

Table of Contents

CRAB configuration file.....	1
CRAB configuration file.....	1
CRAB configuration sections.....	1
Predefined CRAB configuration file with empty skeleton.....	1
CRAB configuration parameters.....	2
Passing CRAB configuration parameters from the command line.....	7
Converting a CRAB2 configuration file into a CRAB3 configuration file.....	9

CRAB configuration file

Complete:  [Go to SWGuideCrab](#)

CRAB configuration file

For convenience, we suggest to place the CRAB configuration file in the same directory as the CMSSW parameter-set file to be used by CRAB.

The expected default name of the CRAB configuration file is `crabConfig.py`, but of course one can give it any name (respecting always the filename extension `.py` and not adding other dots in the filename), as long as one specifies the name when required (e.g. when issuing the CRAB submission command).

In CRAB3 the configuration file is in Python language. It consists of creating a `Configuration` object imported from the `WMCore` library:

```
import CRABClient
from WMCore.Configuration import Configuration
config = Configuration()
```

Once the `Configuration` object is created, it is possible to add new sections to it with corresponding parameters. This is done using the following syntax:

```
config.section_("<section-name>")
config.<section-name>.<parameter-name> = <parameter-value>
```

CRAB configuration sections

The table below shows what are the sections currently available for CRAB configuration.

Section	Description
General	In this section, the user specifies generic parameters about the request (e.g. request name).
JobType	This section aims to contain all the parameters of the user job type and related configurables (e.g. CMSSW parameter-set configuration file, additional input files, etc.).
Data	This section contains all the parameters related to the data to be analyzed, including the splitting parameters.
Site	Grid site parameters are defined in this section, including the stage out information (e.g. stage out destination site, white/black lists, etc.).
User	This section is dedicated to all the information relative to the user (e.g. voms information).
Debug	For experts use only.

Predefined CRAB configuration file with empty skeleton

To simplify life a bit, CRAB provides a function `config` that returns a `Configuration` object with pre-defined sections. The function is in the `CRABClient.UserUtilities` module. Users can import and use the function in their CRAB configuration file:

```
from CRABClient.UserUtilities import config
config = config()
```

which, from the point of view of the `Configuration` instance, is equivalent to:

```
from WMCore.Configuration import Configuration
config = Configuration()
```

```

config.section_("General")
config.section_("JobType")
config.section_("Data")
config.section_("Site")
config.section_("User")
config.section_("Debug")

```

CRAB configuration parameters

The table below provides a list of all the available CRAB configuration parameters (organized by sections), including a short description. Mandatory parameters are marked with two stars (**). Other important parameters are marked with one star (*).

Parameter	Type	Description
Section General		
requestName (*)	string	A name the user gives to it's request/task. In particular, it is used to create a project directory (named <code>crab_<requestName></code>) where corresponding to this particular task will be stored. Defaults to <code><time-stamp></code> , where the time stamp is of the form <code><YYYYMMDD></code> and corresponds to the submission time. The maximum allowed length is 100 characters, according to the format in <code>RX_TASKNAME</code> . The submission will fail with "Incorrect 'workflow' parameter" if other characters are used.
workArea (*)	string	The area (full or relative path) where to create the CRAB project. If the area doesn't exist, CRAB will try to create it using the <code>mkdir</code> command. Defaults to the current working directory.
transferOutputs (*)	boolean	Whether or not to transfer the output files to the storage site. If set to <code>False</code> , the output files are discarded and the user can not recover them. Defaults to <code>True</code> .
transferLogs (*)	boolean	Whether or not to copy the jobs log files to the storage site. If set to <code>False</code> , the log files are discarded and the user can not recover them. Note however that a short version of the log files containing the first 1000 and the last 3000 lines are still available through the monitoring web interface. Defaults to <code>False</code> .
failureLimit	integer	The number of jobs that may fail permanently before the entire task is cancelled. Disabled by default. Note: a very dangerous parameter for expert use, do not touch it unless you are sure of what you are doing.
instance (**)	string	The CRAB server instance where to submit the task. For users please use <code>'prod'</code> .
activity	string	The activity name used when reporting to Dashboard. For experts only.
Section JobType		
pluginName (**)	string	Specifies if this task is running an analysis ('Analysis') on an existing dataset or is running MC event generation ('PrivateMC').
psetName (*)	string	The name of the CMSSW parameter-set configuration file that should be run via <code>cmsRun</code> . Defaults to <code>'pset.py'</code> .
generator	string	This parameter should be set to <code>'lhe'</code> when running MC generation of LHE files. Automatically set if an <code>LHESource</code> is present in the parameter-set.
pyCfgParams	list of strings	List of parameters to pass to the CMSSW parameter-set configuration file as explained here. For example, if set to <code>['myOption', 'param1=value1', 'param2=value2']</code> , then the job will execute <code>cmsRun JobType.psetName myOption param1=value1 param2=value2</code> .

