

MULTIPLE-POINT STATISTICS TO ASSESS COMPLEX SPATIAL UNCERTAINTY

Philippe Renard

MascotNum, 25th of April 2014, Zürich

Stochastic Hydrogeology Group
University of Neuchâtel
Switzerland



J. Kerrou, G. Mariethoz,
J. Straubhaar, A. Comunian,
G. Pirot, F. Oriani



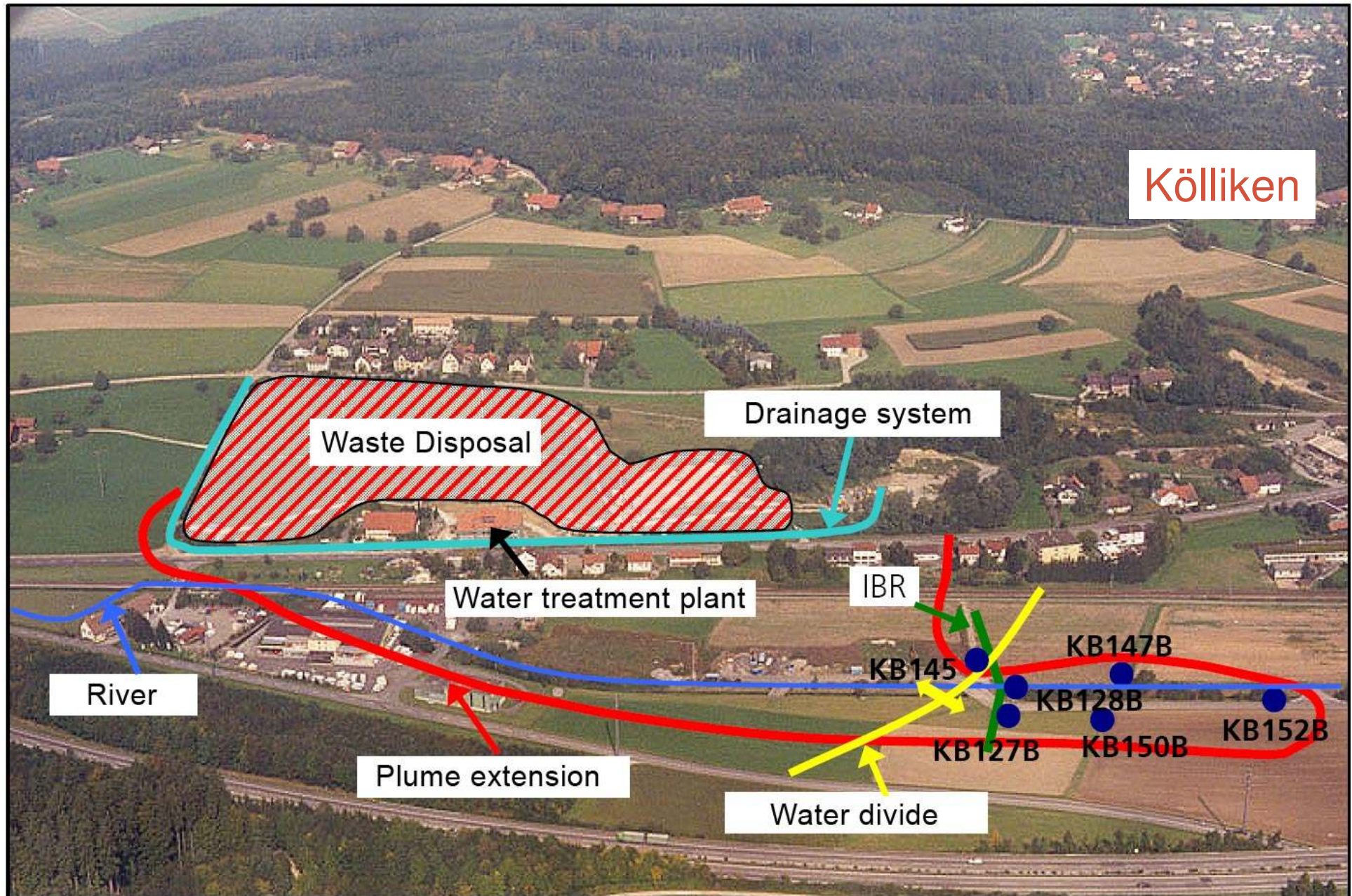
SWISS NATIONAL SCIENCE FOUNDATION



MOTIVATION FOR MPS

What is our problem / limits of current approaches

Will the contamination reach drinking water supply?



Confining building

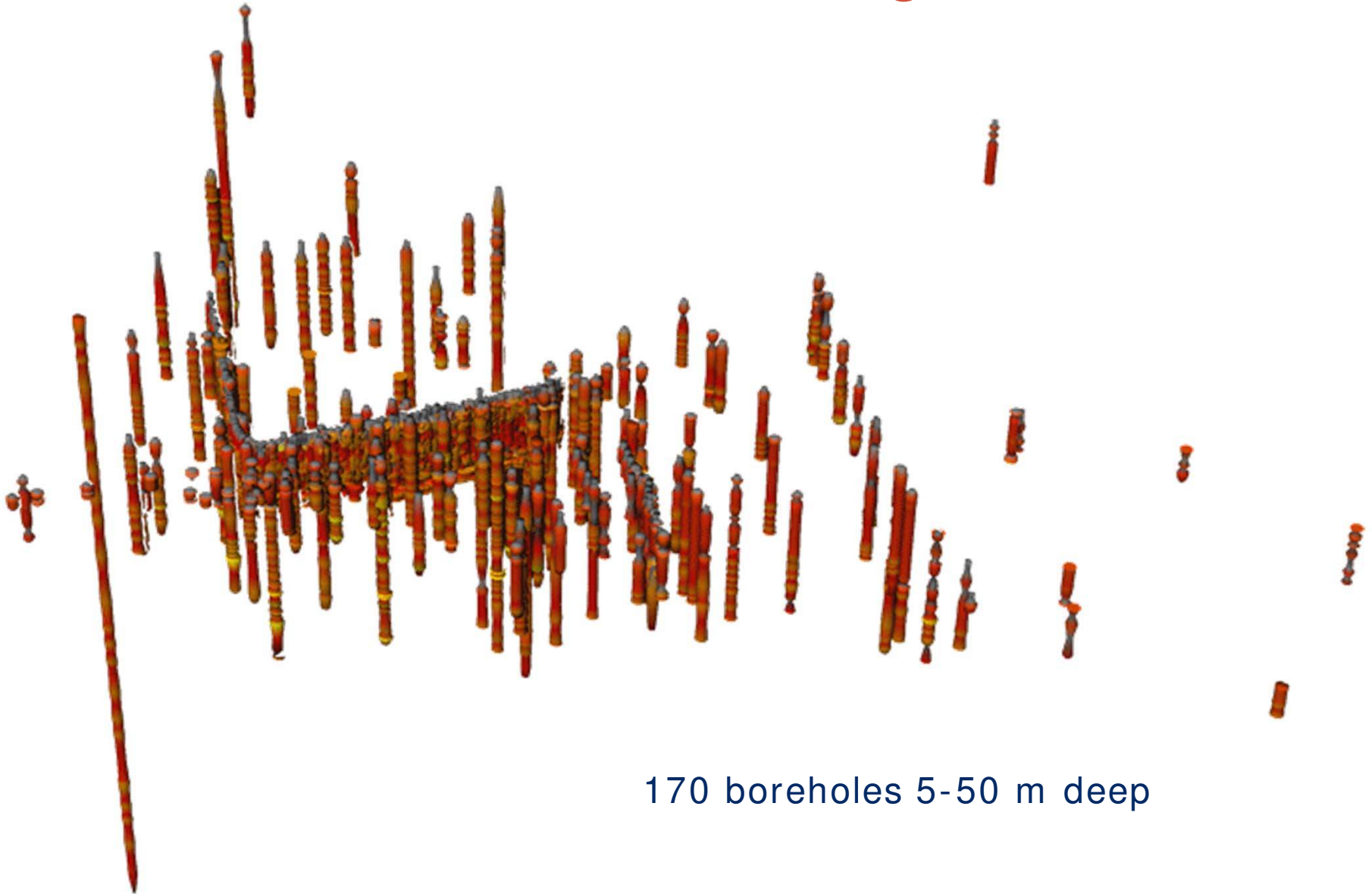


Waste excavation



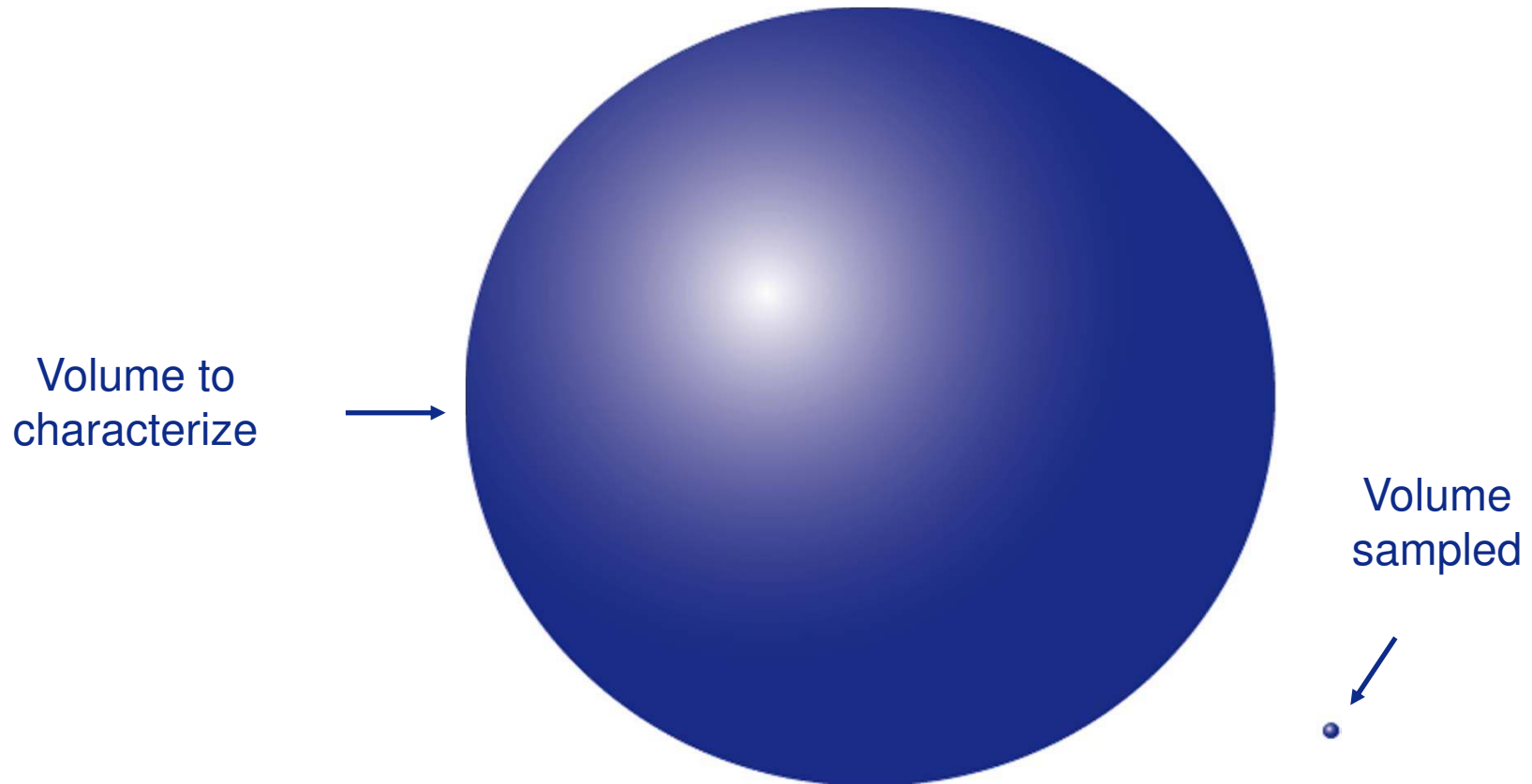
Total cost > 770 M. CHF

Extensive investigations



170 boreholes 5-50 m deep

Characterization issue



Volume to characterize : $850 \times 400 \times 70 = 24 \text{ millions m}^3$

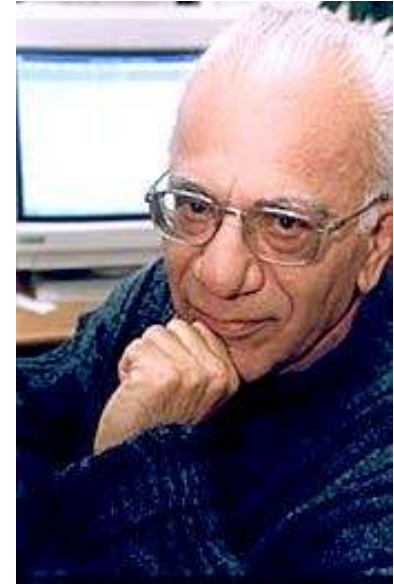
Volume sampled by the boreholes: $6700 \times 0.01 \times \pi = 200 \text{ m}^3$

Understanding groundwater flow

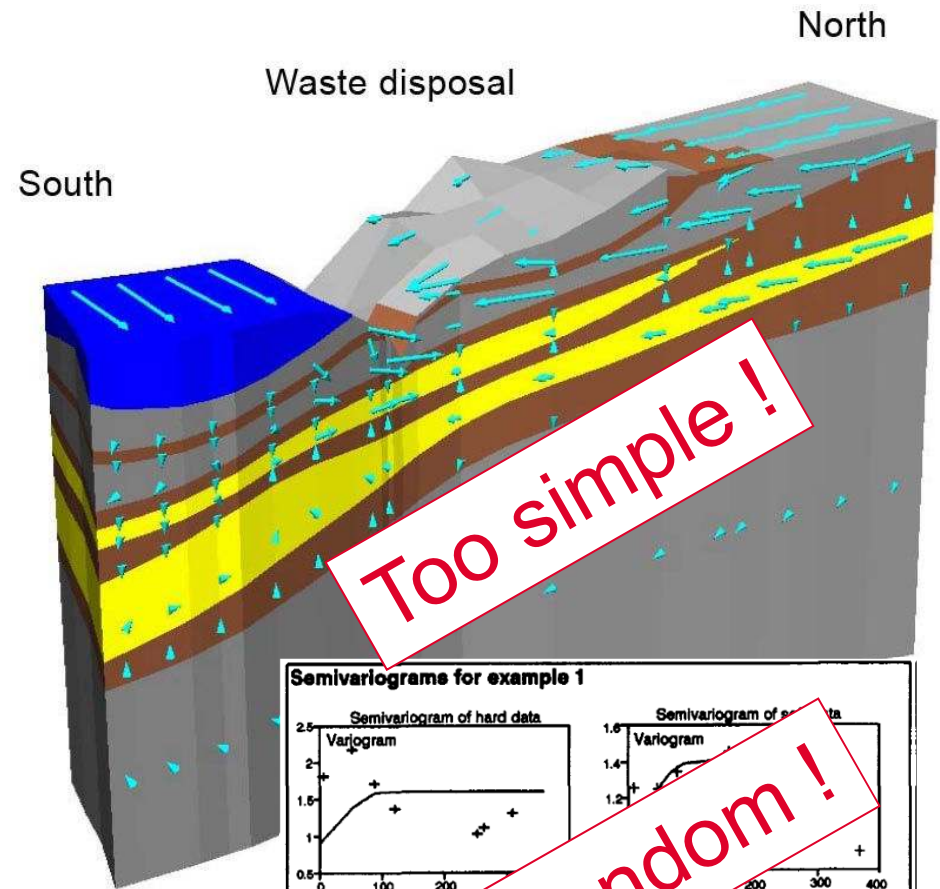
- Well established PDE / Numerical models
- **Huge uncertainty** due to
 - rock **heterogeneity** + lack of data (field of parameter)
 - badly controlled boundary conditions / source terms
- Long tradition (>30 years) to use Gaussian random fields
 - Interpolating parameters
 - Understanding the physics of heterogeneous materials
 - Estimating uncertainty
- Heavy numerical forward models (e.g.CO2 sequestration)
- Today → how to build random fields of input parameter

Gedeon Dagan (2002)

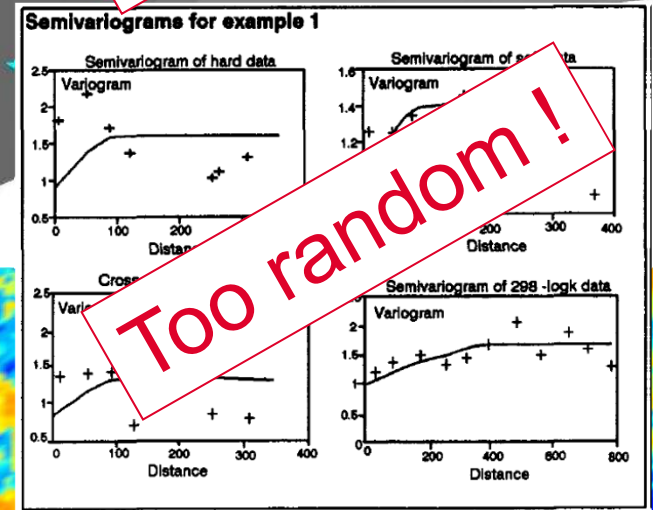
- “The stochastic modeling of groundwater has developed considerably ... [but it] hasn't yet become a routine tool”
- Debate in the community
- Situation is more subtle
- Various issues
 - Education: lack of people
 - Structural issues: consulting market
 - Relevance of the models



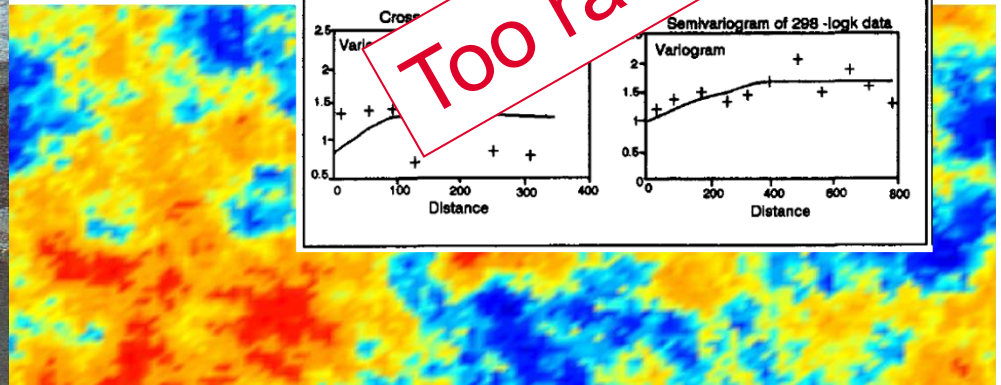
One fundamental reasons



Too simple!

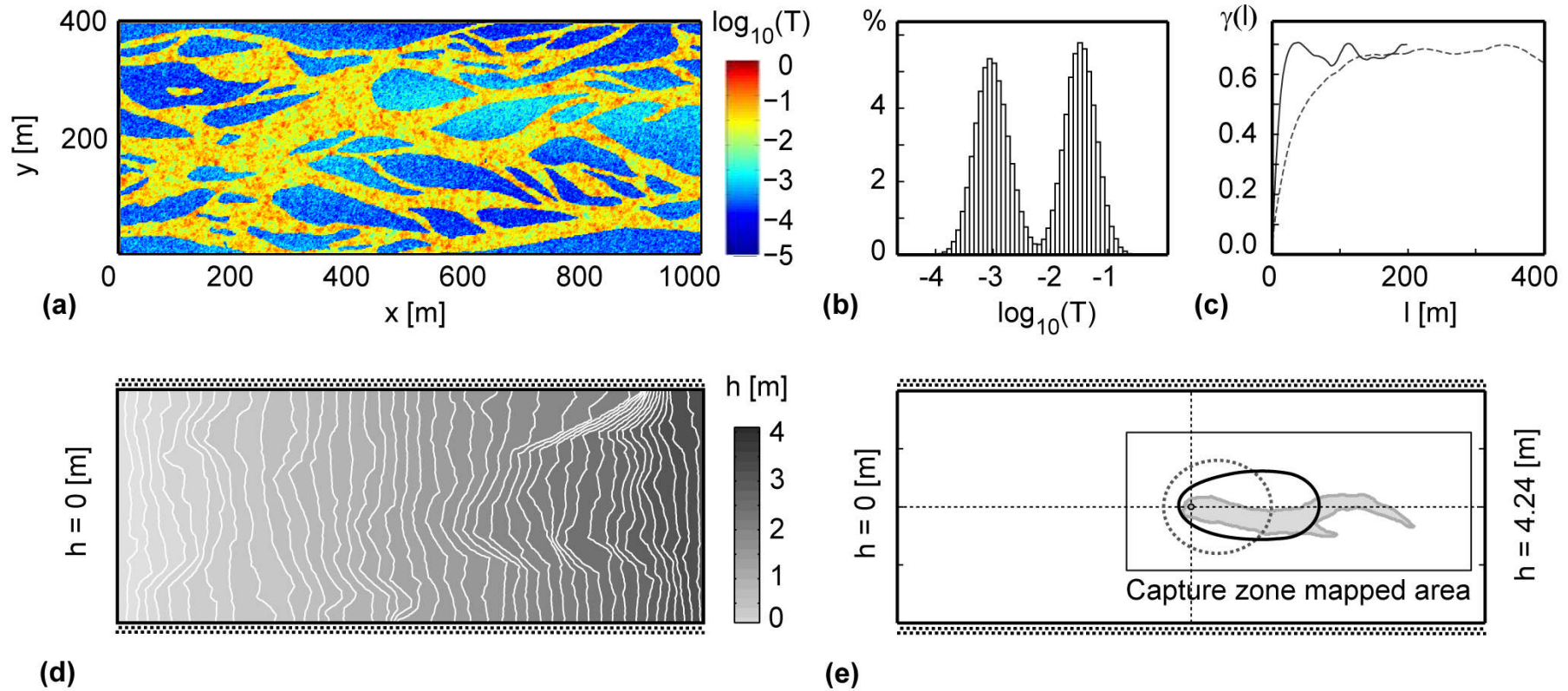


Too random!



