



# Improving Model Identification: Reconciling Theory with Observations & The Problem of Sufficient Statistics

Hoshin Gupta

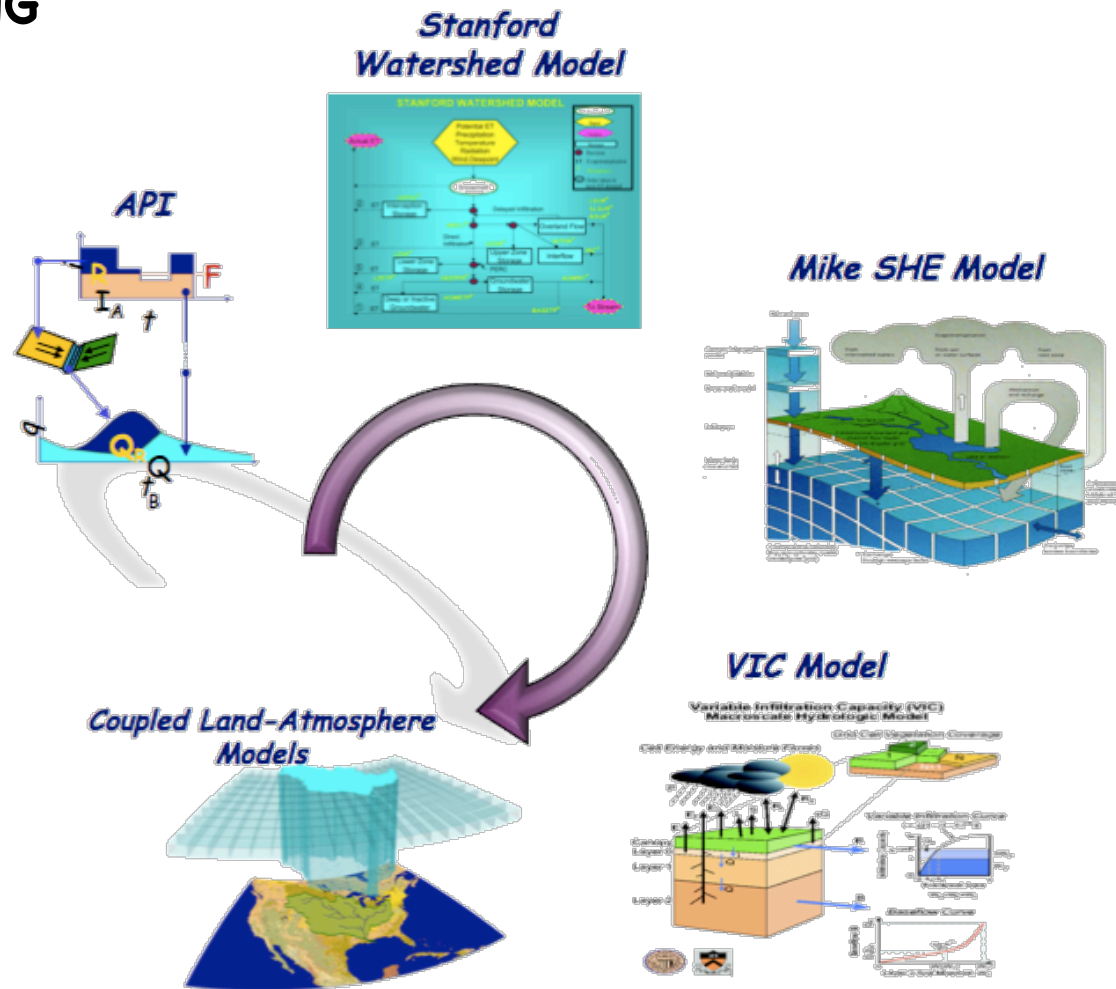
University of Arizona

*Session HS1.6 on  
Metrics and the Use of Data in Hydrology to support  
Model Structure Improvement  
EGU Meeting  
Vienna, Austria*



Hoshin Gupta

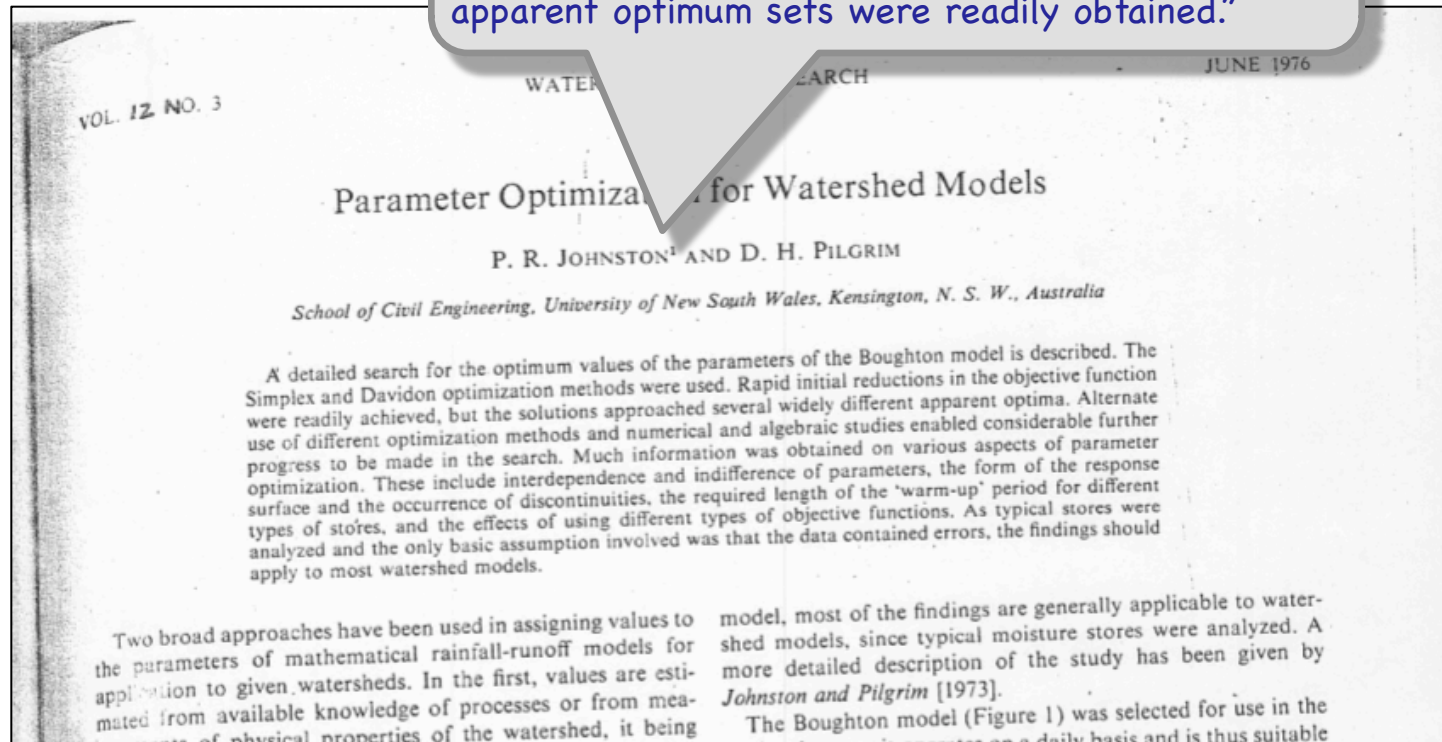
# FIVE DECADES OF COMPUTER-BASED HYDROLOGIC MODELING



## FIVE DECADES OF COMPUTER-BASED HYDROLOGIC

### MANY REPORTS OF DIFFICULTIES IN MODEL IDENTIFICATION

"A true optimum set of (parameter) values was not found in over 2 years of full-time work concentrated on one watershed, although many apparent optimum sets were readily obtained."