CRAB3ConfigurationFile < CMSPublic < TWiki

		param2=value2. NOTE1: No blanks allowed in 'param-value'. double dashes break things, i.e. this works JobType.pyCfgParams=["arg1=1", "arg2=2"] , but this fails JobType.pyCfgParams=["--arg1=1" "--arg2=2"].
inputFiles	list of strings	List of private input files (and/or directories) needed by the jobs. be added to the input sandbox. The input sandbox can not exceed The input sandbox is shipped with each job. The input files will b in the working directory where the users' application (e.g. cmsRu launched regardless of a possible path indicated in this parameter the file name at right of last / is relevant). Directories are tarred a subtree structure is preserved. Please check the FAQ for more det how these files are handled.
disableAutomaticOutputCollection	boolean	Whether to disable or not the automatic recognition of output file produced by PoolOutputModule or TFileService in the CMSSW parameter-set configuration. If set to True, it becomes the user's responsibility to specify in the JobType.outputFiles parameter output files that need to be collected. Defaults to False.
outputFiles	list of strings	List of output files that need to be collected. If disableAutomaticOutputCollection = False (the default), ou produced by PoolOutputModule or TFileService in the CMSSW parameter-set configuration are automatically recognized by CRA don't need to be included in this parameter.
eventsPerLumi	integer	When JobType.pluginName = 'PrivateMC', this parameter spe many events should a luminosity section contain. Note that every with a fresh luminosity section, which may lead to unevenly sized luminosity sections if Data.unitsPerJob is not a multiple of this parameter. Defaults to 100.
allowUndistributedCMSSW	boolean	Whether to allow or not using a CMSSW release possibly not ava sites. Defaults to False.
maxMemoryMB	integer	Maximum amount of memory (in MB) a job is allowed to use. De 2000.
maxJobRuntimeMin	integer	The maximum runtime (in minutes) per job. Jobs running longer amount of time will be removed. Defaults to 1315 (21 hours 55 m see the note about maxJobRuntimeMin below. Not compatible with Automatic splitting.
numCores	integer	Number of requested cores per job. Defaults to 1. If you increase to run multi-threaded cmsRun, you may need to increase maxMem well. In the CMSSW parameter-set configuration you may require also number of streams to be larger than one per thread, which affects memory consumption too.
priority	integer	Task priority among the user's own tasks. Higher priority tasks w processed before lower priority. Two tasks of equal priority will b jobs start in an undefined order. The first five jobs in a task are g priority boost of 10. Defaults to 10.
scriptExe	string	A user script that should be run on the worker node instead of the cmsRun. It is up to the user to setup the script properly to run on a node enviroment. CRAB guarantees that the CMSSW environme (e.g. scram is in the path) and that the modified CMSSW paramet configuration file will be placed in the working directory with nam PSet.py. The user must ensure that a properly named framework file will be written; this can be done e.g. by calling cmsRun withi script as cmsRun -j FrameworkJobReport.xml -p PSet.py. Th itself will be added automatically to the input sandbox. Output fil