A synthetic example

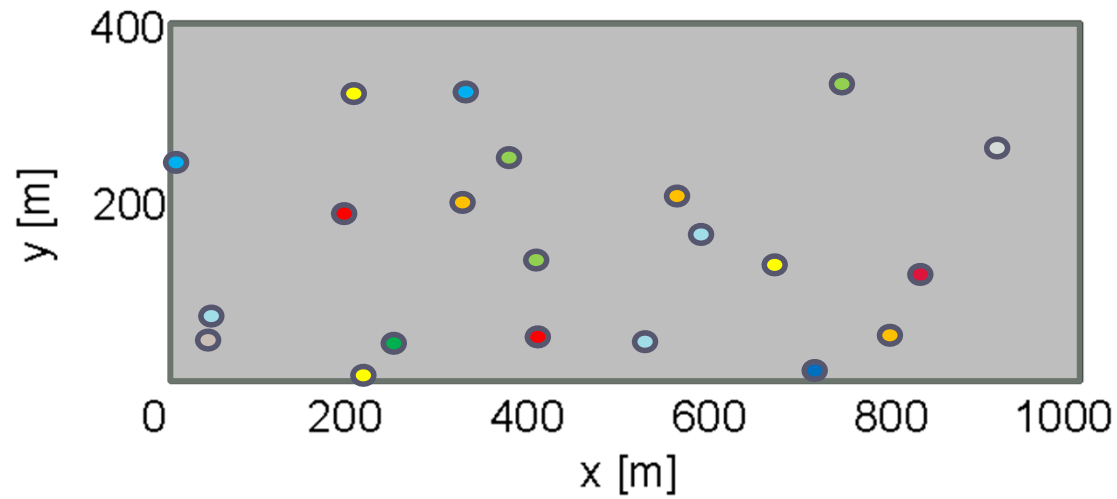
Transmissivity field



Characterization data set

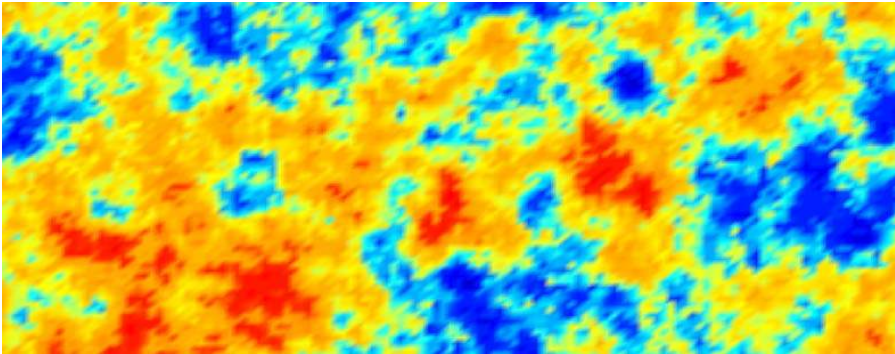
Forecasts data set

Sampling the reference

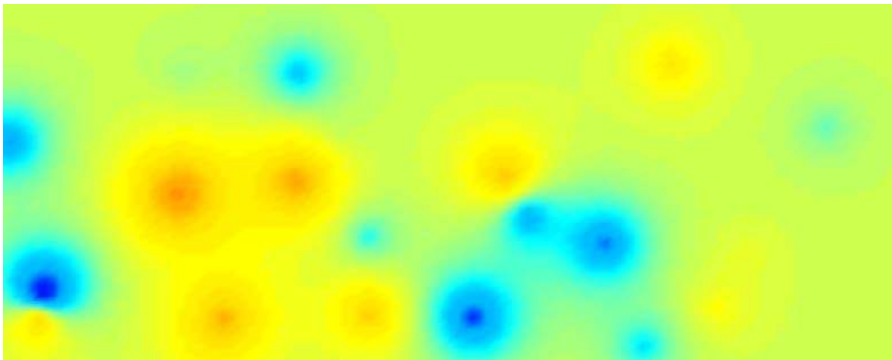


Heterogeneity characterization

one simulation of $\log_{10}(T)$



ensemble average of $\log_{10}(T)$

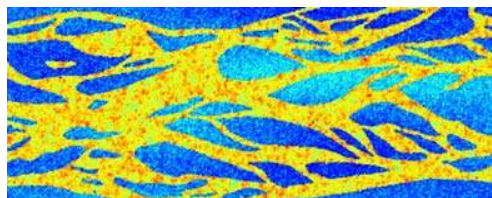


Simulation of N
conditional to 21 values

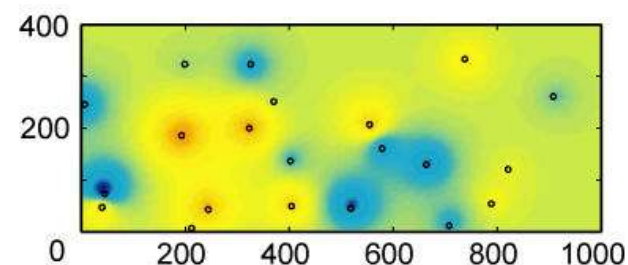
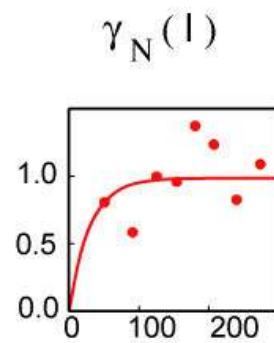
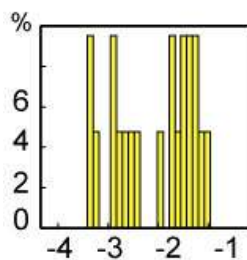
Back transform to get Y
= $\log_{10}(T)$

Turning bands method

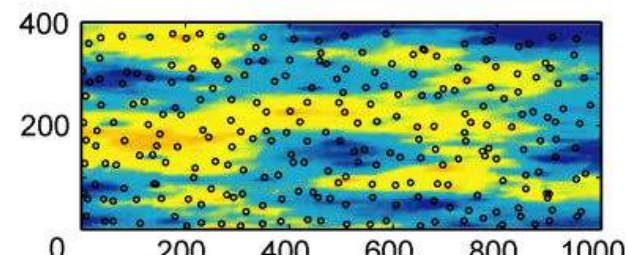
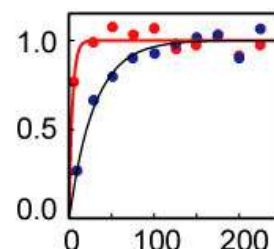
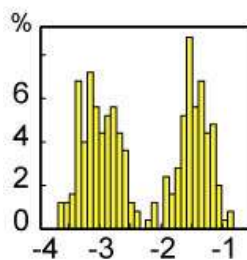
Adding information T measurements



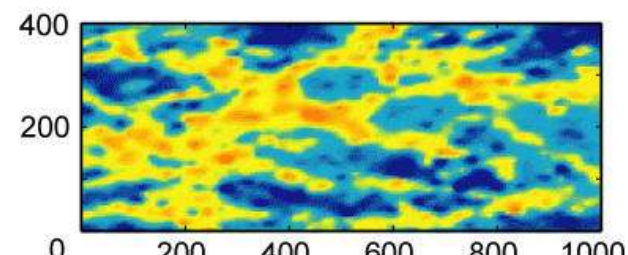
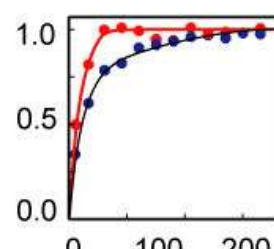
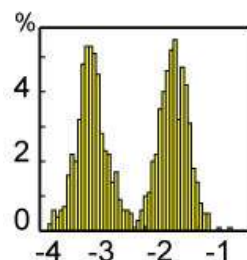
21 T



250 T



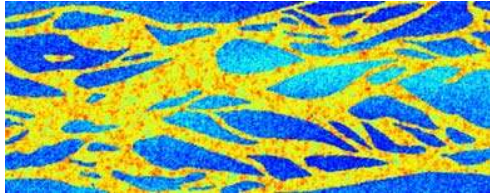
1000 T



$\log_{10}(T)$

I [m]

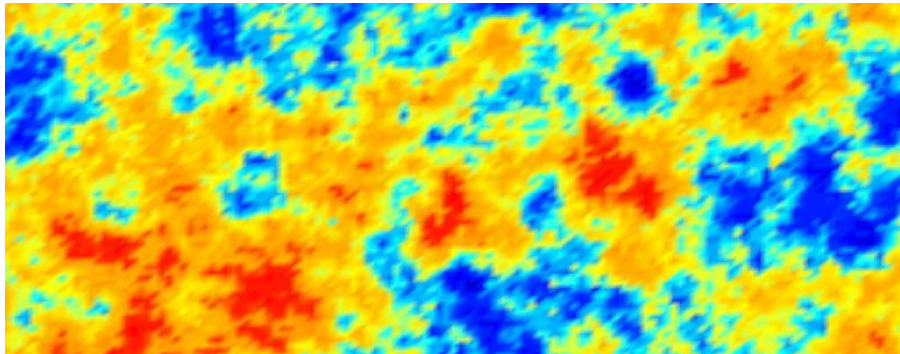
x [m]



Adding head data

Sequential self calibration method
inverto code (Hendricks-Franssen, 2001)

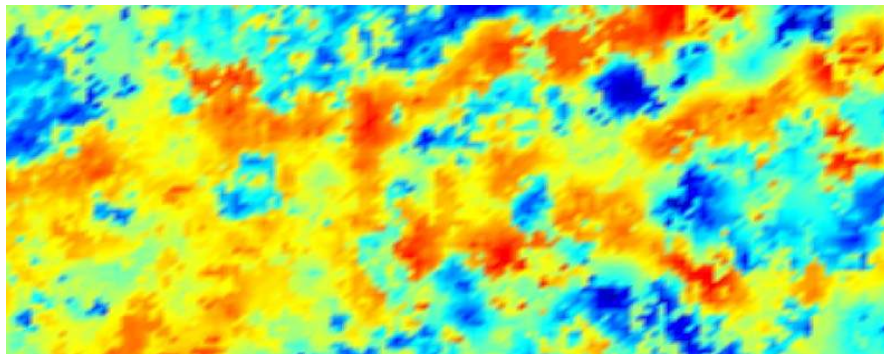
Variogram from the data – 2 master blocks per correlation length



21 T, 0 heads

21 T, 21 heads

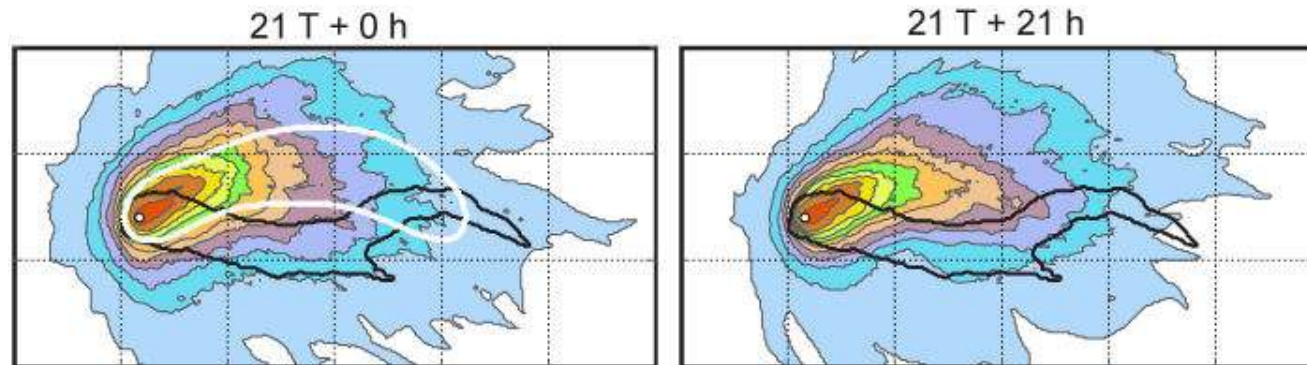
21 T, 250 heads



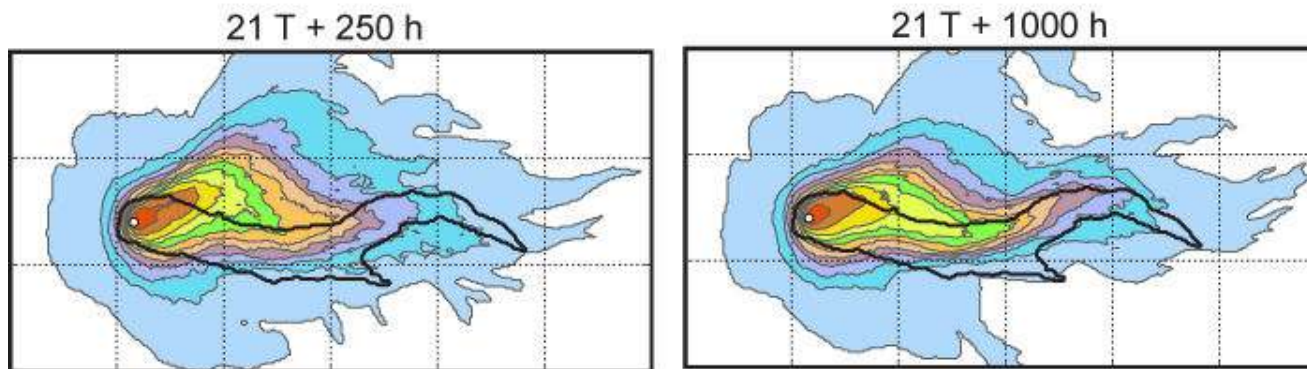
21 T, 1000 heads

100 simulations

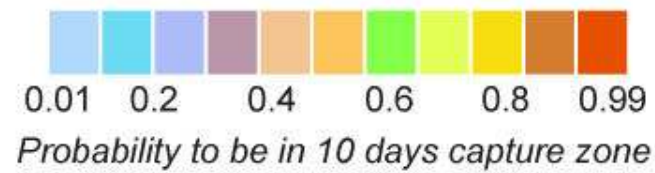
Reliability of transport forecasts



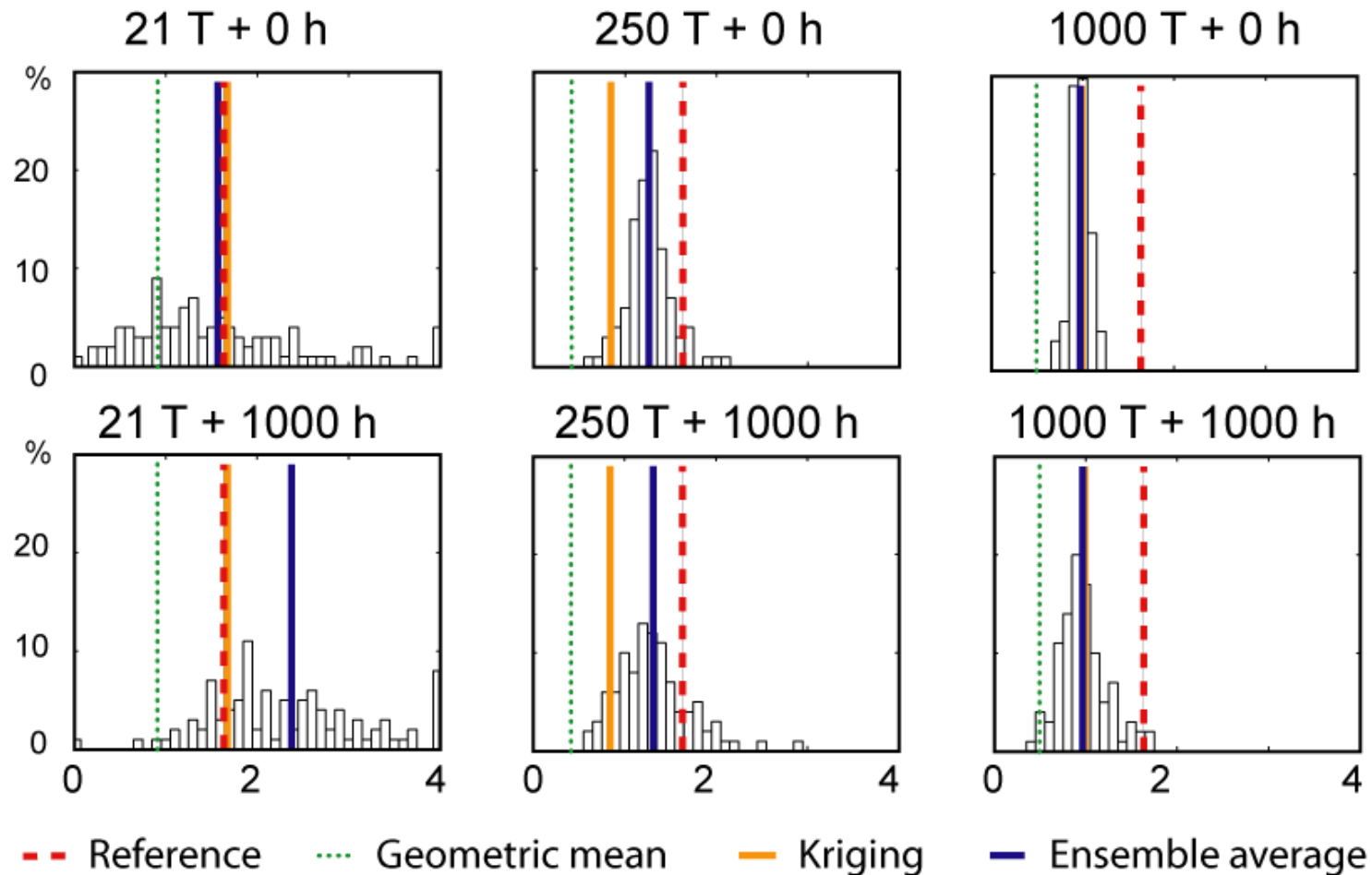
True value belongs to the ensemble forecasts



