

## Towards an ORCHIDEE evaluation chain

Patricia CADULE<sup>1</sup>   Sebastián RIFFO<sup>2</sup>

Institut Pierre-Simon Laplace (IPSL), CNRS<sup>1</sup> and Sorbonne Université<sup>2</sup>

Monday 10<sup>th</sup> February, 2025

TASK FORCE (2020-22)

## TASK FORCE (2020-22)

- ▶ Methodology (variables, gathering data, scripting)

## TASK FORCE (2020-22)

- ▶ Methodology (variables, gathering data, scripting)
- ▶ Hands-on approach

## TASK FORCE (2020-22)

- ▶ Methodology (variables, gathering data, scripting)
- ▶ Hands-on approach
- ▶ Visualization expectations

## TASK FORCE (2020-22)

- ▶ Methodology (variables, gathering data, scripting)
- ▶ Hands-on approach
- ▶ Visualization expectations

## ACCOMPLISHMENTS

- ▶ Evaluation dataset (at February 2025: 13/24 variables)

## TASK FORCE (2020-22)

- ▶ Methodology (variables, gathering data, scripting)
- ▶ Hands-on approach
- ▶ Visualization expectations

## ACCOMPLISHMENTS

- ▶ Evaluation dataset (at February 2025: 13/24 variables)

## CHALLENGES

- ▶ Setting unambiguous and objective criteria for evaluation

## TASK FORCE (2020-22)

- ▶ Methodology (variables, gathering data, scripting)
- ▶ Hands-on approach
- ▶ Visualization expectations

## ACCOMPLISHMENTS

- ▶ Evaluation dataset (at February 2025: 13/24 variables)

## CHALLENGES

- ▶ Setting unambiguous and objective criteria for evaluation
- ▶ Having a reliable dataset to contrast against

## TASK FORCE (2020-22)

- ▶ Methodology (variables, gathering data, scripting)
- ▶ Hands-on approach
- ▶ Visualization expectations

## ACCOMPLISHMENTS

- ▶ Evaluation dataset (at February 2025: 13/24 variables)

## CHALLENGES

- ▶ Setting unambiguous and objective criteria for evaluation
- ▶ Having a reliable dataset to contrast against
- ▶ Development of a open, intuitive, and robust tool

# EVALUATION CHAIN

EVALUATION  
FRAMEWORK

QUALITY  
ASSESSMENT

DATA  
CONSIDERATIONS

## EVALUATION FRAMEWORK

Against what we  
contrast and why?

## QUALITY ASSESSMENT

How do we define  
*being better*?

## DATA CONSIDERATIONS

Setting standards  
for its intended use

The *notion of model evaluation* varies from team to team, in accordance to their own requirements and criteria

The *notion of model evaluation* varies from team to team, in accordance to their own requirements and criteria

### PLUMBER PROTOCOL (ABRAMOWITZ ET AL)

- ▶ *evaluation*  
determine errors with respect to observations, using statistical methods
- ▶ *comparison*  
viewing a model in relation to others, with the purpose of establishing differences in their behavior
- ▶ *benchmarking*  
identifying, from the contrast with a reference (defined in advance), what a model requires to improve

## VISUALIZATION

Useful to convey information to a wide public, yet an explanation of these results relies in the observer's experience, rather than a *settled criterion*

## VISUALIZATION

Useful to convey information to a wide public, yet an explanation of these results relies in the observer's experience, rather than a *settled criterion*

## QUANTIFICATION

Delivering a *measurable* standard depends on a shared space-time domain

- ▶ *metrics*

measures (statistics), indicators (process-oriented), scores (as ILAMB)

## VISUALIZATION

Useful to convey information to a wide public, yet an explanation of these results relies in the observer's experience, rather than a *settled criterion*

## QUANTIFICATION

Delivering a *measurable* standard depends on a shared space-time domain

- ▶ *metrics*

measures (statistics), indicators (process-oriented), scores (as ILAMB)

Decision-making then based on

- ▶ *orders of magnitude*

- ▶ *sensitivity*

to recognize relevant variables in the model's behavior

Data's nature and its processing have an *overlooked impact* in model quality assessment

Data's nature and its processing have an *overlooked impact* in model quality assessment

