Norwegian
Meteorological
Institute

# JSON schema validation of experiment configurations

Joint 30th Workshop All-Staff Meeting 2020

Roel Stappers
Met Norway
roels@met.no

02.04.2020

# Background/Motivation

- In Harmonie-Arome we currently validate only a fraction of our configuration options, e.g.
  - TSTEP needs to be a divisor of 3600
  - NLON, NLAT needs to be of the form $2^a \, 3^b \, 5^c$

- For many variables the allowed values are only mentioned in comments
- We override values based on other configuration options.
- Waste of SBU's by invalid configurations.

- See also RWP: SY3 clean up of Harmonie scripting system
- In addition we have a growing number of configurations in our systems (EPS, reanalysis, climate, high-res, nowcasting )

We need a more structured way to handle and validate configuration settings in Harmonie-Arome.

Norwegian
Meteorological
Institute

# Configuration in Harmonie

Very roughly speaking the Harmonie scripting system consists of 2 parts

1. Scripts called by ecflow tasks (in ecf/ and scr/ directory (CY43))
2. data files to handle configuration *,
   - config_exp.h
   - include.ass
   - Harmonie.pm
   - harmonie_namelists.pm
   - ECFlow Suite definitions (tdf files)
   - Configuration for submission (submit.LinuxPC submit.ecgb, etc. )
   - Configuration for HPC environments  (config.ecgb etc.)
   - Configuration for VarBC
   - Information on ECMWF cycles (MARS parameters etc)
   - Configuration for blacklisting
   - Harmonie_domains.pm
   - Harmonie_configurations.pm
   - Harmonie_testbed.pl (test_defs part)

* These files not pure data files but also contain Perl code.

# Current status.

- Work has started to use a common format for all configuration files.
- Language independent (TOML). Possible alternatives would be YAML or JSON.
- config_exp.toml is finished (for deterministic).
- Next step will be include.ass and harmonie.pm

- Configurations are validated using JSON Schema (next slide)
  - Implementation exist for a wide range of languages (C/C++, Go, Java, PHP, Javascript, Python,Perl, etc)
  - Supported by several editors/IDEs
  - Automatic creation of GUIs

# config_exp.toml

```toml
# **** Model geometry ****
[Geometry]
DOMAIN='DKCOEXP'                    # See definitions in scr/Harmonie_domains.pm
TOPO_SOURCE='gmted2010'            # Input source for orography. Available are (gtopo30|gmted2010)
                                   # For usage of gmted2010 check the documentation first
GRID_TYPE='LINEAR'                 # Type of grid (LINEAR|QUADRATIC|CUBIC)
VLEV='65'                          # Vertical level definition.
                                   # HIRLAM_60, MF_60,HIRLAM_40, or
                                   # BOUNDARIES = same number of levs as on boundary file.
                                   # See the other choices from scr/Vertical_levels.pl

VERT_DISC='vfd'                    # Discretization in the vertical (vfd,vfe)
LGRADSP='yes'                      # Apply Wedi/Hortal vorticity dealiasing (yes|no)

# **** High level forecast options ****
[Physics]
DYNAMICS="nh"                      # Hydrostatic or non-hydrostatic dynamics (h|nh)
PHYSICS="arome"                    # Main model physics flag (arome|alaro)
MASS_FLUX_SCHEME='edmfm'           # Version of EDMF scheme (edkf|edmfm)
                                   # Only applicable if PHYSICS=arome
                                   # edkf is the AROME-MF version
                                   # edmfm is the KNMI implementation of Eddy Diffusivity Mass Flux scheme for Meso-
STATNW="yes"                       # Switch for new set up cloud sscheme (yes|no)
HARATU="yes"                       # Switch for HARATU turbulence scheme (yes|no)
ALARO_VERSION=0                    # Alaro version (1|0)
XRIMAX=0.0                         # Maximum allowed Richardson number in the surface layer (cy40h default was 0.0)
```

