



HAL
open science

Assessing the performance and robustness of two conceptual rainfall-runoff models on a worldwide sample of watersheds

Thibault Mathevet, Hoshin Gupta, Charles Perrin, Vazken Andréassian,
Nicolas Le Moine

► To cite this version:

Thibault Mathevet, Hoshin Gupta, Charles Perrin, Vazken Andréassian, Nicolas Le Moine. Assessing the performance and robustness of two conceptual rainfall-runoff models on a worldwide sample of watersheds. 4es Rencontres HydroGR, Dec 2021, Antony, France. hal-03537150

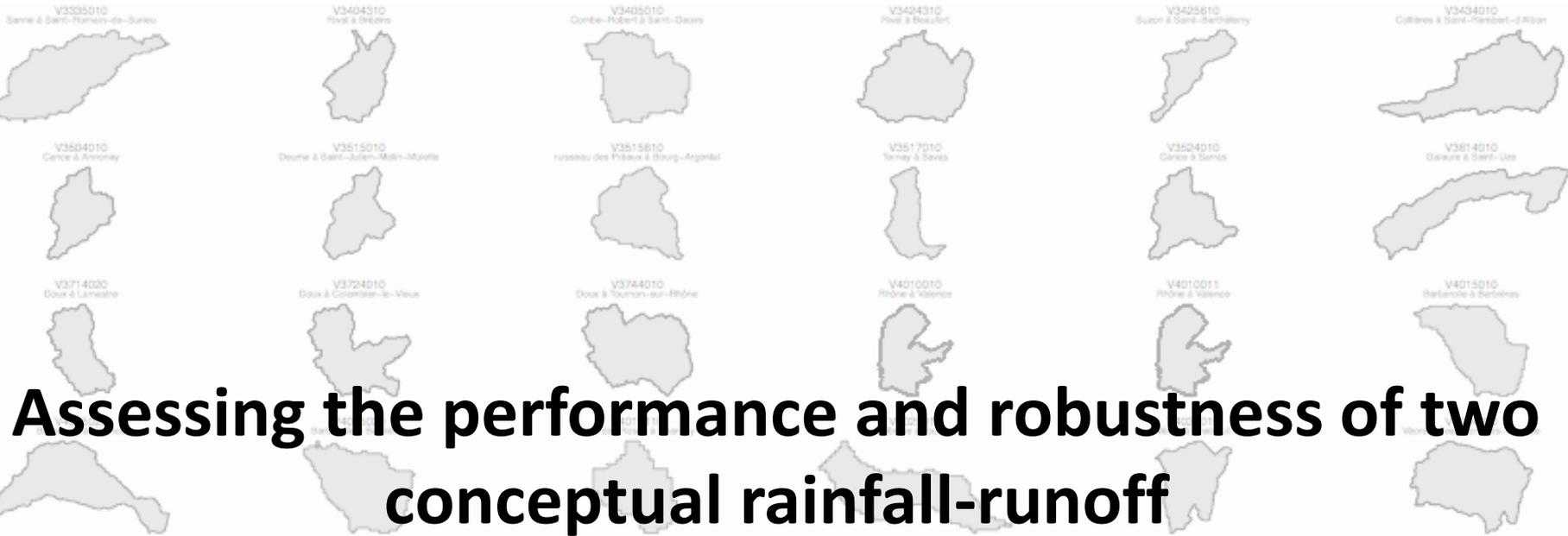
HAL Id: hal-03537150

<https://hal.inrae.fr/hal-03537150>

Submitted on 20 Jan 2022

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Assessing the performance and robustness of two conceptual rainfall-runoff models on a worldwide sample of watersheds

Thibault MATHEVET¹, Hoshin GUPTA², Charles PERRIN³, Vazken ANDRÉASSIAN³, Nicolas LE MOINE⁴
thibault.mathevet@gmail.com

¹ Visiting scientist at Hydrology & Water Resources Dep., Univ. of Arizona (2013/2014)

Modélisation économique des aménagements hydroélectriques, EDF

² Hydrology & Atmospheric sciences Dep., Univ. of Arizona

³ Université Paris-Saclay, INRAE, UR HYCAR, Antony, France

⁴ SorbonneUniversité, UMR Metis, Paris, France



Acknowledgments

Vazken Andr assian
(IRSTEA)

Hoshin Gupta



MOPEX 2004, Paris

Charles Perrin
(IRSTEA)



Background of this study

« How different are different models ? »

« I sometimes trust more *my* model than the observations »

« If *my* model can't make it, almost no model can make it »

• **Question 1: How statistically comparable (based on a detailed evaluation procedure) are the simulation performances of two models?**

• **Question 2: Is the simulation performance of the models essentially identical when provided with the same observational information?**

• **Question 3: Are differences in model performance dependent on watershed characteristics or on hydrometeorological processes?**



Research papers
Assessing the performance and robustness of two conceptual rainfall-runoff models on a worldwide sample of watersheds

Thibault Mathevet^{a,*}, Hoshin Gupta^b, Charles Perrin^c, Vazken Andréassian^d, Nicolas Le Moine^d

^a EDF ERDF (Bureau de Paris), 134 chemin de l'Étang, 92850 Saint Martin de Valmy, France
^b Department of Hydrology and Atmospheric Sciences, University of Arizona, Tucson, AZ, USA
^c Institut Paris-Saclay, IRM4, UR STIC4, Antony, France
^d Sorbonne Université, CNRS Météo, Paris, France

ARTICLE INFO

This manuscript was handled by A. Hanboley, Editor in Chief, with the assistance of Ingrid Moussa, Associate Editor
Keywords:
Hydrological modeling
Large sample hydrology
Calibration
Evaluation
Diagnostic
Kling Gupta efficiency

ABSTRACT

To assess the predictive performance, robustness and generality of watershed-scale hydrological models, we conducted a detailed multi-objective evaluation of two conceptual rainfall-runoff models (the GR2 model, based on the GR4 model, and the MIX model, based on the MORDDOR model), of differing complexity (with respectively, 5 and 11 free parameters in the rainfall-runoff module, and 4 and 11 free parameters in the snow module). These models were compared on a large sample of 2050 watersheds worldwide. Our results, based on the three components of the Kling Gupta Efficiency metric (KGE), indicate that both models provide (on average) similar levels of performance in evaluation when calibrated with KGE, for water balance (mean bias lower than 2%), time-series variability (mean variability bias lower than 2%) and temporal correlation (mean correlation around 0.8). Further, both models clearly suffer from lack of robustness when simulating water balance, with a significant increase of the proportion of biased simulations over the evaluation period (absolute bias lower than 2% in calibration and lower than 20% in evaluation for 80% of the watersheds). Simulation performance depend more on the hydro-meteorological conditions of a given period than on the complexity of the model structure. We also show that long-term aggregate statistics (computed on the overall period) can fail to reveal considerable sub-period variability in model performance, thereby providing inaccurate diagnostic assessments of the predictive model performance. Typically the median absolute bias is lower than 8% in evaluation, but the median maximum bias can be as high as 50% within a sub-period, for both models, when calibrated with KGE.

1. Introduction

Rainfall-runoff (RR) models are widely used for a broad range of research and operational objectives, from hypothesis testing to improving process understanding to streamflow prediction for flood design. Whatever the application, hydrologists and modelers share a particular interest in: (i) the efficiency, robustness and realism of model structures (and their consequent simulations); (ii) the generality (transferability) of model structures across locations (i.e. ability to be efficient in a variety of hydroclimatic contexts); and (iii) methods for parameter identification (Gupta et al., 2014). To achieve these objectives, a variety of strategies for model development and specification have been pursued, ranging from detailed site-specific investigations to more general studies. The term robustness is often used to describe some expected model properties in a broad sense. Here, robustness is understood as the capability of a model to hold a certain level of

performance in changing conditions, i.e. independently from the input/output information used for calibration. Robustness is usually assessed by comparing the difference of evaluation metrics under changing conditions (typically from calibration to evaluation periods, but also from dry to wet conditions, etc.).

The investigations discussed in this paper are rooted in the past experience of the authors with RR model intercomparison studies (Perrin et al., 2001; Perrin et al., 2003; Perrin et al., 2006; Le Moine et al., 2007; Pushpalatha et al., 2011, 2012; Coron et al., 2012, 2014), as well as investigations into diagnostic model identification procedures (Gupta et al. 2008, 2009; Gupta et al., 2012; Vilmarz et al., 2008; Martinez and Gupta, 2010, 2011; de Vos et al., 2010; Pohlken et al., 2012).

* Corresponding author.
E-mail address: thibault.mathevet@edf.fr (T. Mathevet).
¹ Visiting research scholar at Hydrology and Atmospheric Sciences, University of Arizona, in 2014.
<https://doi.org/10.1016/j.jhydrol.2020.124698>
Received 24 October 2019; Received in revised form 22 January 2020; Accepted 14 February 2020
Available online 19 February 2020
0022-1694/© 2020 Published by Elsevier B.V.

Why large sample hydrology?

- **Improving understanding:**
more rigorous testing and comparison of competing model hypotheses and structures on common grounds;
- **Improving the robustness of generalizations:**
allowing statistical analyses of model performances and avoid giving too much weight to outliers;
- **Facilitating classification, regionalization and model transfer:**
gathering a wide diversity of hydrometeorological contexts, enabling testing classification and regionalisation strategies;
- **Supporting the estimation of uncertainties:**
establishing the predictive capabilities and performance of hydrological models on a variety of hydrometeorological contexts.

Hydro. Earth Syst. Sci., 18, 463–477, 2014
www.hydro-earth-syst-sci.net/18/463/2014/
doi:10.5194/hess-18-463-2014
© Author(s) 2014. CC Attribution 3.0 License.



Large-sample hydrology: a need to balance depth with breadth

H. V. Gupta¹, G. Perrin², G. Blöschl³, A. Montanari⁴, R. Kumar⁵, M. Clark⁶, and V. Andreoloni⁷
¹Department of Hydrology and Water Resources, The University of Arizona, Tucson, AZ, USA
²Inria, Hydrosystems and Negresses Research Unit (H2AN), Antony, France
³Institute of Hydraulic Engineering and Water Resources Management, Vienna University of Technology, Vienna, Austria
⁴Department DICAM, University of Bologna, Bologna, Italy
⁵131Z – Hydrology Centre for Environmental Research, Leipzig, Germany
⁶Hydro-meteorological Applications Program, Research Applications Laboratory, Boulder, CO, USA

Correspondence to: H. V. Gupta (hgupta@arizona.edu)

Received: 26 June 2013 – Published in Hydro. Earth Syst. Sci. Discuss.: 12 July 2013
Revised: 18 December 2013 – Accepted: 26 December 2013 – Published: 6 February 2014

Abstract. A holy grail of hydrology is to understand catchment processes well enough that models can provide detailed simulations across a variety of hydrologic settings at multiple spatiotemporal scales, and under changing environmental conditions. Clearly, this cannot be achieved only through intensive place-based investigation at a small number of heavily instrumented catchments, or by empirical methods that do not fully exploit our understanding of hydrology. In this opinion paper, we discuss the need to actively promote and pursue the use of a “large catchment sample” approach to enabling the “multi-scale” process, thereby balancing depth with breadth. We summarize the history of such investigations, discuss the benefits (improved process understanding resulting in robustness of prediction at ungauged locations and under change), examine some practical challenges to implementation and, finally, provide perspectives on issues that need to be taken into account as we move forward. Ultimately, our objective is to provide further discussion and participation, and to promote a potentially important theme for the upcoming Scientific Decade of the International Association of Hydrological Sciences entitled *Fans Data*.

and that each trial cover a period of many years.” (Linsley, 1982).

1.1 Motivations for developing large-sample hydrology

A holy grail of hydrological science is to achieve a degree of process understanding that enables construction of models that are capable of providing detailed and physically realistic simulations across a variety of different hydrologic environments, and at multiple spatial and temporal scales (Osh and Srinivasan, 1970; Klemes, 1986a; Michel et al., 2006). With its focus on reducing predictive uncertainty, the Prediction in Ungauged Basins (PUB) initiative of the International Association of Hydrological Sciences (IAHS) has helped move the culture of hydrologic science closer to this objective, and away from a reliance upon universal models applied via re-calibration at each study location (Hochrainer et al., 2013). This move has deep roots (see Linsley, 1982), but has become considerably stronger since the 1999 IAHS meeting in Strasbourg where the idea was extensively discussed. It has helped drive the search for improved understanding of the hydrological cycle, and for modeling approaches that

1. Introduction
- “Because almost any model with sufficient free parameters can yield good results when applied to a short sample from a single catchment, effective testing requires that the model be tried on many catchments of widely differing characteristics,
- a. achieve the three R’s (reliability, robustness and realism);
- b. have greater generality and transparency; and for which
- c. the parameters can be more easily specified from data.