CRAB3ConfigurationFile < CMSPublic < TWiki

		produced by PoolOutputModule or TFileService in the CMSSW parameter-set configuration file will be automatically collected (CRAB will look in the framework job report). The user needs to specify the output files to be collected in the <code>JobType.outputFiles</code> parameter. See CRAB3AdvancedTopic#Running_a_user_script_with_CRAB for more information.
<code>scriptArgs</code>	list of strings	Additional arguments (in the form <code>param=value</code>) to be passed to the script specified in the <code>JobType.scriptExe</code> parameter. The first argument to the script is always the job number
<code>sendPythonFolder</code>	boolean	Determine if the 'python' folder in the CMSSW release (<code>\$CMSSW_BASE/python</code>) is included in the sandbox or not. Defaults to <code>False</code> .
<code>sendExternalFolder</code>	boolean	Determine if the 'external' folder in the CMSSW release (<code>\$CMSSW_BASE/external</code>) is included in the sandbox or not. See https://hypernews.cern.ch/HyperNews/CMS/get/computing-tools . Defaults to <code>False</code> .
<code>externalPluginFile</code>	string	Name of a plug-in provided by the user and which should be run instead of the standard CRAB plug-in <code>Analysis</code> or <code>PrivateMC</code> . Can not be used together with <code>pluginName</code> ; is either one or the other. Not supported.
Section Data		
<code>inputDataset (*)</code>	string	When running an analysis over a dataset registered in DBS, this parameter specifies the name of the dataset. The dataset can be an official CRAB dataset or a dataset produced by a user.
<code>allowNonValidInputDataset</code>	boolean	Allow CRAB to run over (the valid files of) the input dataset given by <code>Data.inputDataset</code> even if its status in DBS is not <code>VALID</code> . Defaults to <code>False</code> .
<code>outputPrimaryDataset (*)</code>	string	When running an analysis over private input files or running MC generation, this parameter specifies the primary dataset name that will be used in the LFN of the output/log files and in the publication data (see Data handling in CRAB).
<code>inputDBS (*)</code>	string	The URL of the DBS reader instance where the input dataset is produced. The URL is of the form <code>'https://cmsweb.cern.ch/dbs/prod/<instance>/DBSReader'</code> . The <code>instance</code> can be <code>global</code> , <code>phys01</code> , <code>phys02</code> or <code>phys03</code> . The default is <code>global</code> . The aliases <code>global</code> , <code>phys01</code> , <code>phys02</code> and <code>phys03</code> in place of whole URLs are also supported (and indeed recommended to avoid typos). For datasets that are not of <code>USER</code> tier, CRAB only allows to read from <code>global</code> DBS.
<code>splitting (*)</code>	string	Mode to use to split the task in jobs. When <code>JobType.pluginName = 'Analysis'</code> , the splitting mode can either be <code>'Automatic'</code> (the default, please read the dedicated FAQ), <code>'FileBased'</code> , <code>'LumiBased'</code> , or <code>'EventAwareLumiBased'</code> (for Data the recommended mode is <code>'Automatic'</code> or <code>'LumiBased'</code>). For <code>'EventAwareLumiBased'</code> , CRAB will split the task by luminosity sections, where each job will contain a certain number of luminosity sections such that the number of events analyzed in each job is roughly <code>unitsPerJob</code> . When <code>JobType.pluginName = 'PrivateMC'</code> , the splitting mode can only be <code>'EventBased'</code> .
<code>unitsPerJob (*)</code>	integer	Mandatory when <code>Data.splitting</code> is not <code>'Automatic'</code> , suggests (and imposes) how many units (i.e. files, luminosity sections or events - depending on the splitting mode - see the note about <code>Data.splitting</code> below) to include in each job. When <code>Data.splitting = 'Automatic'</code> , <code>unitsPerJob</code> represents the jobs target runtime in minutes and its minimum allowed value is 180 (i.e. 3 hours).