# JSON Schema (example 1)

```json
"CISBA": {
    "type": "string",
    "description": "ISBA scheme",
    "enum": [
        "3-L",
        "2-L",
        "DIF"
    ],
    "default": "3-L",
    "links" : [
        {
        "rel" : "ISBA documentation",
        "href" : "https://www.umr-cnrm.fr/isbadoc/model.html"
        }
    ]
}
```

# JSON Schema (example 2)

```json
"CROUGH": {
    "type": "string",
    "description": "type of orographic roughness length",
    "enum": [
        "NONE",
        "Z01D",
        "BE04"
    ],
    "options" : {
        "enum_titles" : [
            "NONE | no orographic treatment",
            "Z01D | orographic roughness length does not depend on wind direction",
            "BE04 | Beljaars 2004 orographic drag"
        ]
    },
    "default": "NONE",
    "links": [
        {
            "rel": "CROUGH Surfex documentation",
            "href": "http://www.umr-cnrm.fr/surfex/spip.php?article126"
        }
    ]
}
```

Norwegian
Meteorological
Institute

# JSON Schema.

Combining Schema's. anyOf

```
{
  "anyOf": [
    { "type": "string", "maxLength": 5 },
    { "type": "number", "minimum": 0 }
  ]
}
```

"short" ✔️

"too long" ❌

12 ✔️

-5 ❌

# JSON Schema (Example)

To validate against `oneOf`, the given data must be valid against exactly one of the given subschemas.

```
{
  "oneOf": [
    { "type": "number", "multipleOf": 5 },
    { "type": "number", "multipleOf": 3 }
  ]
}
```

| 10 | ✔ |

| 9 | ✔ |

Not a multiple of either 5 or 3.

| 2 | ✘ |

Multiple of *both* 5 and 3 is rejected.

| 15 | ✘ |

Norwegian
Meteorological
Institute

# JSON Schema

assimilation.schema.json

```json
{
  "type": "object",
  "oneOf": [
    {
      "title": "3DVAR",
      "$ref" : "3dvar.schema.json"
    },
    {
      "title" : "4DVAR",
      "$ref" : "4dvar.schema.json"
    },
    {
      "title": "blending",
      "$ref" : "blending.schema.json"
    },
    {
      "title": "none",
      "$ref" : "none.schema.json"
    }
  ]
}
```

# Validation using JSON Schema

See https://json-schema.org/ for full specification

```
roels@pc4523:/media/roels/_disk2/git/Harmonie.jl/docs/schema (master)$ tree
.
├── archiving
│   └── archiving.schema.json
├── assimilation
│   ├── 3dvar.schema.json
│   ├── 4dvar.schema.json
│   ├── anasurf_mode.schema.json
│   ├── anasurf.schema.json
│   ├── assimilation.schema.json
│   ├── blending.schema.json
│   ├── ilres.schema.json
│   ├── inco.schema.json
│   ├── incv.schema.json
│   ├── lsmixbc.schema.json
│   ├── none.schema.json
│   ├── nouterloop.schema.json
│   └── tstep4d.schema.json
├── aux
│   └── aux.schema.json
├── build
│   └── build.schema.json
├── dfi
│   └── dfi.schema.json
├── eda
│   └── eda.schema.json
├── geometry
│   ├── domain_name.schema.json
│   ├── domain.schema.json
│   ├── enum_for_nlon_nlat.jl
│   ├── geometry.schema.json
│   └── nlon_nlat.schema.json
├── main
│   ├── branches.json
│   ├── branches.schema.json
│   ├── date.schema.json
│   ├── emails.schema.json
│   ├── main.schema.json
│   ├── paths.schema.json
│   └── timelists.schema.json
├── nesting
│   ├── mars.schema.json
│   └── nesting.schema.json
├── observations
│   ├── liste_loc.schema.json
│   ├── observations.schema.json
│   └── odb
│       ├── codetype.schema.json
│       ├── ec2keyval.jq
│       ├── ecfilter.jq
│       ├── extract_from_ecmwf
│       ├── obstype2.json
│       ├── obstype.json
│       ├── obstype.schema.json
│       ├── reporttype.schema.json
│       └── varno.schema.json
├── physics
│   ├── alaro.schema.json
│   ├── arome.schema.json
│   ├── dynamics.schema.json
│   └── physics.schema.json
├── postprocessing
│   └── postprocessing.schema.json
├── surfex
│   ├── namelist
│   │   ├── nam_sson.schema.json
│   │   ├── nam_teb.schema.json
│   │   └── surfex_namelist.schema.json
│   └── surfex.schema.json
└── system
    ├── config
    │   └── ecgb-cca.json
    ├── config.schema.json
    ├── hostdescriptions.schema.json
    ├── linuxpc.json
    └── system.schema.json

17 directories, 57 files
roels@pc4523:/media/roels/_disk2/git/Harmonie.jl/docs/schema (master)$ 
```

