



# Towards a new validation system - DAVAÏ

*Alexandre Mary*  
in coll. w/ E. Arbogast & F. Suzat

Code Training Days - Sept. 12th, 2019

## 1 Why

- Building a new cycle
- From Mitrailllette to...

## 2 DAVAÏ

- Scripting and organisation of tests
- Automated comparisons
- Graphical interface
- Open discussion

# Outline

## 1 Why

- Building a new cycle
- From Mitrallette to...

## 2 DAVAT

- Scripting and organisation of tests
- Automated comparisons
- Graphical interface
- Open discussion

Building a new cycle

## Building a new cycle : process

## Building a new cycle

## Building a new cycle : process

- ① merge N user branches (or T+R branches) → CY##\_t1.01

## Building a new cycle

## Building a new cycle : process

- ① merge N user branches (or T+R branches) → **CY##\_t1.01**
- ② phasers & GMAP staff : run Mitrallette tests

## Building a new cycle

## Building a new cycle : process

- ① merge N user branches (or T+R branches) → **CY##\_t1.01**
- ② phasers & GMAP staff : run Mitrallette tests
- ③ analyse crashes/outputs ⇒ report(s)

## Building a new cycle

## Building a new cycle : process

- ① merge N user branches (or T+R branches) → **CY##\_t1.01**
- ② phasers & GMAP staff : run Mitrallette tests
- ③ analyse crashes/outputs ⇒ report(s)
- ④ search for crashes & causes of non-expected numerical differences

## Building a new cycle

## Building a new cycle : process

- ① merge N user branches (or T+R branches) → **CY##\_t1.01**
- ② phasers & GMAP staff : run Mitrallette tests
- ③ analyse crashes/outputs ⇒ report(s)
- ④ search for crashes & causes of non-expected numerical differences
- ⑤ collect bugfixes → **CY##\_t1.02**

## Building a new cycle

## Building a new cycle : process

- ① merge N user branches (or T+R branches) → **CY##\_t1.01**
- ② phasers & GMAP staff : run Mitrallette tests
- ③ analyse crashes/outputs ⇒ report(s)
- ④ search for crashes & causes of non-expected numerical differences
- ⑤ collect bugfixes → **CY##\_t1.02**
- ⑥ ... repeat (3-4-5) as long as necessary ...

## Building a new cycle

## Building a new cycle : process

- ① merge N user branches (or T+R branches) → **CY##\_t1.01**
- ② phasers & GMAP staff : run Mitrallette tests
- ③ analyse crashes/outputs ⇒ report(s)
- ④ search for crashes & causes of non-expected numerical differences
- ⑤ collect bugfixes → **CY##\_t1.02**
- ⑥ ... repeat (3-4-5) as long as necessary ...
- ⑦ when happy with results, declare **CY##\_t1.nn** ⇒ **CY##T1**

## Why building a new cycle is difficult

## Building a new cycle

## Why building a new cycle is difficult

- user branches are ***insufficiently validated***  
⇒ need for scientists/developers to ***systematically*** validate their branches against a wide range of model configurations and algorithms

## Why building a new cycle is difficult

- user branches are ***insufficiently validated***  
⇒ need for scientists/developers to ***systematically*** validate their branches against a wide range of model configurations and algorithms
- much of the ***evaluation of the tests is manual*** (though ciboulette brought some help on that), whereas it could be automatised  
⇒ *save time* spent on comparing listings or fields, to have more time to dive *in the code* (looking for crashes or numerical differences)

## Why building a new cycle is difficult

- user branches are ***insufficiently validated***  
⇒ need for scientists/developers to ***systematically*** validate their branches against a wide range of model configurations and algorithms
- much of the ***evaluation of the tests is manual*** (though ciboulette brought some help on that), whereas it could be automatised  
⇒ *save time* spent on comparing listings or fields, to have more time to dive *in the code* (looking for crashes or numerical differences)
- validation of pre-cycle is split between GMAP staff & ALADIN phasers, each with its working habits :  
→ status of validation of a pre-cycle is spread between e-mails, different formatting and degrees of precision...  
⇒ need for a *clear & complete dashboard*

## Why building a new cycle is difficult

- user branches are ***insufficiently validated***  
⇒ need for scientists/developers to ***systematically*** validate their branches against a wide range of model configurations and algorithms
- much of the ***evaluation of the tests is manual*** (though ciboulette brought some help on that), whereas it could be automatised  
⇒ *save time* spent on comparing listings or fields, to have more time to dive *in the code* (looking for crashes or numerical differences)
- validation of pre-cycle is split between GMAP staff & ALADIN phasers, each with its working habits :  
→ status of validation of a pre-cycle is spread between e-mails, different formatting and degrees of precision...  
⇒ need for a ***clear & complete dashboard***

not to mention Data Assimilation, postponed by months (when it is done)

## Flaws of Mitrallette

- design not very flexible
- code duplication
- not adapted to Data Assimilation validation (dependency between tasks, ODB handling)

From Mitrallette to...