CRAB3ConfigurationFile < CMSPublic < TWiki

totalUnits (*)	integer	Mandatory when <code>JobType.pluginName = 'PrivateMC'</code> , in which parameter tells how many events to generate in total. When <code>JobType.pluginName = 'Analysis'</code> , this parameter tells how many luminosity sections (when <code>Data.splitting = 'FileBased'</code>), luminosity sections (when <code>Data.splitting = 'LumiBased'</code>) or events (when <code>Data.splitting = 'EventAwareLumiBased'</code> or <code>Data.splitting = 'Automatic'</code> - note about "Data.splitting" below) to analyze (after applying the lumi and/or run range filters).
useParent	boolean	Adds corresponding parent dataset in DBS as secondary input source. Allows to gain access to more data tiers than present in the current dataset. This will not check for parent dataset availability; jobs may fail with errors or due to missing dataset access. Defaults to <code>False</code> .
secondaryInputDataset	string	An extension of the <code>Data.useParent</code> parameter. Allows to specify grandparent dataset in DBS (same instance as the primary dataset) as secondary input source. CRAB will internally set this dataset as the secondary input and will set <code>Data.useParent = True</code> . Therefore, <code>Data.useParent</code> and <code>Data.secondaryInputDataset</code> can not be used together <i>a priori</i> .
lumiMask (*)	string	A lumi-mask to apply to the input dataset before analysis. Can either be a URL address or the path to a JSON file on disk. Default to an empty string (no lumi-sections filter).
runRange (*)	string	The runs and/or run ranges to process (e.g. <code>'193093-193999,198050,199564'</code>). It can be used together with <code>lumiMask</code> . Defaults to an empty string (no run filter).
outLFNDirBase (*)	string	The first part of the LFN of the output files (see Data handling in CRAB). Accepted values are <code>/store/user/<username>[/<subdir>*</code>] (the <code>/</code> after <code><username></code> can not be omitted if a <code>subdir</code> is not given) and <code>/store/group/<groupname>[/<subgroupname>*</code>] (and <code>/store/local/<dir>[/<subdir>*</code>] if <code>Data.publication = False</code>). Defaults to <code>/store/user/<username>/</code> . CRAB creates the <code>outLFNDirBase</code> path on the storage site if needed, do not create it yourself otherwise stage-out may fail due to permissions inconsistency.
publication (*)	boolean	Whether to publish or not the EDM output files (i.e. output files produced by <code>PoolOutputModule</code>) in DBS. Notice that for publication to be successful, the corresponding output files have to be transferred to the permanent storage element. Defaults to <code>True</code> .
publishDBS (*)	string	The URL of the DBS writer instance where to publish. The URL has the form <code>'https://cmsweb.cern.ch/dbs/prod/<instance>/DBSWriter'</code> where <code>instance</code> can so far only be <code>phys03</code> , and therefore it is set to <code>phys03</code> by default, so the user doesn't have to specify this parameter. The alias <code>DBS</code> in place of the whole URL is also supported.
outputDatasetTag (*)	string	A custom string used in both, the LFN of the output files (even if <code>Data.publication = False</code>) and the publication dataset name (if <code>Data.publication = True</code>) (see Data handling in CRAB).
publishWithGroupName	boolean	If <code>Data.outLFNDirBase</code> starts with <code>/store/group/<groupname></code> then <code>groupname</code> instead of the username in the publication dataset name is used (see Data handling in CRAB). This feature allows different users running the same workflow to publish in the same dataset. Defaults to <code>False</code> .
ignoreLocality	boolean	Defaults to <code>False</code> . Set to <code>True</code> to allow the jobs to run at sites regardless of where the input dataset is hosted (this parameter has effect only when <code>Data.inputDataset</code> is used). Remote file access is done using Xrootd. The parameter <code>Site.whitelist</code> is mandatory and <code>Site.blacklist</code> can be used and it is respected. This parameter is useful to allow the jobs to run on other sites when for example a dataset is hosted only on sites which are not running CRAB jobs. When set to <code>True</code> , the user must provide also