# Implementation

https://github.com/Hirlam/Harmonie.jl

```
-------------------------------------------------------------------
Language                      files        blank      comment       code
-------------------------------------------------------------------

JSON                             56           98            0       5062
TOML                            115          219          418       2273
YAML                              4           14            0        387
Julia                             7           71           29        169
HTML                              1           20           22         49
Markdown                          1           12            0         32
-------------------------------------------------------------------

SUM:                            184          434          469       7972
-------------------------------------------------------------------
```

Note 97% of SLOC is language independent JSON/TOML/YAML files.

HTML used for prepIFS like gui (see later slides)

Julia used to

- Convert TOML to config_exp.h format
- Unit-tests (next slide)
- Validate TOML files.

https://julialang.org/blog/2012/02/why-we-created-julia

# Unit tests

Any push to github automatically runs travis-CI to validate testbed configurations (work in progress).

Testbed configurations are created by merging toml files.

```julia
# AROME 4DVAR
@testset "AROME_4DVAR" begin
    config_exp  = TOML.parsefile("config/config_exp.toml")
    arome_4dvar = TOML.parsefile("harmonie_configurations/arome_4dvar.toml")
    config_exp_with_arome_4dvar = merge(merge, config_exp, arome_4dvar)
    @test Harmonie.isvalid(config_exp_with_arome_4dvar)
    merged_arome_4dvar = TOML.parsefile("testbed_configurations/arome_4dvar.toml")
    testbedconfig = merge(merge, config_exp_with_arome_4dvar, merged_arome_4dvar)
    @test Harmonie.isvalid(testbedconfig)
end
```

Norwegian
Meteorological
Institute

# Harmonie GUI

# LISTE_LOC example



**Harmonie experiment configuration** ✎ JSON ✎ Properties

See Surfex and Observations tabs for examples. Observations.LISTE_LOC[].obstype shows how enum_titles can be used to give meaningful names to integers. Geometry.DOMAIN shows how we can insert value directly in the config_exp.json. Surfex.Namelist shows how namelist information can be included in the gui. All of this could help to reduce the dependency on perl in the scripts.

Geometry    Nesting    Assimilation    DFI    Physics    Surfex    EDA    Postprocessing    Archiving    System    Paths    Aux    Times    **Observations**

## Observations ✎ JSON ✎ Properties

Example of using the grid layout style. Current on/off style based on include.ass. Better to use booleans in which can we can use select box. Alternative is to use a single array with multi-select. We need url's for each obstype here

| **Synop** | **Mode-S** | **Bouy** | **Temp** |
|---|---|---|---|
| 1 | 0 | 1 | 1 |
| | Mode S Enhanced Surveillance | | TEMP, TEMPSHIP |

## LISTE_LOC ⌄  [+ item] [✗ Last item] [✗ All]

This needs more work. Included here as an example of how arrays work, e.g. to be used in EPS? Click add item a few times.

**Array must have unique items**

E 10 40 3

E 1 3 3

E 1 3 3

### E 1 3 3 ⌄  [✗ Item] ↑  ✎ JSON  ✎ Properties

**action**

Exclude

**obstype**

1 Land SYNOP and SHIP reports

specifies the observation subtype (BUFR code) or the satellite channel for images.