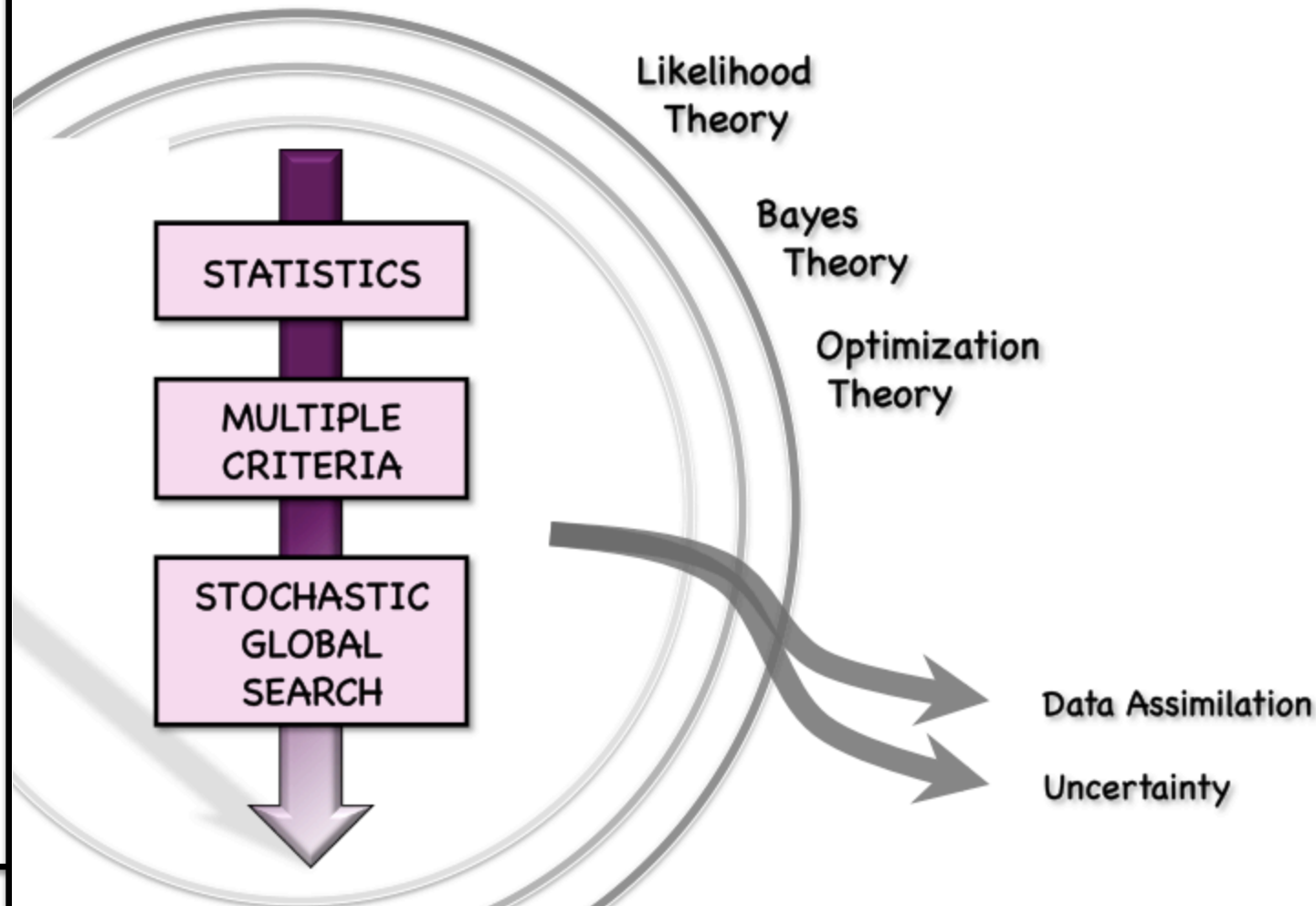


# FIVE DECADES OF COMPUTER-BASED HYDROLOGIC

MANY REPORTS OF DIFFICULTIES IN MODEL

I

MANY ATTEMPTS TO IMPROVE MODEL IDENTIFICATION





# FIVE DECADES OF COMPUTER-BASED HYDROLOGIC

MANY REPORTS OF DIFFICULTIES IN MODEL

MANY ATTEMPTS TO IMPROVE MODEL IDENTIFICATION

NOW RECOGNIZED BY “NSF” AS A “GRAND CHALLENGE”

## GRAND CHALLENGES OF THE FUTURE FOR ENVIRONMENTAL MODELING

Report of the NSF Project (Award # 0630367)  
May 2006 - May 2008

DRAFT

M B Beck

Warnell School of Forestry and Natural Resources  
University of Georgia, Athens, Georgia 30602-2152

October, 2008

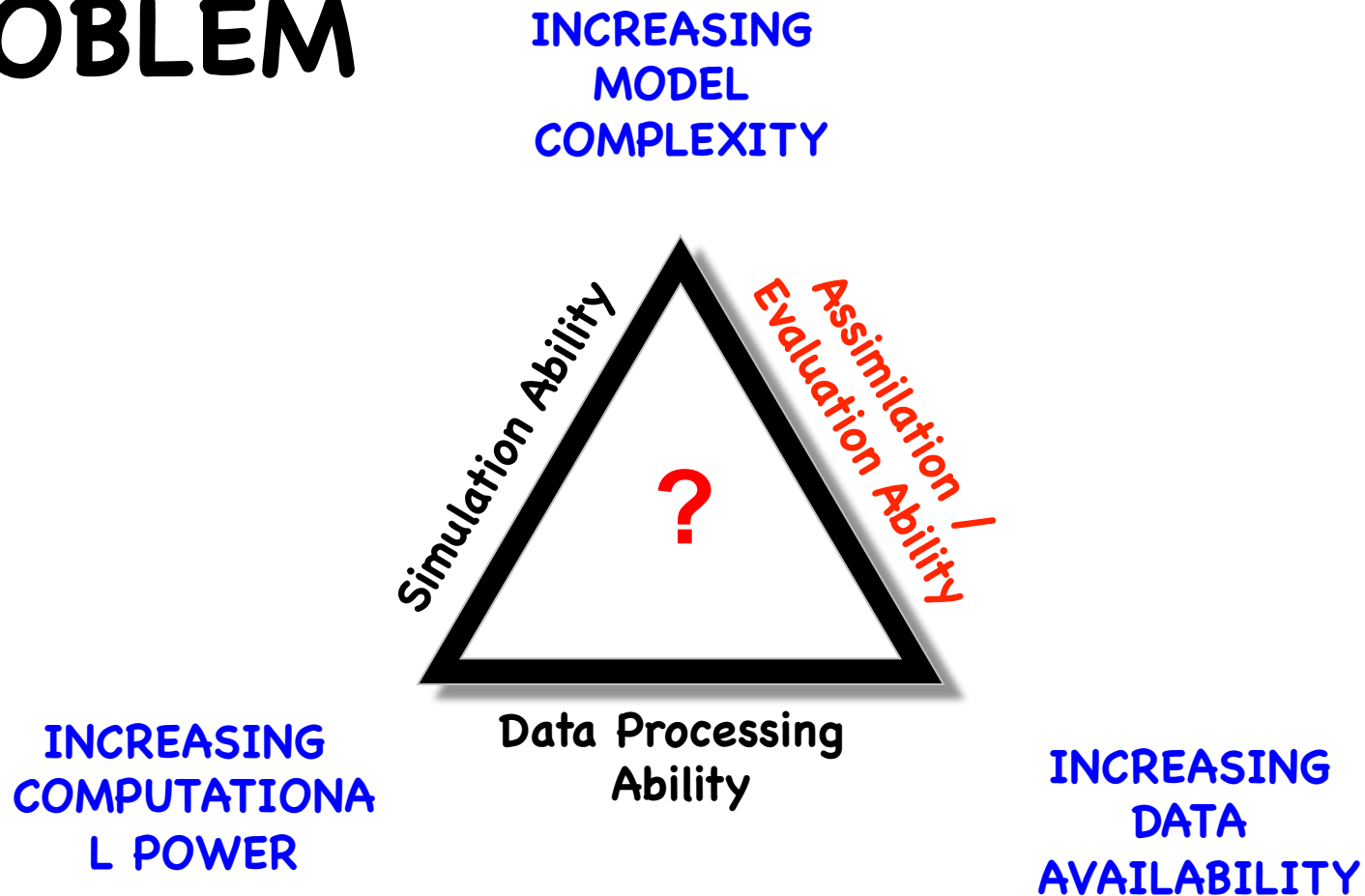
*Beck et al (2008)*

We do a poor job at reconciling  
complex models with field data.  
How do we improve this?

Challenge # 7:

*Under the expectation of massive expansion in the scope and volume of field observations generated by the Environmental Observatories, coupled and integrated with the prospect of equally massive expansion in data processing and scientific visualization enabled by the future environmental cyber-infrastructure, what radically novel procedures and algorithms are needed to rectify the chronic, historical deficit of the past four decades in engaging complex models (VHOMs) systematically and successfully with field data for the purposes of learning and discovery and, thereby, enhancing the growth of environmental knowledge?*

# THE PROBLEM

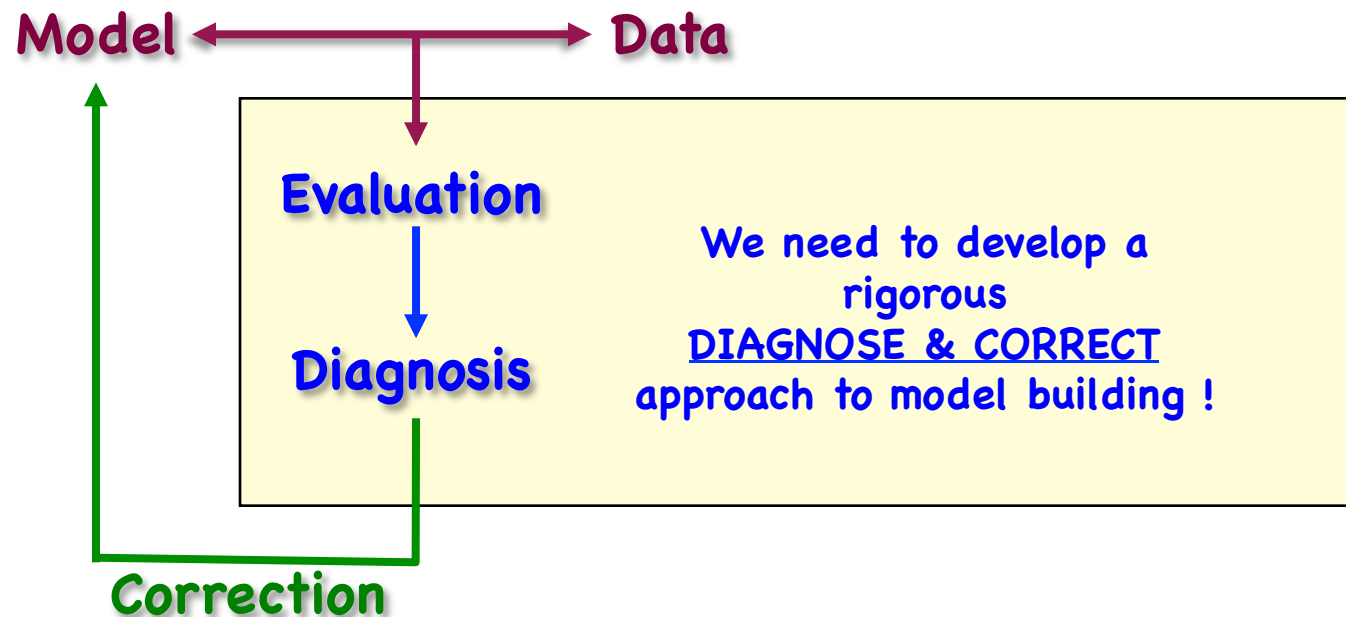


# THE SOLUTION



## WE NEED A THEORY OF DIAGNOSTIC EVALUATION

The theory should enable us to link what we “see” in the data to what is “right” and “wrong” with our models.