Published by Copernicus Publications on behalf of the European Geosciences Union.

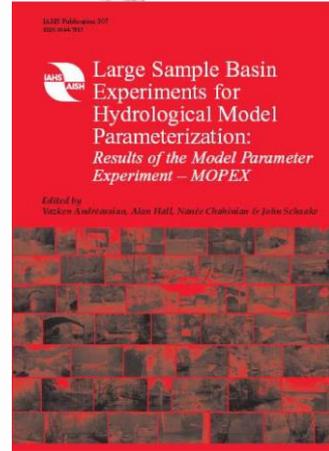
Why large sample hydrology ?

Improving understanding:

- 1 – What are the respective performances of different RR model structures ?
- 2 – Are the performances of RR structure dependant of watershed characteristics, climatological or hydrological processes ?

Improving the robustness of generalizations:

- How to properly compare two (n) RR model structures ?
- How can I state than two (n) RR structures are different ?



A bounded version of the Nash-Sutcliffe criterion for better model assessment on large sets of basins

THEBAULT MATHIERET¹, CLAUDE MICHEL¹, VALERIE ANNEHEIMAN² & CHARLES FERRIS¹

¹Université de Lyon, UMR 5175 Hydrosciences, France
²Université de Bourgogne, UMR 5076 Hydrosciences, France

Abstract: Hydrological model results are useful tools for hydrological research, water engineering and environmental applications. Given the large number of available hydrological model results, model comparison studies have been done to compare model performances. In 1970, the Nash-Sutcliffe criterion was proposed as a simple and effective way to compare model performances. This paper requires large samples of test catchments. However, large test samples may be scarce in some basins or specific model performances. This paper shows that the use of large test sets and the bounded version of the Nash-Sutcliffe criterion can be used to compare model performances. This paper also shows that the bounded version of the Nash-Sutcliffe criterion can be used to compare model performances. This paper also shows that the bounded version of the Nash-Sutcliffe criterion can be used to compare model performances.

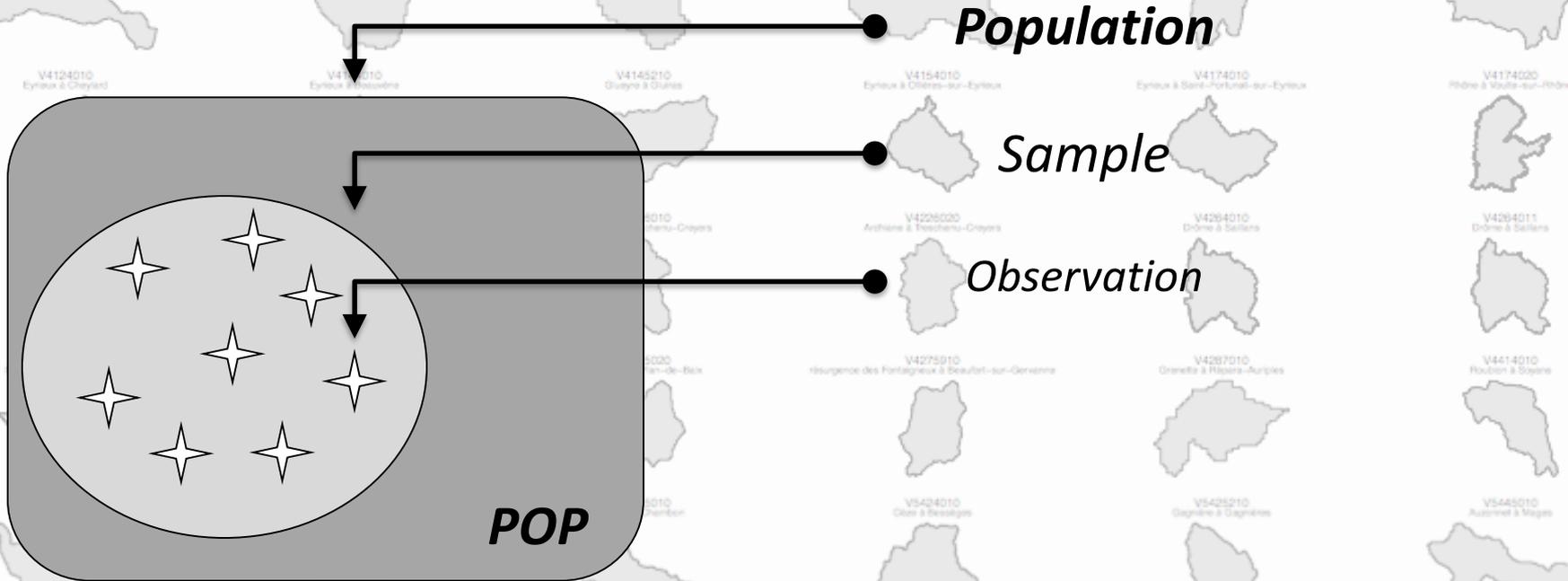
INTRODUCTION

Since the 1960s, hydrologists have developed a large number of more or less complex rainfall-runoff models. In a consequence of the proliferation, the need for comparative studies appeared quite early (Cohn, 1976; Linsley (1982) suggested that "because almost any model with sufficient free parameters can yield good results when applied to a short sample from a single basin, effective testing requires that models be tested on many basins of widely differing characteristics, and that each test cover a period of many years". Few modelers have, however, followed these recommendations, and most of the RR models studies reported in the literature present the performances of one RR model on a single basin (or on a small number of similar basins). If studies include two or three basins, the validity of these conclusions will be limited to the hydro-climatic domain of the test sample. Conversely, a large set of basins provides a global overview of the efficiency of one or several models, in a wide range of hydro-climatological conditions and watershed physical characteristics (topology, soil, vegetation, topography, land use, etc.). Model performance becomes complex since a complete distribution of results is obtained, often over a large range of performance. In such cases, one must find ways to compare flow distributions, or to summarize them into proper statistics, provided that the formulation of the selected criteria does not bias the results.

In this paper we propose a criterion formulation suitable for comparing model performances on large basin samples. An application is made on a sample of 313 basins. We also show the usefulness of large basin sets for model assessment.

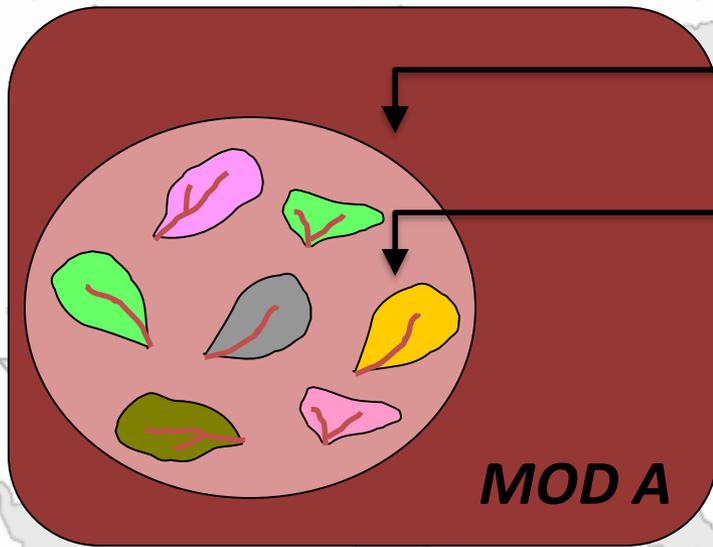
Why large sample hydrology ?

- Allowing statistical comparison of RR model structures
 - Infer the properties of a population from a sample of observations



Why large sample hydrology ?

- Allowing statistical comparison of RR model structures
 - Infer the properties of a population from a sample of observations



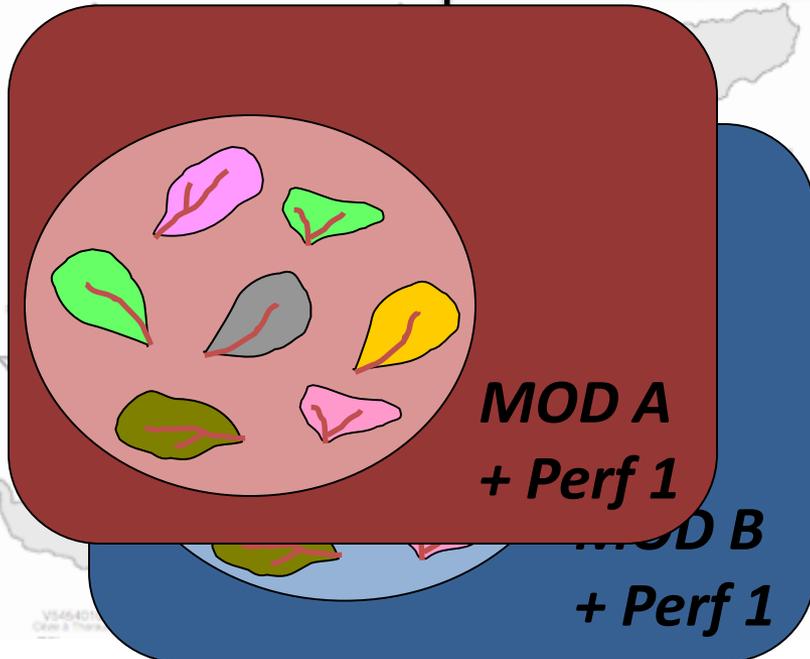
• **RR Model A performance**

Distribution of perf. on a sample of n watersheds

• *Performance on Watershed i*

Why large sample hydrology ?

- Allowing statistical comparison of RR model structures
 - Infer the properties of a population from a sample of observations

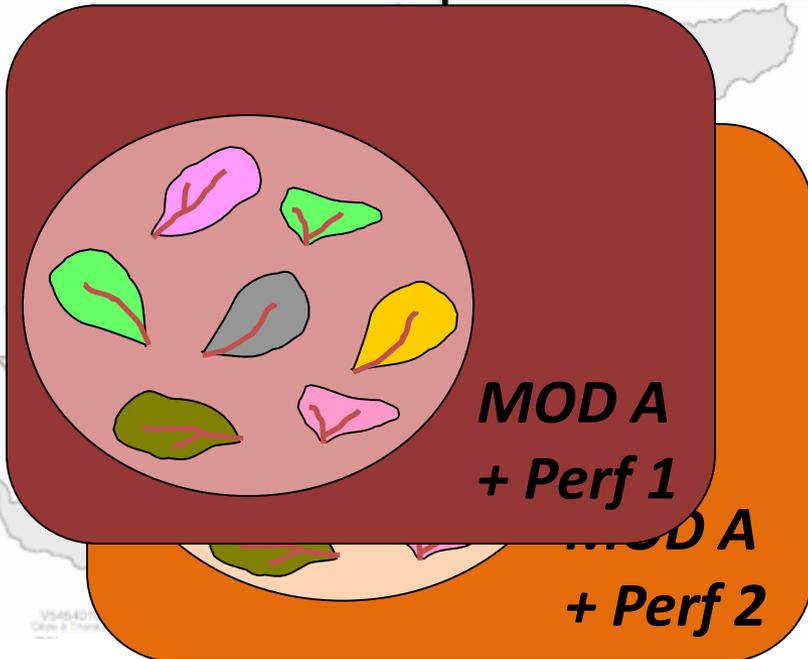


Why large sample hydrology ?