		white-list of sites where the jobs could run that are physically close to the site that hosts the input dataset. For example, if the input dataset is hosted by a site in the USA, a reasonable site white-list would be ['T2_US_*'] while ['T2*'] is NOT reasonable and will most likely result in more jobs than gain
userInputFiles	list of strings	This parameter serves to run an analysis over a set of (private) input files as opposed to run over an input dataset from DBS. One has to provide in this parameter the list of input files: <code>Data.userInputFiles = ['file1', 'file2', 'etc']</code> , where 'fileN' can be an LFN, a PFN or even a Xrootd redirector. One could also have a local text file containing the list of input files (one file per line; don't include quotation marks nor comments) and then specify in this parameter the following: <code>Data.userInputFiles = open('/path/to/local/file.txt').readlines()</code> . When this parameter is used, the only allowed splitting mode is 'FileBased'. Also, since there is no input dataset from where to extract the primary dataset name, one should use the parameter <code>Data.outputPrimaryDataset</code> to define it, otherwise CRAB will use 'CRAB_UserFiles' as the primary dataset name. This parameter can not be used together with <code>Data.inputDataset</code> . CRAB will not do any data discovery, meaning that most probably jobs will be run at the sites where the input files are hosted (and therefore they will be accessed via Xrootd). But since it is in general more efficient to run jobs at the sites where the input files are hosted, it is strongly recommended that the user forces the jobs to be submitted to these sites using the <code>Site.whitelist</code> parameter.
Section Site		
storageSite (**)	string	Site where the output files should be permanently copied to. See about storageSite below.
whitelist	list of strings	A user-specified list of sites where the jobs can run. For example: ['T2_CH_CERN', 'T2_IT_Bari', ...]. Jobs will not be assigned to sites that is not in the white list. Note that at times this list may not be respected, see this FAQ
blacklist	list of strings	A user-specified list of sites where the jobs should not run. Useful to prevent jobs to run on a site where the user knows they will fail (e.g. because of temporary problems with the site). Note that at times this list may not be respected, see this FAQ
ignoreGlobalBlacklist	boolean	Whether or not to ignore the global site blacklist provided by the CMS Computing Board. Should only be used in special cases with a custom whitelist and blacklist to make sure the jobs land on the intended sites.
Section User		
voGroup	string	The VO group that should be used with the proxy and under which the jobs should be submitted.
voRole	string	The VO role that should be used with the proxy and under which the jobs should be submitted.
Section Debug		
oneEventMode	boolean	For experts use only.
asoConfig	list of dictionaries	For experts use only.
scheddName	string	For experts use only.
extraJDL	list of strings	For experts use only.
collector	string	For experts use only.

 **Note for** `Data.splitting = 'EventAwareLumiBased'`

When CRAB does data discovery of the input dataset in DBS, the number of events is only known per input

file (because that's the information available on DBS) and not per luminosity section. CRAB can therefore only estimate the number of events per luminosity section in a given input file as the number of events in the file divided by the number of luminosity sections in the file. Because of that, `Data.unitsPerJob` and `Data.totalUnits` should not be considered by the user as rigorous limits, but as limits applicable on average.

📌 **Note for** `maxJobRuntimeMin`

We strongly encourage every user to tune their splitting parameters aiming for jobs to run for a few hours, ideally 8-10 hours, and set the `maxJobRuntimeMin` accordingly. Having many jobs increases the chance of failure, since the number of problems is roughly proportional to the number of run jobs. Moreover, short jobs suffer of start/end overheads resulting in poor CPU/Wall-clock ratio, which impacts negatively CMS and makes it harder to secure additional resources. The large default value for `maxJobRuntimeMin` prevents jobs to be scheduled within shorter site queues which are polled faster.

📌 **Note for** `storageSite`

In CRAB3 the output files of each job are transferred first to a temporary storage element in the site where the job ran and later from there to a permanent storage element in a destination site. The transfer to the permanent storage element is done asynchronously by a service called AsyncStageOut (ASO). The destination site must be specified in the `Site.storageSite` parameter in the form '`Tx_yy_zzzzz`' (e.g. '`T2_IT_Bari`', '`T2_US_Nebraska`', etc.). The official names of CMS sites can be found in the CRIC [web page](#). The user **MUST** have write permission in the storage site.

Passing CRAB configuration parameters from the command line

It is possible to define/overwrite CRAB configuration parameters by passing them through the command line when the `crab submit` command is executed. Parameters can be set with the convention `<parameter-name>=<parameter-value>` and can be sequentially listed separating them with a blank space. Here is an example on how one would pass the request name and the publication name:

```
crab submit -c my_crab_config_file.py General.requestName=my_request_name Data.outputDatasetTag=m
```

📌 **Note:** Currently it is only possible to overwrite the parameters that take as value a string, an integer, a float or a boolean. Parameters that take a list can not be overwritten this way.