## Flaws of Mitrallette

- design not very flexible
- code duplication
- not adapted to Data Assimilation validation (dependency between tasks, ODB handling)

## Opportunities

- Karim leave
- **Vortex** : now used for **all** our operational and research configurations, and contains much **wrapping utilities** for DA runs
- **OOPS** : thanks to its object-oriented design, enables to test observation operators **separately** from algorithms (hence in parallel to the model)

# Outline

## 1 Why

- Building a new cycle
- From Mitraillette to...

## 2 DAVAÏ

- Scripting and organisation of tests
- Automated comparisons
- Graphical interface
- Open discussion

## Scripting a Task

- 1 fetch input resources for task
- 2 run executable
- 3 save output : cache/archive

## Scripting a Test

- 1 fetch input resources for task
- 1' *fetch reference output (for comparison)*
- 2 run executable
- 2' *compare output to reference, and feed DAVAI graphical interface (database) with the comparison*
- 3 save output : cache/archive (+ *metadata of the task, e.g. mem/CPU*)

## Classification of the tests

- model :
    - forecast
    - fullpos
    - algorithmic components around forecast (DFI, IAU, IOserv, restart...)
    - TL & AD
- upon a wide variety of geometries (LAM/global) and model options  
(largely imported from *Mitrallette*)

## Classification of the tests

- model :
  - forecast
  - fullpos
  - algorithmic components around forecast (DFI, IAU, IOserv, restart...)
  - TL & AD
- upon a wide variety of geometries (LAM/global) and model options  
(largely imported from *Mitrallette*)
- obs preparation

## Classification of the tests

- model :
  - forecast
  - fullpos
  - algorithmic components around forecast (DFI, IAU, IObserv, restart...)
  - TL & AD
- upon a wide variety of geometries (LAM/global) and model options  
(largely imported from *Mitrallette*)
- obs preparation
- obs operators (obstype per obstype)

## Classification of the tests

- model :
  - forecast
  - fullpos
  - algorithmic components around forecast (DFI, IAU, IOserv, restart...)
  - TL & AD
- upon a wide variety of geometries (LAM/global) and model options  
(largely imported from *Mitrallette*)
- obs preparation
- obs operators (obstype per obstype)
- variational

## Classification of the tests

- model :
  - forecast
  - fullpos
  - algorithmic components around forecast (DFI, IAU, IOserv, restart...)
  - TL & AD
- upon a wide variety of geometries (LAM/global) and model options  
(largely imported from *Mitrallette*)
- obs preparation
- obs operators (obstype per obstype)
- variational
- ensemble

## Classification of the tests

- model :
  - forecast
  - fullpos
  - algorithmic components around forecast (DFI, IAU, IOserv, restart...)
  - TL & AD
- upon a wide variety of geometries (LAM/global) and model options  
(largely imported from *Mitrallette*)
- obs preparation
- obs operators (obstype per obstype)
- variational
- ensemble
- physiography

## Classification of the tests

- model :
  - forecast
  - fullpos
  - algorithmic components around forecast (DFI, IAU, IOserv, restart...)
  - TL & AD
- upon a wide variety of geometries (LAM/global) and model options  
(largely imported from *Mitrallette*)
- obs preparation
- obs operators (obstype per obstype)
- variational
- ensemble
- physiography
- miscellaneous

## Classification of the tests

- model :

- forecast
- fullpos
- algorithmic components around forecast (DFI, IAU, IObserv, restart...)
- TL & AD

→ upon a wide variety of geometries (LAM/global) and model options  
(largely imported from *Mitrallette*)

- obs preparation

- obs operators (obstype per obstype)

- variational

- ensemble

- physiography

- miscellaneous

→ on small truncation/grids, for **efficiency** (a few minutes or less per test, max)  
→ and ± small number of obs (**memory** vs. representativity)

## Comparison to reference

Each task has its specific outputs to be checked.

⇒ Attach one or several dedicated "*experts*" to each task, an expert being in charge of comparing two similar outputs.

## Comparison to reference

Each task has its specific outputs to be checked.

⇒ Attach one or several dedicated "*experts*" to each task, an expert being in charge of comparing two similar outputs.

So far, following experts have been developed :

## Comparison to reference

Each task has its specific outputs to be checked.

⇒ Attach one or several dedicated "*experts*" to each task, an expert being in charge of comparing two similar outputs.

So far, following experts have been developed :

- **Norms** : from listings, parse norms and compares numbers of different digits, step by step, field by field

## Comparison to reference

Each task has its specific outputs to be checked.