### OPERATIONAL DATASET

Setting standards around

- ▶ *scale*  
original resolution, spatial and temporal coverage
- ▶ *regridding*  
upsampling (interpolations) or downscaling (e.g. averaging) effects
- ▶ *reliability*  
accounts for variability and uncertainty, using an ensemble of realizations or sources

## EVALUATION FRAMEWORK

- ▶ evaluation
- ▶ comparison
- ▶ benchmarking

## QUALITY ASSESSMENT

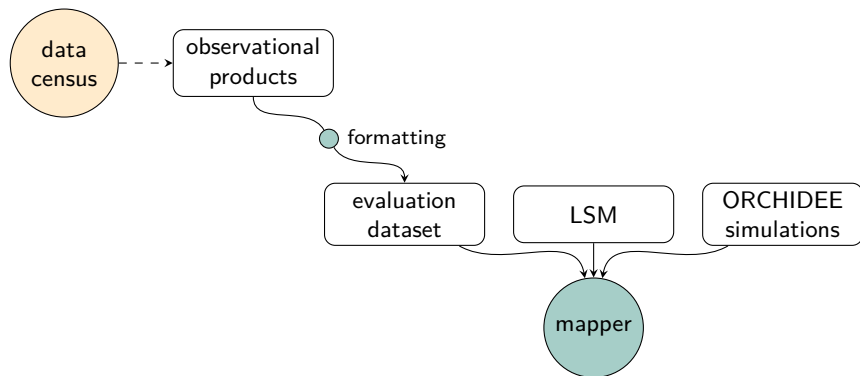
- ▶ visualization
- ▶ metrics
- ▶ orders of magnitude
- ▶ sensitivity

## DATA CONSIDERATIONS

- ▶ scale
- ▶ regridding
- ▶ reliability

*The tool should « be open (...), documented enough (...), easy [exporting the figures for possible publication]; not require (...) to have specific knowledge of where datasets are hosted (...); provide a basic set of typical evaluation techniques (...) »*

Task Force (2021)



## MOMA (METRICS FOR ORCHIDEE MODEL ANALYSIS)

- ▶ Web application running on spirit
- ▶ Interactive and dynamic visualization
- ▶ Real-time catalog
- ▶ Statistical measures

## MOMA (METRICS FOR ORCHIDEE MODEL ANALYSIS)

- ▶ Web application running on spirit
- ▶ Interactive and dynamic visualization
- ▶ Real-time catalog
- ▶ Statistical measures

## TECHNICAL STACK

- ▶ Front end  
JavaScript (D3, Gridstack, DataTables)
- ▶ Back end  
Python (Flask; Xarray, Numpy, Dask)
- ▶ Version control  
Git