- Allowing statistical comparison of RR model structures
 - Infer the properties of a population from a sample of observations

RR Model A perf. 1

RR Model A perf. 2



Insights from previous studies (1/2)

- **Perrin et al. 2001 :**

- 20 RR model structures, +400 watersheds, daily time-step, NSE ;
- Complex models suffers from a lack of robustness and 4-6 free parameters seems sufficient to give the « best » results ;

- **Mathevet et al., 2006 :**

- 4 RR model structures, +300 watersheds, hourly time-step, NSE + modification ;
- NSE do not allow robust statistical comparisons ;
- Framework to state if two RR structures performances are significantly different or not ;

- **Coron et al. 2011 :**

- 3 RR model structures, +200 watersheds, daily time-step, 2 performance metrics ;
- RR model are extremely dependent to climatic conditions during calibration and have a strong lack of robustness when evaluated on contrasted climatic periods ;

Insights from previous studies (2/2)

- **Fenizia et al. 2011, Kavetski et al. 2011 :**

- SUPERFLEX : flexible modeling framework, with a collection of conceptual structures and constitutive functions ;
- Hypothese of a better representation of underlying « true » hydrological processes ;

- **van Esse et al. 2013 (including Perrin & Fenizia) :**

- 30 RR model structures, +200 watersheds, hourly time-step, 4 performance metrics ;
- **Allmost no difference between a flexible modeling (SUPERFLEX) and a fixed modeling (GR4H) framework ;**

- **Gupta et al. 2009, Gupta & Kling 2011:**

- Nash-Sutcliffe Efficiency is not an accurate objective function for RR model calibration ;
- Bias on the water balance and the variability of streamflows ;
- Introduction of the Kling-Gupta Efficiency (KGE) ;

- **To be updated**

Experimental design (1/4)

- **A (very) large sample of watersheds :**

- Collect samples already used in litterature (*Chiew et al., 2000; Duan et al., 2006; Le Moine et al., 2008 ; Vaze et al. 2010; Coron et al., 2011 ; Valery et al., 2009 & 2010; Nicolle et al. 2014 ; Top-Down modeling working group*);
- French national projects (PEMHYCE : *Nicolle et al. 2014*; R2D2 : *Kuentz, 2013*);
- « My » sample at EDF ;

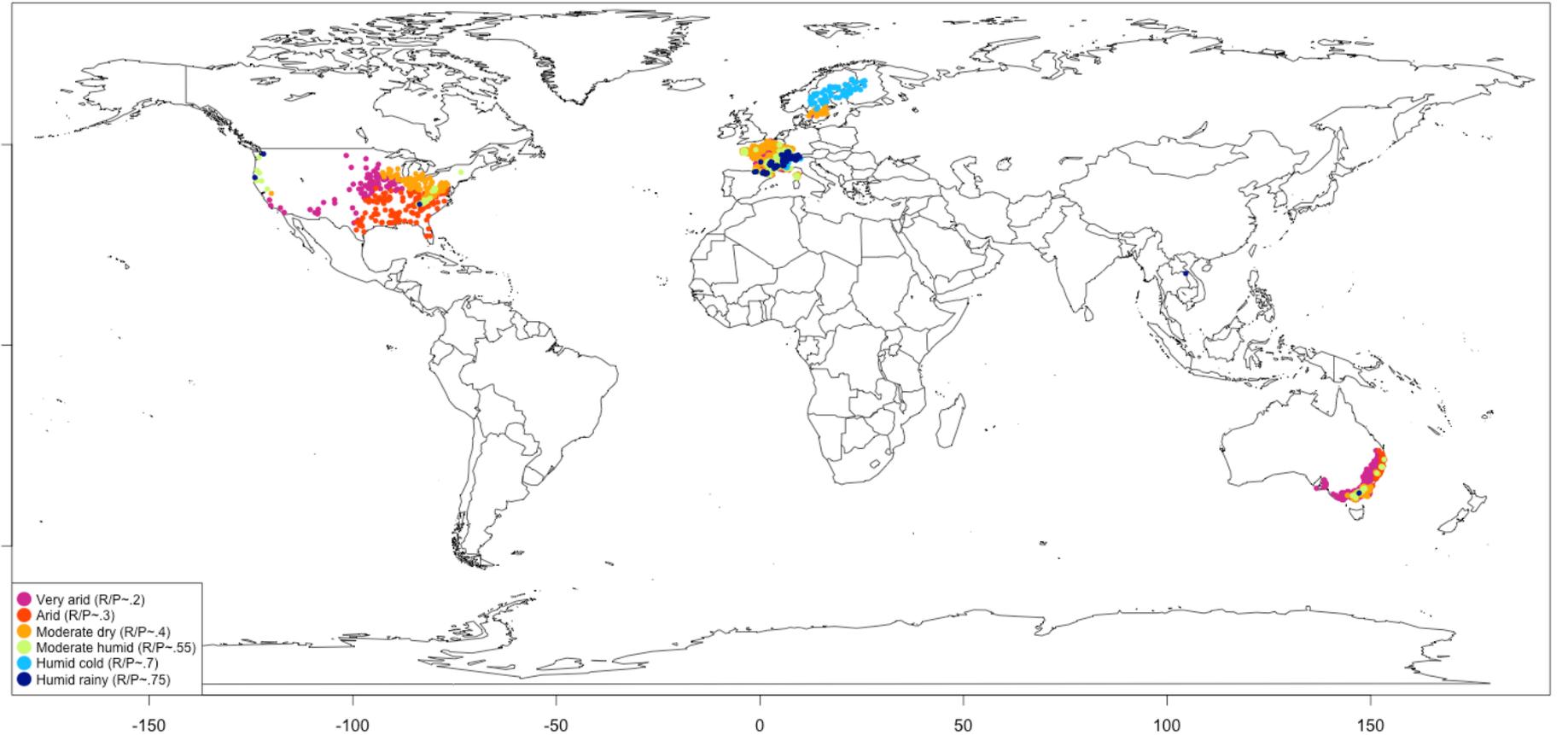
- **2050 watersheds worldwide (+ ~200 not used):**

- France, USA, Australia (80%);
- Switzerland, Sweden, UK, Laos, Italy (20%) ;

- **Since this study :**

- Many open-source & unified hydrometeorological samples ;
- **Camels** initiatives largely supported by N. Addor & colleagues (USA, UK, NZ, Chile, Brasil, Australia, etc.) ;

Experimental design (1/4)



Experimental design (2/4)

- **2 Rainfall-Runoff model structures :**

- Used in many different comparative studies since 2004 ;
- Statistically the most efficient among 20 different RR on hundreds of watersheds;

GRX (IRSTEA/Cemagref, Paris)

- Empirical development on 100 to 1000 of watersheds worldwide
- 2 buckets
- 5 free parameters
- Undergroud exchanges function
- PET based on Tair and extra-terrestrial radiation
- Snow : 2 buckets & 4 free param.

MRX (EDF / Grenoble)

- Conceptual development on <10 watersheds in the Alps
- 4 buckets
- 11 free parameter
- No Undergroud exchanges function
- « optimised » PET
- Snow : 2 buckets & 11 free param.

Experimental design (3/4)

- Evaluation metrics :

- **NSE : Nash-Sutcliffe efficiency** (Nash & Sutcliffe., 1970)

$$NSE_Q = 1 - \frac{MSE}{\sigma_Q^2} = 1 - \frac{\frac{1}{n} \sum_{i=1}^n (Q_i - \hat{Q}_i)^2}{\frac{1}{n} \sum_{i=1}^n (Q_i - \bar{Q})^2}$$

- **KGE : Kling-Gupta efficiency** (Gupta et al., 2009)

$$KGE_Q = 1 - \sqrt{(\beta - 1)^2 + (\alpha - 1)^2 + (r - 1)^2}$$

$\beta = \frac{\hat{Q}}{Q}$ $\alpha = \frac{\hat{\sigma}_Q}{\sigma_Q}$ $r = \text{Linear correlation}$

And also (Kling et al., 2012): $\gamma = \frac{\hat{\sigma}_Q}{\sigma_Q} / \frac{\hat{Q}}{Q} = \frac{\alpha}{\beta}$

Experimental design (4/4) :

- **Classical Split-sample test** (Klemes, 1986) : 2 periods of calibration and 2 periods of validation

Calibration (P1)

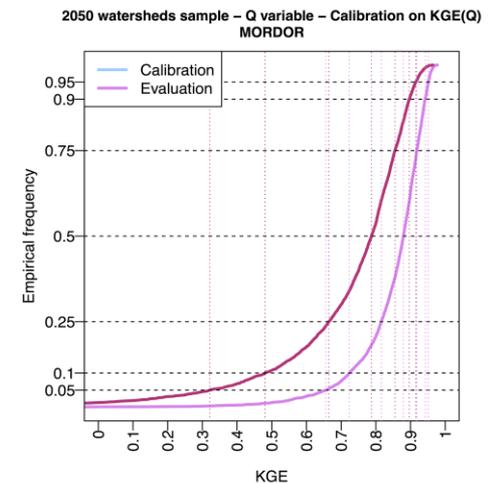
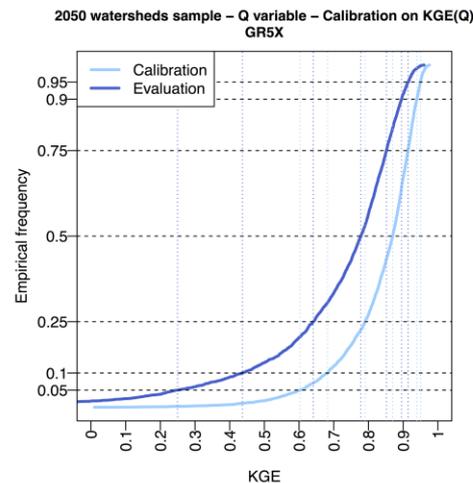
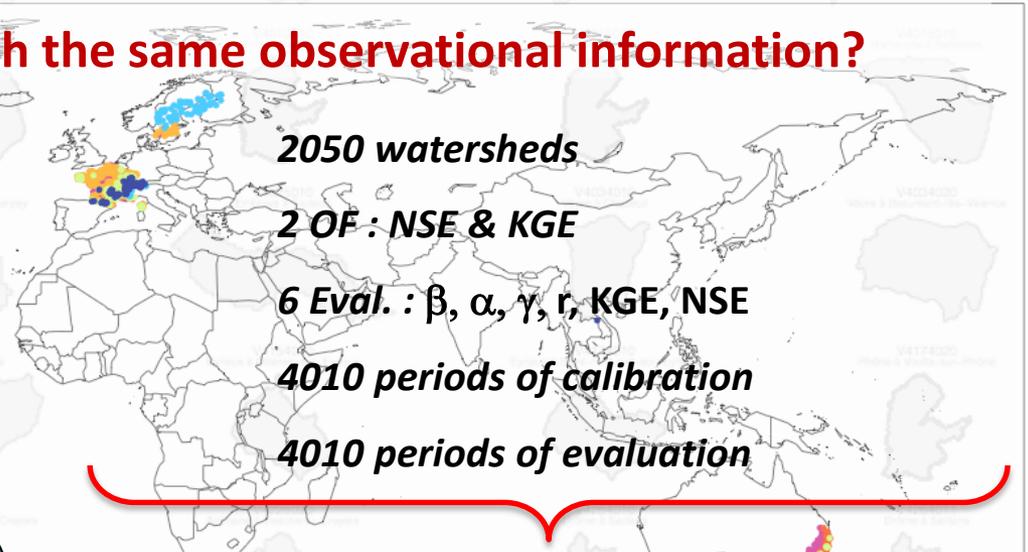
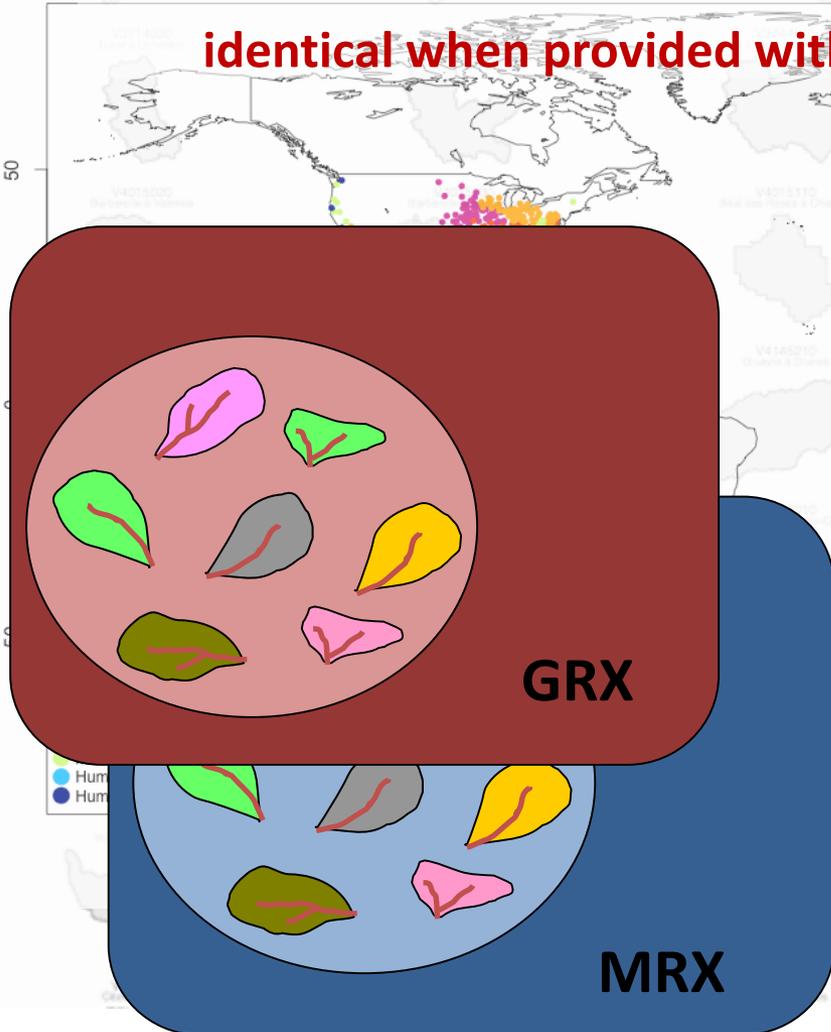
Evaluation (P2)