Head conditioning reduces uncertainty and increases accuracy

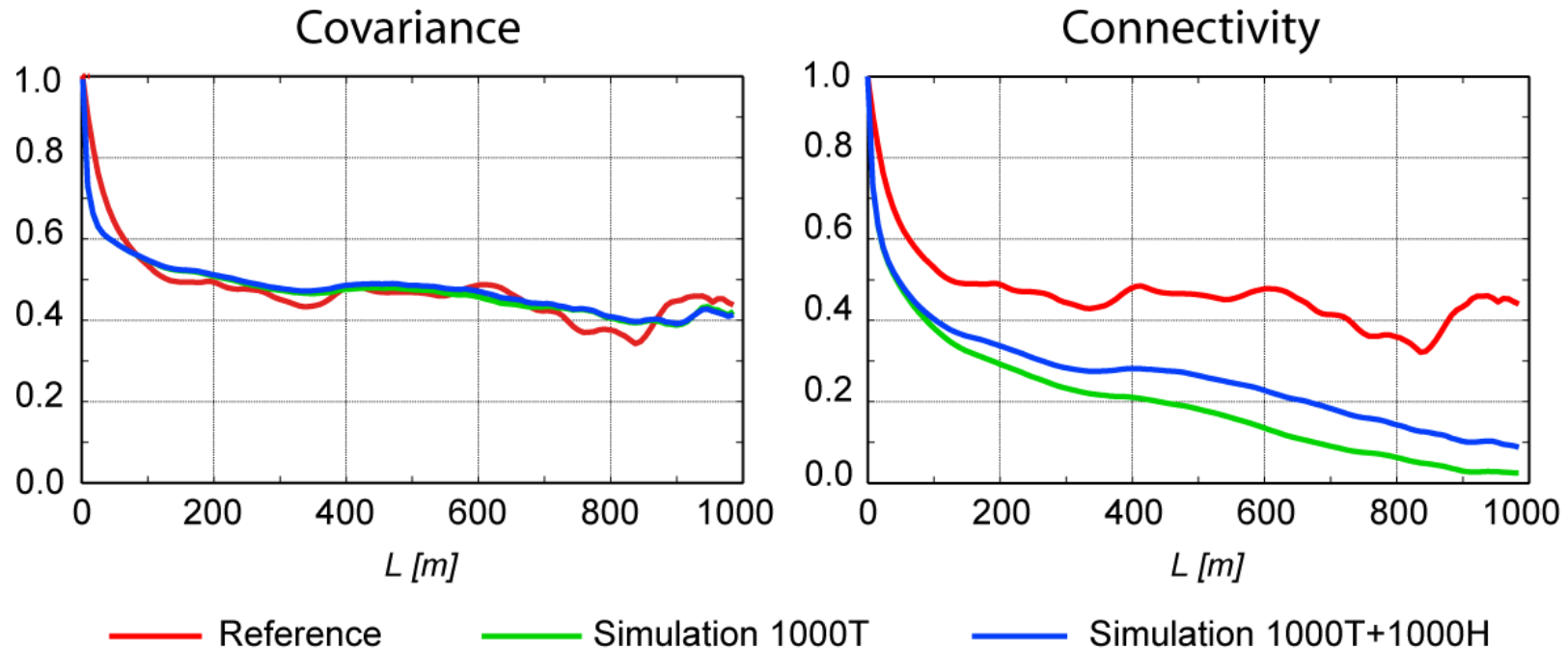


Limits of the multi-Gaussian approach



Adding transmissivities reduces uncertainty and increases bias
The bias is partly compensated by conditioning to head

Multi-Gaussian approach is insufficient



Covariances are well reproduced by the simulations
Connectivity are not
Conditioning to head improves connectivity

Same observations with field data

An Evaluation of Conditioning Data for Solute Transport Prediction

by Timothy D. Scheibe^{1,2} and Yi-Ju Chien¹

ground
water

The results show that **conditioning to a large number of small-scale measurements** does not significantly improve model predictions, and **may lead to biased or overly confident predictions.**

However, **conditioning to geophysical interpretations with larger spatial support significantly improves the accuracy and precision of model predictions.**

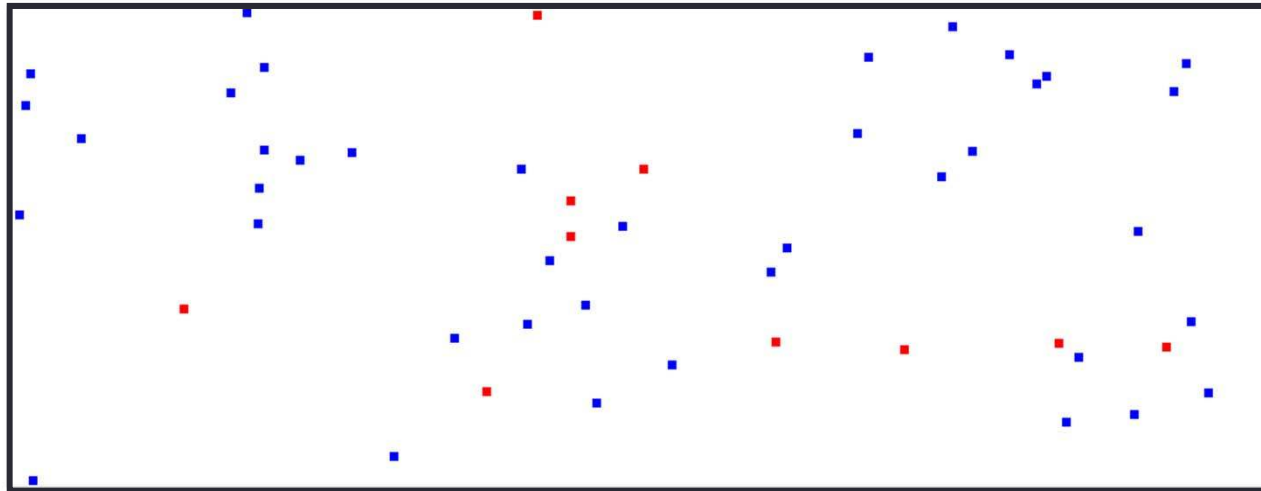
all cases, the effects of model error appear to be significant in relation to parameter uncertainty.

MULTIPLE POINT STATISTICS

What is it? / Principle of the method

Data set

Map of geological samples

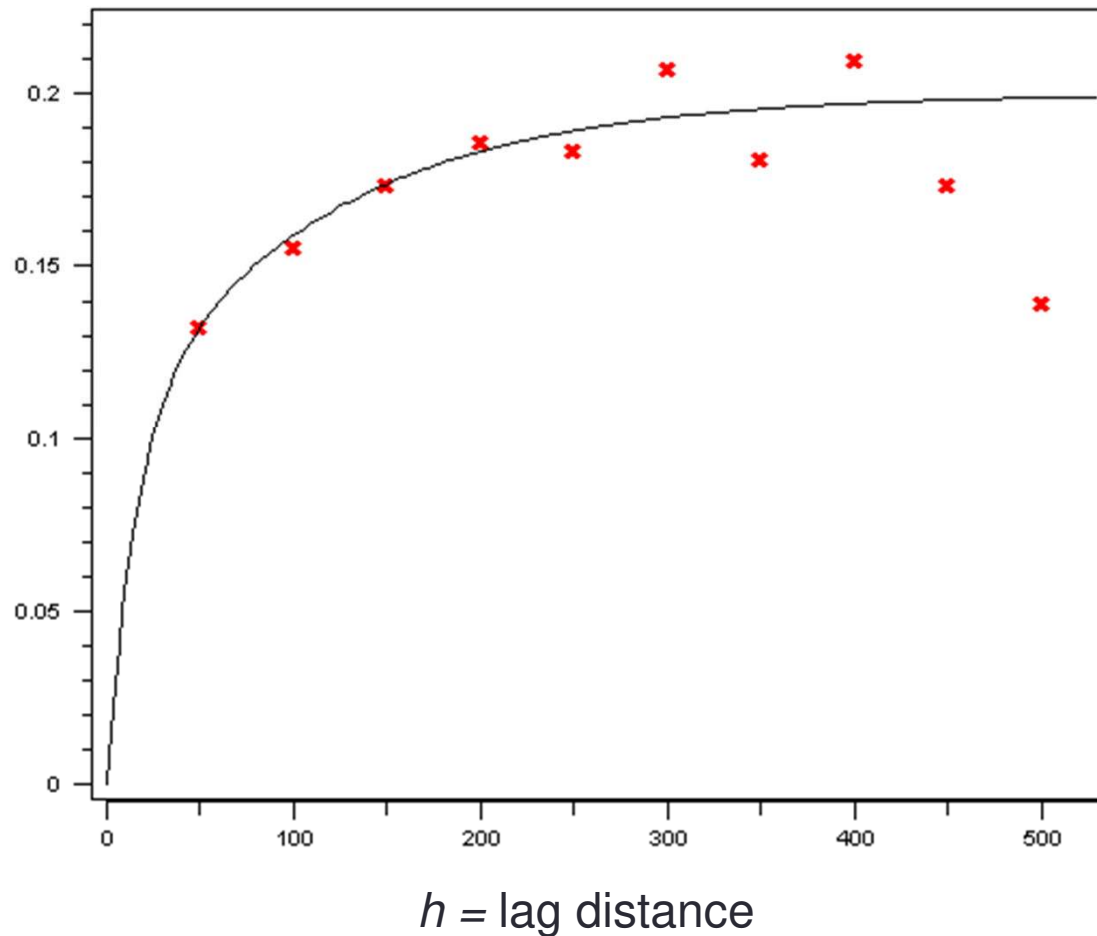


■ $I(x) = 0$ \iff $x \in \text{Clay}$

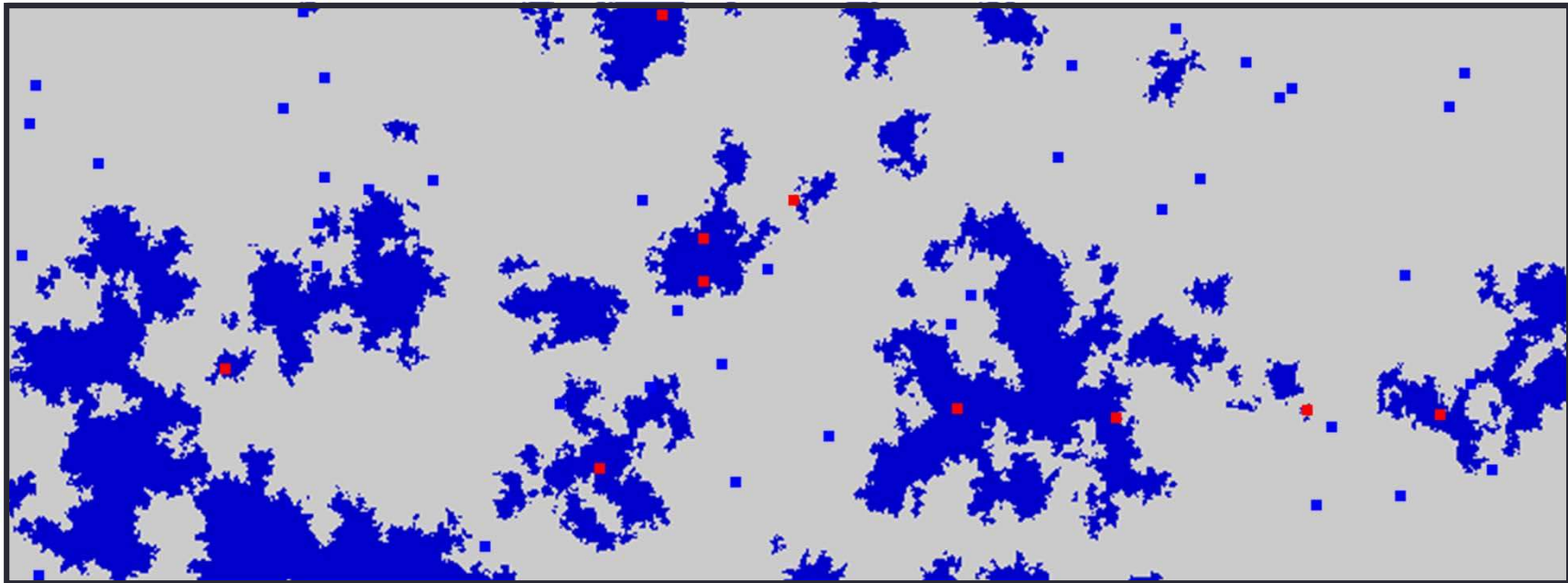
■ $I(x) = 1$ \iff $x \in \text{Sand}$

Indicator variogram

$$\gamma_i(h) = \frac{1}{2} E[(I_i(x+h) - I_i(x))^2]$$

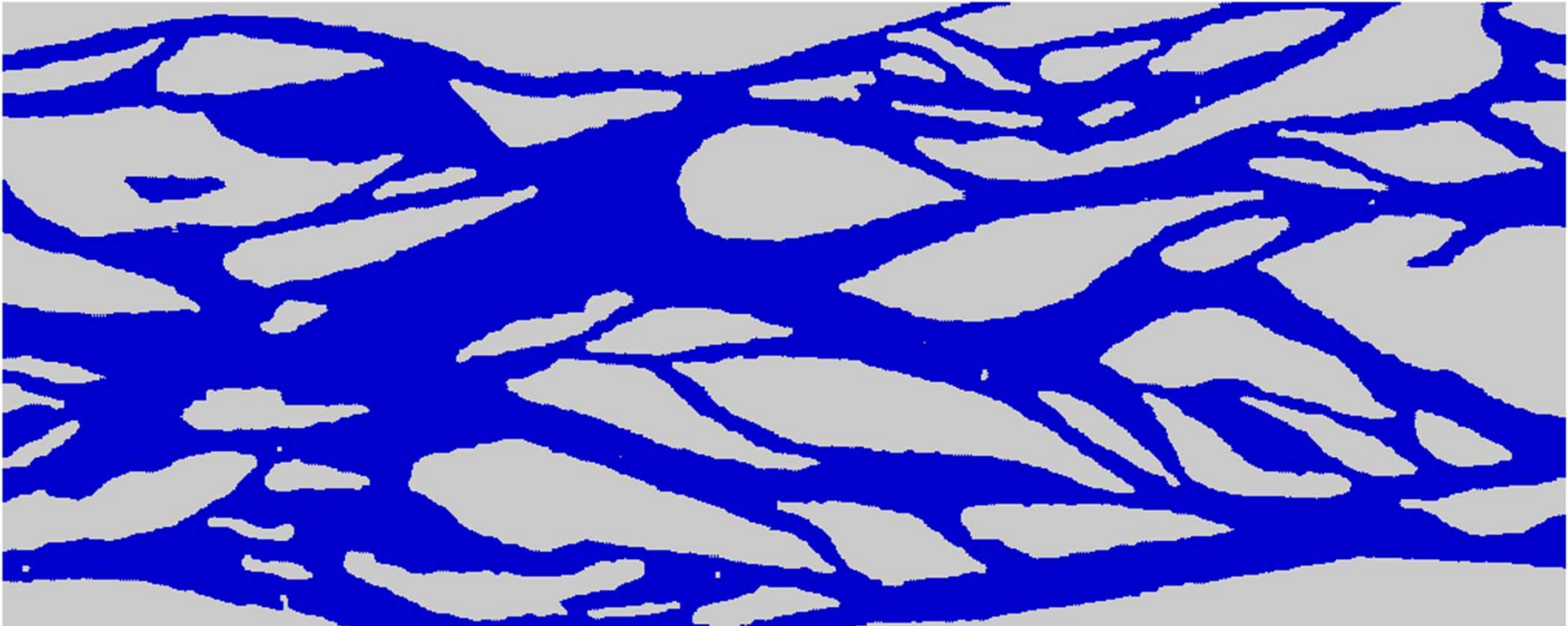


Sequential indicator simulation



Honors the variogram and proportions of the data

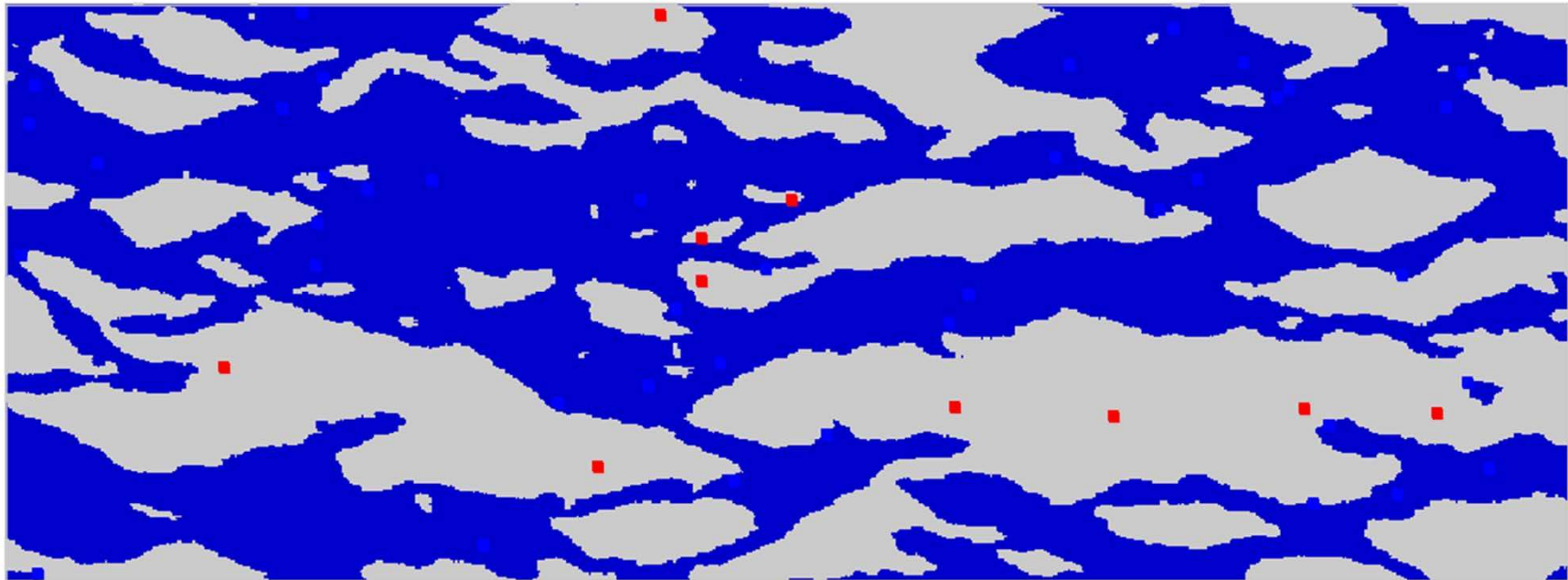
Braided channel



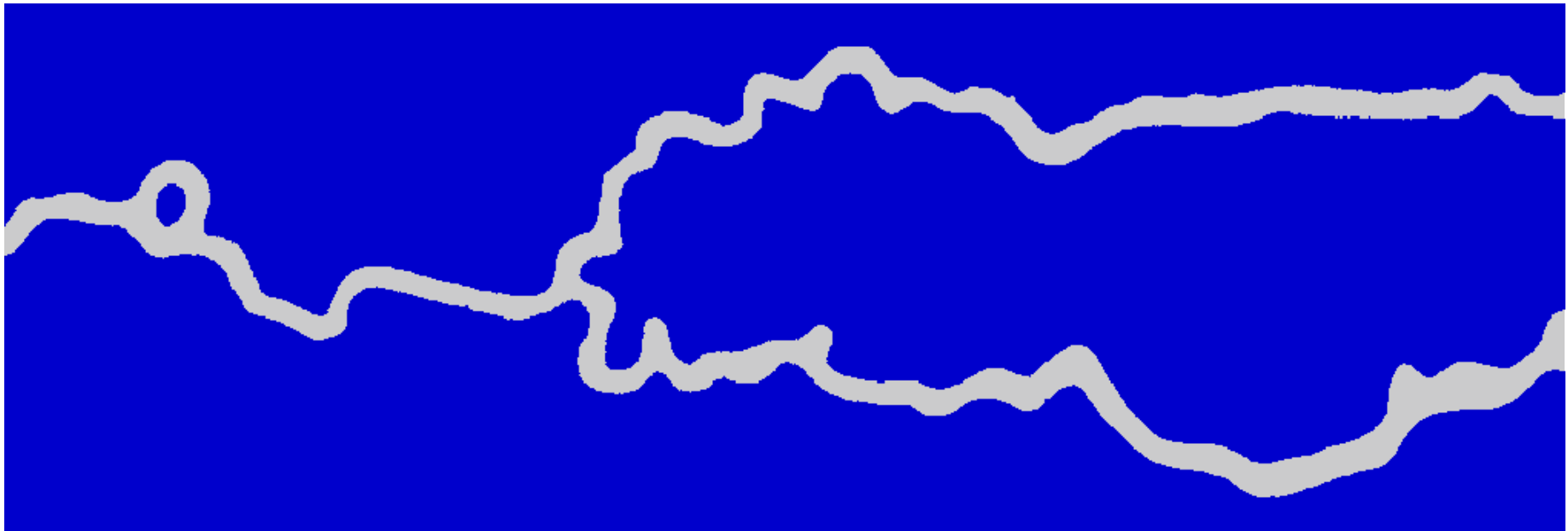
Ohau river, New Zealand.