# METRICS FOR ORCHIDEE MODEL ANALYSIS

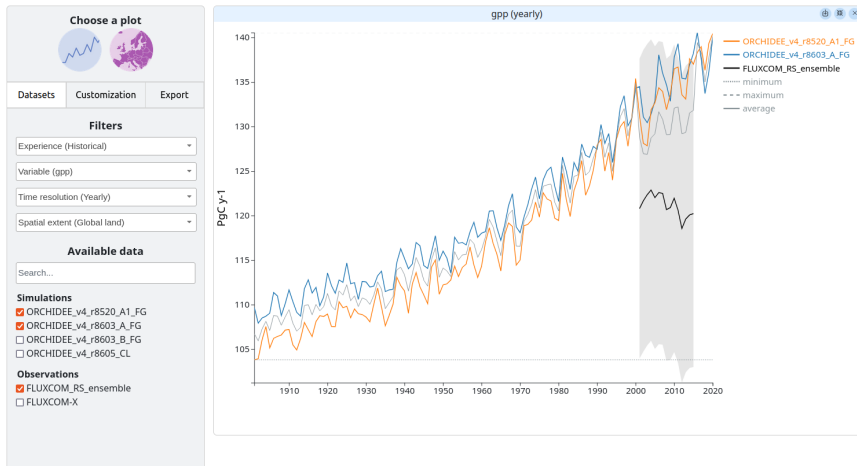


Figure 1: Front end

# METRICS FOR ORCHIDEE MODEL ANALYSIS

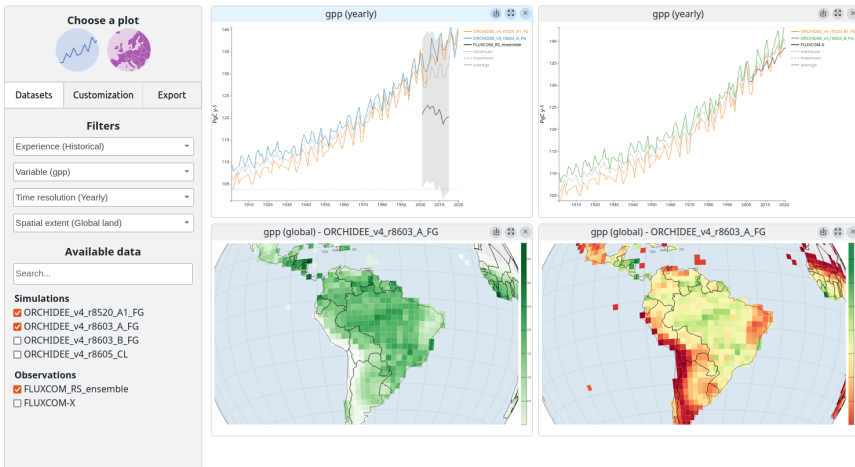
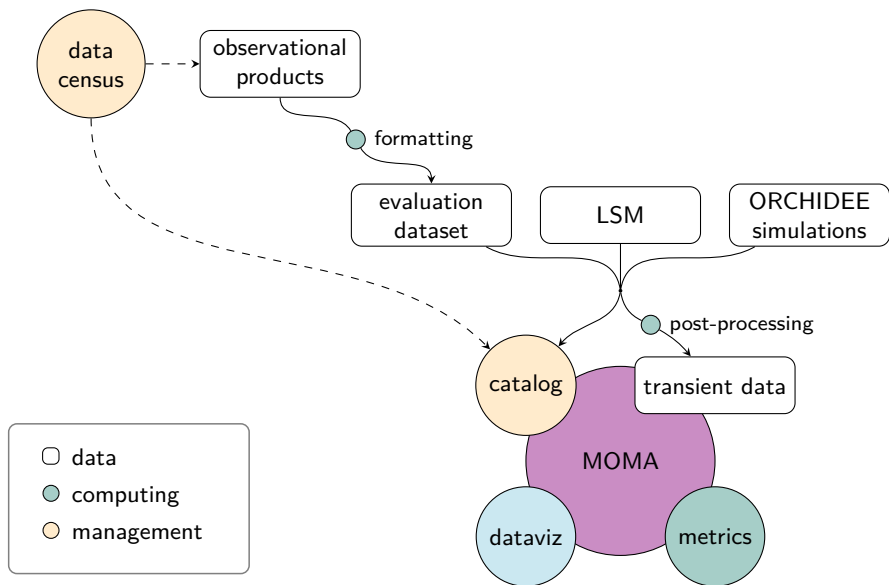
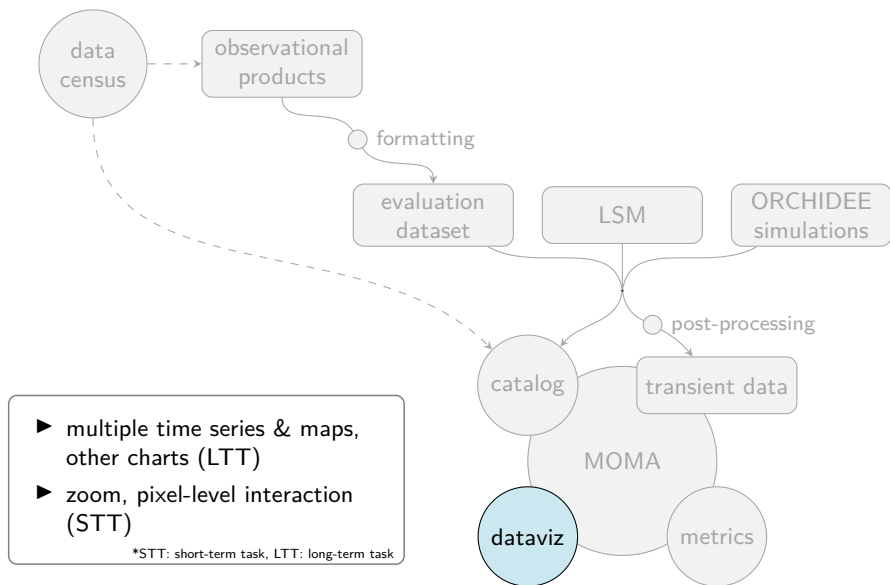