Converting a CRAB2 configuration file into a CRAB3 configuration file

CRAB3 is essentially new compared to CRAB2; it is not just a re-write. As a consequence, the configuration is different and there is no direct trivial translation that can be done automatically for every CRAB2 configuration file into a CRAB3 one. There is only a basic CRAB3 utility, called `crab2cfgTocrab3py`, meant to help the user to convert an existing CRAB2 configuration file into a CRAB3 configuration file template. The user has to provide the name of the CRAB2 configuration file he/she wants to convert and the name he/she wants to give to the CRAB3 configuration file (both arguments have default values; `crab.cfg` and `crabConfig.py` respectively).

```
crab2cfgTocrab3py [crab2configName.cfg] [crab3configName.py]
```

Instead of blindly taking the produced CRAB3 configuration file and run it, the user should always inspect the produced file, understand what each parameter means, edit them and add other parameters that might be needed, etc.

Here we give a usage example. Suppose we have the following CRAB2 configuration file with the default name `crab.cfg`:

CRAB3ConfigurationFile < CMSPublic < TWiki

```
[CRAB]
jobtype           = cmssw
scheduler         = remoteGlidein
use_server        = 0

[CMSSW]
datasetpath       = /GenericTTbar/HC-CMSSW_5_3_1_START53_V5-v1/GEN-SIM-RECO
dbs_url           = global
pset              = my_CMSSW_config.py
number_of_jobs    = 100
events_per_job    = 20
output_file       = output.root

[GRID]
se_white_list     = T2_IT_Bari
se_black_list     = T2_IT_Legnaro
data_location_override = T2_IT_Bari

[USER]
ui_working_dir    = my_CRAB_project_directory
return_data       = 0
copy_data         = 1
storage_element   = T2_IT_Legnaro
user_remote_dir   = my_remote_directory
publish_data      = 1
publish_data_name = my_publication_name
dbs_url_for_publication = phys03
```

If we run the tool without specifying any input parameters:

```
crab2cfgTOcrab3py
```

it will create a file `crabConfig.py` with the following content:

```
from WMCore.Configuration import Configuration
config = Configuration()
config.section_('General')
config.General.transferOutputs = True
config.General.requestName = 'my_CRAB_project_directory'
config.section_('JobType')
config.JobType.psetName = 'my_CMSSW_config.py'
config.JobType.pluginName = 'Analysis'
config.JobType.outputFiles = ['output.root']
config.section_('Data')
config.Data.inputDataset = '/GenericTTbar/HC-CMSSW_5_3_1_START53_V5-v1/GEN-SIM-RECO'
config.Data.publication = True
config.Data.unitsPerJob = 20
config.Data.publishDBS = 'https://cmsweb.cern.ch/dbs/prod/phys03/DBSWriter/'
config.Data.splitting = 'EventBased'
config.Data.inputDBS = 'https://cmsweb.cern.ch/dbs/prod/global/DBSReader/'
config.Data.outputDatasetTag = 'my_publication_name'
config.section_('Site')
config.Site.blacklist = ['T2_IT_Legnaro']
config.Site.whitelist = ['T2_IT_Bari']
config.Site.storageSite = 'T2_IT_Legnaro'
```

and it will show the following screen output:

```
Conversion done!
crab2cfgTOcrab3py report:
CRAB2 parameters not YET supported in CRAB3:
  data_location_override, user_remote_dir
CRAB2 parameters obsolete in CRAB3:
  return_data, jobtype, scheduler, use_server
```

CRAB3ConfigurationFile < CMSPublic < TWiki

As we already emphasized, the template configuration file produced by the `crab2cfgTocrab3py` utility should not be used before carefully looking into its content. Along this line, one can see for example that the parameter `JobType.outputFiles` was set to `['output.root']`. If `output.root` is defined in the CMSSW parameter-set configuration file in an output module, then it doesn't have to be included in the `JobType.outputFiles` list (although it doesn't harm).

-- AndresTanasijczuk - 02 Oct 2014

This topic: CMSPublic > CRAB3ConfigurationFile

Topic revision: r77 - 2020-08-28 - StefanoBelforte



Copyright &© 2008-2020 by the contributing authors. All material on this collaboration platform is the property of the contributing authors.
or Ideas, requests, problems regarding TWiki? use [Discourse](#) or [Send feedback](#)