ECMWF documentation
**codetype**

3, Radar Rain Rates

ECMWF documentation
**varno**

3, upper air u component

ECMWF documentation

# GUI

- The GUI should remain lightweight, and low maintenance
- Use of the GUI should be  optional.
- Relation between GUI and toml should be transparent
- GUI will be similar to prepIFS
- Need to a solution to handle the configuration for the linking of files in scripts (similar to Olive/Vortex)

# Editor support (JSON/YAML only)

# Code structure (repositories)

- In Harmonie scripts and configuration have always been in a single repository.
- Currently Harmonie.jl is a separate repository with configuration for config_exp.toml. This is similar to how prepIFS is used at ECMWF. Is it preferable to have a single repository with scripts and configuration or are there benefits to split this in two repositories ?
- Automatic generation of scripts (like Olive/Vortex)?

# Summary/To do list

A prototype for the handling and validating configuration data has been presented.

- Based on language independent format (TOML/YAML/JSON)
- Validation using JSON-Schema
- Works in Harmonie `/perm/ms/no/fars/worktrees/jsonschema`
- Testbed runs successfully
- Configuration now in a separate git-repository

To do:

- Needs further extension for 1) namelists, 2) param_bator.cfg etc., 3) LISTE_LOC, LISTE_NOIRE_DIAP, 4) codetype obstype, 5) submit.ecgb-cca, submit.LinuxPC, 6) job submission, 7) compiler options, 8) VarBC predictors, 9) MARS request files, 10) ECMWF cycles, 11) ecflow tasks/families 12) testbed
- Grouping of variables.
- Scripts should start extracting from config_exp.toml directly without needing the `export` statements in config_exp.h

Norwegian
Meteorological
Institute

# Script simplification

- Use TOML config to move "business logic" out of scripts

E.g. in Climate

## roelstappers / **Harmonie.jl**

<> Code    ⓘ Issues **0**    ⌥ Pull requests **0**    ☰ Projects **0**    ☷ Wiki    ⛉ Security

Branch: master ▾    **Harmonie.jl** / test / config / GRID_TYPE / **CUBIC.toml**

🔘 **Roel Stappers** Update CUBIC and LINEAR grid

**0** contributors

12 lines (8 sloc) | 103 Bytes

```
1    TRUNC = 4
2
3    [NAMRIP]
4    NFOST = 6
5
6    [NAMDYN]
7    LBOUND_D3 = true
8
9    # From Climate
10   [NAMCLA]
11   LSPSMORO = false
```

```
31  #############################################
32  # Determine use of smoothing or not
33  case $GRID_TYPE in
34      "LINEAR" )
35        LSPSMORO=.TRUE.
36        TRUNC=2
37      ;;
38      "QUADRATIC" )
39        LSPSMORO=.FALSE.
40        TRUNC=3
41      ;;
42      "CUBIC" )
43        LSPSMORO=.FALSE.
44        TRUNC=4
45      ;;
46      "CUSTOM" )
47        LSPSMORO=.FALSE.
48        TRUNC=2.4
49      ;;
50      *)
51      echo "Wrong grid type"$GRID_TYPE
52      exit 1
53      ;;
54  esac
55
56  # Redefining the spectral truncation and C+I
57  # if not given by user
58  if [ $LNMSMAX -eq 0 ] ; then
```
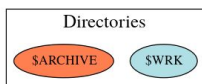
# default values in Namelist schema's

JSON-schema allows specification of default values. They are not part of the validation but tools can use this to fill in missing values.

