern.ch/lenzip/mkshapesrdf
4            entrypoint: ["/bin/sh", "-c"]
5
6    test:
7        tags:
8            - cvmfs
9
10       before_script:
11           - source /code/start.sh
12           - source /code/fix_xrdfs.sh
13
14       script:
15           - . .gitlab/init_infn_AF_token.sh
16           - ls /ca.crt
17           - condor_q
18           - condor_q -debug
19           - printf $proxy | base64 -d > myproxy
20           - export X509_USER_PROXY=$(pwd)/myproxy
21           - export X509_CERT_DIR=/cvmfs/cms.cern.ch/grid/etc/grid-security/certificates/
22           - source /code/fix_xrdfs.sh
23           - echo $X509_USER_PROXY
24           - echo $X509_CERT_DIR
25           #- xrdfs root://eoscms.cern.ch ls /eos/cms/store/group/phys_higgs/cmshww/amassiro/HWWNano/Summer20UL18_106x_nAOD
26           - voms-proxy-info
27           #- root -l -q root://eoscms.cern.ch//store/group/phys_higgs/cmshww/amassiro/HWWNano/Run2018_UL2018_nAODv9_Full2
28           - which checkCondor
29           - ls -a
30           - cd Full2017_v9
31           - ls -a
32           - mkShapesRDF -c 1
33           - ls -a
34           - condor_q
```

# The AF token

- this script is used to get the token needed to access the condor workers on the AF machine

```
1   IAM_TOKEN_ENDPOINT=https://cms-auth.web.cern.ch/token
2
3   #IAM_USER=dciangot
4
5   result=$(curl -s -L \
6     -d client_id=${IAM_CLIENT_ID} \
7     -d client_secret=${IAM_CLIENT_SECRET} \
8     -d grant_type=client_credentials \
9     -d username=${IAM_CLIENT_ID} \
10    -d password=${IAM_CLIENT_SECRET} \
11    -d scope="openid profile offline_access wlcg" \
12    ${IAM_TOKEN_ENDPOINT})
13
14  if [[ $? != 0 ]]; then
15    echo "Error!"
16    echo $result
17    exit 1
18  fi
19
20
21  access_token=$(echo $result | jq -r .access_token)
22  refresh_token=$(echo $result | jq -r .refresh_token)
23
24  echo $access_token > my_access_token
25
26
27  export _condor_SCHEDD_NAME=131.154.96.124.myip.cloud.infn.it
28  export _condor_SCHEDD_HOST=131.154.96.124.myip.cloud.infn.it
29  export _condor_COLLECTOR_HOST=131.154.96.124.myip.cloud.infn.it:30618
30  export _condor_SCITOKENS_FILE=$(pwd)/my_access_token
31  export _condor_AUTH_SSL_CLIENT_CAFILE=/ca.crt
32  export _condor_SEC_DEFAULT_AUTHENTICATION_METHODS=SCITOKENS
33  export _condor_TOOL_DEBUG=D_FULLDEBUG,D_SECURITY
```

# The proxy

- In order to access the files on eos a token is needed
- The token is personal, has a limited validity, and is generated with the command:
  - voms-proxy-init –rfc –voms cms -valid 192:00
  - base64 -w0 X509USERPROXY
  - go to the gitlab page containing your project, create a variable inside Settings → CI/CD → Variables
  - copy paste the content of x509_user_proxy to that variable
- If everything goes smoothly, you should see the following result and be finally able to submit your jobs to condor via the CI

```
33  $ condor_q
34  -- Schedd: 131.154.96.124.myip.cloud.infn.it : <131.154.96.124:31618?... @ 11/14/
    23 16:00:23
35  OWNER     BATCH_NAME   SUBMITTED   DONE   RUN    IDLE  TOTAL JOB_IDS
36  cmscat ID: 136305  11/14 14:50    _      _      _     411 136305.0-410
37  cmscat ID: 136307  11/14 15:32    2      _      _     411 136307.0-410
38  Total for query: 820 jobs; 820 completed, 0 removed, 0 idle, 0 running, 0 held, 0
    suspended
39  Total for cmscat: 820 jobs; 820 completed, 0 removed, 0 idle, 0 running, 0 held,
    0 suspended
40  Total for all users: 898 jobs; 873 completed, 0 removed, 0 idle, 25 running, 0 he
    ld, 0 suspended
```

# Submitting jobs to condor from the CI

- When submitting jobs to Condor from the CI the jdl file should contain a line pointing to the docker image that will be used by the workers:
  - +SingularityImage: /cvmfs/unpacked.cern.ch/gitlab-registry.cern.ch/lenzip/mkshapesrdf:latest/
- Once the jobs have been submitted a scripts is used to check the status of the condor jobs every n seconds to keep the CI busy
- Once all the jobs are done running the script will exit the loop and all data are transferred back to the CI runner, merged together and added to artifacts
- The full pipeline is here:
  `https://gitlab.cern.ch/cms-analysis/smp/wpwmjj_`
  `polarizations/analysis_code/-/jobs/33957811`

# Final result

# Conclusions

- We have been setting up a CI pipeline running a full latino analysis on an INFN analysis facility
- We overcame the initial struggles with authentication and tokens
- Detailed instructions will be thoroughly documented and made available
- Job submission is entirely based on condor at the moment, but we plan to start experimenting soon the use of dask to improve handling and merging of the full dataset