# BLUEPRINT FOR SUCH A THEORY

HYDROLOGICAL PROCESSES

*Hydrol. Process.* (2008)

Published online in Wiley InterScience

(www.interscience.wiley.com) DOI: 10.1002/hyp.6989

## Reconciling theory with observations: elements of a diagnostic approach to model evaluation

Hoshin V. Gupta,<sup>1\*</sup> Thorsten Wagener<sup>2</sup> and Yuqiong Liu<sup>1</sup>

<sup>1</sup> SAHRA, Department of Hydrology & Water Resources, The University of Arizona, Tucson AZ 85721

<sup>2</sup> Department of Civil and Environmental Engineering, Pennsylvania State University, University Park, PA 16802

This paper discusses the need for a well-cons has clear and compelling diagnostic power. Th 'Predictions in Ungaged Basins' initiative and initiative, among others. It is suggested that observational data are inadequate in the face environmental science, and steps are proposed. This paper presents the concept of a diagnostic signature indices that measure theoretically rele issue of degree of system complexity resolvable facilitate uncertainty analysis, and can be readily in ungaged basins. Copyright © 2008 John Wil

KEY WORDS model identification; information;

Received 26 June 2007; Accepted 15 December

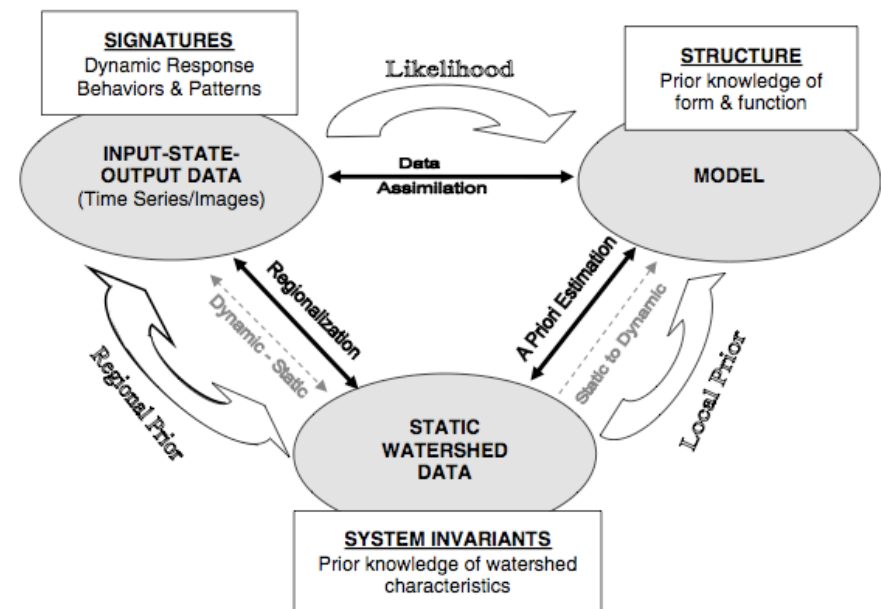
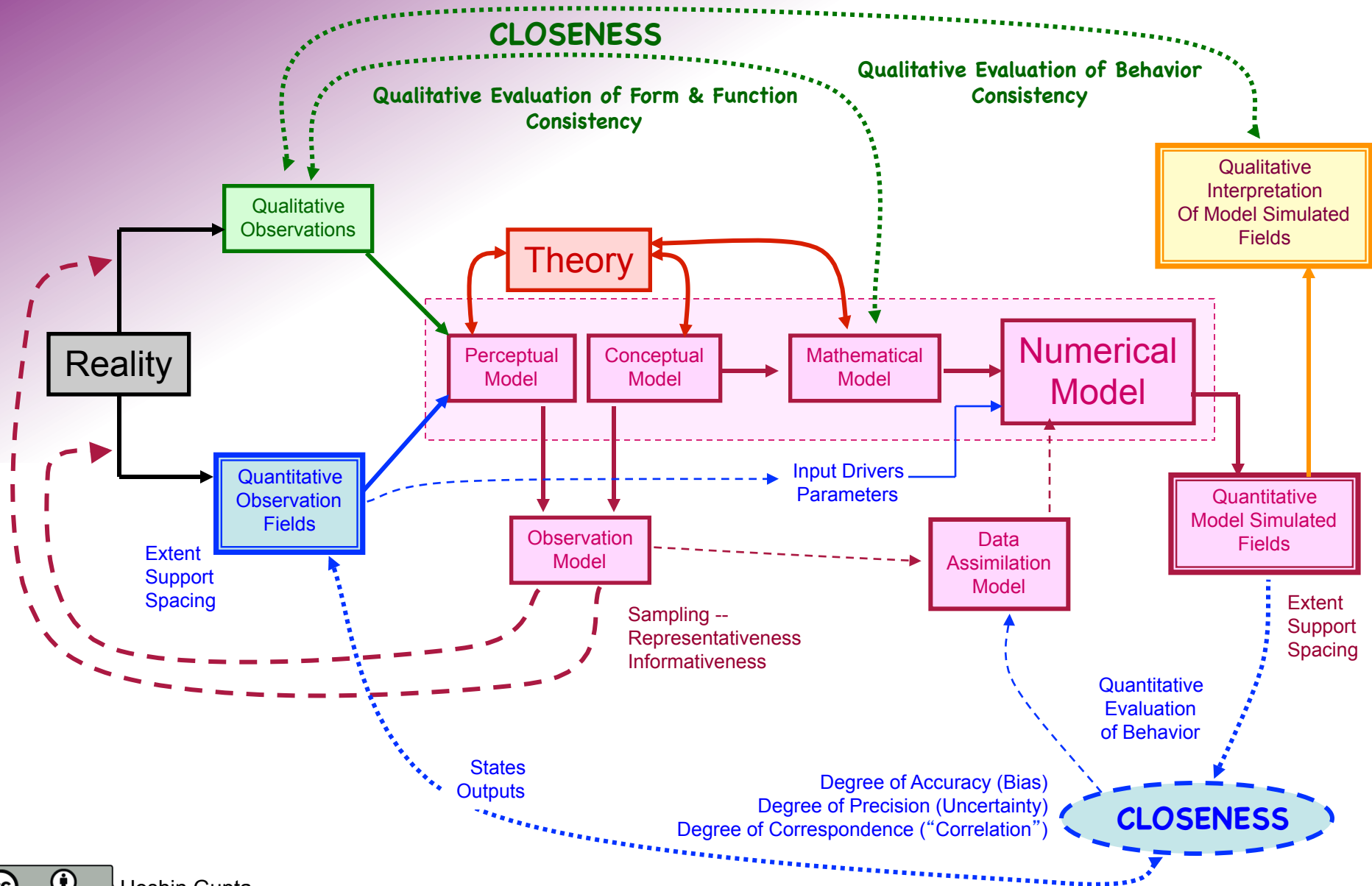
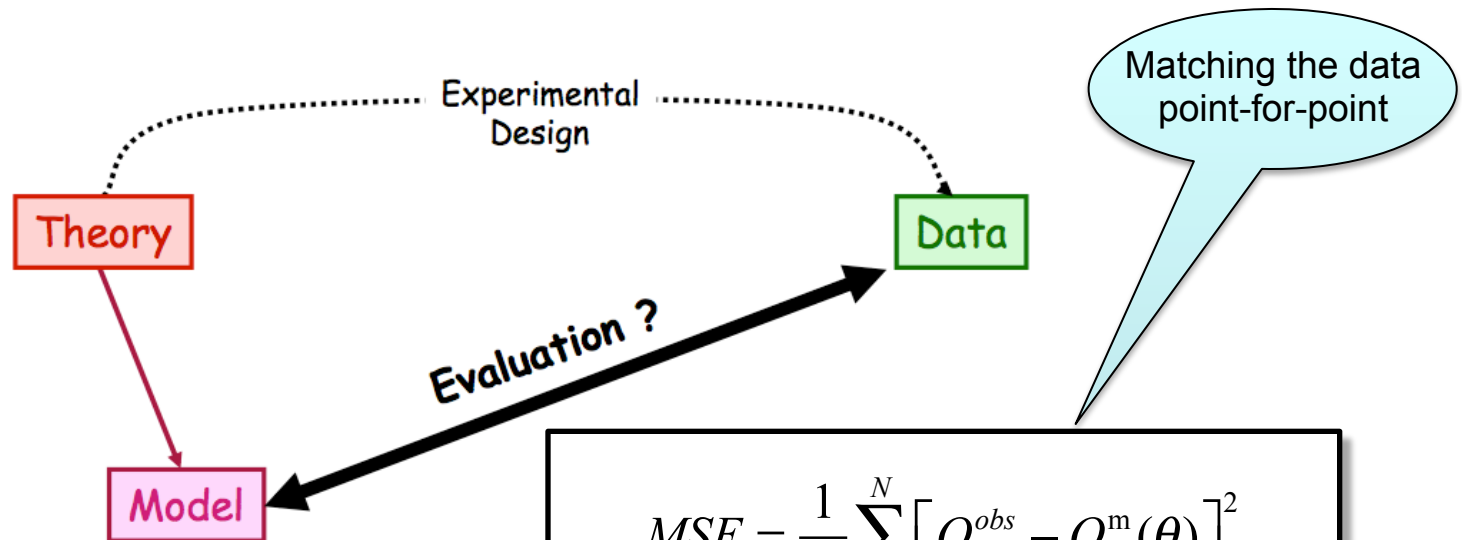