Renard (2007) Groundwater, 45(5): 531-541

Multiple-points simulation



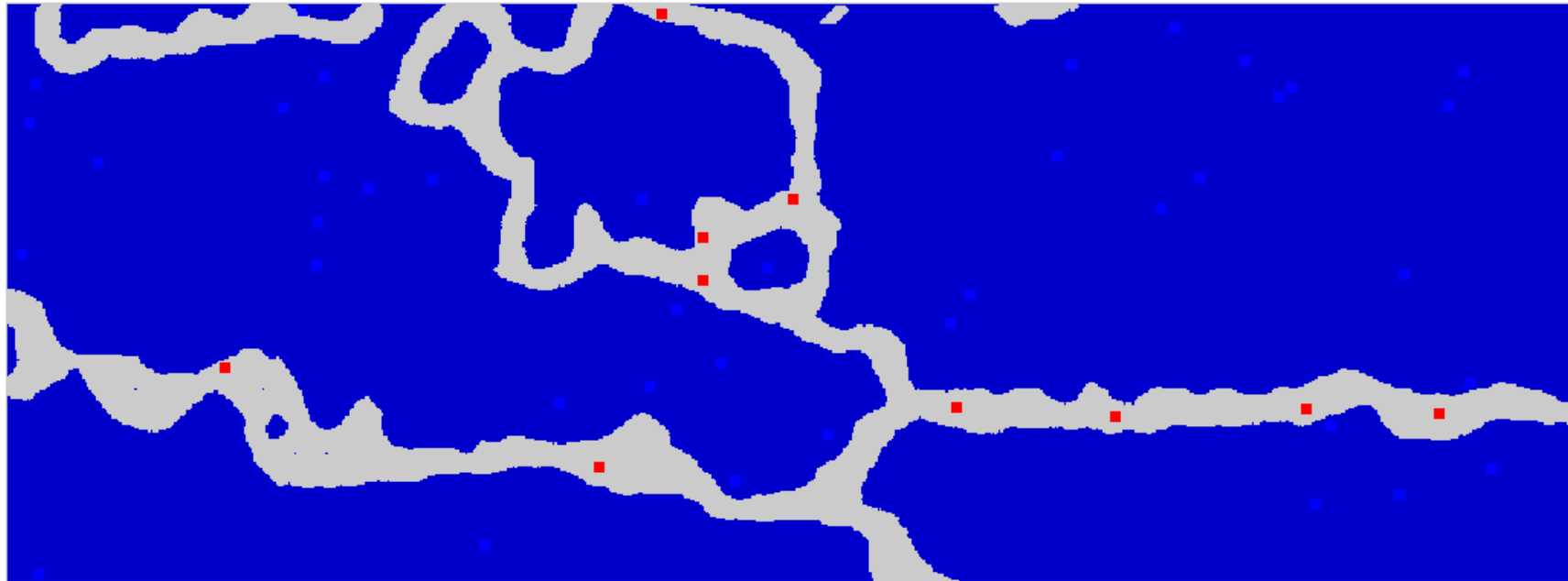
Meanders



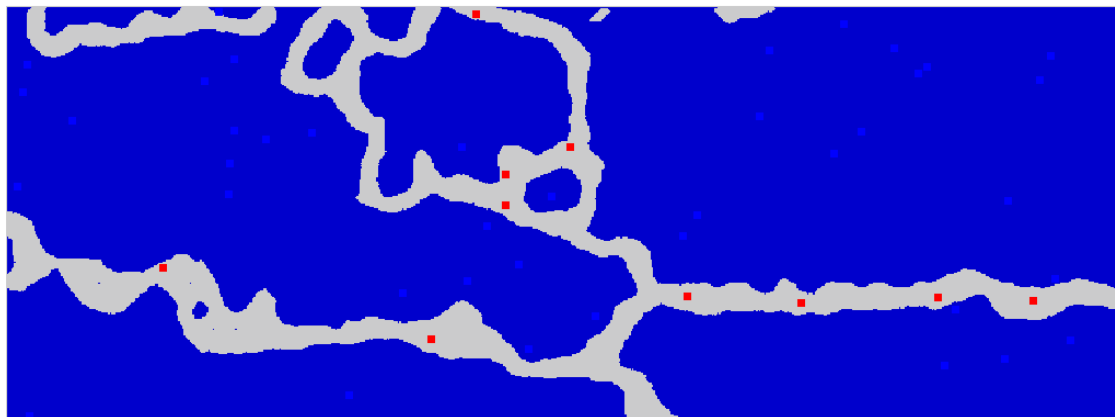
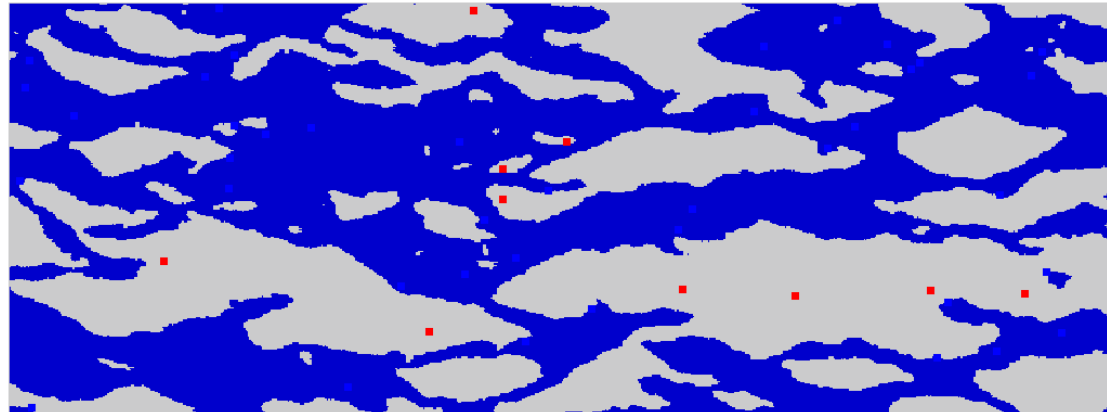
Citronelle oil field, Alabama

Renard (2007) Groundwater, 45(5): 531-541

Multiple-points simulation



Importance of the conceptual model



Renard (2007) Groundwater, 45(5): 531-541

3 innovations

- Field data are not sufficient:

Training Image (TI)

- Two point statistics are not sufficient:

Multiple-point statistics (MPS)

- Analytical statistical model not tractable:

non parametric approach

Principle of the method

Domain to model

Data: geological observations

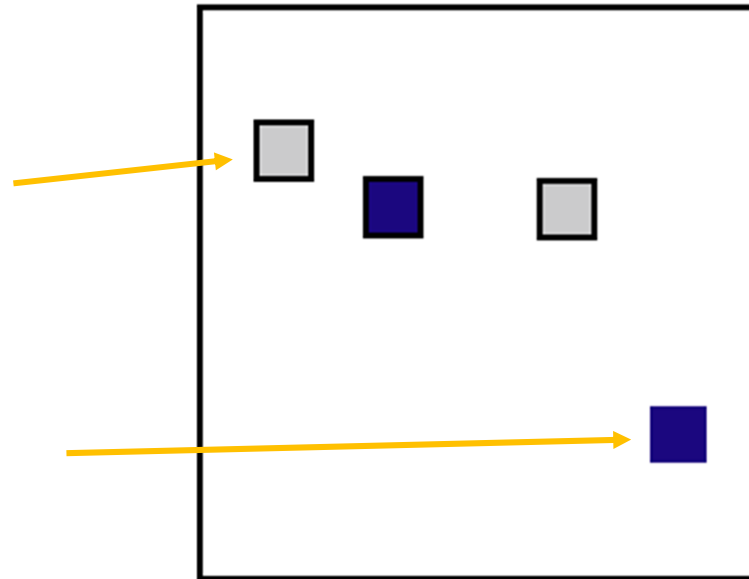
Indicator variable $I(x)$

$$I(x) = 1$$

Sand

$$I(x) = 0$$

Clay



Principle of the method

Sequential simulation method

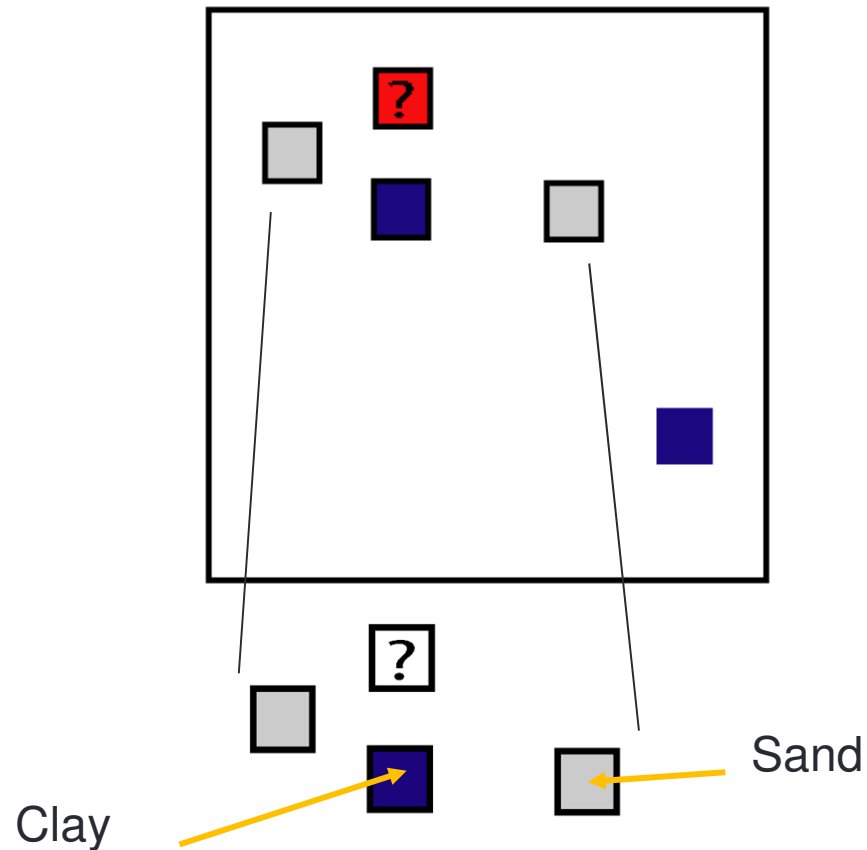
Choice of a location x

We want to find $I(x)$

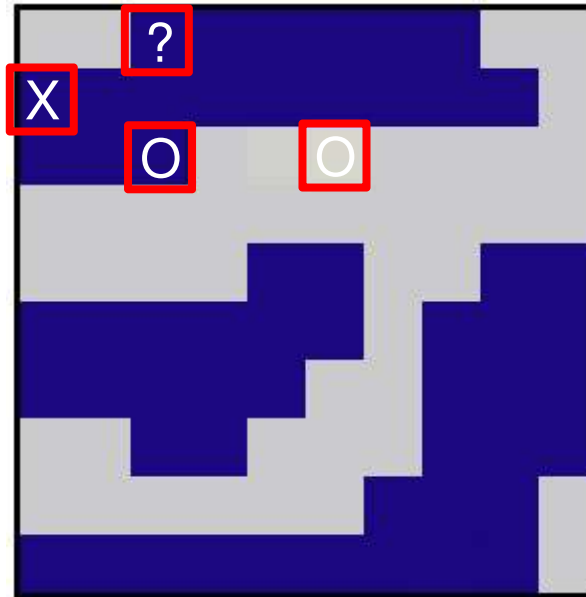
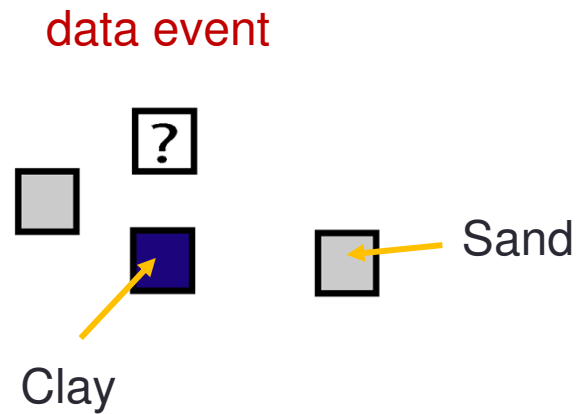
Find the n closest nodes

Find the corresponding values of I

data event

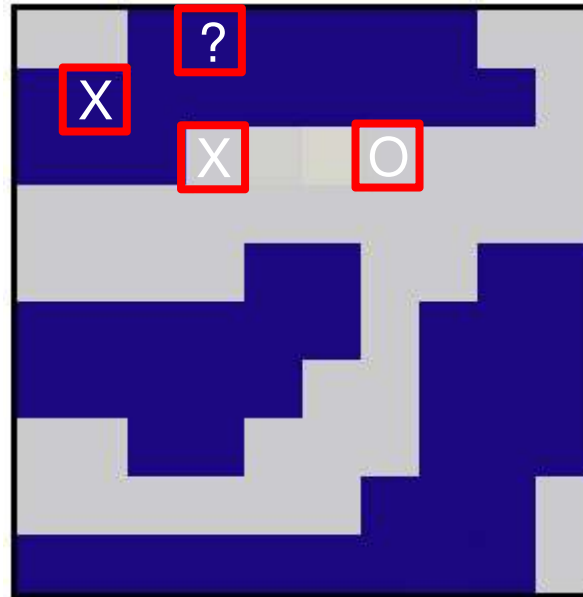
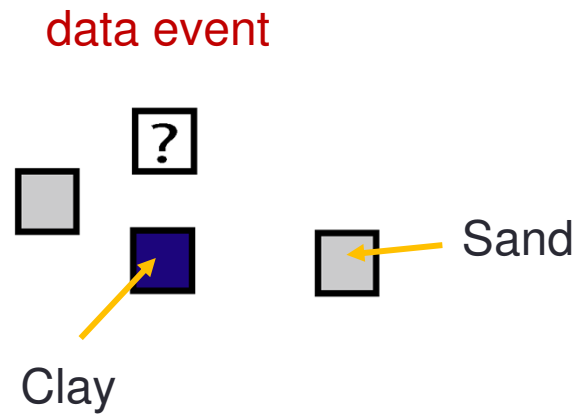


Analysis of the Training image



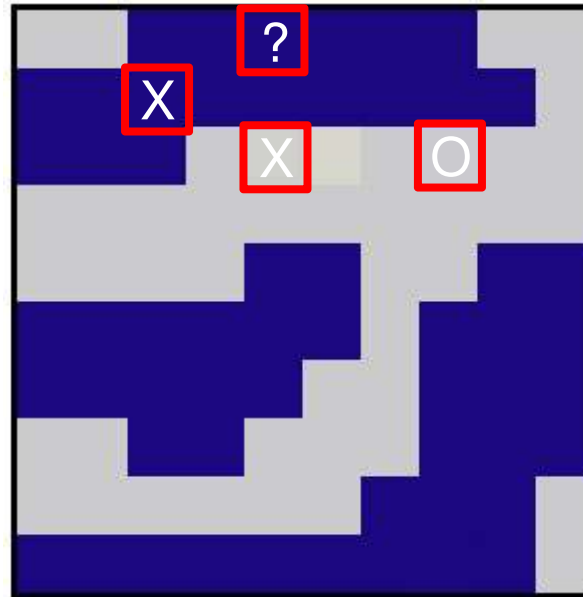
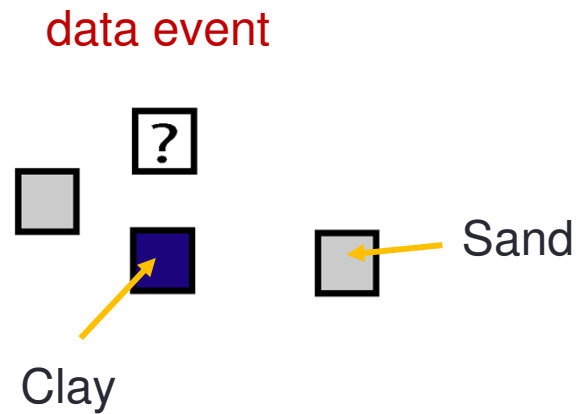
Event		
Counter	0	0

Analysis of the Training image



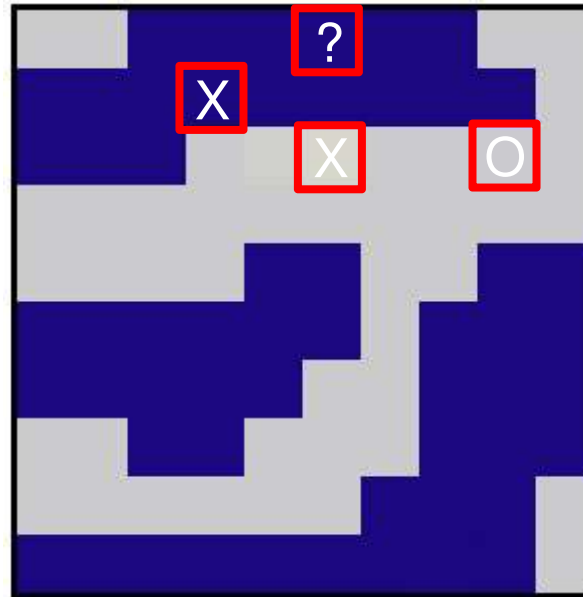
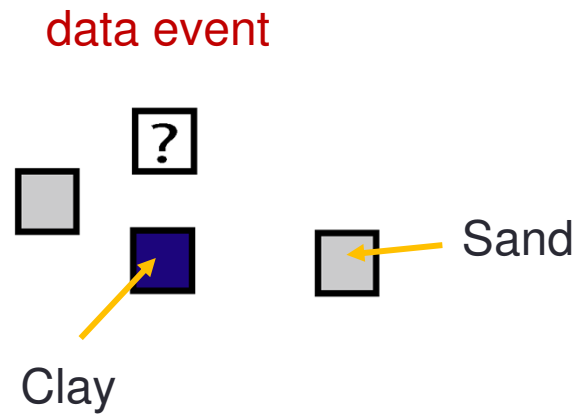
Event		
Counter	0	0

Analysis of the Training image



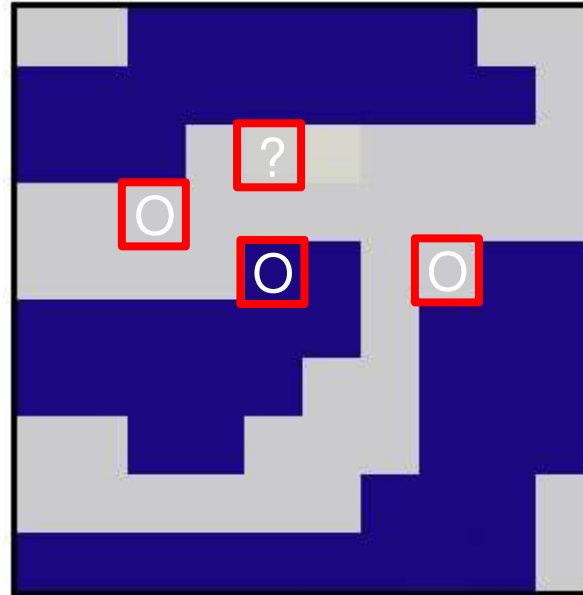
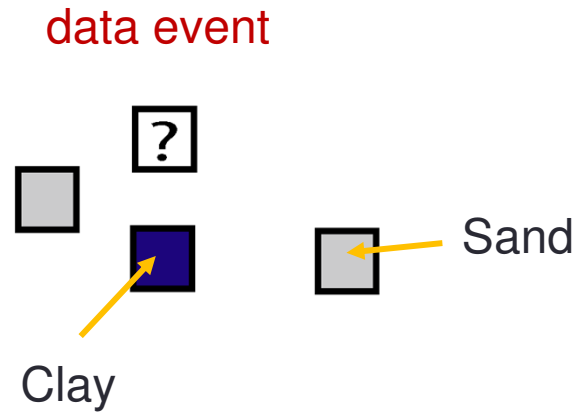
Event		
Counter	0	0

Analysis of the Training image



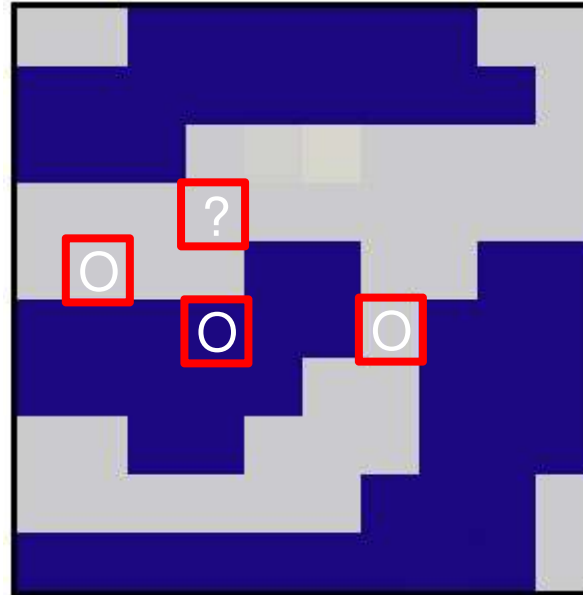
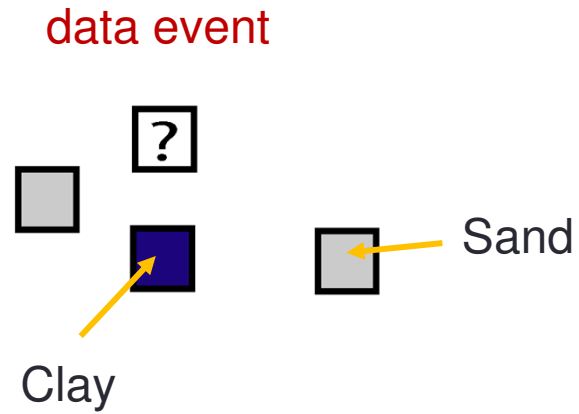
Event		
Counter	0	0