Calibration (P2)

Evaluation (P1)

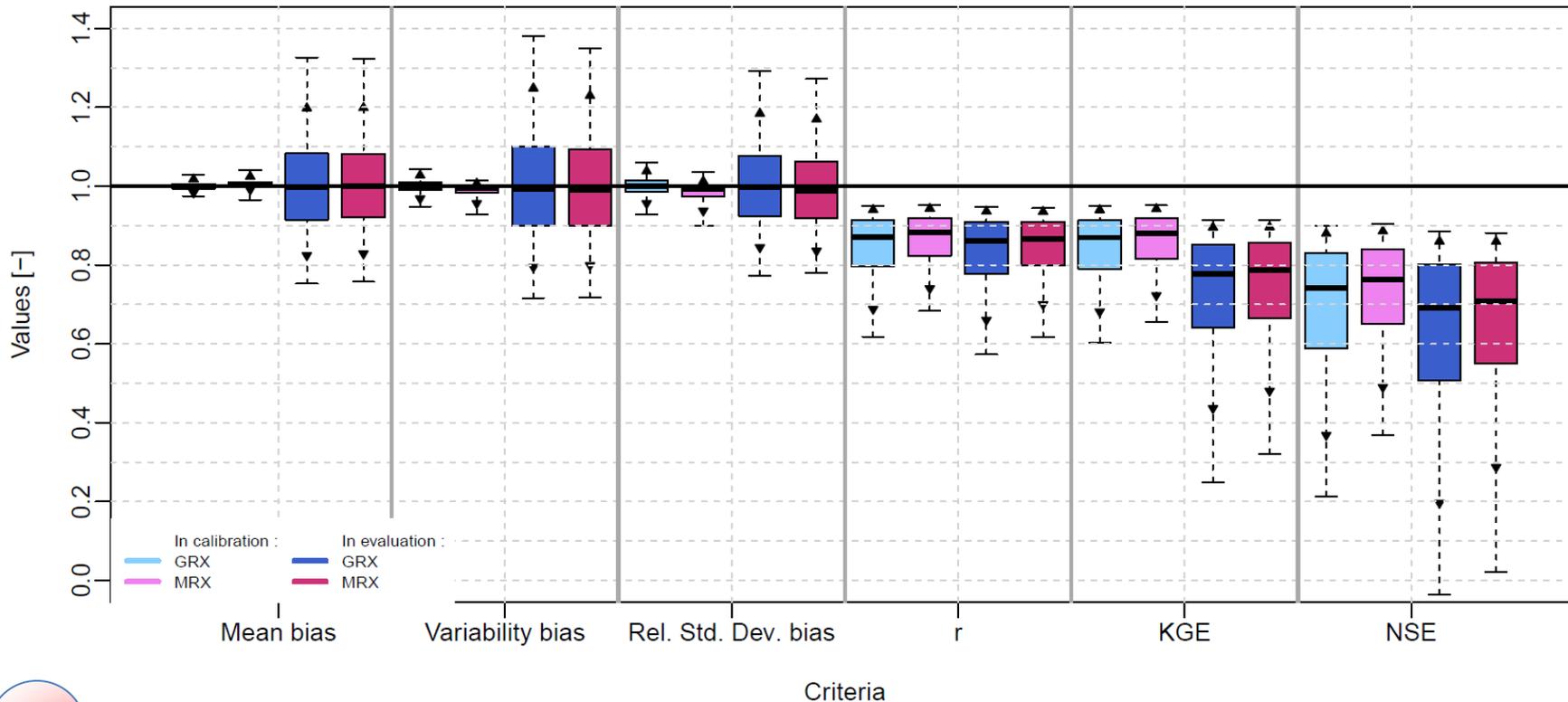


Question 1: How statistically comparable (based on a detailed evaluation procedure) are the simulation performances of two models? & Question 2: Is the simulation performance of the models essentially identical when provided with the same observational information?



Results : Boxplots

Calibration on KGE(Q)
Calibration & Evaluation

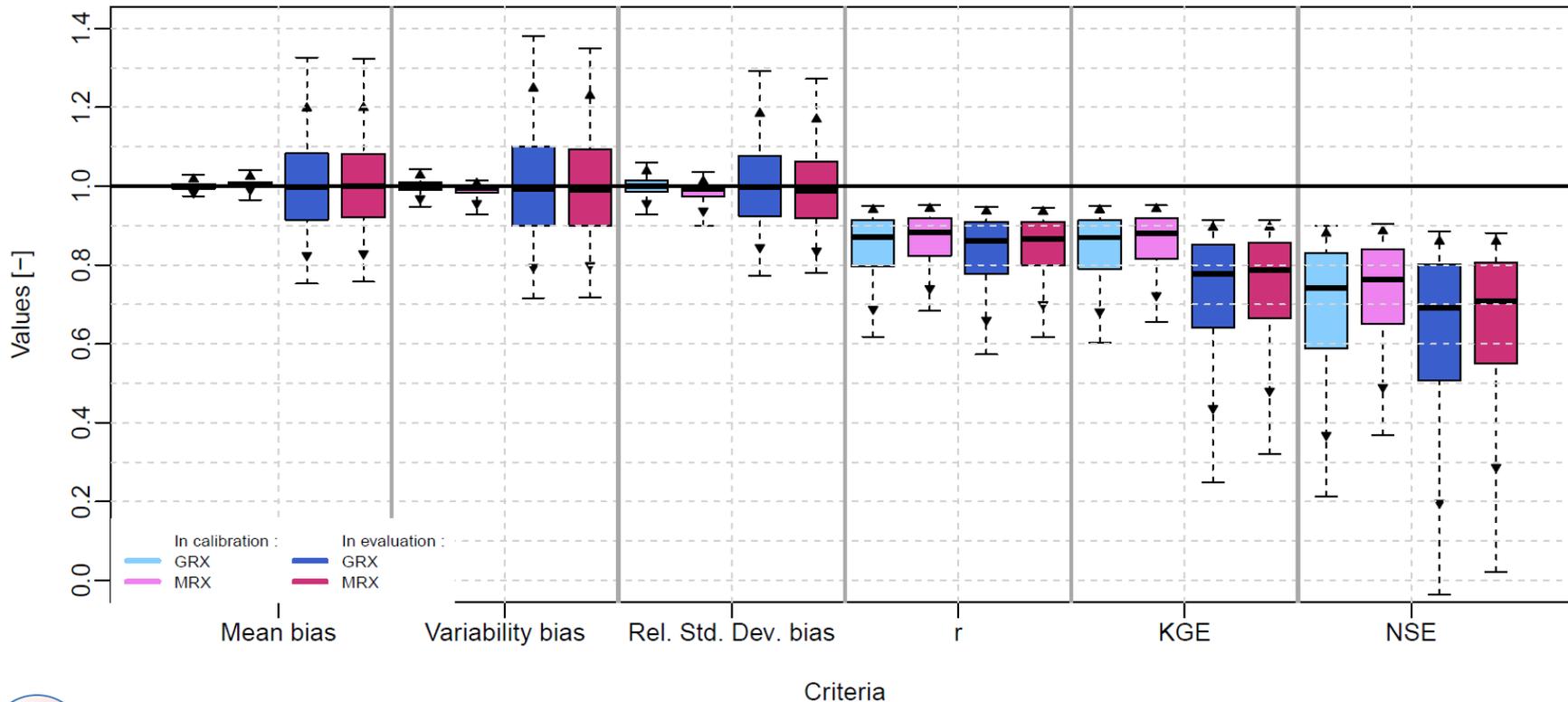


0

How to read these boxplots ?

Results : Boxplots

Calibration on KGE(Q)
Calibration & Evaluation



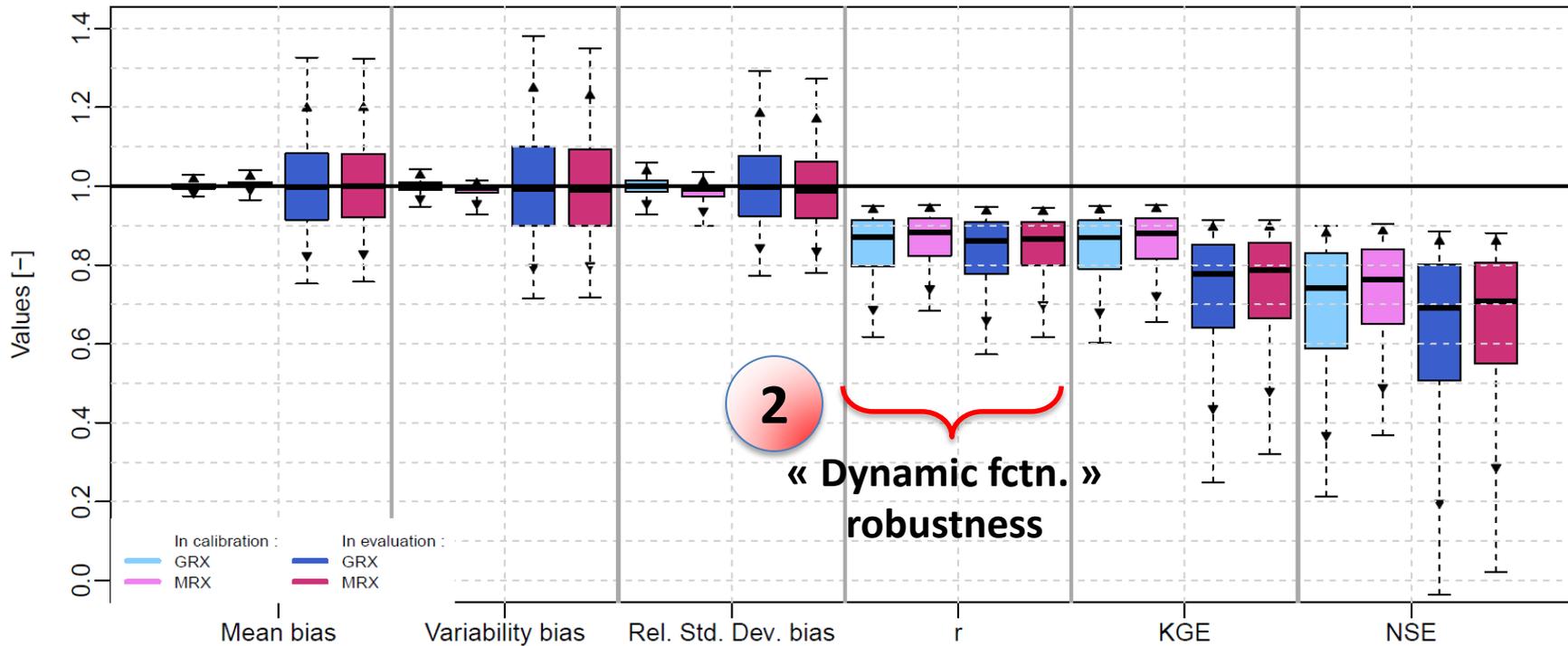
1

WB & Var. lack of robustness

Results : Boxplots



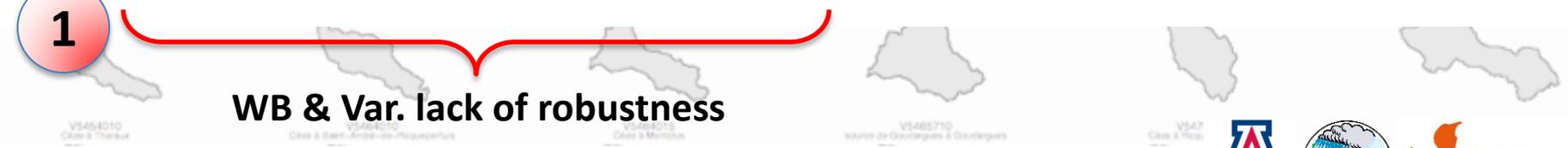
Calibration on KGE(Q)
Calibration & Evaluation