Figure 6. The three kinds of information used to constrain the predictive model

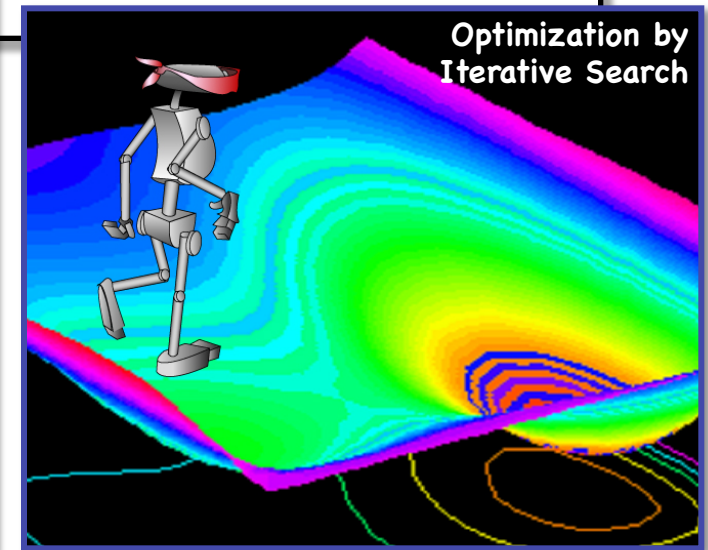
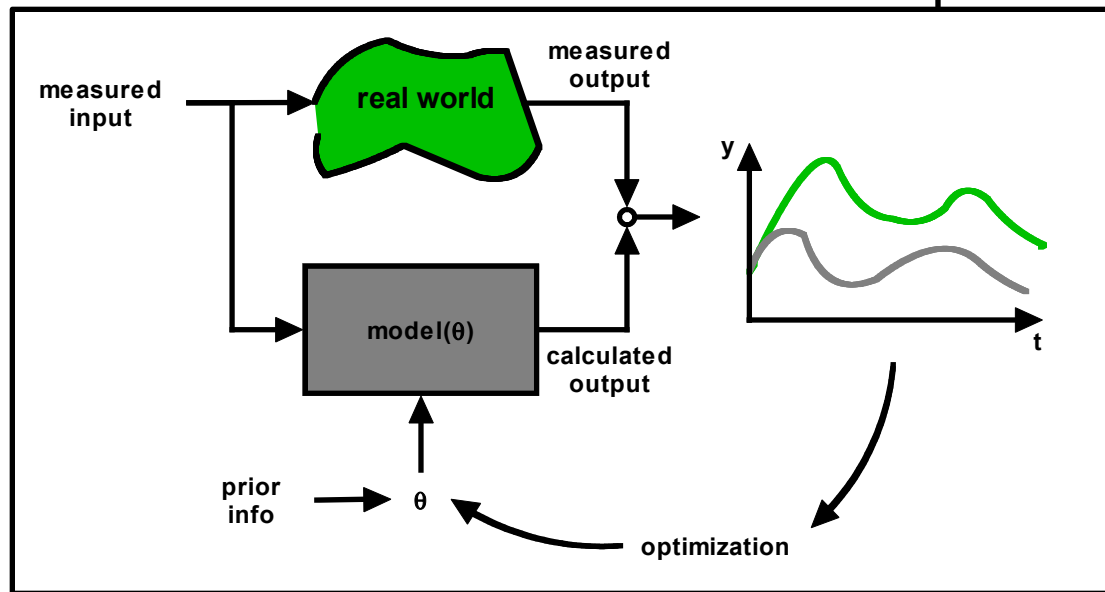
# The Model Building & Evaluation Process



# NEO-CLASSICAL APPROACH TO MODEL EVALUATION



$$MSE = \frac{1}{N} \sum_{t=1}^N [Q_t^{obs} - Q_t^m(\theta)]^2$$



**BUT ...**  
**'MSE' DECOMPOSES\***  
**INTO THREE STATISTICS**  
**OF MODEL PERFORMANCE**

$$MSE = \frac{1}{N} \sum_{t=1}^N [Q_t^{obs} - Q_t^m(\theta)]^2$$

$$MSE = (\mu_m - \mu_{obs})^2 + (\sigma_m - \sigma_{obs})^2 + 2\sigma_m \sigma_{obs} (1 - r)$$

$$MSE = \mathcal{F}(\text{Mean Error}) + \mathcal{F}(\text{Variability Error}) + \mathcal{F}(\text{Linear Correlation})$$

**Water Balance**

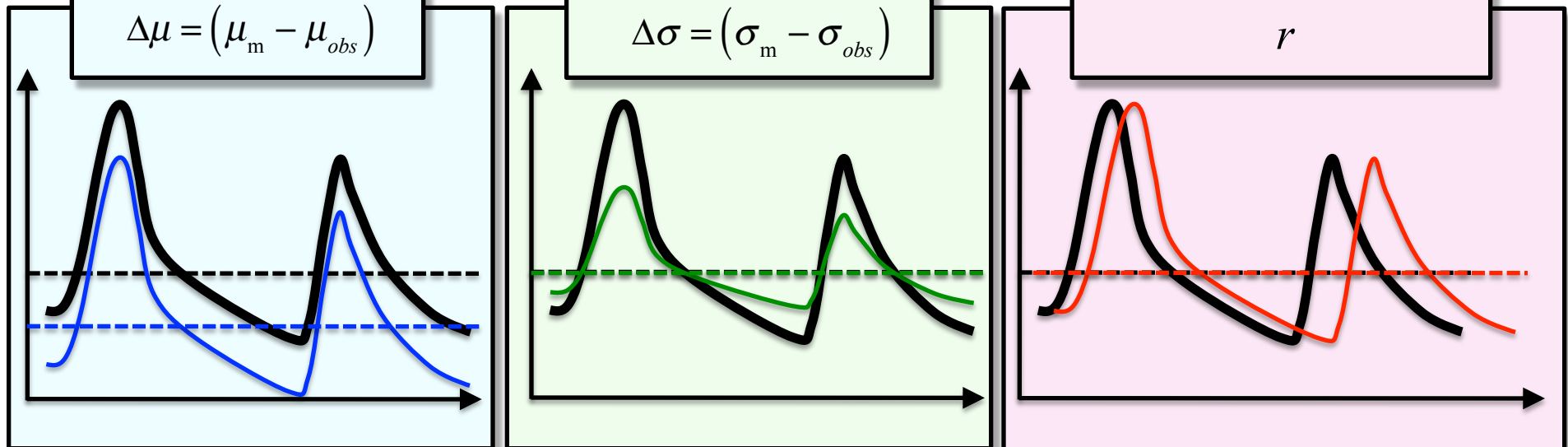
$$\Delta\mu = (\mu_m - \mu_{obs})$$

**Variability**

$$\Delta\sigma = (\sigma_m - \sigma_{obs})$$

**Timing & Shape**

$$r$$



\* Decomposition of the Mean Squared Error & NSE Performance Criteria: Implications for Improving Hydrological Modeling

Gupta, H.V., H. Kling, Y.K. Yilmaz & G.F. Martinez-Baquero, *Manuscript submitted to Journal of Hydrology*, 2009.

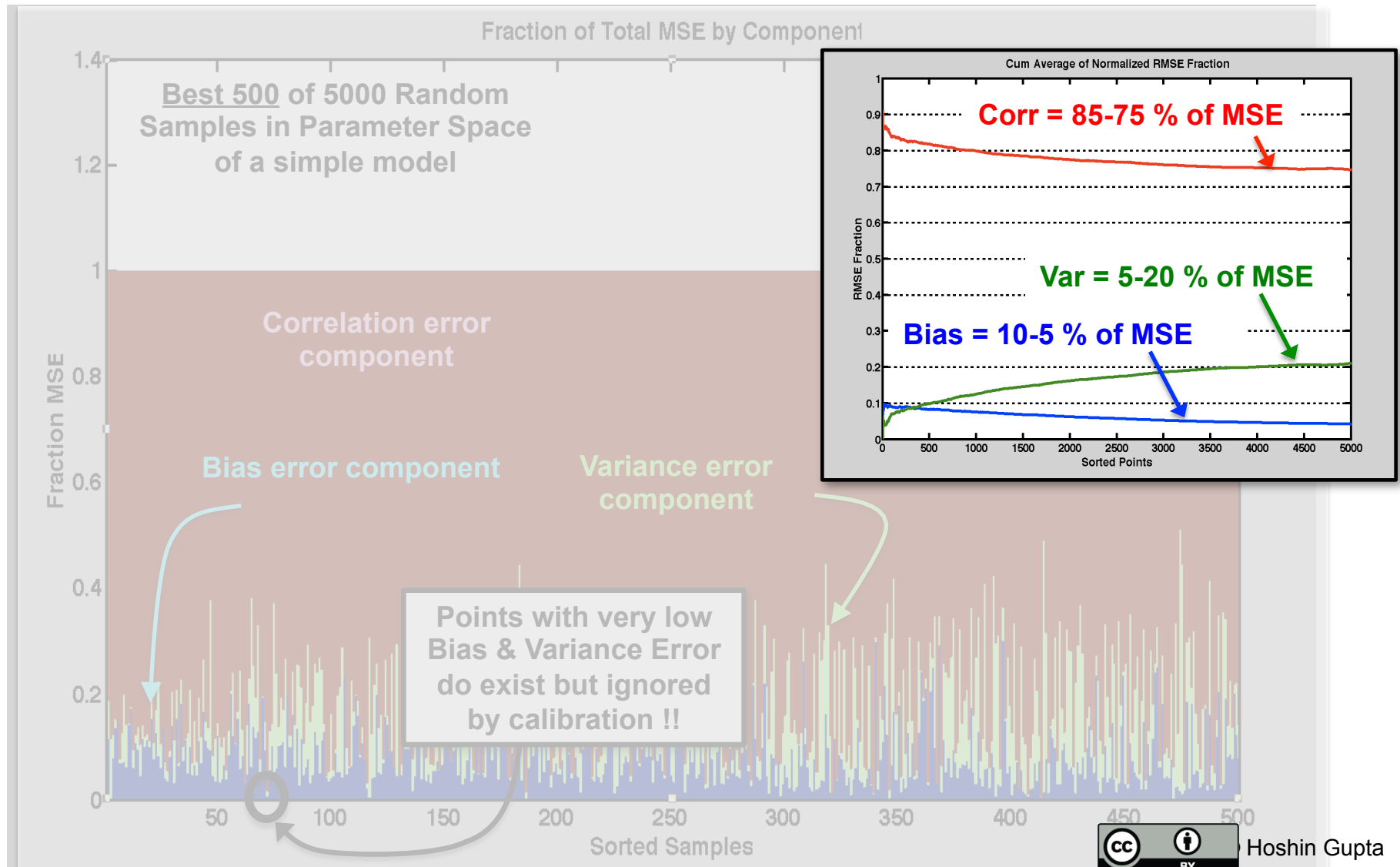


Hoshin Gupta



# PROBLEM 1. CORRELATION COMPONENT DOMINATES

$$MSE = \mathcal{F}(\text{Bias Error}) + \mathcal{F}(\text{Var Error}) + \mathcal{F}(\text{Corr Error})$$



## PROBLEM 2.

### MODEL PERFORMANCE WILL BE SIGNIFICANTLY OVER-ESTIMATED

$$MSE = (\mu_m - \mu_{obs})^2 + (\sigma_m - \sigma_{obs})^2 + 2\sigma_m \sigma_{obs}(1-r)$$