Analysis of the Training image



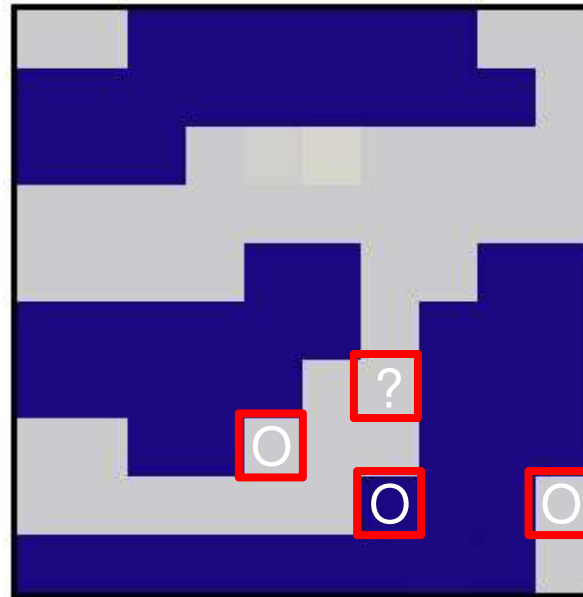
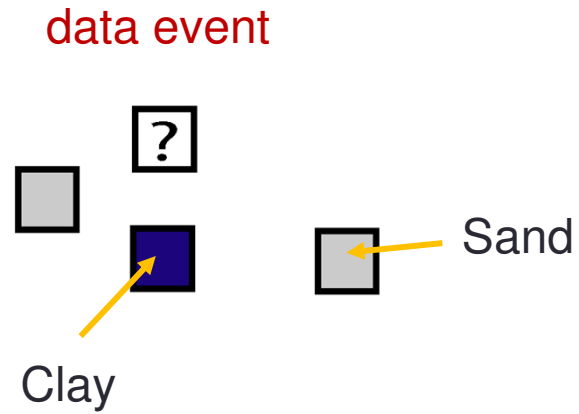
Event		
Counter	0	1

Analysis of the Training image



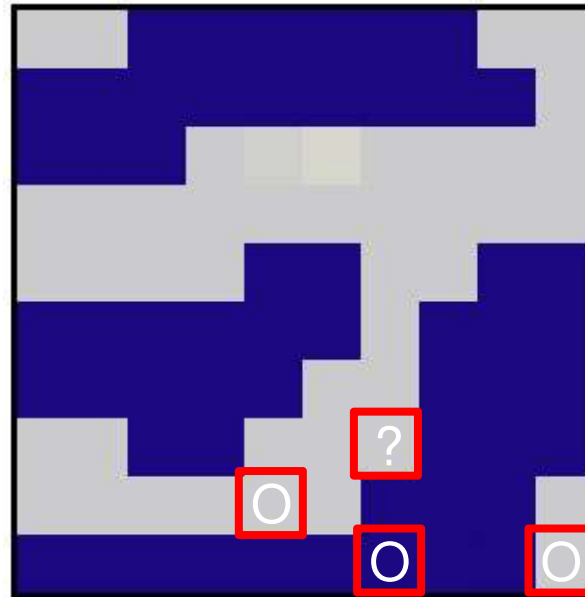
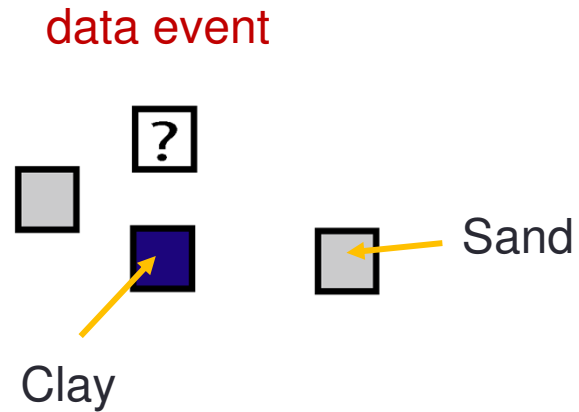
Event		
Counter	0	2

Analysis of the Training image



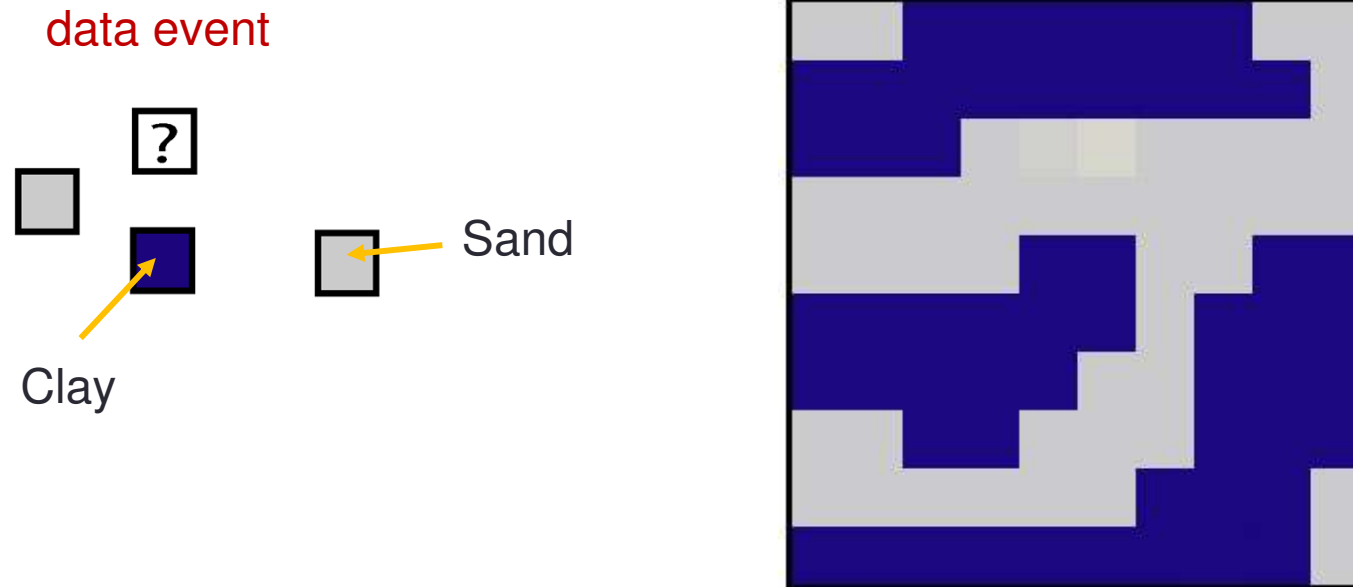
Event		
Counter	0	3

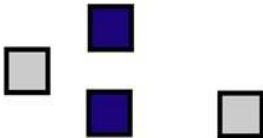
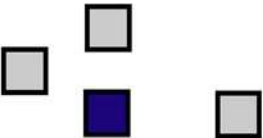
Analysis of the Training image



Event		
Counter	0	4

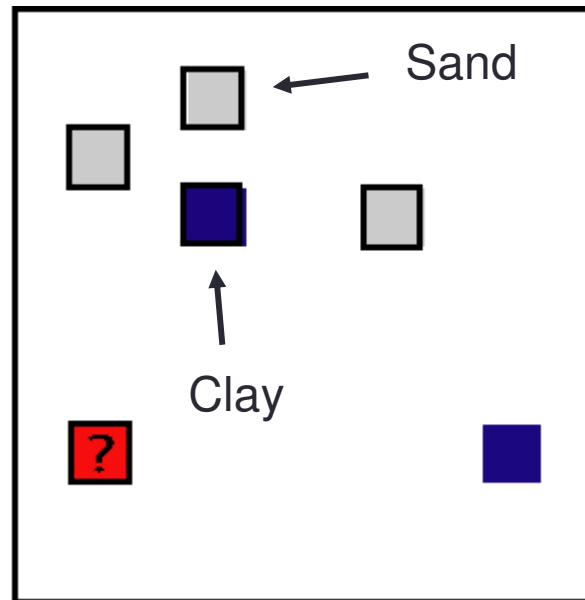
Analysis of the Training image



Event			Total
Counter	0	4	4
Probability	$0 / 4 = 0$	$4 / 4 = 1$	

Sampling the cpdf

$I(x)$ is drawn from the conditional distribution



Another point is randomly selected, simulated, and so on until the whole domain is filled

Technical difficulties

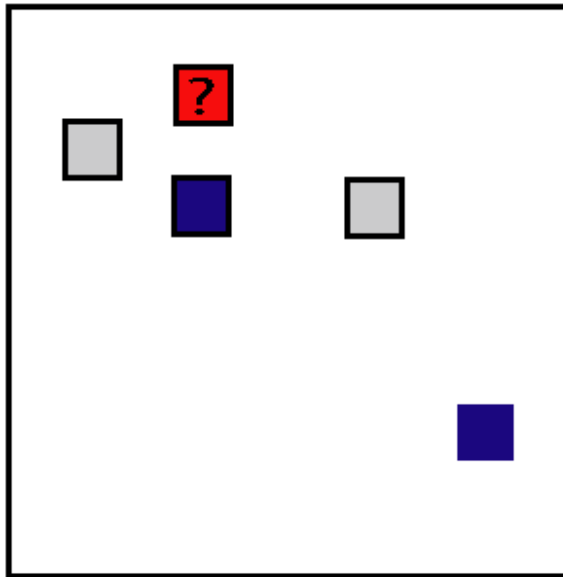
- Scanning the TI for every pixel is inefficient
- Solutions
 - Analyzing the TI and storing the events within a predefined neighborhood (limited dimension)
snesim / impala
Implies additional algorithmic tricks (multigrids, subgrids, data migration, etc)
 - Directly sample the training image

Direct sampling (Deesse)

- Does not use a catalog of patterns
- Allows to extend the technique to continuous and multiple variables
- Allows to get rid of the fixed template size and multigrids
- Our main tool today

Direct sampling

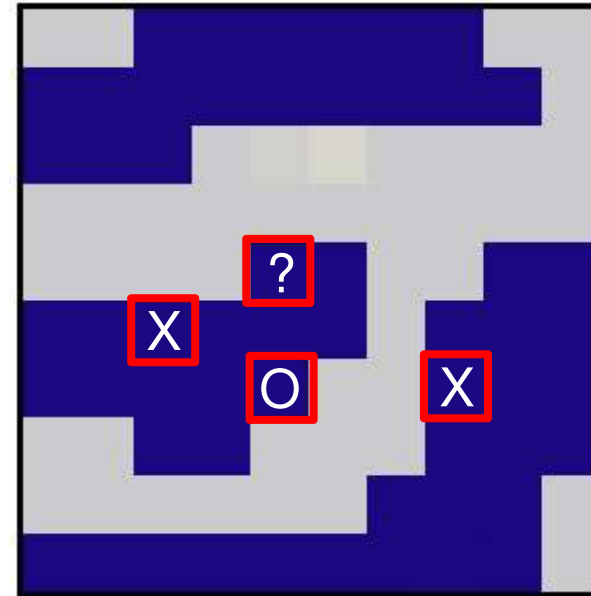
Simulation grid



Conditioning data event

$$\mathbf{d}_n(\mathbf{x}) = \{Z(\mathbf{x} + \mathbf{h}_1), \dots, Z(\mathbf{x} + \mathbf{h}_n)\}$$

Training image



Candidate data event

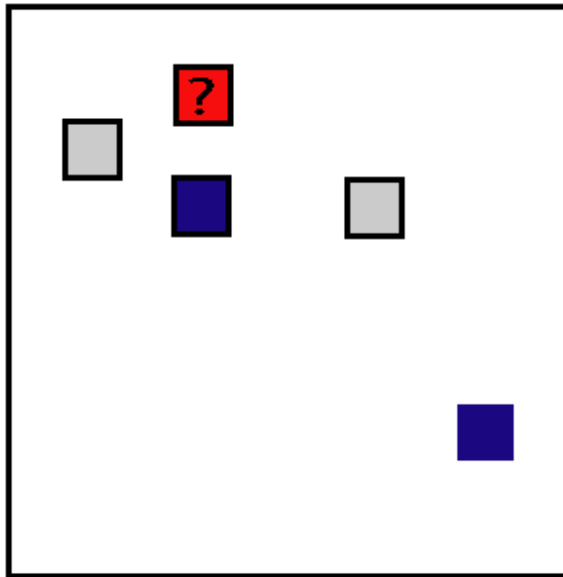
$$\mathbf{d}_n(\mathbf{y}) = \{Z(\mathbf{y} + \mathbf{h}_1), \dots, Z(\mathbf{y} + \mathbf{h}_n)\}$$

$$d\{\mathbf{d}_n(\mathbf{x}), \mathbf{d}_n(\mathbf{y})\} = \frac{1}{n} \sum_{i=1}^n a_i$$

$$a_i = \begin{cases} 0 & \text{if } Z(\mathbf{x} + \mathbf{h}_i) = Z(\mathbf{y} + \mathbf{h}_i) \\ 1 & \text{if } Z(\mathbf{x} + \mathbf{h}_i) \neq Z(\mathbf{y} + \mathbf{h}_i) \end{cases}$$

Direct sampling

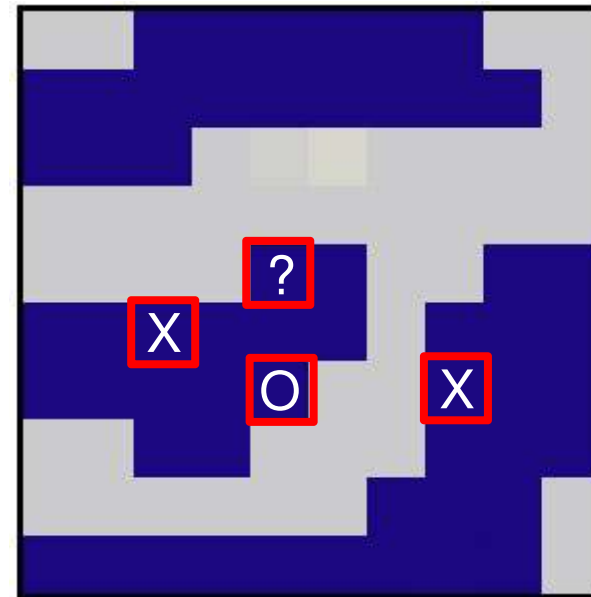
Simulation grid



Conditioning data event

$$\mathbf{d}_n(\mathbf{x}) = \{Z(\mathbf{x} + \mathbf{h}_1), \dots, Z(\mathbf{x} + \mathbf{h}_n)\}$$

Training image



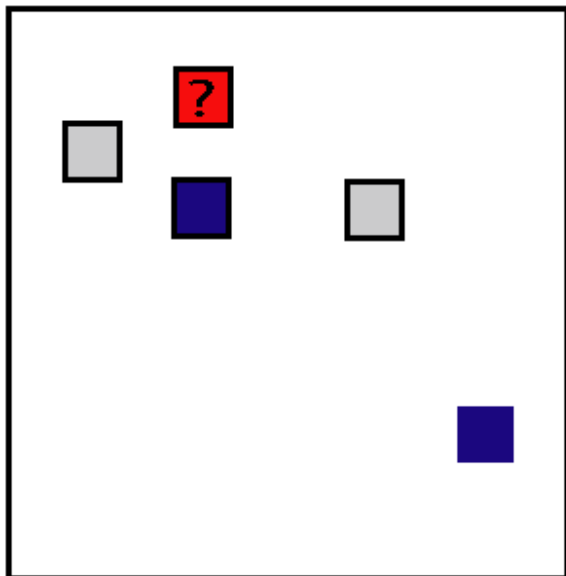
Candidate data event

$$\mathbf{d}_n(\mathbf{y}) = \{Z(\mathbf{y} + \mathbf{h}_1), \dots, Z(\mathbf{y} + \mathbf{h}_n)\}$$

$$d\{\mathbf{d}_n(\mathbf{x}), \mathbf{d}_n(\mathbf{y})\} = \frac{2}{3}$$

Direct sampling

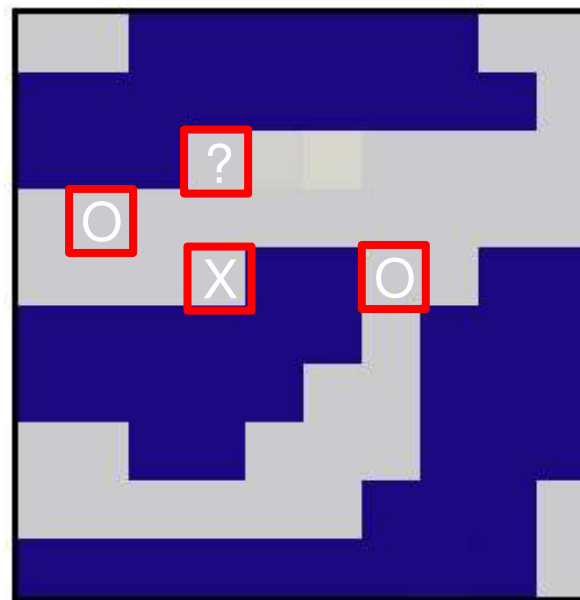
Simulation grid



Conditioning data event

$$\mathbf{d}_n(\mathbf{x}) = \{Z(\mathbf{x} + \mathbf{h}_1), \dots, Z(\mathbf{x} + \mathbf{h}_n)\}$$

Training image



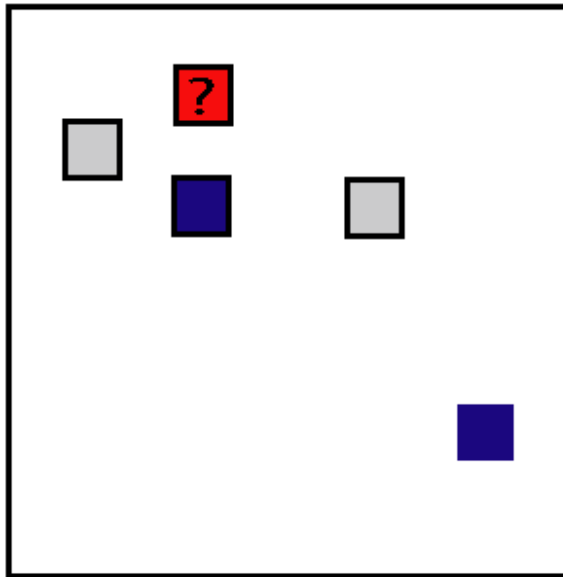
Candidate data event

$$\mathbf{d}_n(\mathbf{y}) = \{Z(\mathbf{y} + \mathbf{h}_1), \dots, Z(\mathbf{y} + \mathbf{h}_n)\}$$

$$d\{\mathbf{d}_n(\mathbf{x}), \mathbf{d}_n(\mathbf{y})\} = \frac{1}{3}$$

Direct sampling

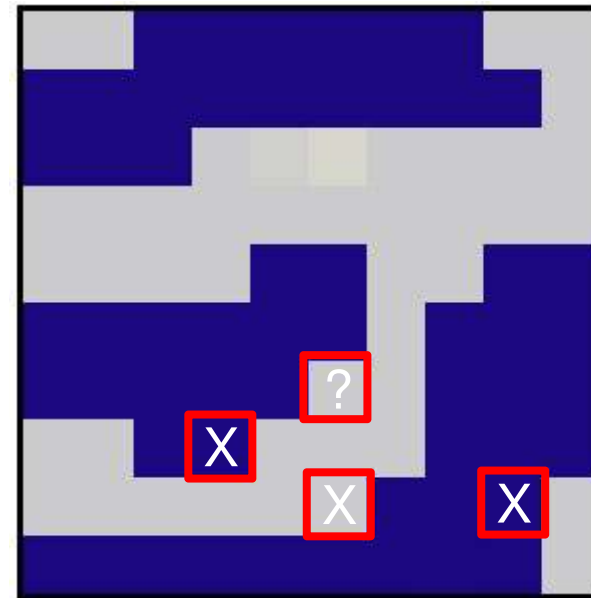
Simulation grid



Conditioning data event

$$\mathbf{d}_n(\mathbf{x}) = \{Z(\mathbf{x} + \mathbf{h}_1), \dots, Z(\mathbf{x} + \mathbf{h}_n)\}$$