1

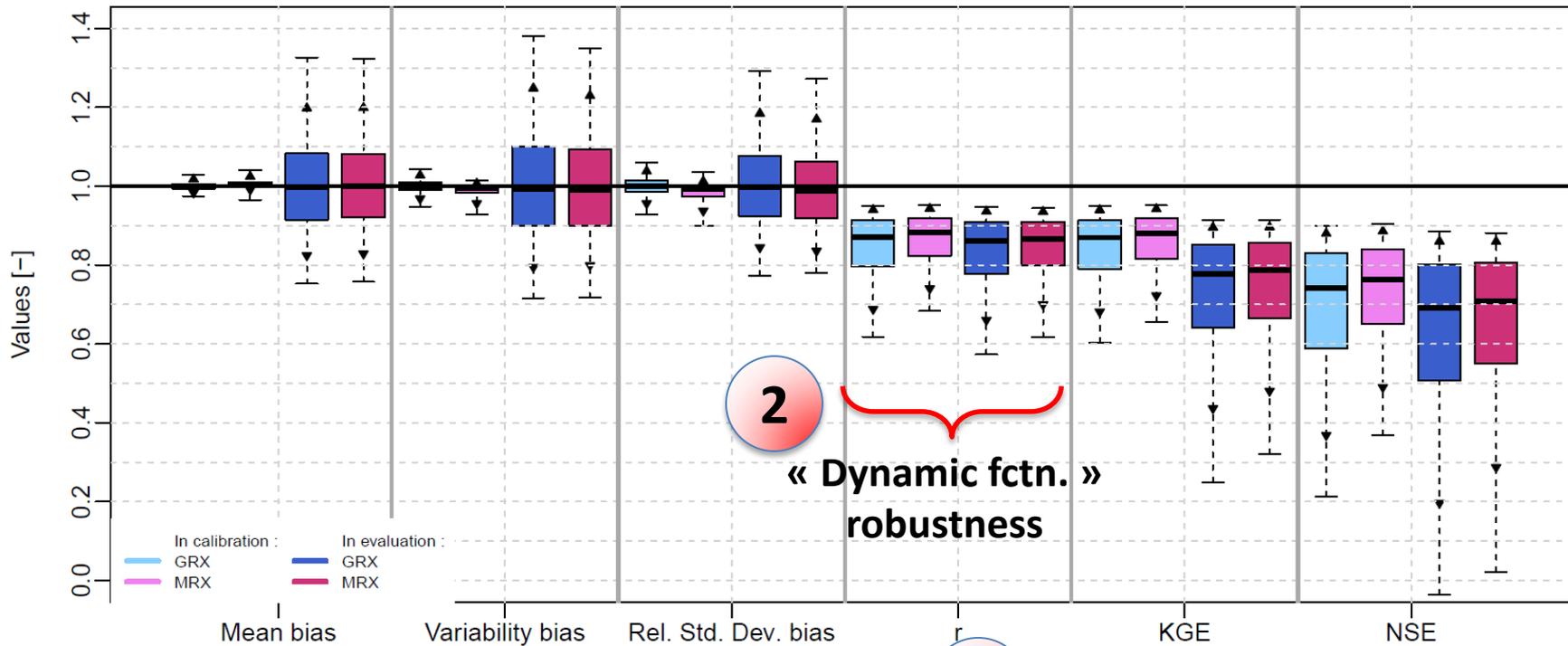
WB & Var. lack of robustness

Criteria



Results : Boxplots

Calibration on KGE(Q)
Calibration & Evaluation



1

WB & Var. lack of robustness

3

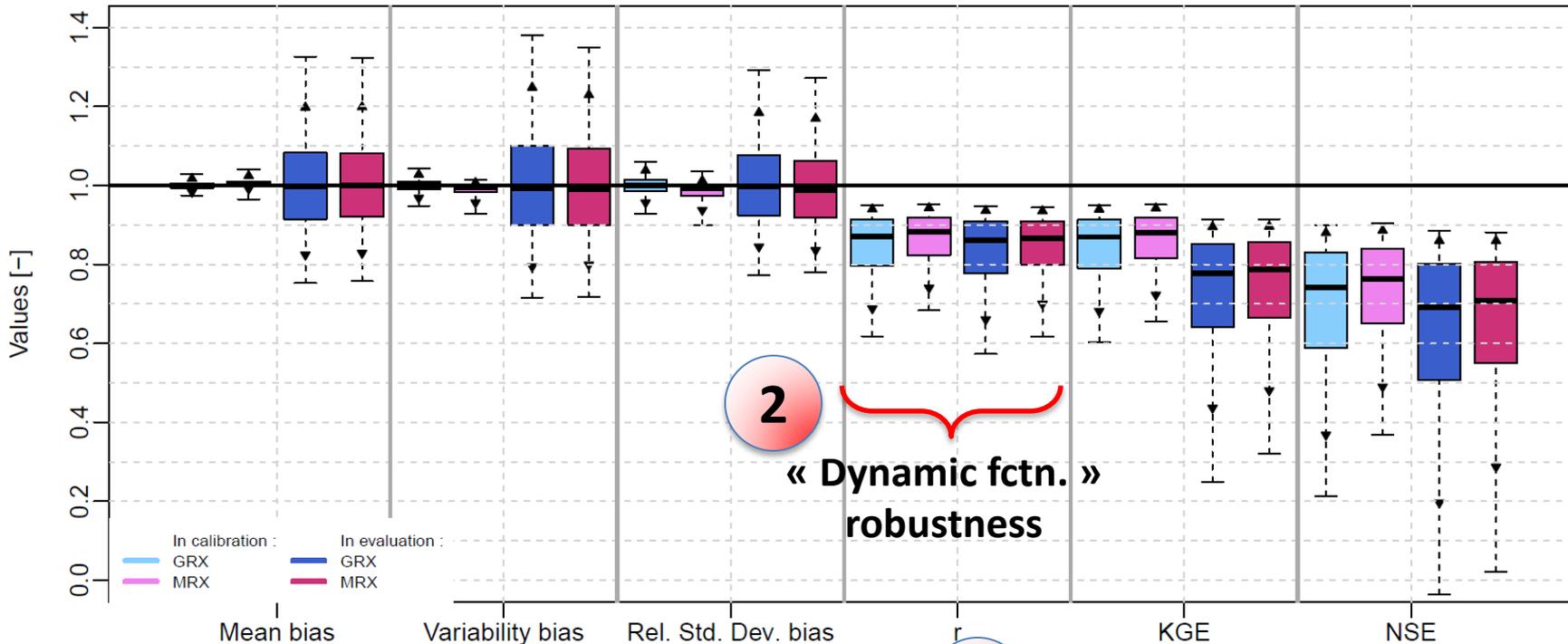
Impact of WB & Var.
lack of robustness

Results : Boxplots

4

Very similar results

Calibration on KGE(Q)
Calibration & Evaluation



2

1

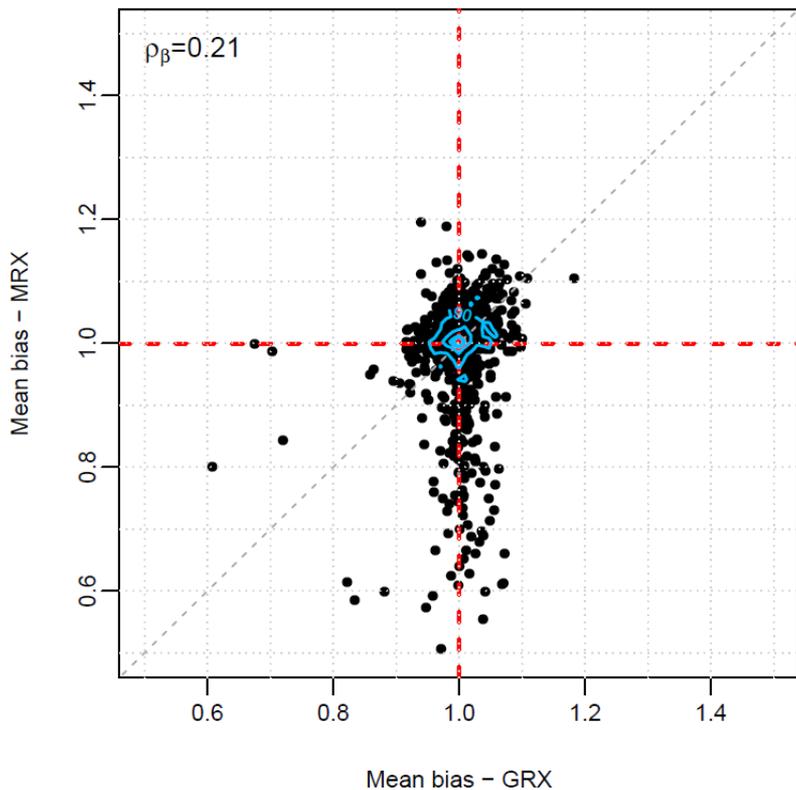
WB & Var. lack of robustness

3

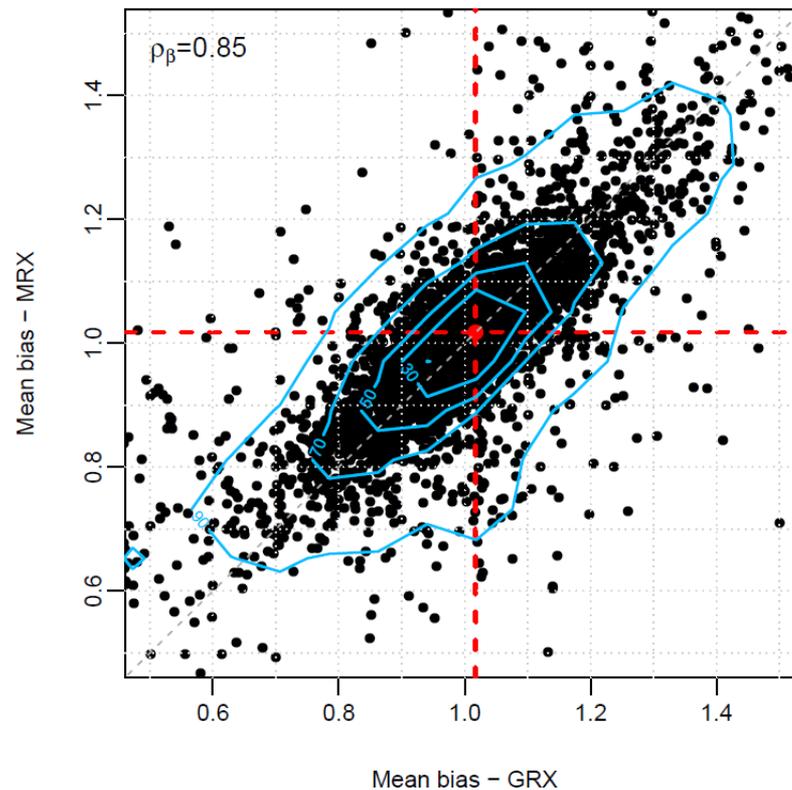
Impact of WB & Var.
lack of robustness

Results : Scatterplots

Calibration on KGE(Q)
Calibration



Calibration on KGE(Q)
Evaluation

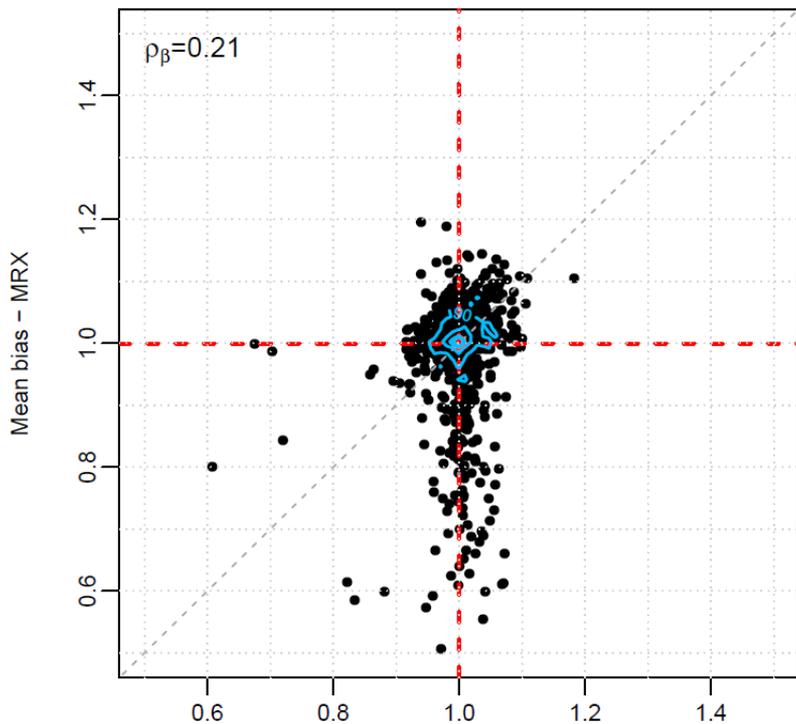


0

How to read these scatterplots ?