#### IDEAL

If we ensure  $\mu_m = \mu_{obs}$  and  $\sigma_m = \sigma_{obs}$   
the expected 'best' value for  $MSE$   
is:

$$\frac{MSE_{ideal}}{\sigma_{obs}^2} = 2(1 - r_{ideal})$$

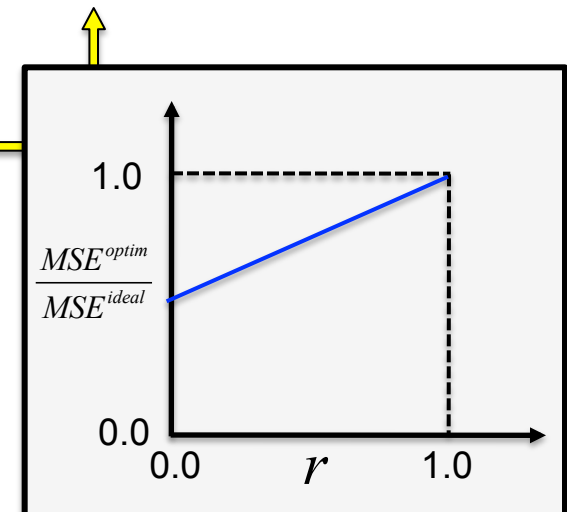
#### OPTIMIZED

But if we optimize on  $MSE$  without  
constraining  $\mu_m$  &  $\sigma_m$  we will get:

$$\frac{MSE_{optim}}{\sigma_{obs}^2} = (1 - r_{optim}^2)$$

$$MSE^{optim} \approx \left( \frac{1+r}{2} \right) MSE^{ideal}$$

$$MSE_{optim} < MSE_{ideal}$$



assumes  $r_{ideal} \sim r_{optim}$

### PROBLEM 3. VARIABILITY WILL BE UNDER- ESTIMATED

*'Optimal'* MSE is achieved when:

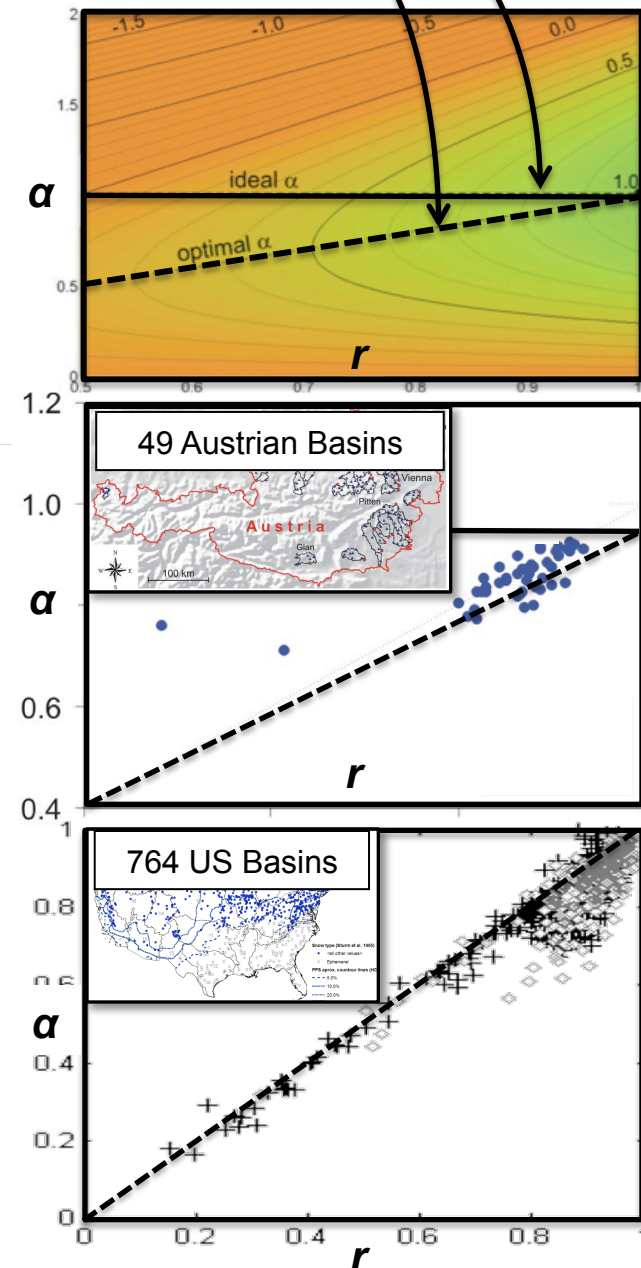
$$\sigma_m = r \cdot \sigma_{obs}$$

In other words:

$$\alpha = \frac{\sigma_m}{\sigma_{obs}} = r < 1.0$$

*'Optimal' model will  
underestimate  
the observed variability  
of the data*

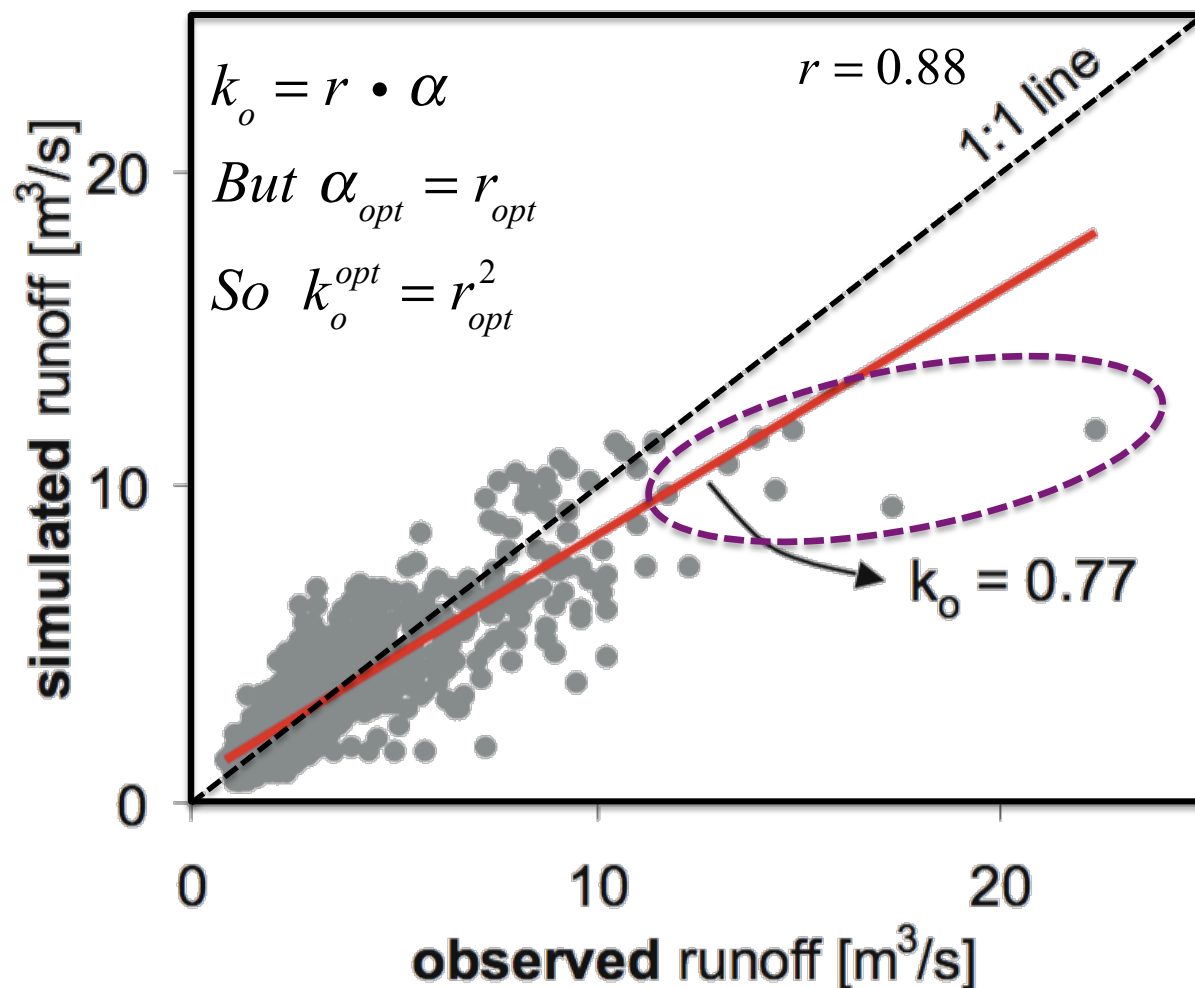
Ideal Value for  $\alpha$   
Optimized Value



**PROBLEM 4.**  
**PEAK FLOWS WILL**  
**BE UNDERESTIMATED**

THIS IS IN SPITE OF THE FACT THAT  
 'MSE' IS SUPPOSED TO GIVE  
 BETTER FIT TO THE LARGE  
 VALUES !!!