Training image



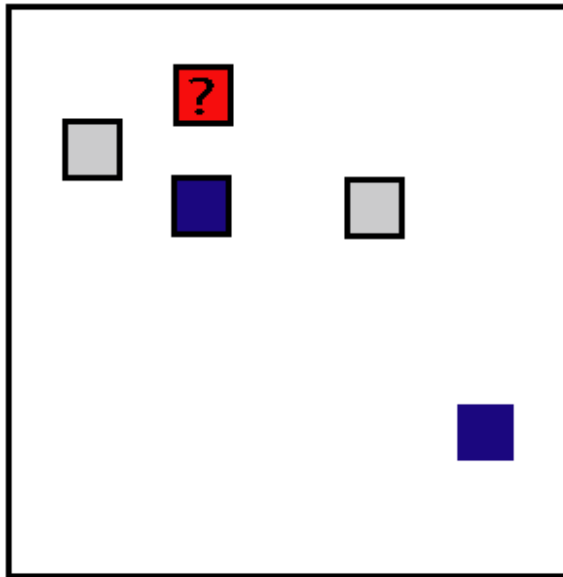
Candidate data event

$$\mathbf{d}_n(\mathbf{y}) = \{Z(\mathbf{y} + \mathbf{h}_1), \dots, Z(\mathbf{y} + \mathbf{h}_n)\}$$

$$d\{\mathbf{d}_n(\mathbf{x}), \mathbf{d}_n(\mathbf{y})\} = \frac{3}{3}$$

Direct sampling

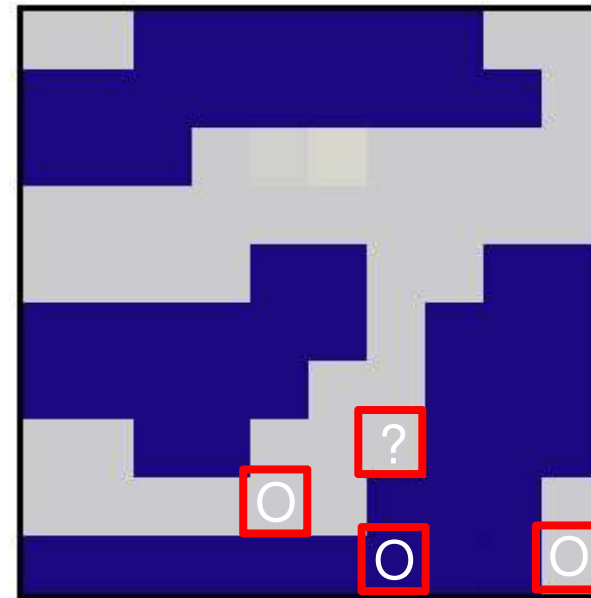
Simulation grid



Conditioning data event

$$\mathbf{d}_n(\mathbf{x}) = \{Z(\mathbf{x} + \mathbf{h}_1), \dots, Z(\mathbf{x} + \mathbf{h}_n)\}$$

Training image



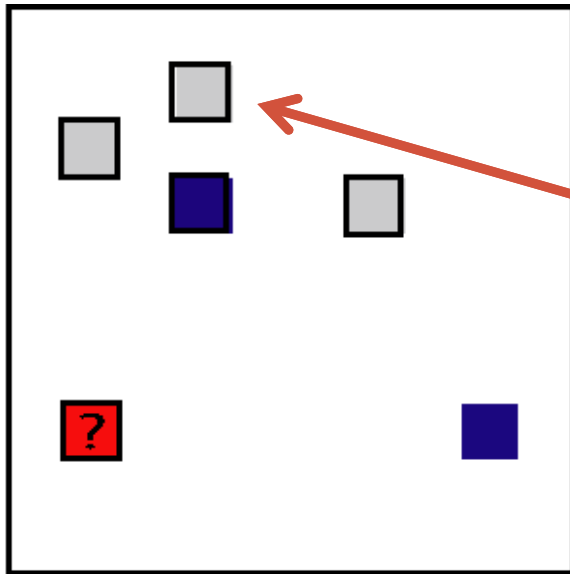
Candidate data event

$$\mathbf{d}_n(\mathbf{y}) = \{Z(\mathbf{y} + \mathbf{h}_1), \dots, Z(\mathbf{y} + \mathbf{h}_n)\}$$

$$d\{\mathbf{d}_n(\mathbf{x}), \mathbf{d}_n(\mathbf{y})\} = \frac{0}{3}$$

Direct sampling

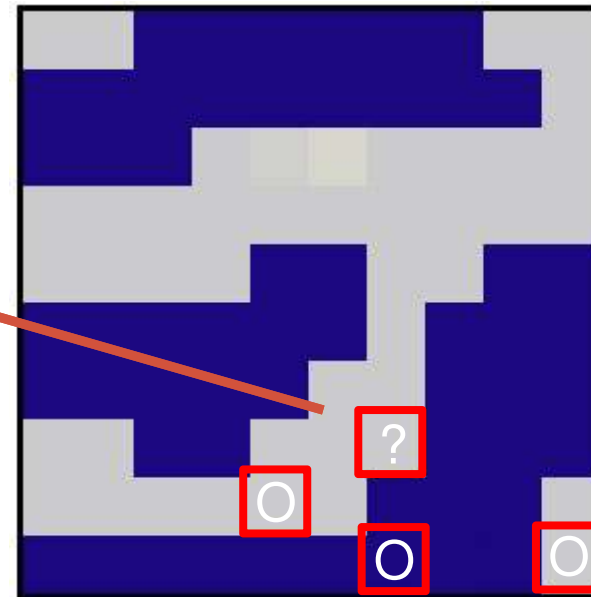
Simulation grid



Conditioning data event

$$\mathbf{d}_n(\mathbf{x}) = \{Z(\mathbf{x} + \mathbf{h}_1), \dots, Z(\mathbf{x} + \mathbf{h}_n)\}$$

Training image



Candidate data event

$$\mathbf{d}_n(\mathbf{y}) = \{Z(\mathbf{y} + \mathbf{h}_1), \dots, Z(\mathbf{y} + \mathbf{h}_n)\}$$

$$d\{\mathbf{d}_n(\mathbf{x}), \mathbf{d}_n(\mathbf{y})\} = \frac{0}{3}$$

Basic Direct Sampling algorithm

- Distance :

$$d\{\mathbf{d}_n(\mathbf{x}), \mathbf{d}_n(\mathbf{y})\} = \frac{1}{n} \sum_{i=1}^n a_i \quad a_i = \begin{cases} 0 & \text{if } Z(\mathbf{x} + \mathbf{h}_i) = Z(\mathbf{y} + \mathbf{h}_i) \\ 1 & \text{if } Z(\mathbf{x} + \mathbf{h}_i) \neq Z(\mathbf{y} + \mathbf{h}_i) \end{cases}$$

- DS algorithm consists in scanning the Training Image
 - Until we find the **first data event** such that

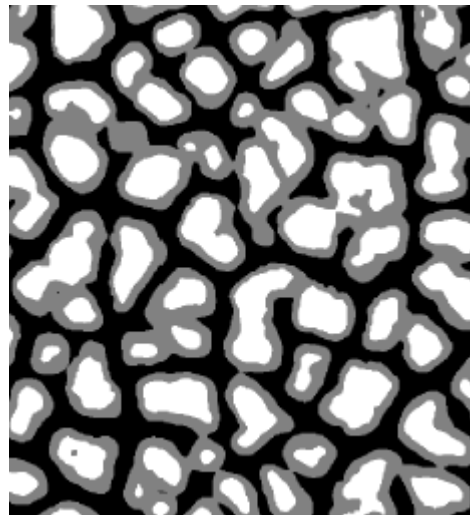
$$d\{\mathbf{d}_n(\mathbf{x}), \mathbf{d}_n(\mathbf{y})\} < t$$

- Or until a certain fraction **f** of the training image has been scanned, then the best data event is selected

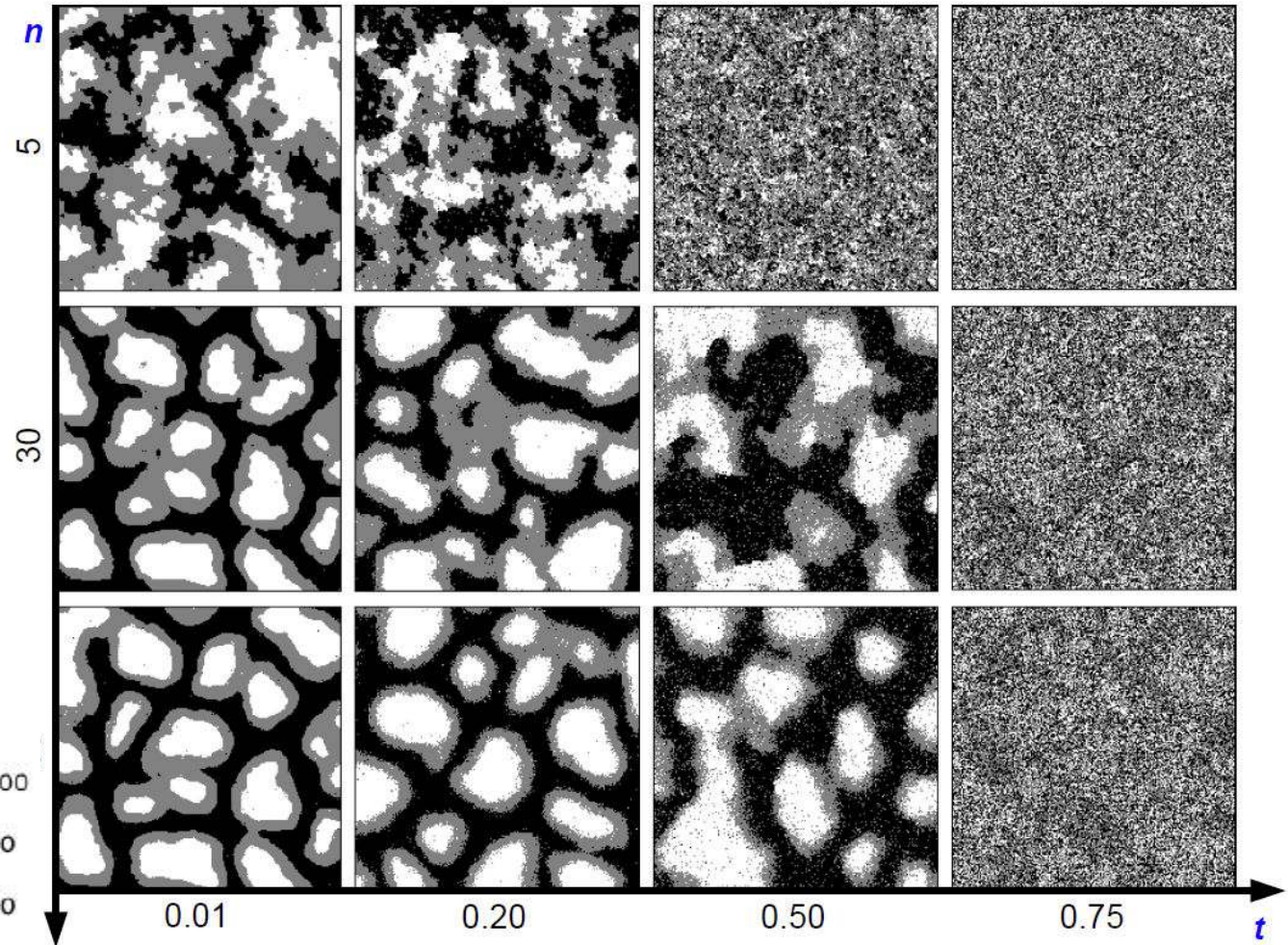
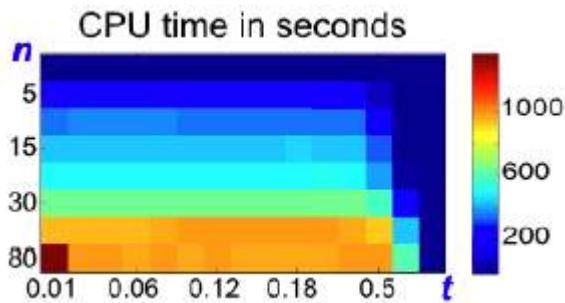
Direct Sampling parameters

- Size of the neighborhood: **n**
 - Depends on space dimension, 10 to 100
- Acceptance threshold for the distance: **t**
 - A value between 0 and 1
- Maximum scan fraction of the Training Image: **f**
 - Between 0.1 and 0.5

Parameter sensitivity



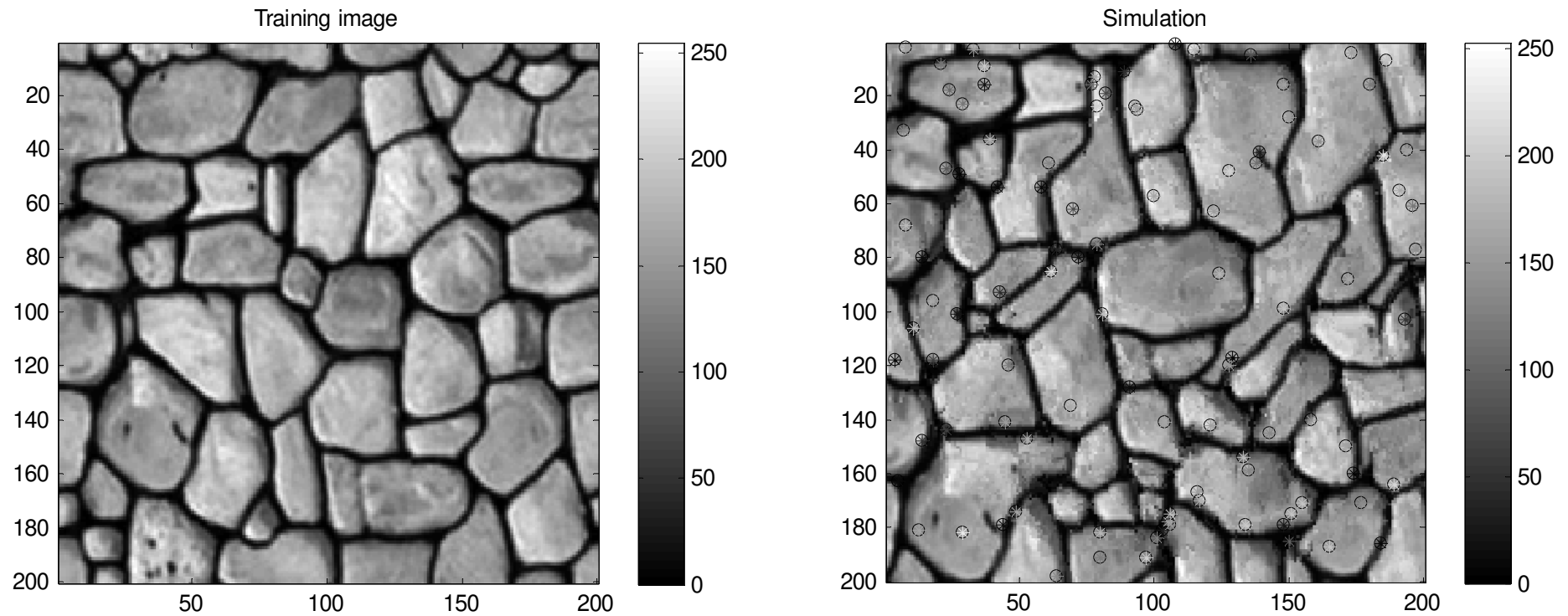
Training image



Simulations with $f=0.3$

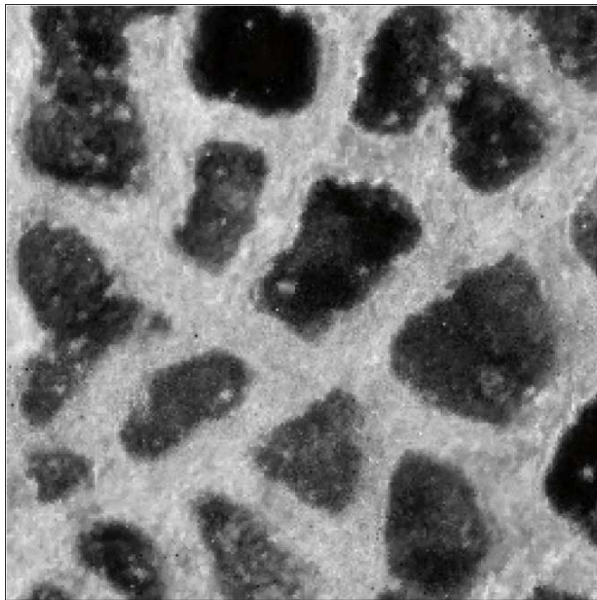
Continuous variable

$$d\{\mathbf{d}_n(\mathbf{x}), \mathbf{d}_n(\mathbf{y})\} \propto \sum_{i=1}^n [Z(\mathbf{x} + \mathbf{h}_i) - Z(\mathbf{y} + \mathbf{h}_i)]^2$$