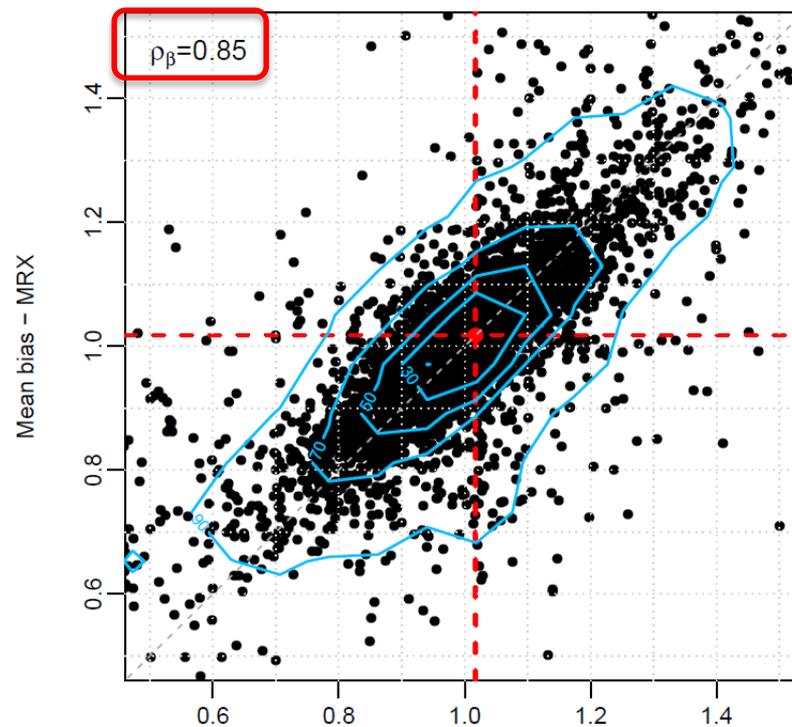
Results : Scatterplots

Calibration on KGE(Q)
Calibration



Mean bias - GRX

Calibration on KGE(Q)
Evaluation



Mean bias - GRX



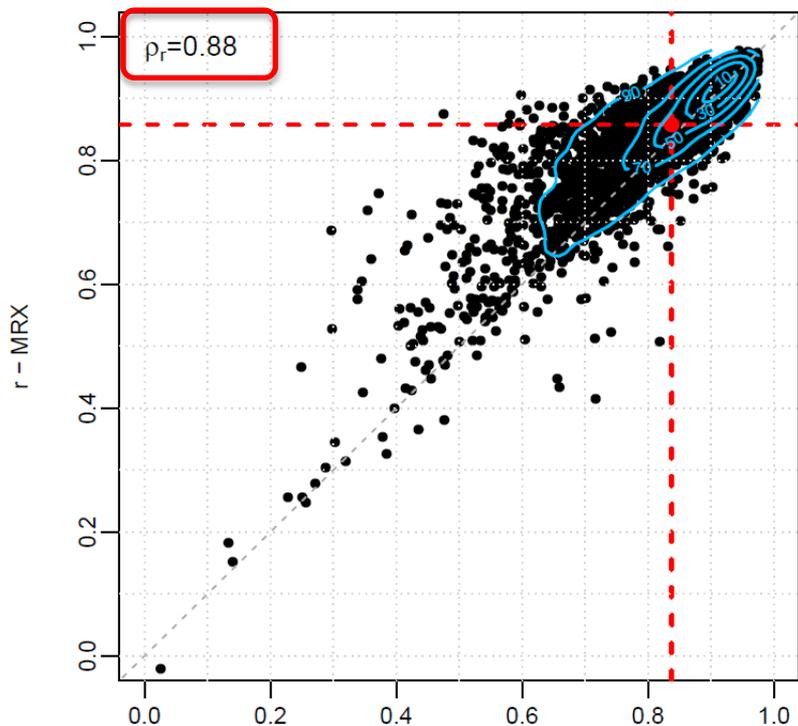
1

Strongly correlated behaviour
for β , α , γ

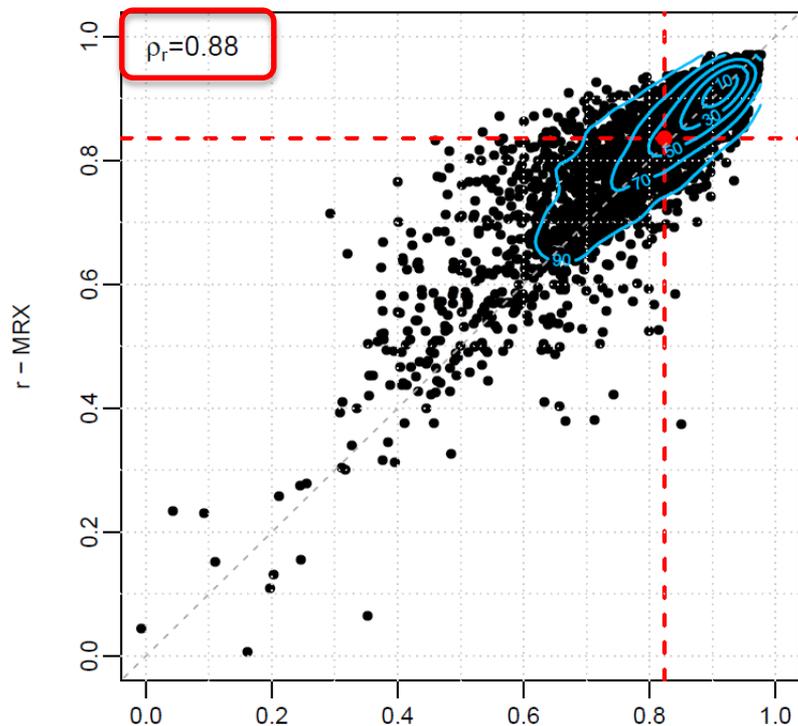


Results : Scatterplots

Calibration on KGE(Q)
Calibration



Calibration on KGE(Q)
Evaluation



$r - \text{GRX}$

$r - \text{GRX}$

2

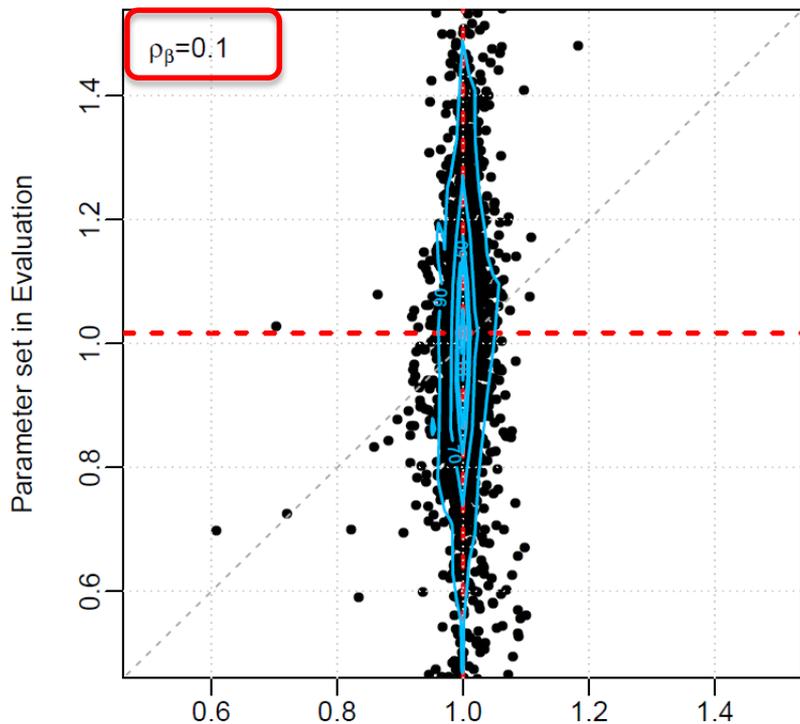
Strongly correlated behaviour
for r , KGE & NSE



Results : Scatterplots

Calibration on KGE(Q)

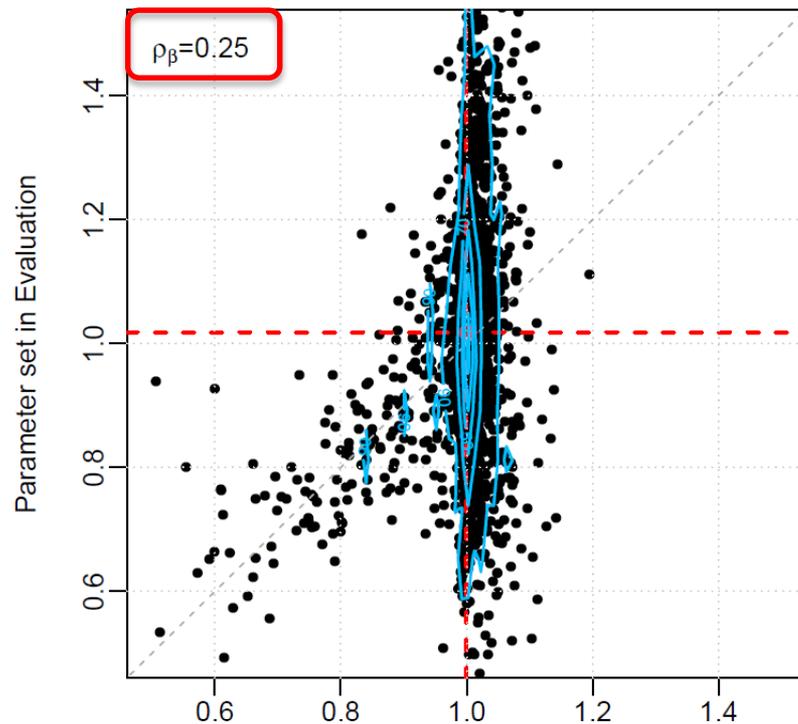
GRX - Comparison of Mean bias in calibration & evaluation



Parameter set in Calibration

Calibration on KGE(Q)