For namelists this would mean we have to pick default values consistent with the IFS/surfex code.
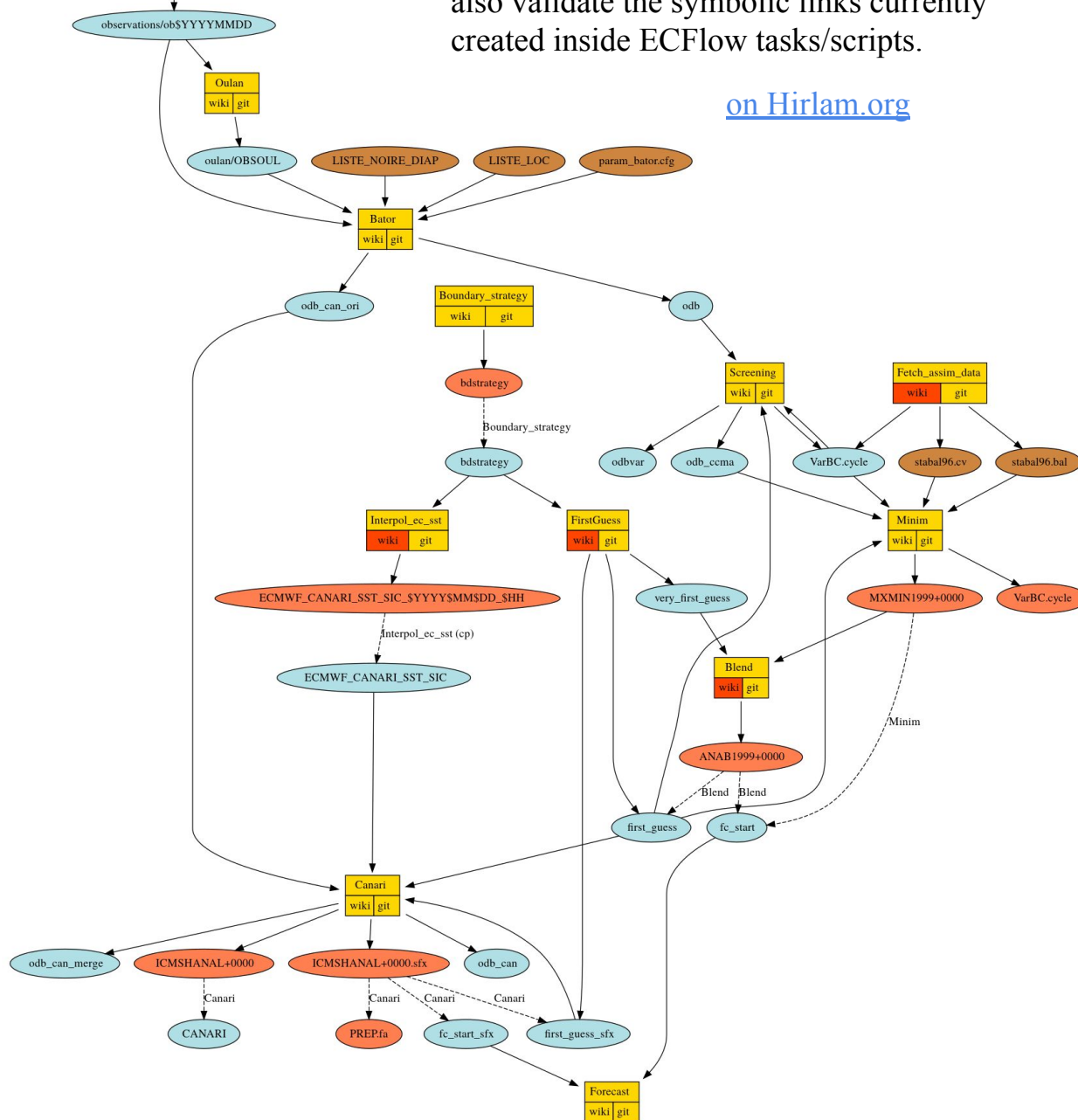
Which would force the use of the default LELAM=false

However we can have `enum : [ true ]` such that not explicitly setting specifying LELAM=false is invalid.

Norwegian
Meteorological
Institute

# Data Assimilation in Harmonie-Arome



Code refactoring needed such that unit tests can also validate the symbolic links currently created inside ECFlow tasks/scripts.