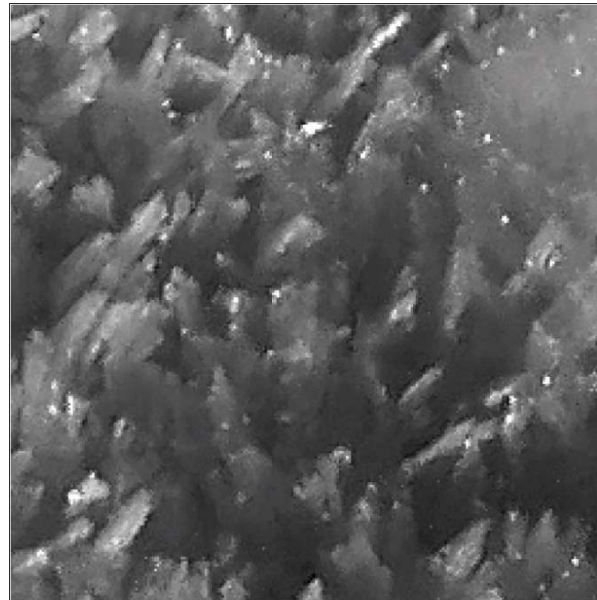


Examples of simulations

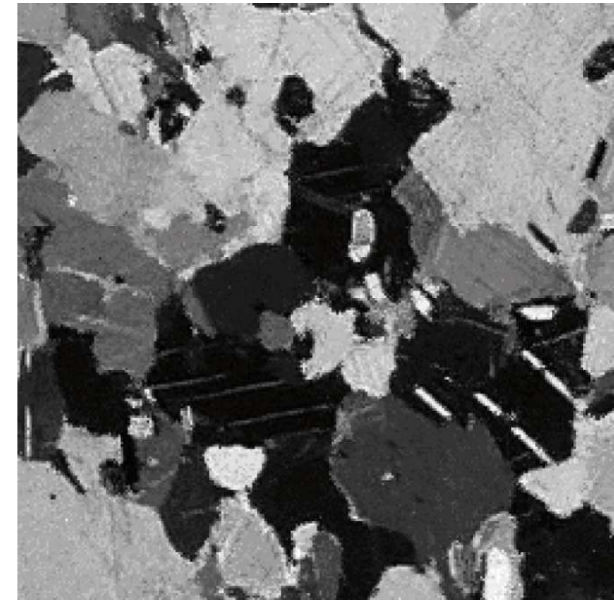
Ice wedge polygonal soil



Snow



Thin slice marble

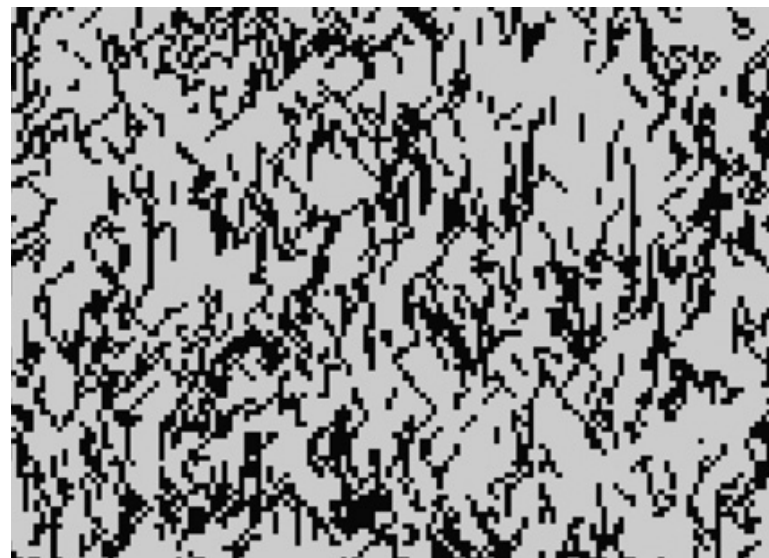


Training image stationarity

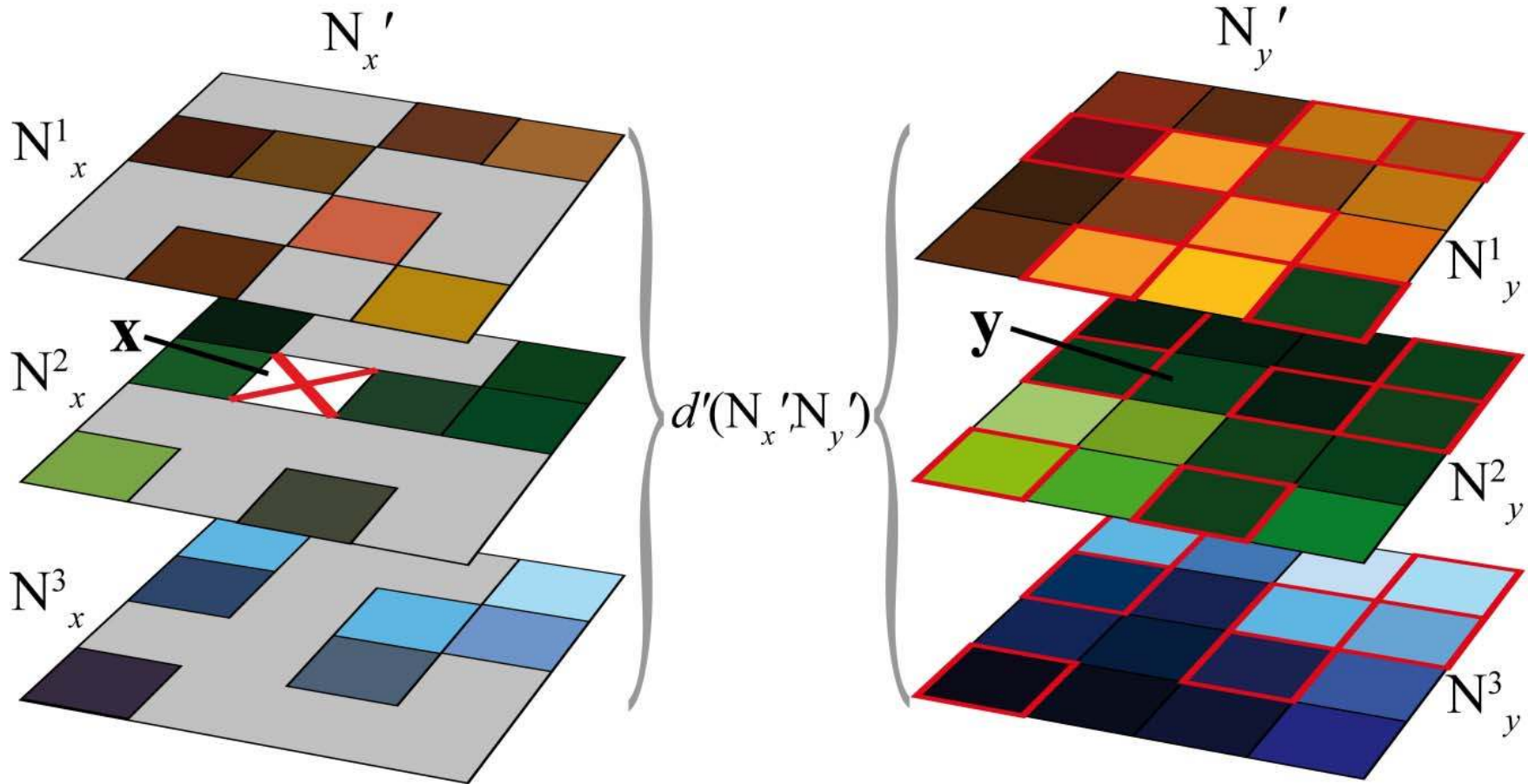
Training image (TI)



Simulation



Multiple variables

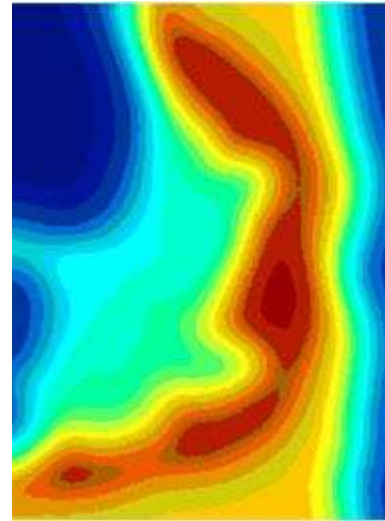


$$d^i \{ \mathbf{d}_n(\mathbf{x}), \mathbf{d}_n(\mathbf{y}) \} < t_i$$

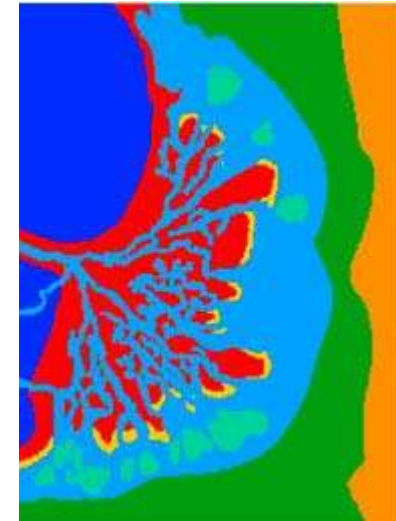
Multivariate simulation

Training data set

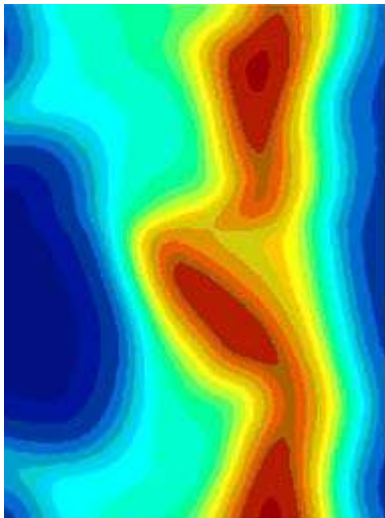
Variable 1



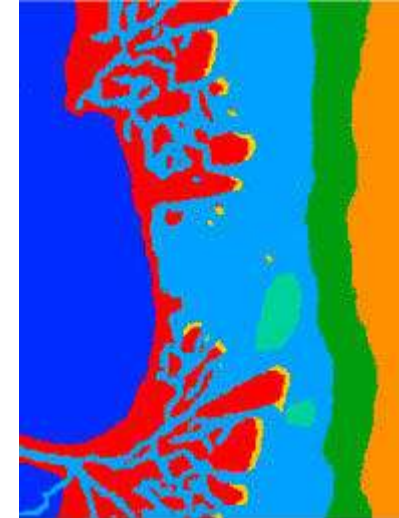
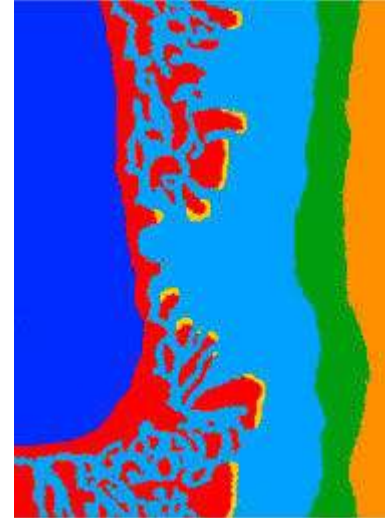
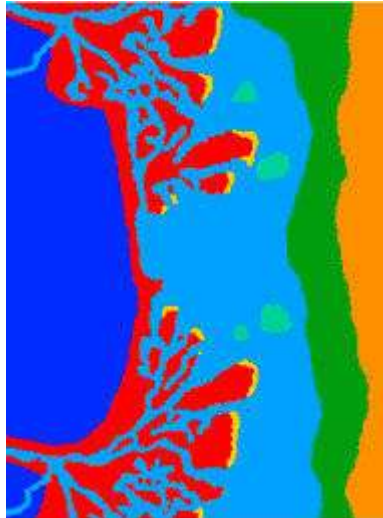
Variable 2



Conditioning variable
Variable 1



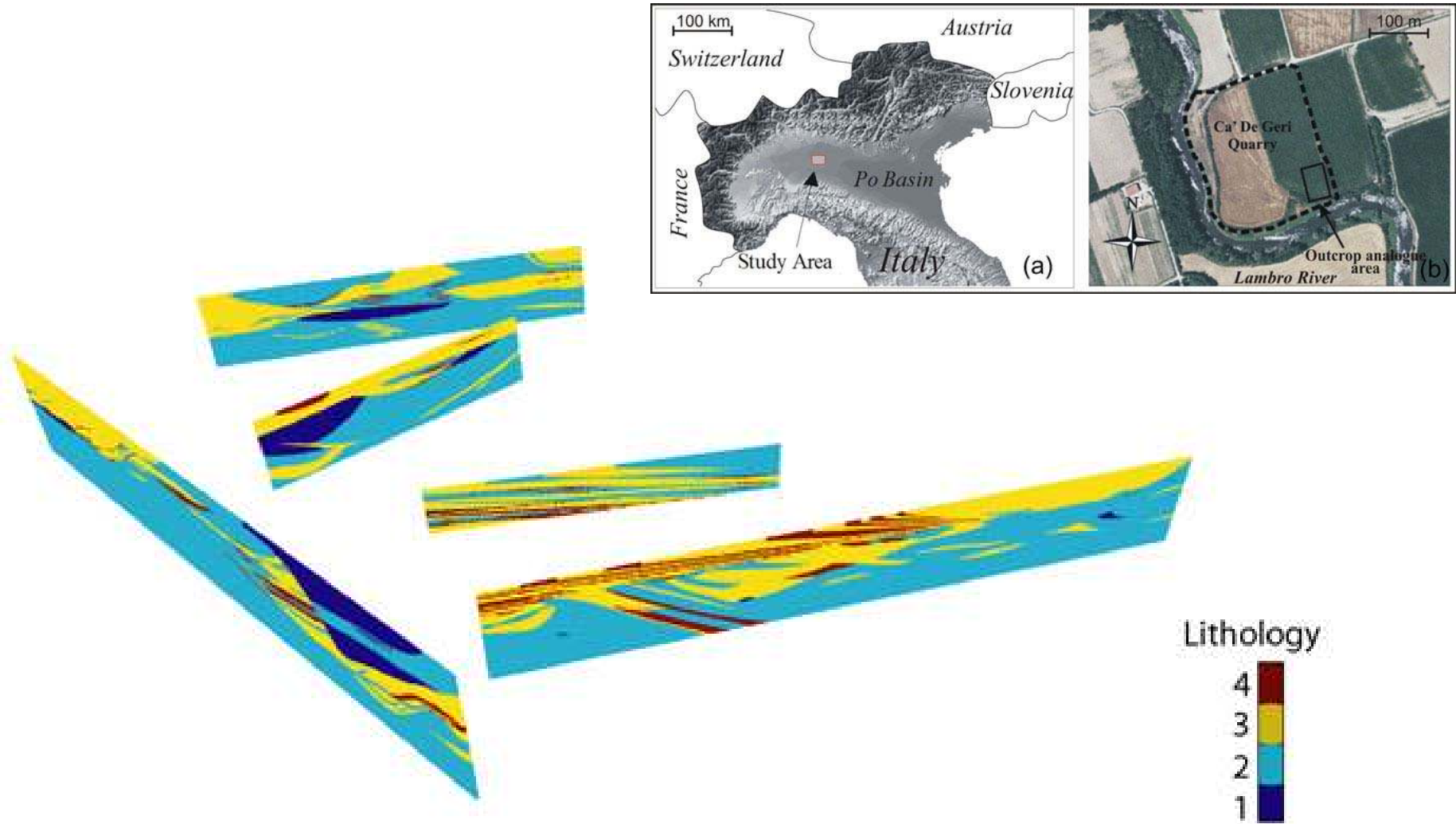
3 Simulations
Variable 2



A FEW EXAMPLES OF APPLICATIONS

3D geology / Rainfall simulation / Reconstruction of missing data / etc.

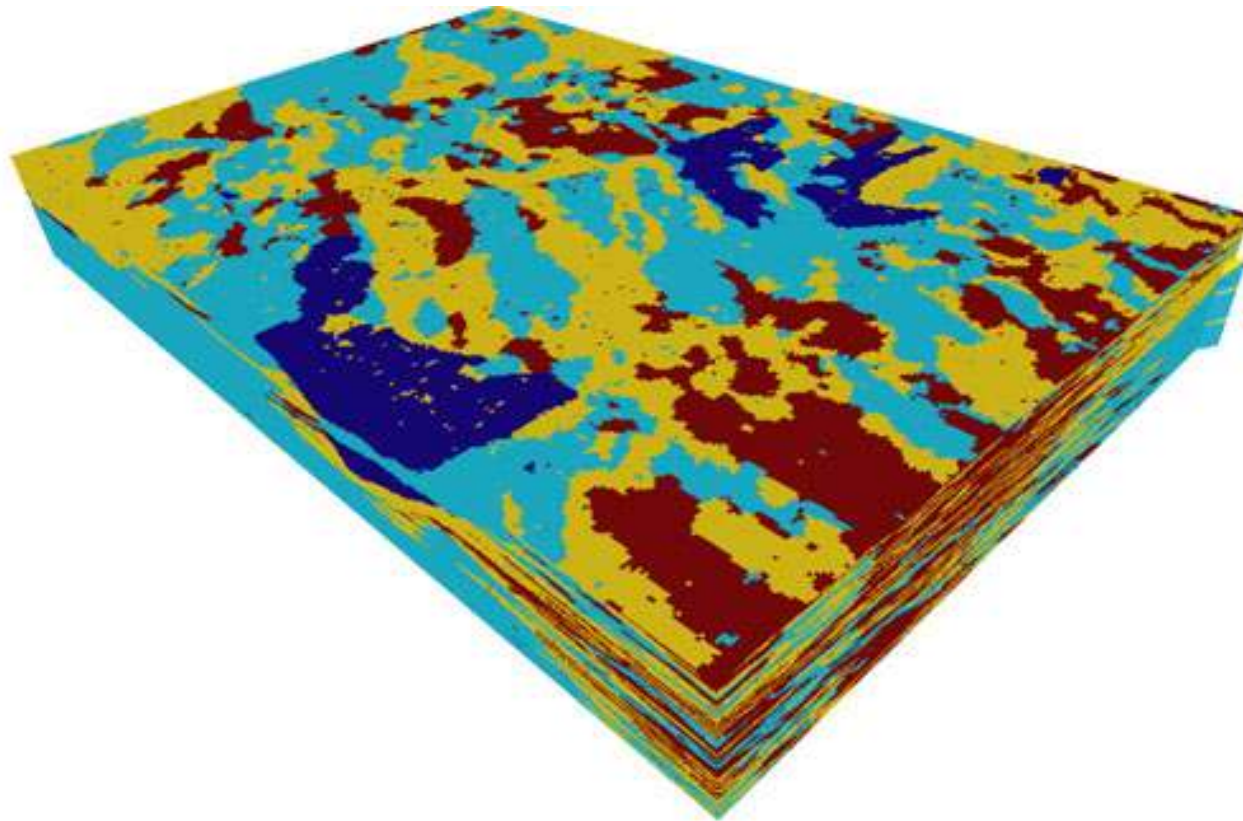
Reconstructions from sections



Data provided by Dell'Arciprete, Felletti, Bersezio

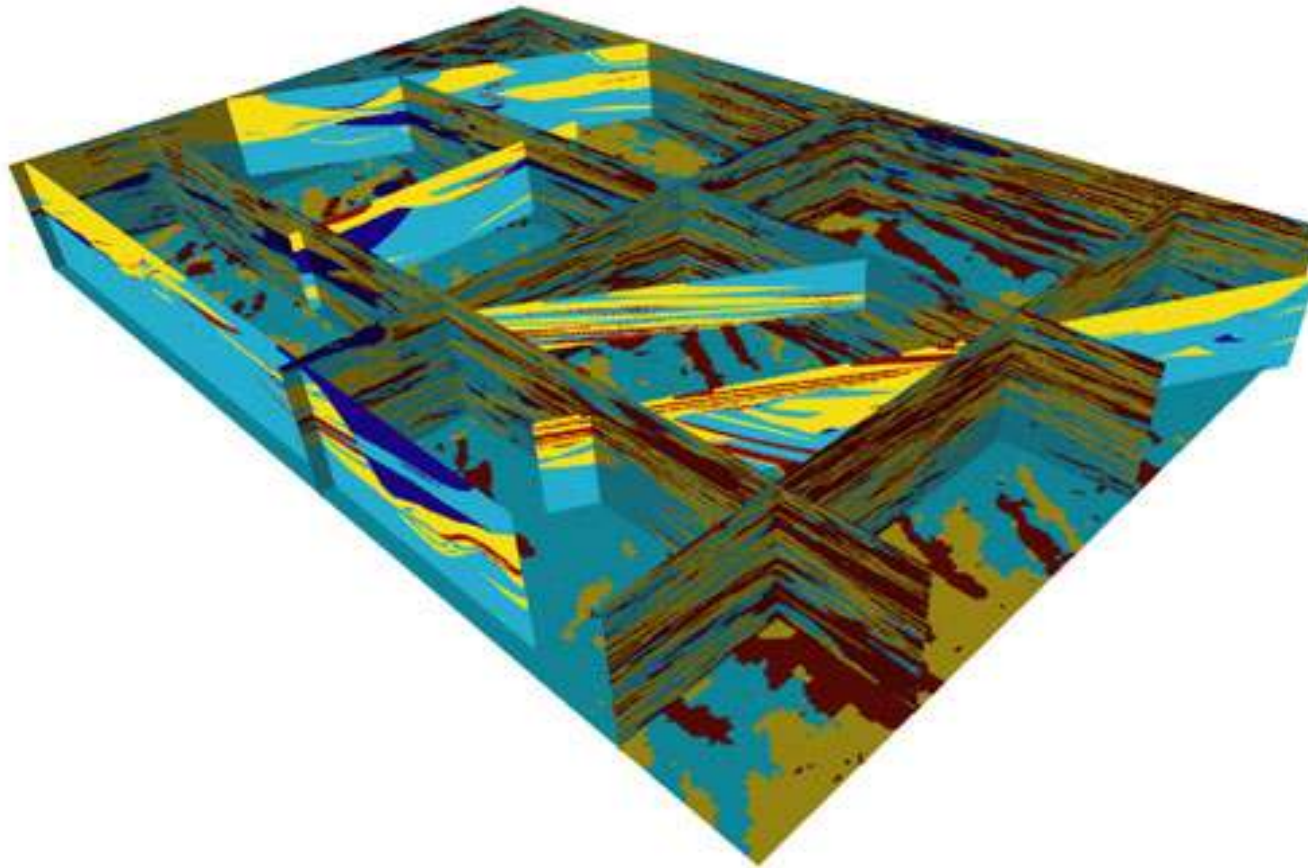
Mariethoz, Renard (2009) *Mathematical Geosciences*. 42(3): 245–268