regression against observed runoff

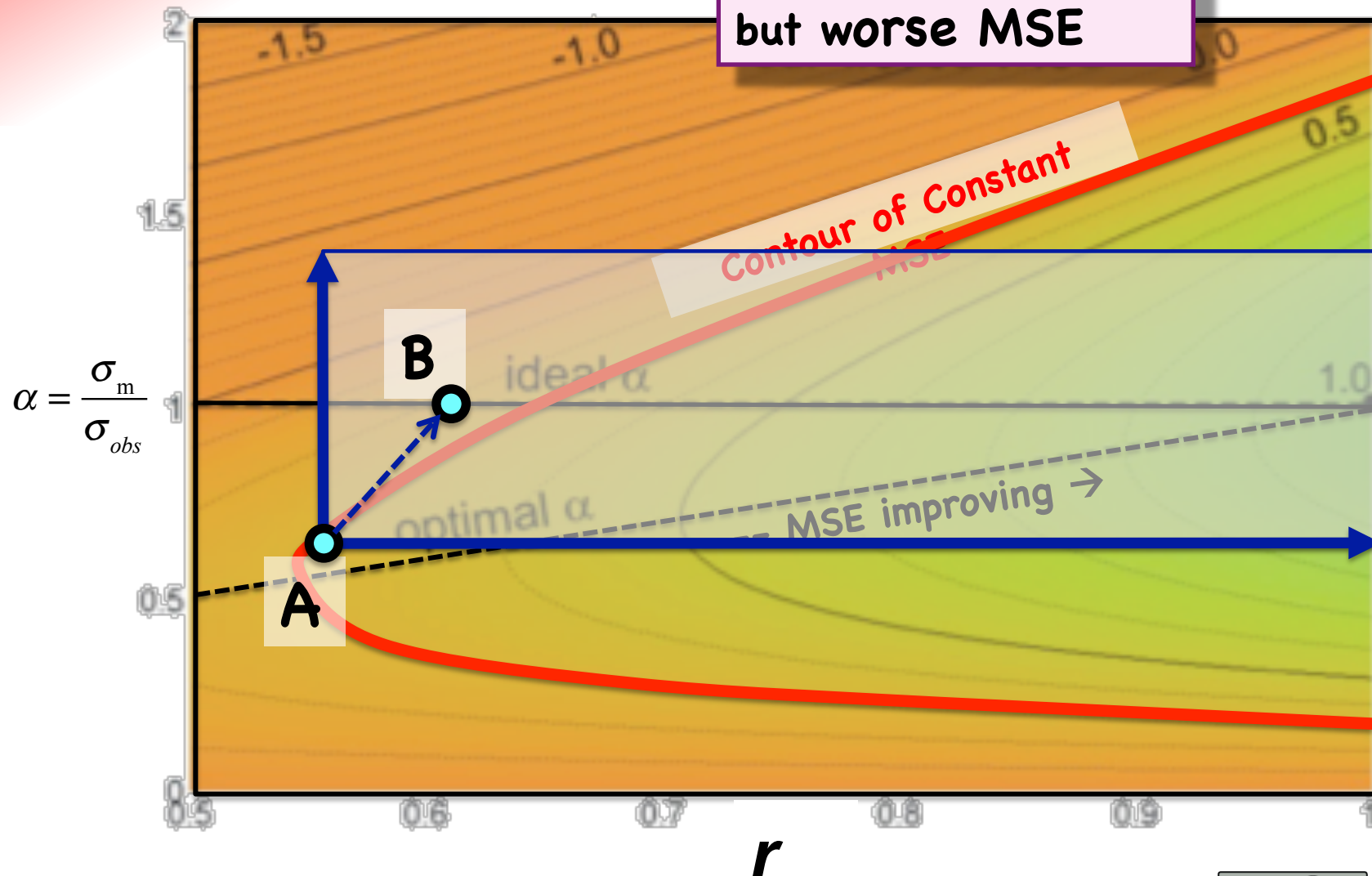


**PROBLEM 5.**

MODEL RESULT CAN BE  
BETTER FOR WORSE VALUE OF  
'MSE'

$$NSE = 1 - \frac{MSE}{\sigma_{obs}^2}$$

Point 'B' has  
better ' $r$ ' and ' $\alpha$ '  
but worse MSE



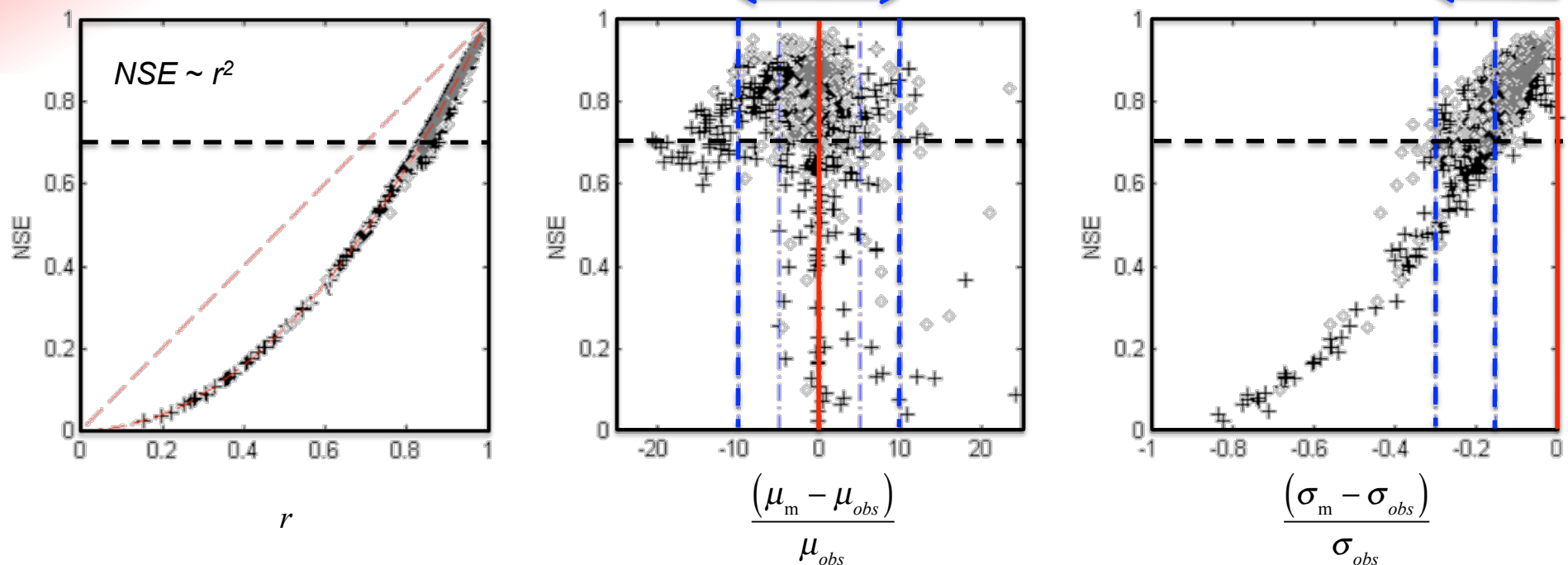
# OVERALL PROBLEM

COMBINED EFFECT IS  
*NOT WELL CONTROLLED*

$$NSE = 1 - \frac{MSE}{\sigma_{obs}^2}$$

$\pm 10\%$   
Volume Balance Error

$- 30\%$   
Variability Error



**Water Balance Model**  
**764 Catchments in Continental USA**  
*Calibrated using NSE measure and SCE optimization algorithm*

\* A Continental Scale Diagnostic Evaluation of the 'abcd' Monthly Water Balance Model for the Conterminous US  
 G.F. Martinez-Baquero & H.V. Gupta, *Manuscript in preparation*, 2009.



Hoshin Gupta

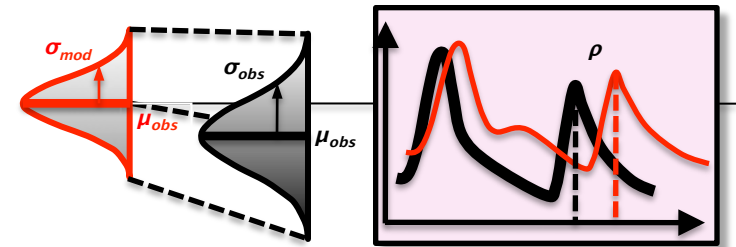
# WHAT CAN WE LEARN FROM THIS?

$$MSE = \frac{1}{N} \sum_{t=1}^N [Q_t^{obs} - Q_t^m(\theta)]^2$$

$$MSE = \Delta\mu^2 + \Delta\sigma^2 + 2\sigma_m \sigma_{obs}(1 - \rho)$$

1. Optimization using  $MSE$  is equivalent to trying to match *THREE* statistical properties of the data

*Data Mean (1<sup>st</sup> moment)*      -  $\mu_{obs}$   
*Data Variance (2<sup>nd</sup> moment)*      -  $\sigma_{obs}^2$   
*Data Correlation structure*      -  $\rho$



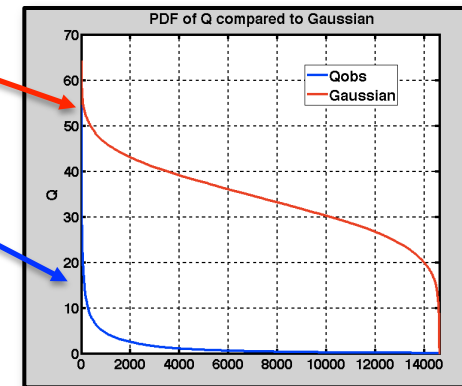
2. Two are properties of the data PDF, and the third is a property of the spatial &/or temporal correlation structure
3. These are combined in a way that emphasizes certain aspects of system behavior ... at the expense of others
4. For catchment modeling this can result in poor Water Balance and under-estimation of Variability – both being important system behaviors we wish to reproduce



## BUT... WHY ONLY *THESE* THREE PROPERTIES ?

Gaussian CDF  
(symmetrical)

Typical CDF of  
Streamflow  
(highly skewed)



1. Data *PDF*'s are very rarely Gaussian !
2. The Model should also reproduce *other* statistical properties of the data – particularly ones with hydrological significance !
3. *Linear correlation* ' $r$ ' aggregates different kinds of information about *spatio-temporal correlation structure* into ONE measure

THIS IS AN INEFFICIENT  
WAY TO EXTRACT  
INFORMATION !