Figure 2: Front end

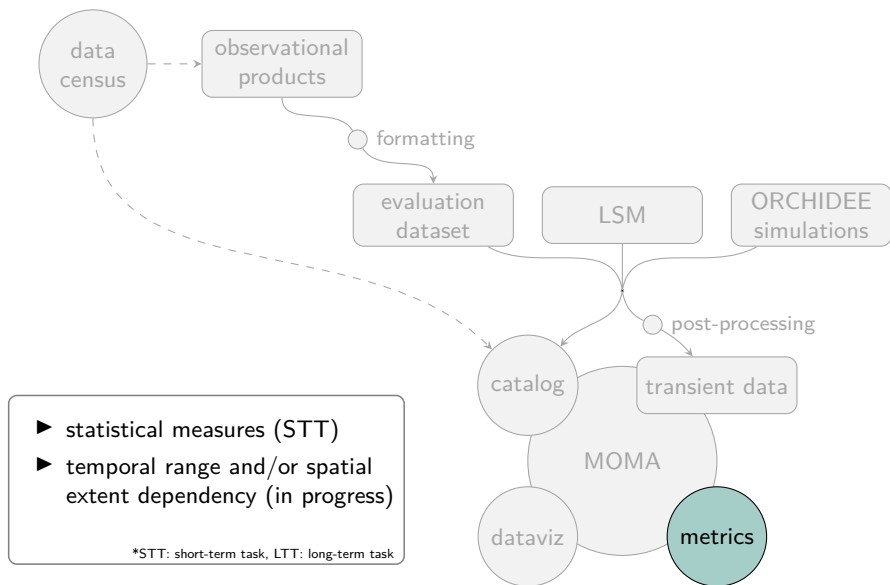
# METRICS FOR ORCHIDEE MODEL ANALYSIS



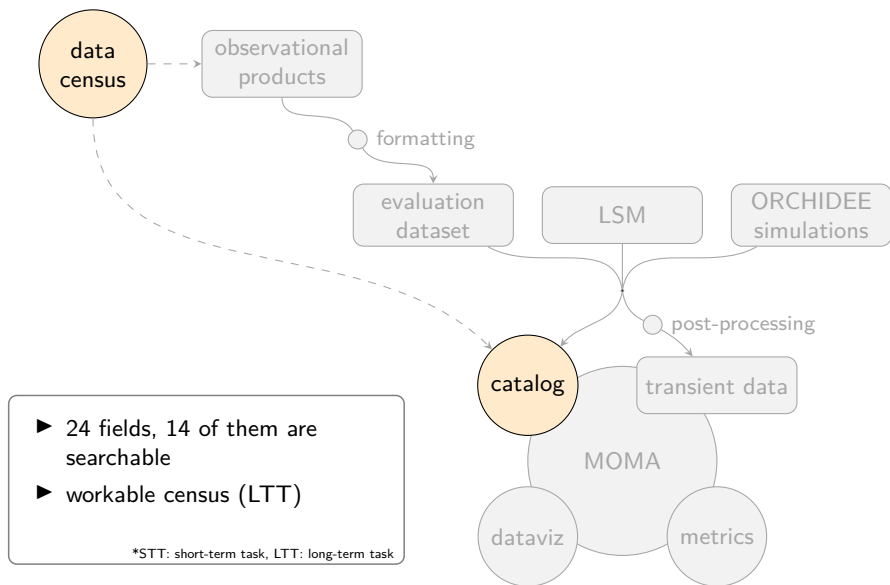
# METRICS FOR ORCHIDEE MODEL ANALYSIS



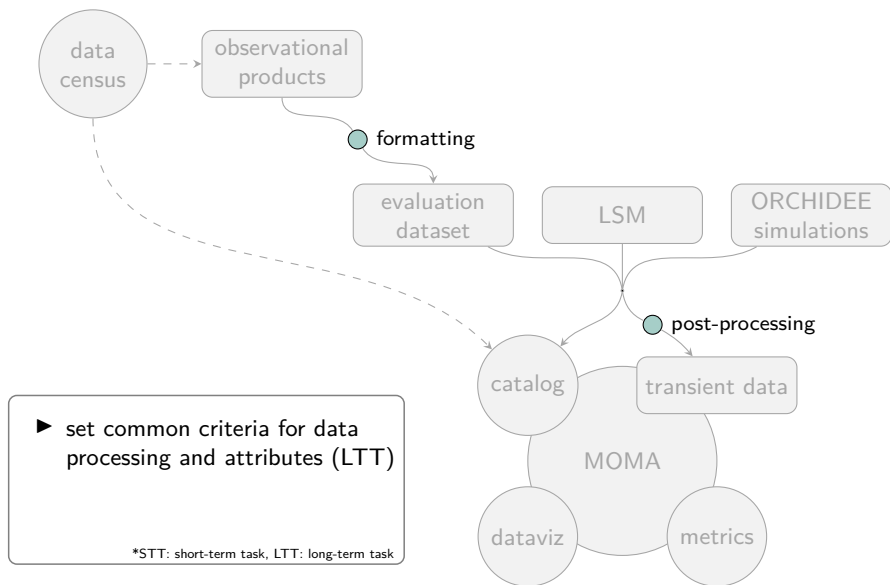
# METRICS FOR ORCHIDEE MODEL ANALYSIS



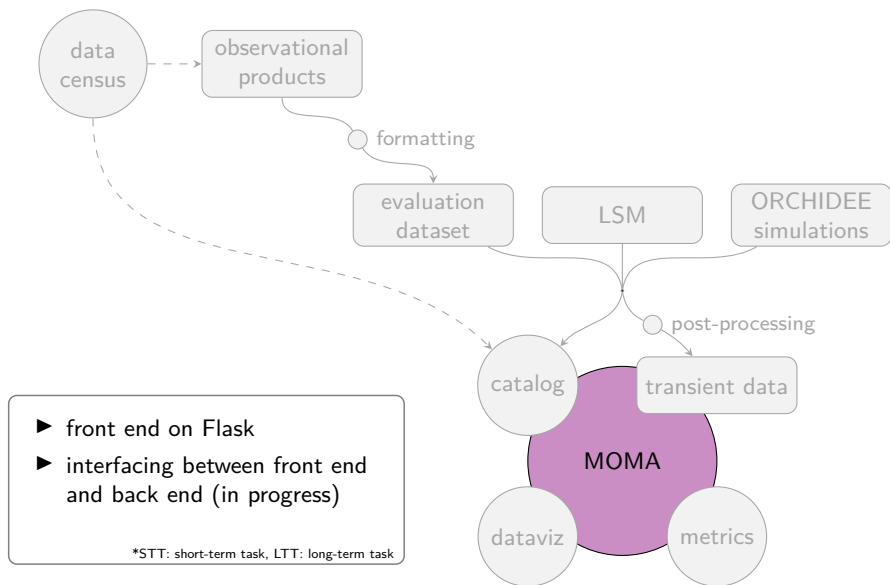
# METRICS FOR ORCHIDEE MODEL ANALYSIS



# METRICS FOR ORCHIDEE MODEL ANALYSIS



# METRICS FOR ORCHIDEE MODEL ANALYSIS



### NATIONWIDE

- ▶ CliMAF & C-ESM-EP (CNRM and IPSL)
- ▶ ILAMB (RUBISCO, US)
- ▶ AMBER (CLASSIC, Canada)
- ▶ PMP (LLNL, US)
- ▶ LVT (NASA, US)

### INTERNATIONAL

- ▶ ESMValTool
- ▶ PLUMBER2 MIP

## WIKI

<https://forge.ipsl.jussieu.fr/orchidee/wiki/Documentation/Validation>

## REPORT (IN PROGRESS)

- ▶ Documentation
- ▶ Publication

Thank you for your attention !