⇒ Attach one or several dedicated "*experts*" to each task, an expert being in charge of comparing two similar outputs.

So far, following experts have been developed :

- **Norms** : from listings, parse norms and compares numbers of different digits, step by step, field by field
- **Jo-Tables** : from listings, parse Jo-Tables and compares Jo obstype by obstype

## Comparison to reference

Each task has its specific outputs to be checked.

⇒ Attach one or several dedicated "*experts*" to each task, an expert being in charge of comparing two similar outputs.

So far, following experts have been developed :

- **Norms** : from listings, parse norms and compares numbers of different digits, step by step, field by field
- **Jo-Tables** : from listings, parse Jo-Tables and compares Jo obstype by obstype
- **Fields in files** : from FA historical or GRIB post-processed files, read fields and compares field by field (normalized bias, std and max difference)

## Comparison to reference

Each task has its specific outputs to be checked.

⇒ Attach one or several dedicated "*experts*" to each task, an expert being in charge of comparing two similar outputs.

So far, following experts have been developed :

- **Norms** : from listings, parse norms and compares numbers of different digits, step by step, field by field
- **Jo-Tables** : from listings, parse Jo-Tables and compares Jo obstype by obstype
- **Fields in files** : from FA historical or GRIB post-processed files, read fields and compares field by field (normalized bias, std and max difference)
- **RSS** : from *stdout*, compares memory use

## Comparison to reference

Each task has its specific outputs to be checked.

⇒ Attach one or several dedicated "*experts*" to each task, an expert being in charge of comparing two similar outputs.

So far, following experts have been developed :

- **Norms** : from listings, parse norms and compares numbers of different digits, step by step, field by field
- **Jo-Tables** : from listings, parse Jo-Tables and compares Jo obstype by obstype
- **Fields in files** : from FA historical or GRIB post-processed files, read fields and compares field by field (normalized bias, std and max difference)
- **RSS** : from *stdout*, compares memory use
- **DrHook** : from DrHook profiling files, compares elapse time

## Comparison to reference

Each task has its specific outputs to be checked.

⇒ Attach one or several dedicated "*experts*" to each task, an expert being in charge of comparing two similar outputs.

So far, following experts have been developed :

- **Norms** : from listings, parse norms and compares numbers of different digits, step by step, field by field
- **Jo-Tables** : from listings, parse Jo-Tables and compares Jo obstype by obstype
- **Fields in files** : from FA historical or GRIB post-processed files, read fields and compares field by field (normalized bias, std and max difference)
- **RSS** : from *stdout*, compares memory use
- **DrHook** : from DrHook profiling files, compares elapse time
- **OOPS-JoTests** : from *stdout*, compares computed Jo or result of adjoint test of Jo

## Ciboulaï

***Ciboulaï*** is the name of the graphical interface :

- macro view :
  - ***all Davaï experiments from all users***
  - filtering tool among experiments
  - summary of each experiment : number and % of tests OK/KO/ ?/crashed
  - → give access to *individual experiment* : ↓

---

a. a task can be compared to a "*continuity*" and/or a "*consistency*" reference

## Ciboulaï

**Ciboulaï** is the name of the graphical interface :

- macro view :
  - ***all Davaï experiments from all users***
  - filtering tool among experiments
  - summary of each experiment : number and % of tests OK/KO/ ?/crashed
  - → give access to *individual experiment* : ↓
- meso view (*individual experiment*) :
  - ***all tasks ran in the experiment***
  - filtering tool among tasks
  - summary of each task : OK/KO/ ?/crashed + RSS/elapse deltas
  - → give access to *individual task* : ↓

---

a. a task can be compared to a "*continuity*" and/or a "*consistency*" reference

## Ciboulaï

**Ciboulaï** is the name of the graphical interface :

- macro view :
  - ***all Davaï experiments from all users***
  - filtering tool among experiments
  - summary of each experiment : number and % of tests OK/KO/ ?/crashed
  - → give access to *individual experiment* : ↓
- meso view (*individual experiment*) :
  - ***all tasks ran in the experiment***
  - filtering tool among tasks
  - summary of each task : OK/KO/ ?/crashed + RSS/elapse deltas
  - → give access to *individual task* : ↓
- micro view (*individual task*) :
  - expertise about the ***task itself*** : metadata parsed by all ***experts*** involved
  - comparison to reference(s)<sup>a</sup>, by all ***experts***

a. a task can be compared to a "continuity" and/or a "consistency" reference

## DAVAI for/from outside MF

- portability :
  - Vortex : ~ OK ?
  - Olive/Vortex works on ECMWF HPC from Météo France SMS servers
- accessibility :
  - of Olive from outside MF ?
  - export of Olive experiments/tasks to standalone Vortex scripts ?

## DAVAI