MRX - Comparison of Mean bias in calibration & evaluation



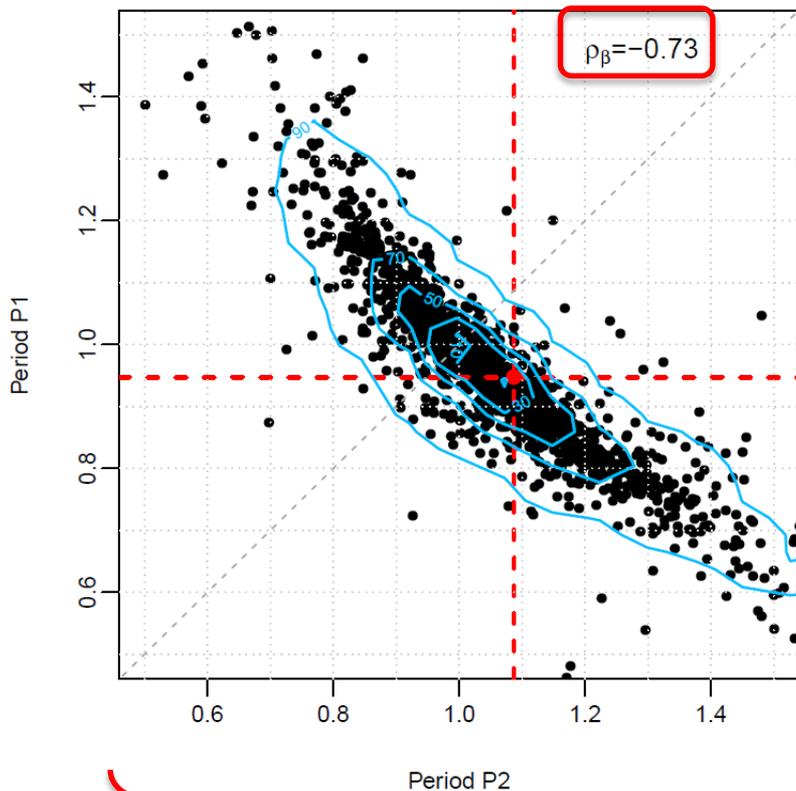
Parameter set in Calibration

3

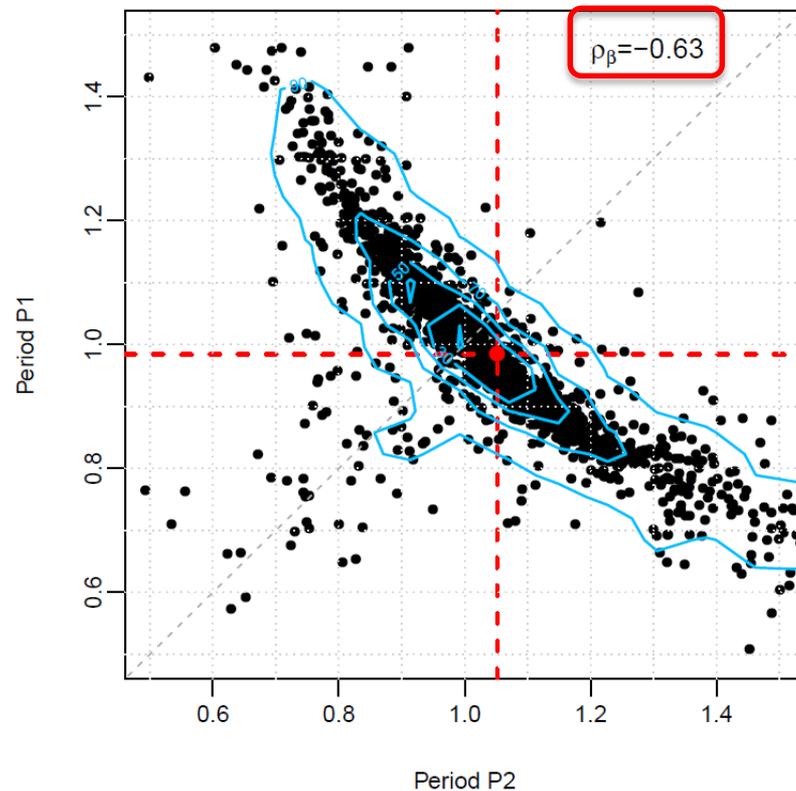
Very poor robustness
for β , α , γ

Results : Scatterplots

Calibration on KGE(Q)
Mean bias in evaluation for GRX



Calibration on KGE(Q)
Mean bias in evaluation for MRX

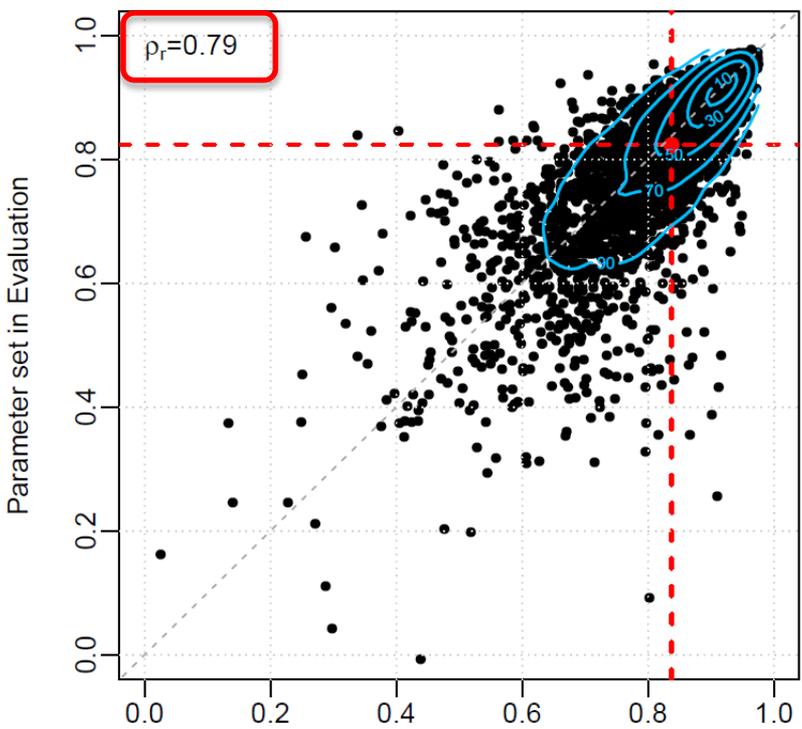


3+ Anti-correlated behaviour
for β on different evaluation periods

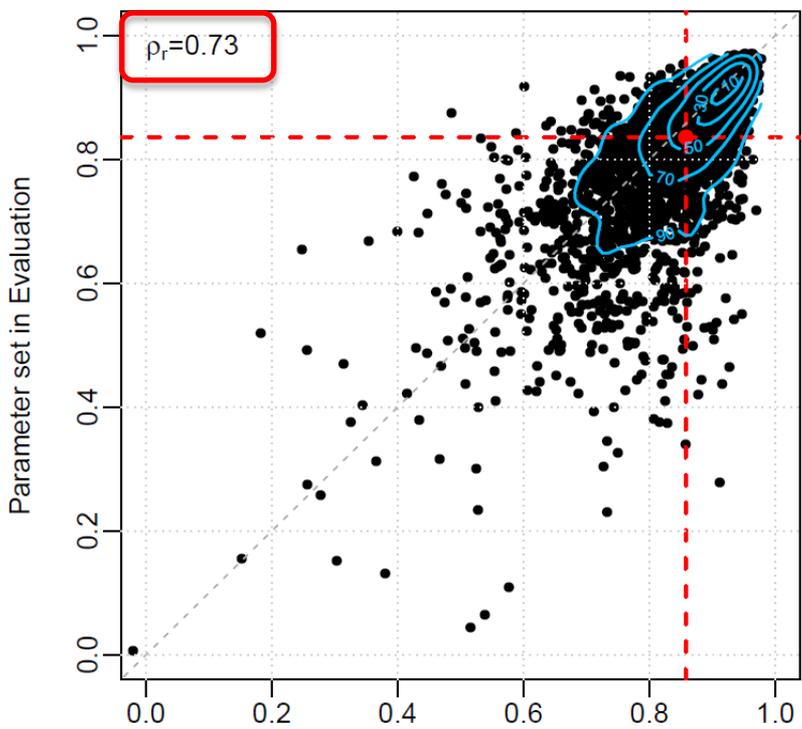
Results : Scatterplots



Calibration on KGE(Q)
GRX - Comparison of r in calibration & evaluation



Calibration on KGE(Q)
MRX - Comparison of r in calibration & evaluation



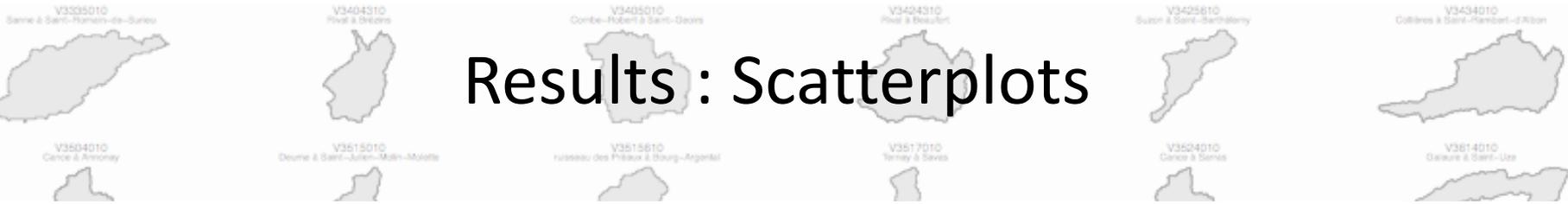
Parameter set in Calibration

Parameter set in Calibration

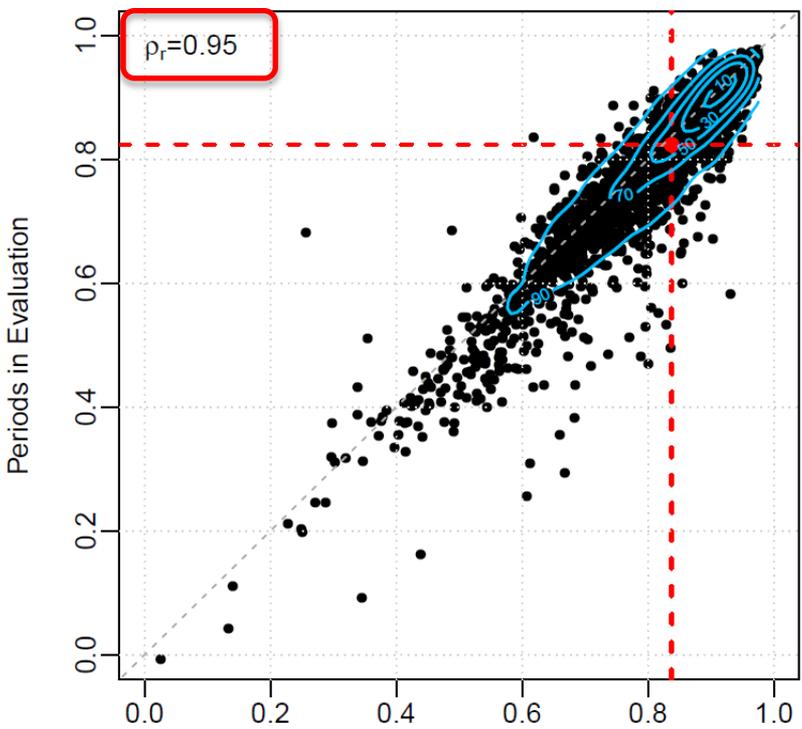
4

Very robust and consistent performances of parameter set (whatever the period)

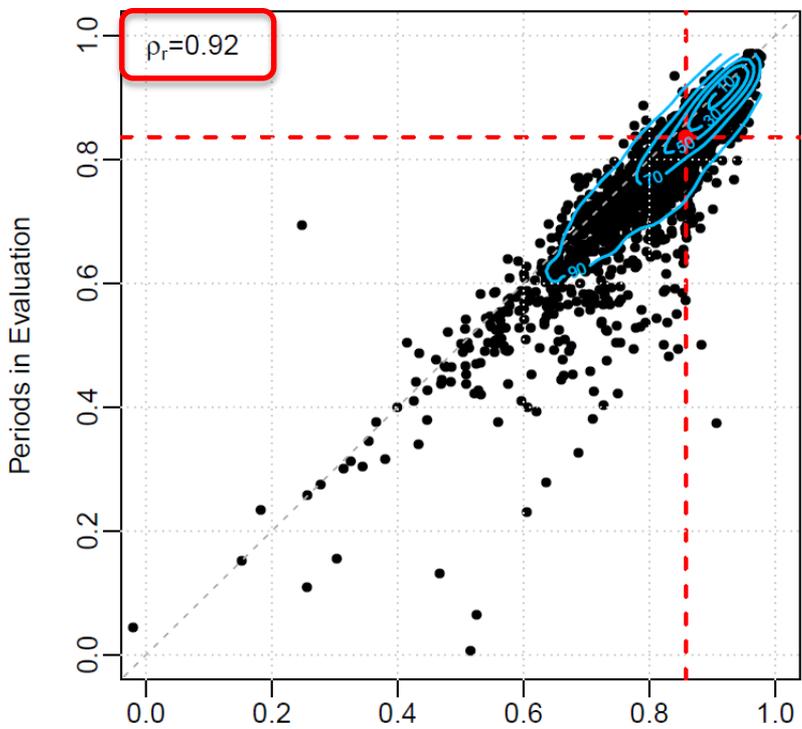
Results : Scatterplots



Calibration on KGE(Q)
GRX - Comparison of r in calibration & evaluation



Calibration on KGE(Q)
MRX - Comparison of r in calibration & evaluation



Periods in Calibration

Periods in Calibration

5

Extremely robust and consistent performances on periods (whatever the parameter set)

Results : Synthesis

1 Both models suffer from a strong lack of robustness in the simulation of water balance and streamflow variability. The water balance bias varies on the range $\pm 10\%$ for 50% of the watersheds, on the range $\pm 20\%$ for 80% of the watersheds. **However, both models are particularly robust concerning the representation of the dynamic functioning of the watershed.**

2 The performance of both models is highly correlated (r ranging from 0.75 to 0.92), despite the strong difference of structure and complexity. This means that model performance correlation (between simulations provided by the two models) is at the same level as the correlation between each of the model simulations and the observations, suggesting that there is no significant difference in overall abilities of the two models across the range of watersheds used for testing.