# THE CHALLENGE OF MODEL EVALUATION

TO DEVELOP  
**“SUFFICIENT STATISTICS”**  
 OF MODEL PERFORMANCE  
 THAT ARE  
**“DIAGNOSTICALLY RELEVANT”**  
 TO THE PROBLEM

Qualitative Evaluation of Behavior  
Consistency

Qualitative  
Interpretation  
Of Model Simulated  
Fields

Numerical  
model

Quantitative  
Model Simulated  
Fields

Extent  
Support  
Spacing

Quantitative  
Evaluation  
of Behavior

CLOSENESS

Degree of Accuracy (Bias)  
Degree of Precision (Uncertainty)  
Degree of Correspondence (“Correlation”)

States  
Outputs

Sampling –  
Representativeness  
Informativeness

Extent  
Support  
Spacing

Reality

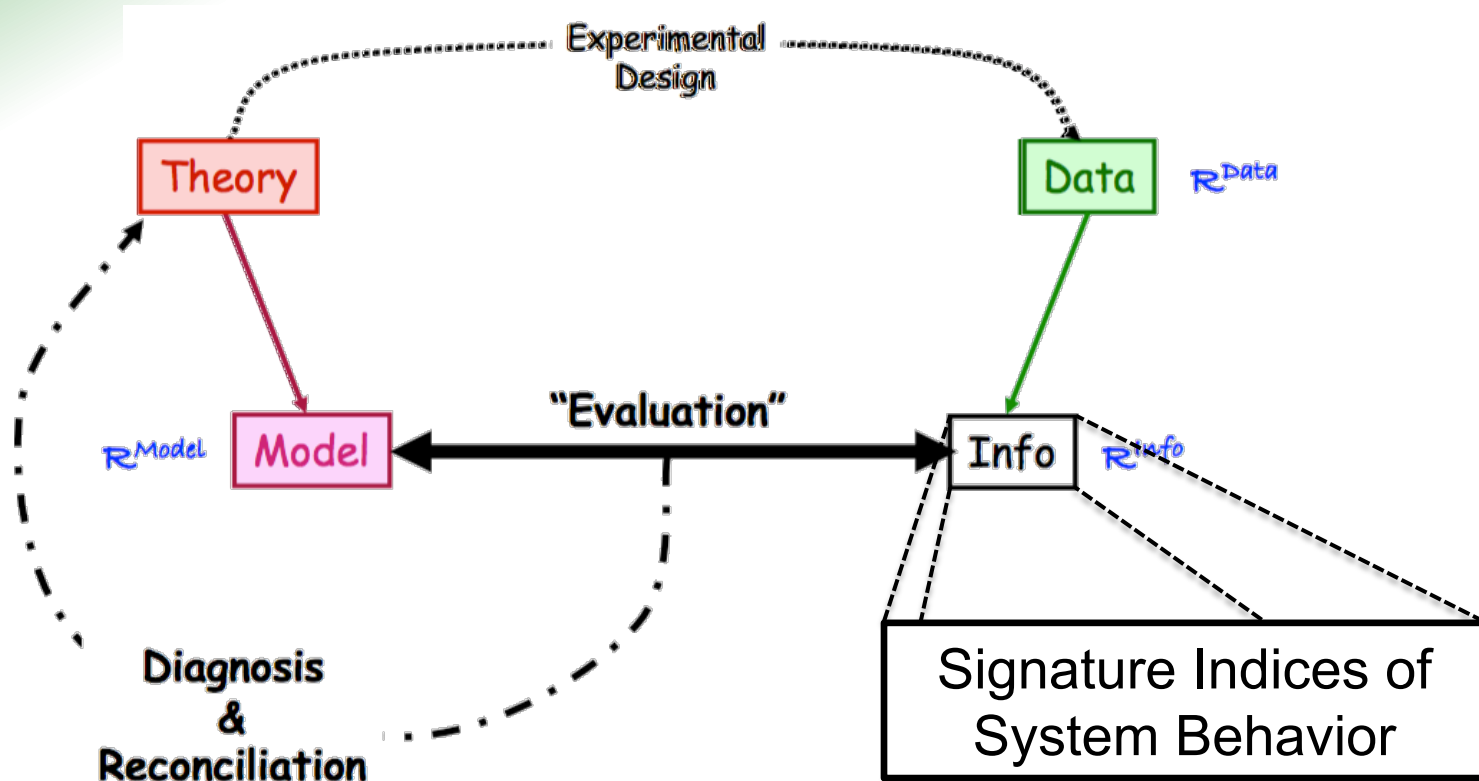
Quant  
Observ  
Field

Model

Data  
Assimilation  
Model

# DIAGNOSTIC APPROACH TO MODEL EVALUATION

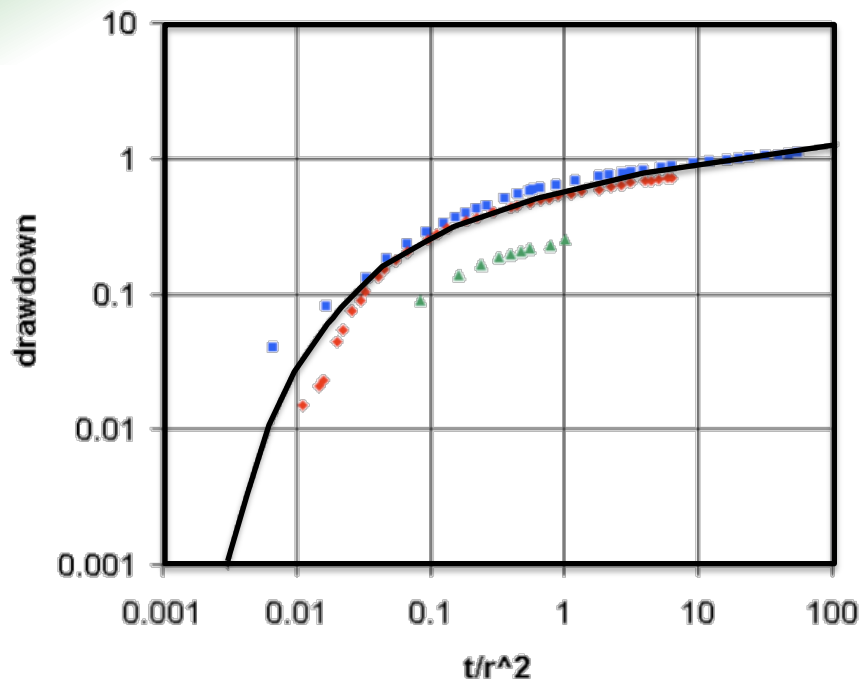
DATA IS NOT  
INFORMATION !!!



MODEL REFERENCED PATTERNS  
SHOULD BE RECONCILED WITH  
DATA REFERENCED PATTERNS

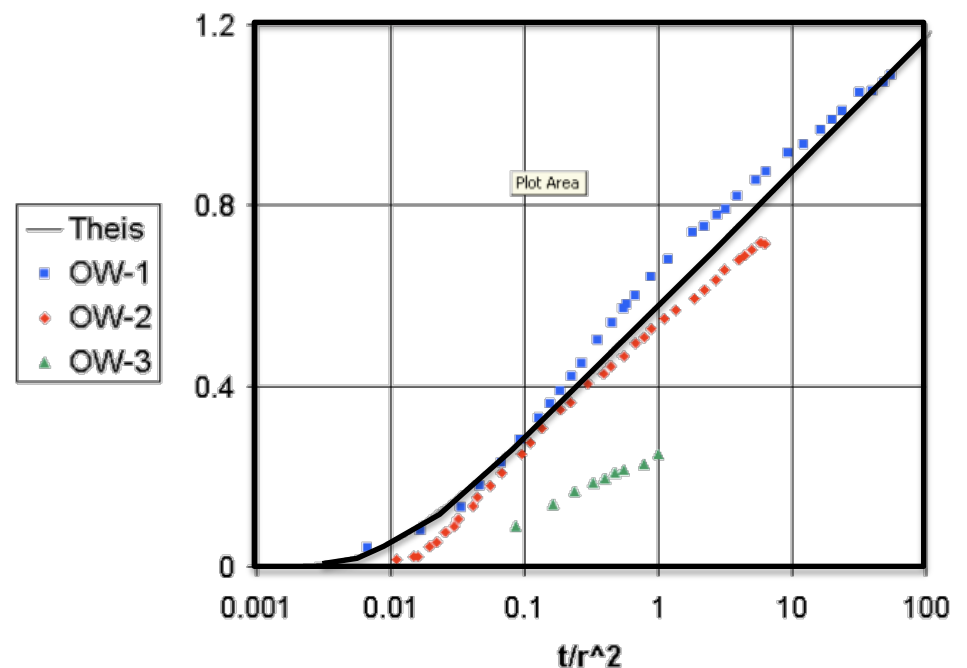
# WHAT CONSTITUTES A SIGNATURE BEHAVIOR ?

Matching the Theis type curve to observed drawdown on log-log plot



## Aquifer Well Test Analysis

Matching the Theis type curve to observed drawdown on semi-log plot



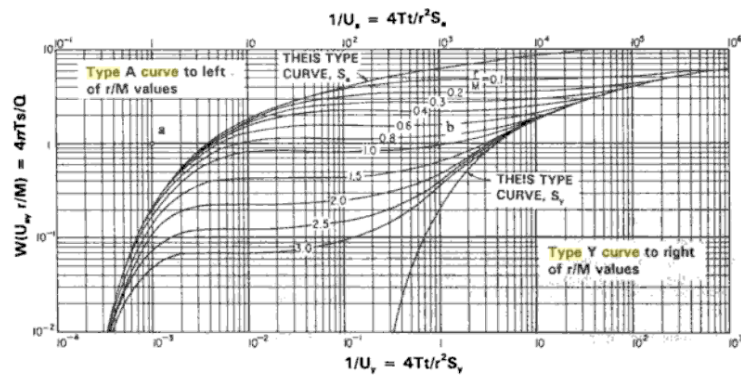
$$D = \frac{Q}{4\pi T} W(u) \quad \text{where} \quad u = \frac{r^2 S}{4Tt}$$