on Hirlam.org

# TOML for Namelists

scr/Mimim

```
194    #--- namelist
195    NAMELIST=$WRK/$WDIR/namelist
196    Get_namelist minimization $NAMELIST
197    sed -e "s/NBPROC/${NPROC}/g"   \
198        -e "s/NPROCX/${NPROCX}/g"  \
199        -e "s/NPROCY/${NPROCY}/g"  \
200        -e "s/ICLOUDFRACTI/$ICLOUDFRACTI/g" \
201        -e "s/LCLOUDFRACTI/$LCLOUDFRACTI/g" \
202        -e "s/NREDNMC/${REDNMC}/"  \
203        -e "s/NBZONVAR_EW=NBZONVAR_EW/NBZONVAR_EW=$NBZONVAR_EW/g" \
204        -e "s/LBVARBC/$LBBVARBC/"  \
205        -e "s/NBUPTRA/${JUPTRA}/g" \
206        -e "s/LBOBS/${JLOBS}/g" \
207        -e "s/LBSKIPMIN/${JLSKIPMIN}/g" \
208        -e "s/LBCHRESINCR/${JLCHRESINCR}/g" \
209        -e "s/LBNHDYN/${JLNHDYN}/g" \
```

```
2372    # Minimization
2373    %minimization=(
2374      NAMCT0=>{
2375      'LFDBOP' => '.FALSE.,',
2376      'NFPOS' => '0,',
2377      'LOBS' => 'LBOBS,',
2378      'LNHDYN'   => 'LBNHDYN,',
2379      'LSIMOB' => '.FALSE.,',
2380      'NCNTVAR' => '2,',
2381      'NFRGDI' => '10000,',
2382      'NFRHIS' => '10000,',
2383      'NFRISP' => '10000,',
2384      'NFRPOS' => '10000,',
```

```
2736      NAMPAR0=>{
2737      'NPRGPEW' => '1,',
2738      'NPRGPNS' => 'NBPROC,',
2739      'NPROC' => 'NBPROC,',
2740      'NPRTRV' => '1,',
2741      'NPRTRW' => 'NBPROC,',
2742      },
```

scr/Get_namelist

```
220
221    minimization)
222      NAMELIST_CONFIG="$DEFAULT $VARBC_NAM minimization ${PHYSICS}_minimization ${EXTRA_FORECAST_OPTIONS} args"
223      ;;
224
```

```
256    # Screening/Minim/4D-Var
257    for $task ('Minim','Screening','4DVtraj','4DVscreen','LETKF','FGerror','ComputeHx'){
258    $job_list{$task}{'TASK_PER_NODE'} = $submit_type.'-l EC_tasks_per_node='.$tasks_high ;
259    $job_list{$task}{'TOTAL_TASKS'}   = $submit_type.'-l EC_total_tasks='.$nproc_high ;
260    $job_list{$task}{'RESOURCES'}     = $submit_type.'-l EC_memory_per_task='.$memory_high.'MB' ;
261    $job_list{$task}{'ZMPPEXEC'}      = 'export MPPEXEC="aprun -n '.$nproc_high.'"' ;
262    $job_list{$task}{'ZNPROCX'}       = 'export NPROCX=1' ;
263    $job_list{$task}{'ZNPROCY'}       = 'export NPROCY='.$nproc_high ;
264    $job_list{$task}{'ZNPROC'}        = 'export NPROC='.$nproc_high ;
265    $job_list{$task}{'ZENSSIZE'}      = 'export ENSSIZE='.$ENV{ENSSIZE} ;
```

Norwegian
Meteorological
Institute

# Domains

## Domains.jl

`build | passing`  `coverage | 70%`  `docs | dev`

## Installation

You can obtain Domains.jl using Julia's Pkg REPL-mode (hitting `]` as the first character of the command prompt):

```
(v1.3) pkg> add https://github.com/Hirlam/Domains.jl
```

## Unit tests

The domains in `src/json/` are validated against the json schema file in `src/jsonschema/domain.schema.json` The schema validates:

- Required fields are present: `TSTEP`, `NLON`, `NLAT`, `LONC`, `LATC`, `LON0`, `LAT0`, `GSIZE`
- `TSTEP` is a divisor of 3600
- `NLON` ( `NLAT` ) are of the form $2^a\ 3^b\ 5^c$ with either $a \geq 1$, $b \geq 0, c \geq 0$ or $a=b=c=0$
- $-180 \leq$ `LON0`, `LONC` $\leq 180$
- $-90 \leq$ `LAT0`, `LATC` $\leq 90$

`EZONE` is not required but currently present in all domains `EZONE=11`

In addition ,for domains that use the Lambert projection. tests validate that the north pole is outside the domain.

Documentation