3 Hence, it seems that differences in hydroclimatic conditions between calibration to evaluation periods play a more important role on the differences in performance from calibration to evaluation than differences in model structures do.

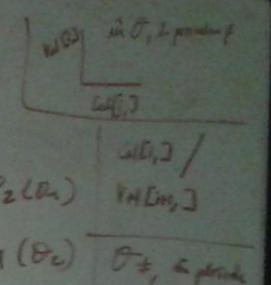
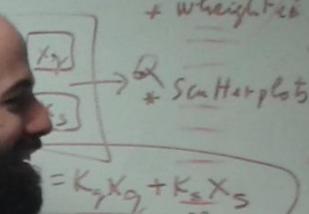


TRiB
 Sample: FRA+EDF+AUS+MPX
 Model: GR+MR
 FO: QSEC+QUCM
 Param: P1, P2, P1+P2

Statistical analysis of PR performance

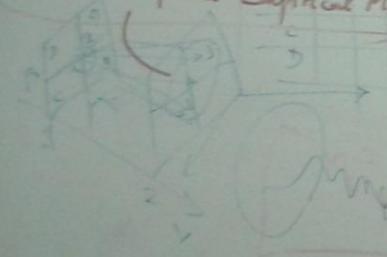
- * Coefficient of Variation
- + BV with problems
- + Spatial correlation of Δ
- + internal states / flux

- * Drunken-dart plots
- + weighted KGE



- 1- BABEL
- 2- GR+MR simpl
- 3- can't analyse within sub-periods
- 4- bootstrap temporal / noise

Hybrid Conceptual-Empirical Mod.



BABEL

- The hunt.
1. Overfitting to data
 2. Service rainfall/errors
 3. Non-stationarity
 4. Model str. instability
- GEF $\begin{matrix} > 0.9 \\ < 0 \end{matrix}$ $\begin{matrix} > 0 \\ < 0 \end{matrix}$
- GR3, GR4, GR5, GR6

$Q = f(X_1, X_2, X_3, X_4, T)$

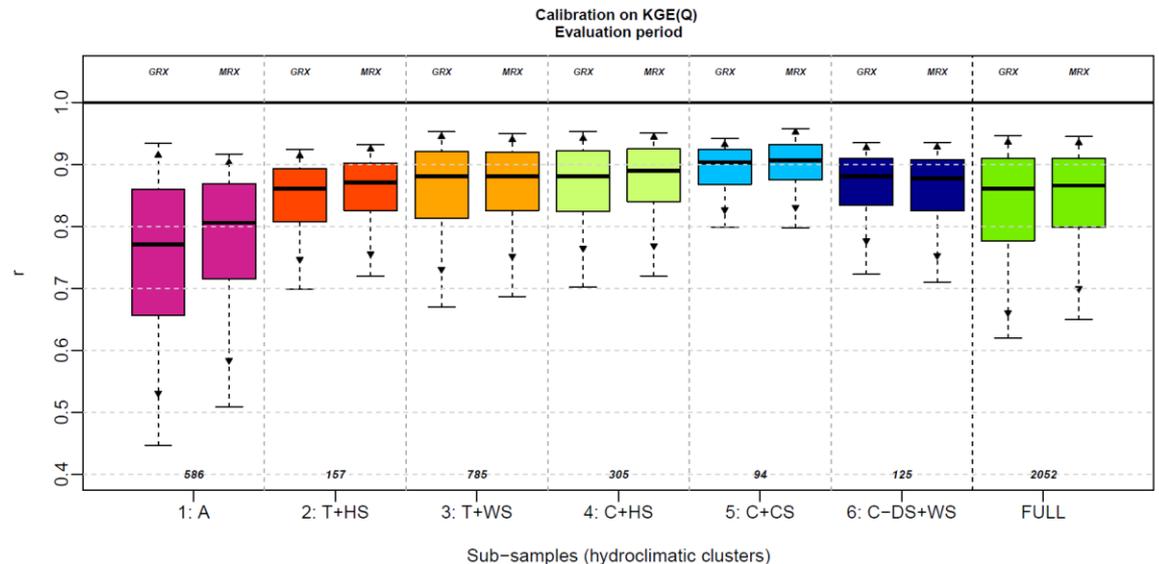
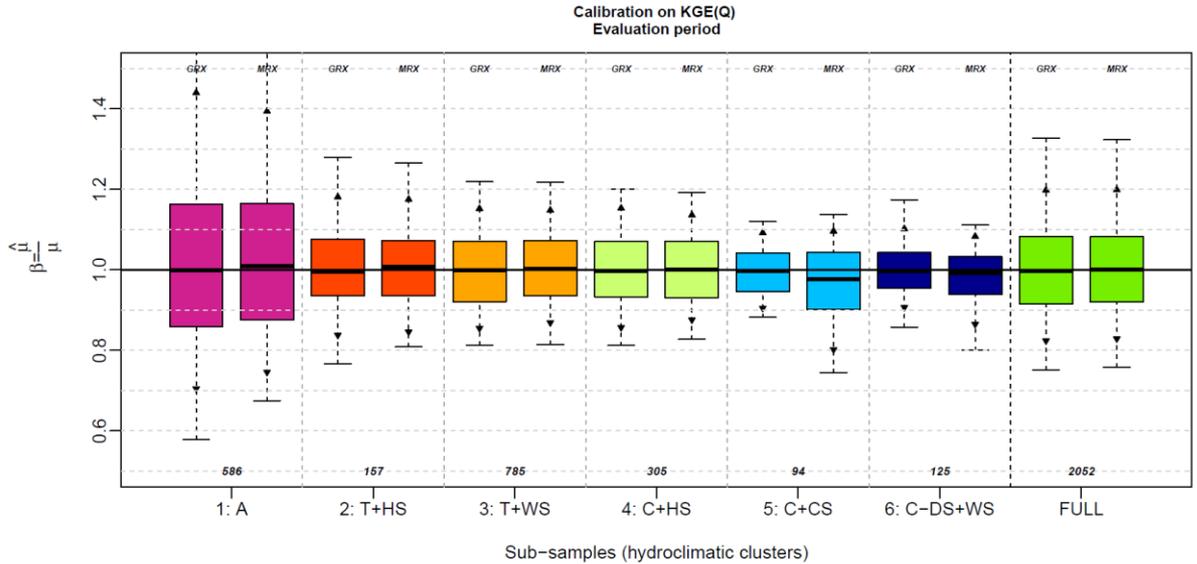
$P \rightarrow Q_{SEC}$
 $Q_{SEC}, Q_{UCM}, Q_{VQA}$

Water Balance ✓
 50%

$Y_t = g(f(x_{t-1}, u_t)) = h(x_{t-1}, u_t)$

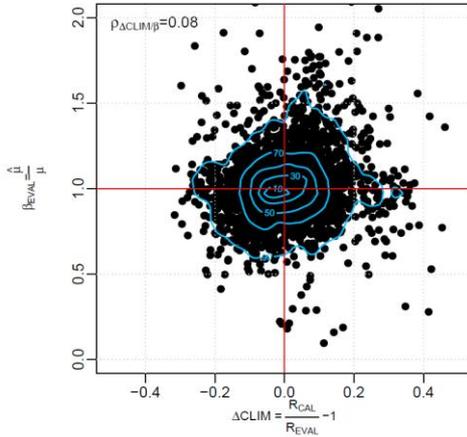
MERCI

Question 3: Are differences in model performance dependent on watershed characteristics or on hydrometeorological processes?

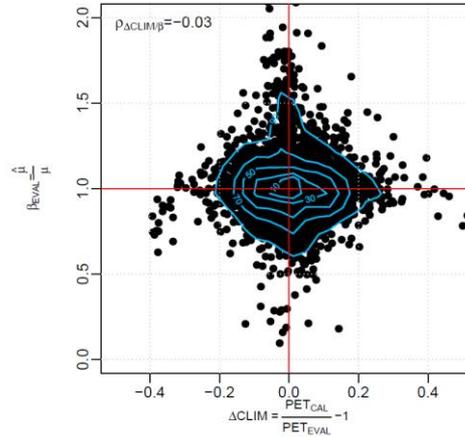


Question 3: Are differences in model performance dependent on watershed characteristics or on hydrometeorological processes?

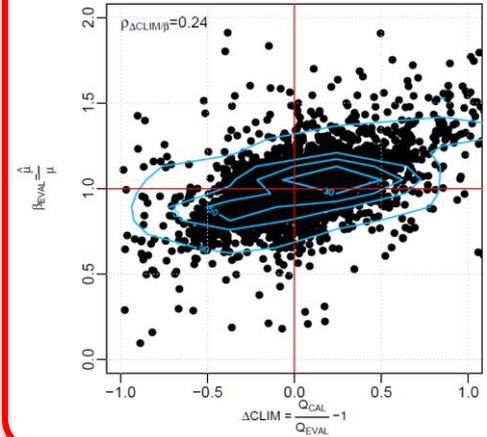
Calibration on KGE(Q)
MRX



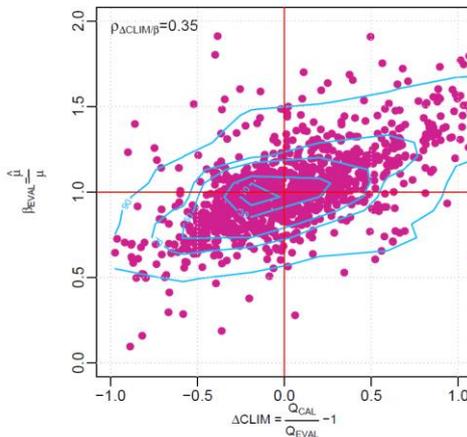
Calibration on KGE(Q)
MRX



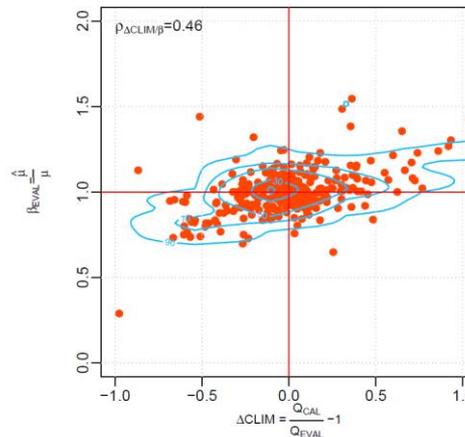
Calibration on KGE(Q)
MRX



Calibration on KGE(Q)
MRX



Calibration on KGE(Q)
MRX



Calibration on KGE(Q)
MRX