D = drawdown  
Q = pumping rate  
T = Transmissivity  
u = dimensionless time

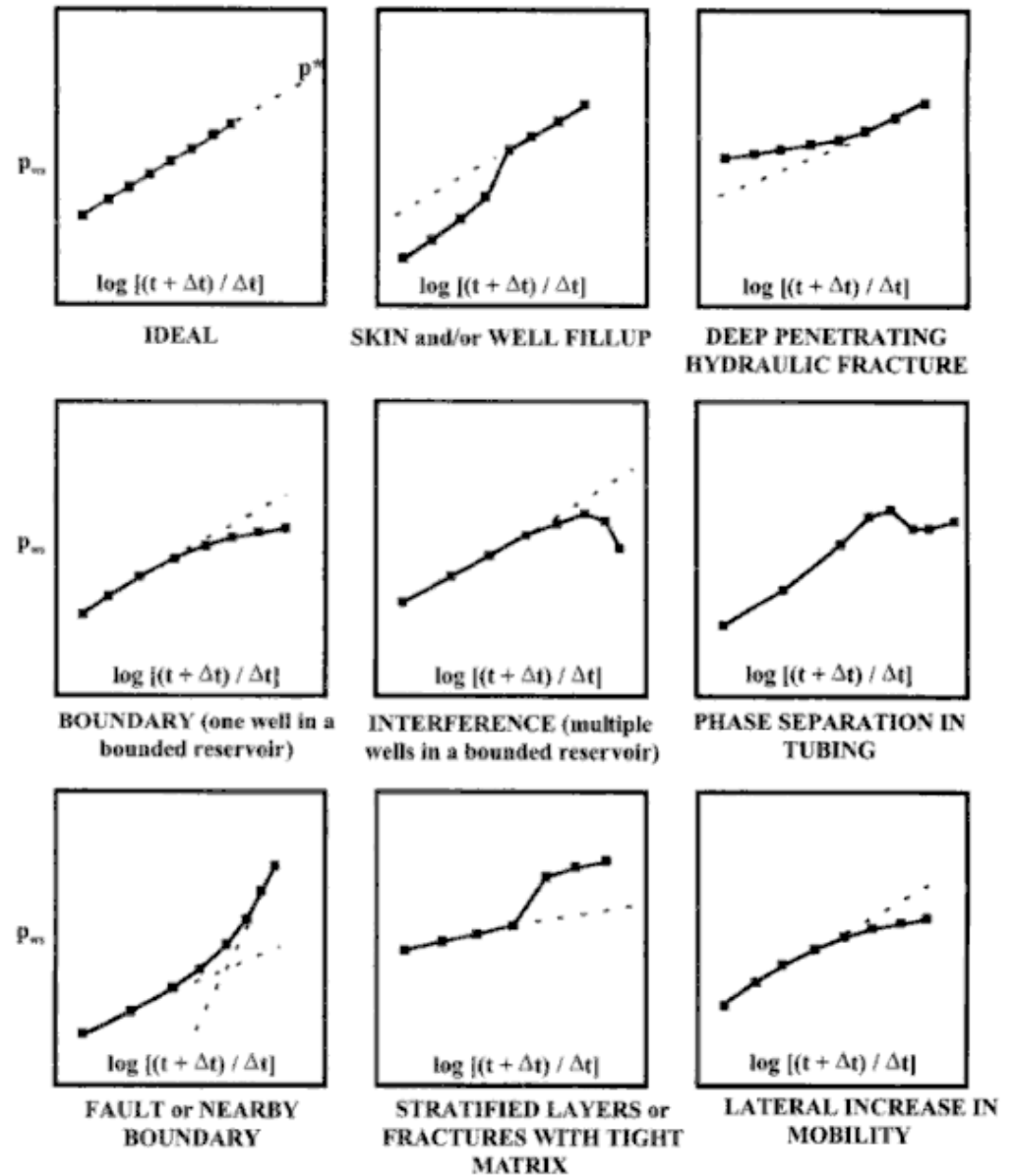
S = Storativity  
r = distance from well  
t = time

# WHAT CONSTITUTES A SIGNATURE BEHAVIOR ?

## Tests in Unconfined Aquifers



## Example Build-up Curves Illustrating Various Effects



From Matthews & Russell<sup>1</sup>



Hoshin Gupta

# IN CONCLUSION

THE PROBLEM OF INFERENCE  
(Reconciling Theory With Obs)



PROBLEM OF DEVELOPING  
DIAGNOSTIC SUFFICIENT  
STATISTICS

**The MODELING problem:**

*To explicitly state*

- a) The Hypothesis to be tested*
- b) The Tests that will unambiguously challenge the hypothesis.*

**The OBSERVATIONAL problem:**

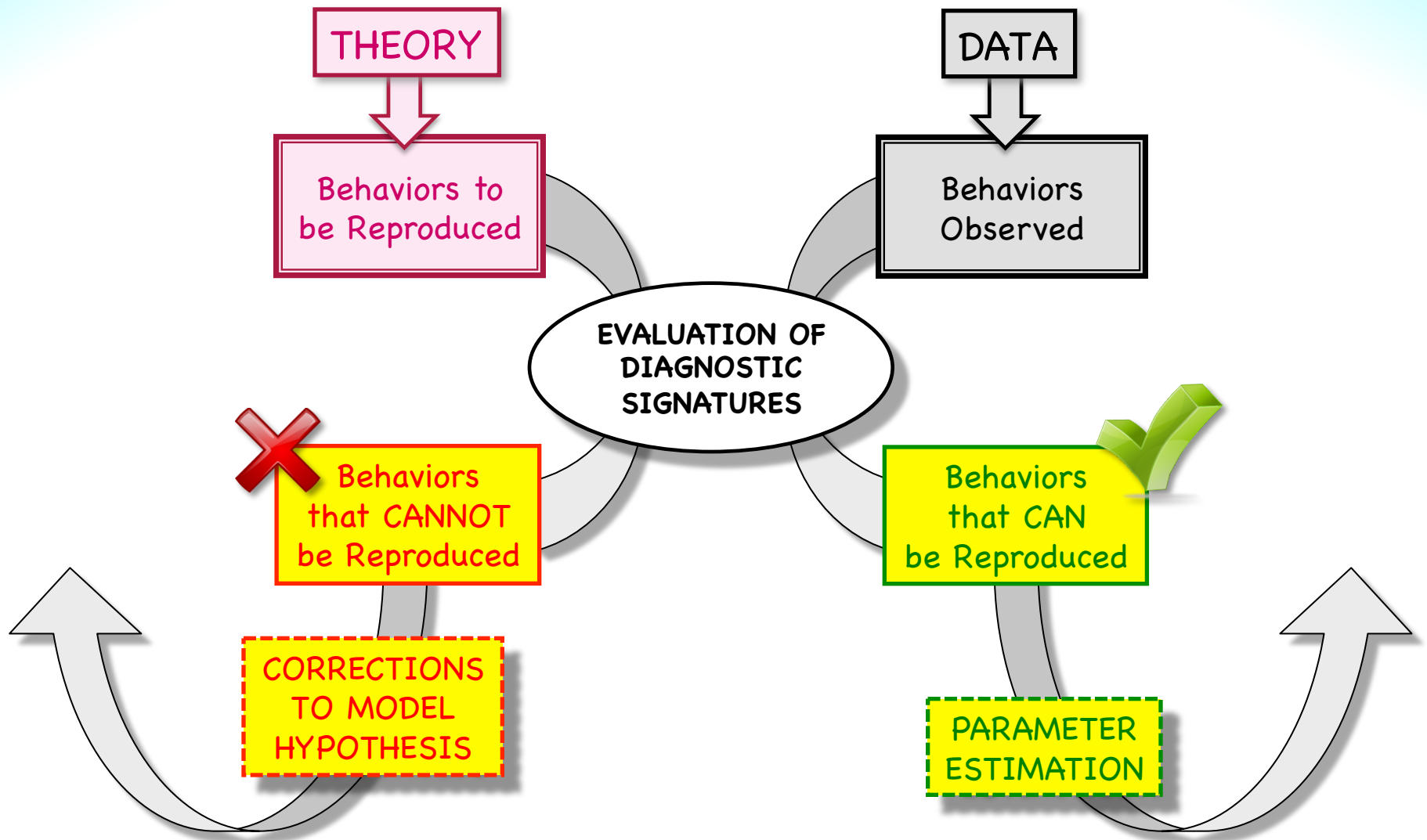
*To extract from the data, INFO that*

- a) Diagnostically characterizes system behavior*
- b) Supports or challenges the model hypothesis*

**The RECONCILIATION problem is to:**

- a) Make robust inferences regarding which aspects of the model hypothesis are (are not) supported by the observations*
- b) Diagnostically guide improvements to the theory (model)*
- c) Suggest improvements in the acquisition of observations*

# MODEL REFERENCED PATTERNS ARE TO BE RECONCILED WITH DATA REFERENCED PATTERNS







# Corrections to Model Hypothesis ... WRR 2009

WATER RESOURCES RESEARCH, VOL. 45, W00B13, doi:10.1029/2007WR006749, 2009

Click  
Here  
for  
Full  
Article

## Estimating the uncertain mathematical structure of a water balance model via Bayesian data assimilation

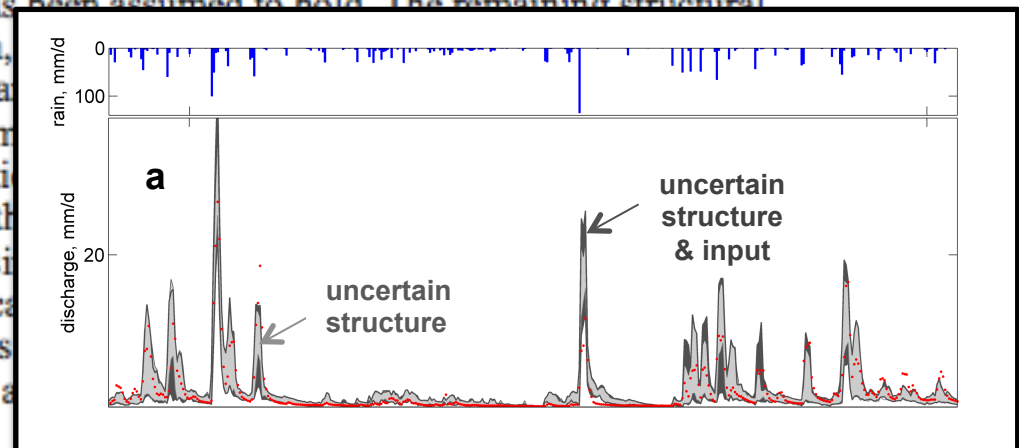
Nataliya Bulygina<sup>1,2</sup> and Hoshin Gupta<sup>1</sup>

Recei

[1]  
struc  
the  
suita

### ... Correcting the Model EQUATIONS ...

determined; that is, the system boundaries have been specified, the important state variables and input and output fluxes to be included have been selected, the major hydrological processes and geometries of their interconnections have been identified, and the continuity equation (mass balance) has been assumed to hold. The remaining structural identification problem that remains, then, dependence of the output on the inputs a model can be constructed for making sim input-state-output behavior. The conventi some fixed (and possibly erroneous) math We show instead how Bayesian data assi (construct) the form of these mathematical consistent with macroscale measurements state variables. The resulting model has a



Hoshin Gupta

**Improving Model Identification:  
Reconciling Theory with Observations &  
The Problem of Sufficient Statistics**

Evaluation should enable us to link  
what we "see" in the data to  
what is "right" and "wrong" with our models.

This task will require the active  
collaboration of Process Scientists,  
Modelers & Systems Theorists.