One simulation

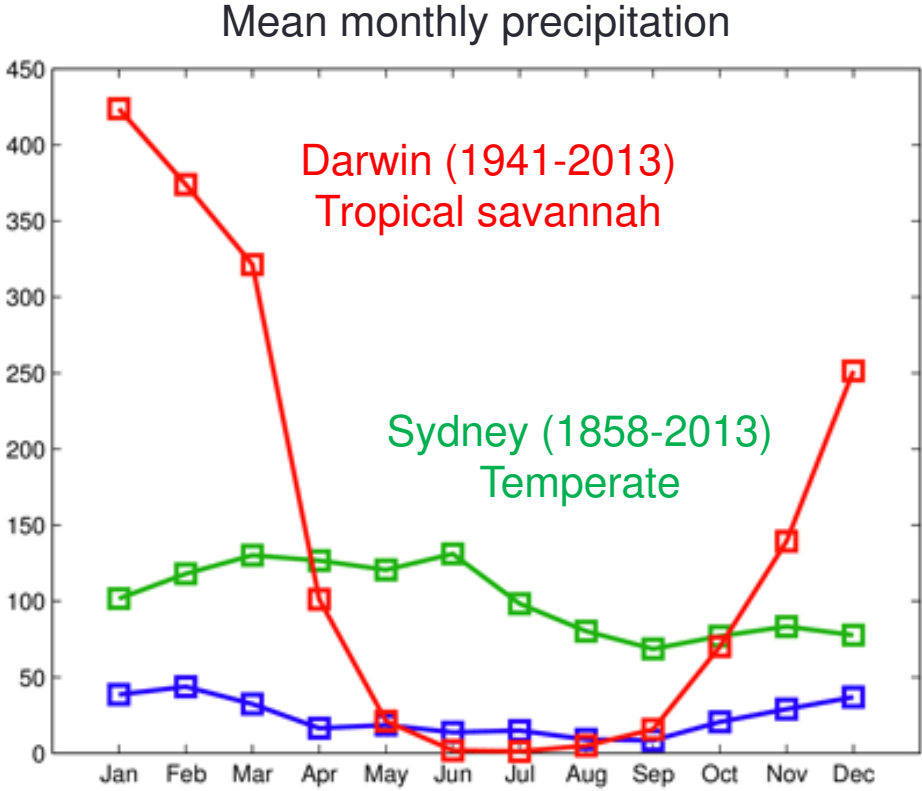
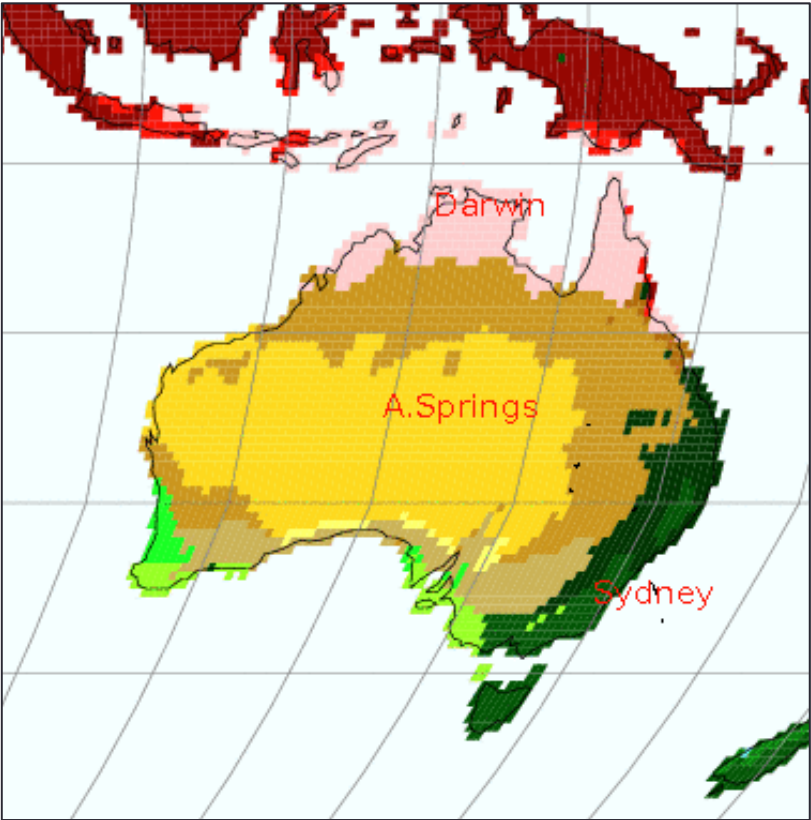


Mariethoz, Renard (2009) *Mathematical Geosciences*. 42(3): 245–268

Sections in the simulation

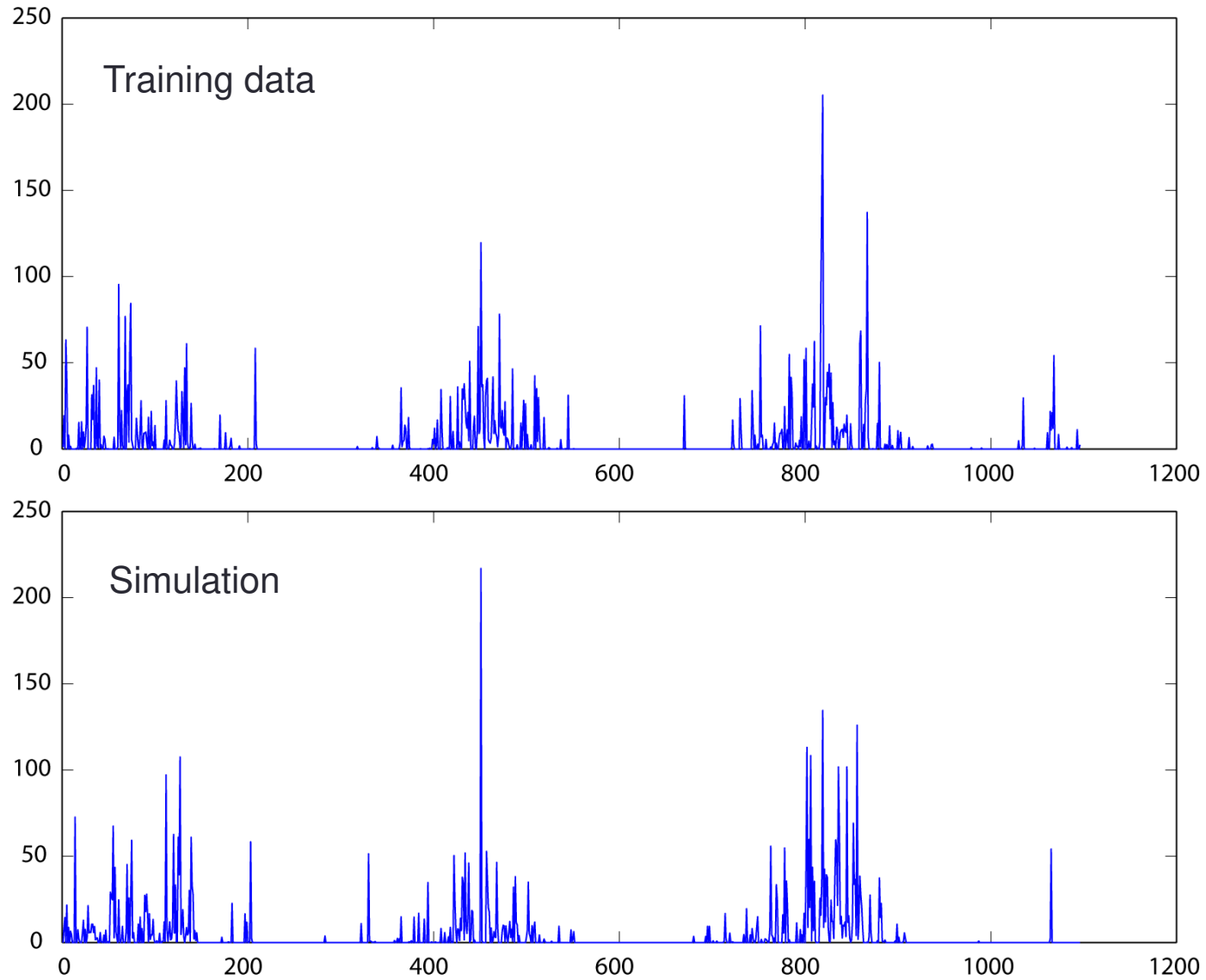


Rainfall simulation



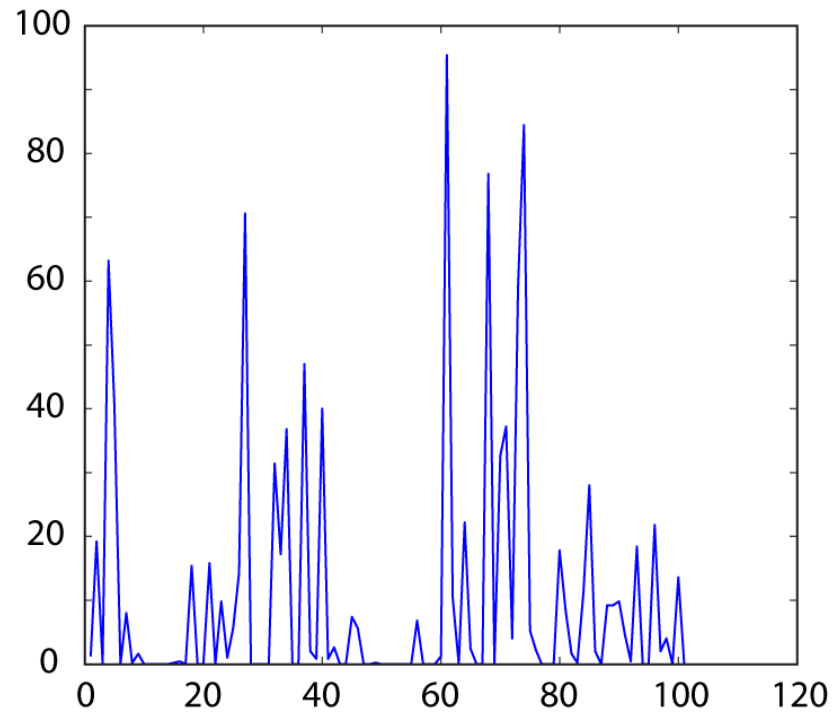
Alice Spring (1941-2013)
Hot desert

Rainfall simulation (Darwin)

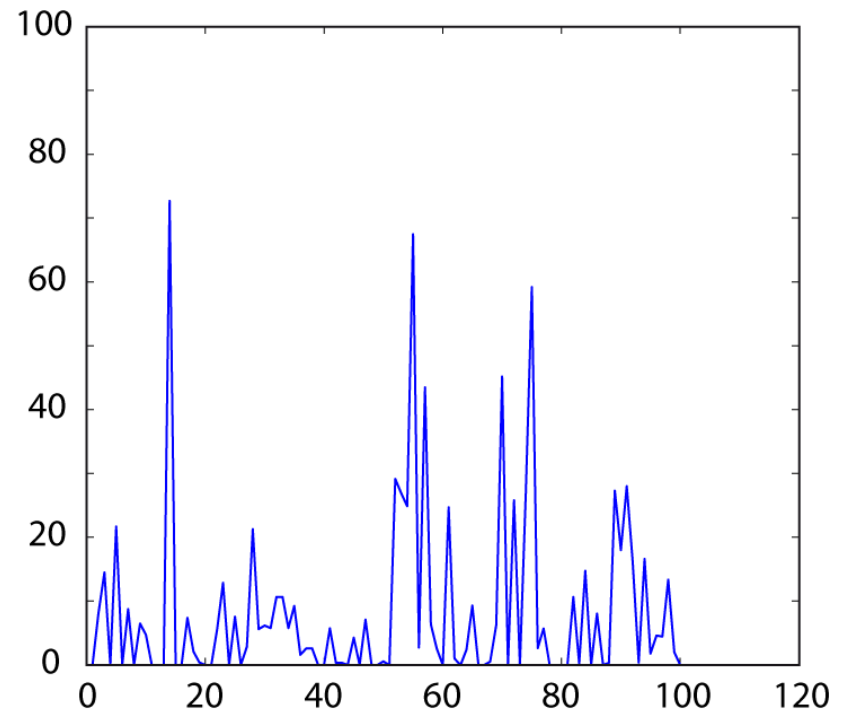


Rainfall simulation (Darwin)

Training data

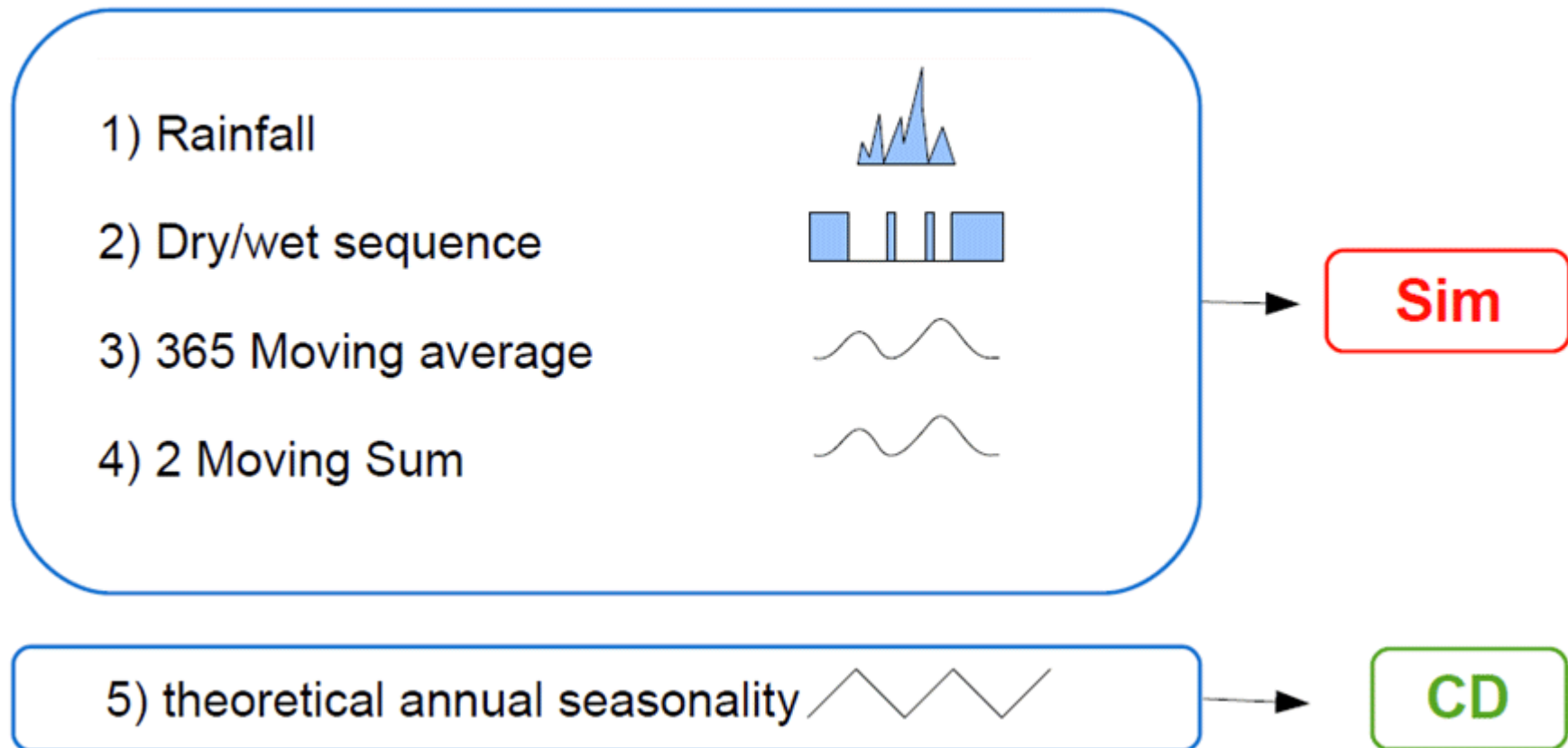


Simulation



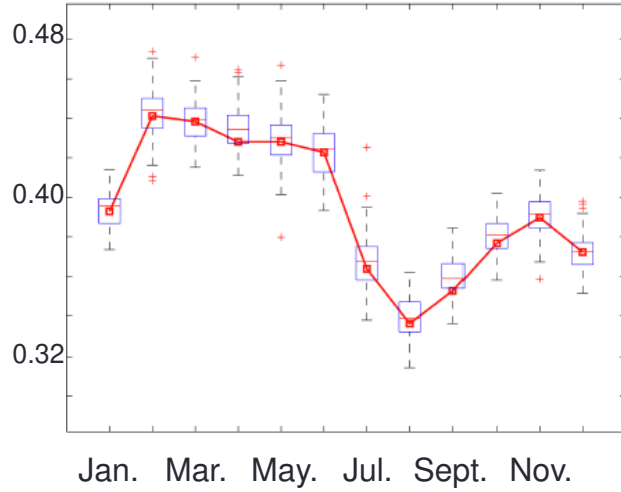
Rainfall simulation procedure

TI

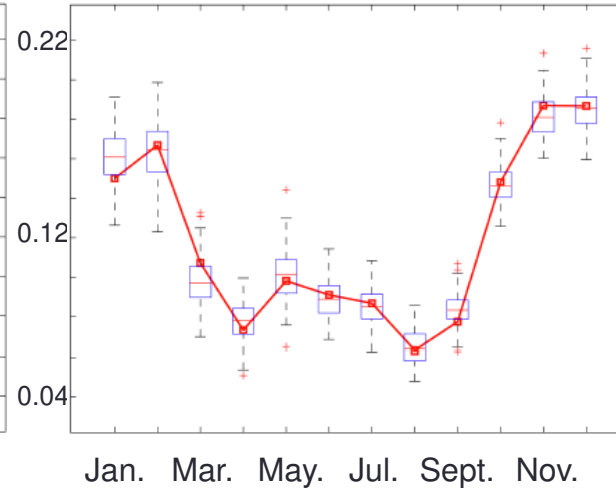


Wet days probability

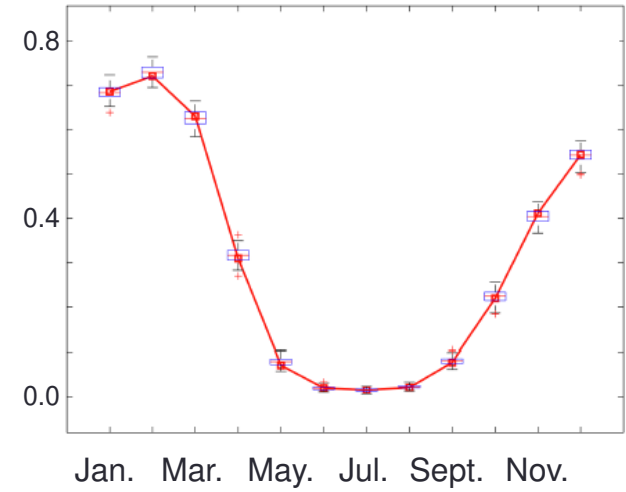
Sydney



Alice Springs

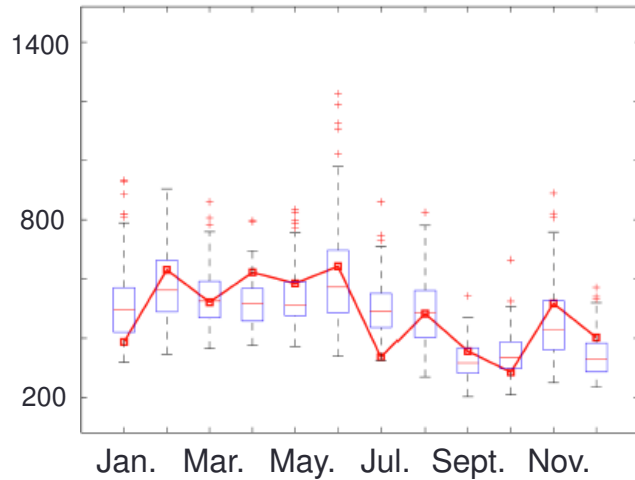


Darwin

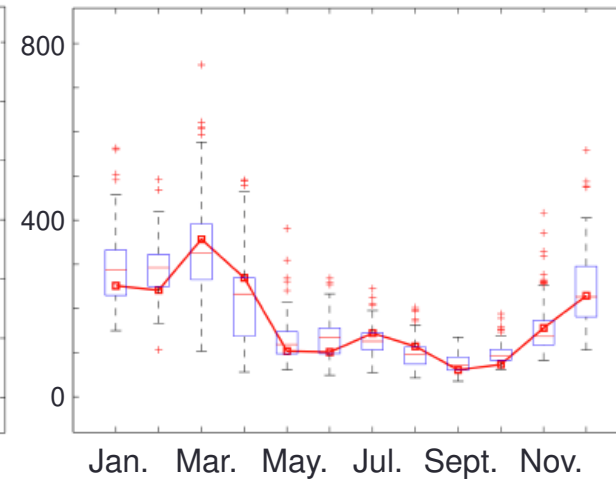


Max cumulated monthly rainfall [mm]

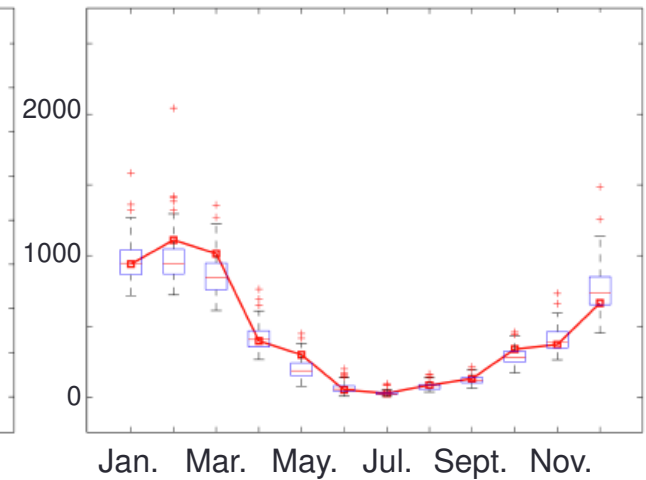
Sydney



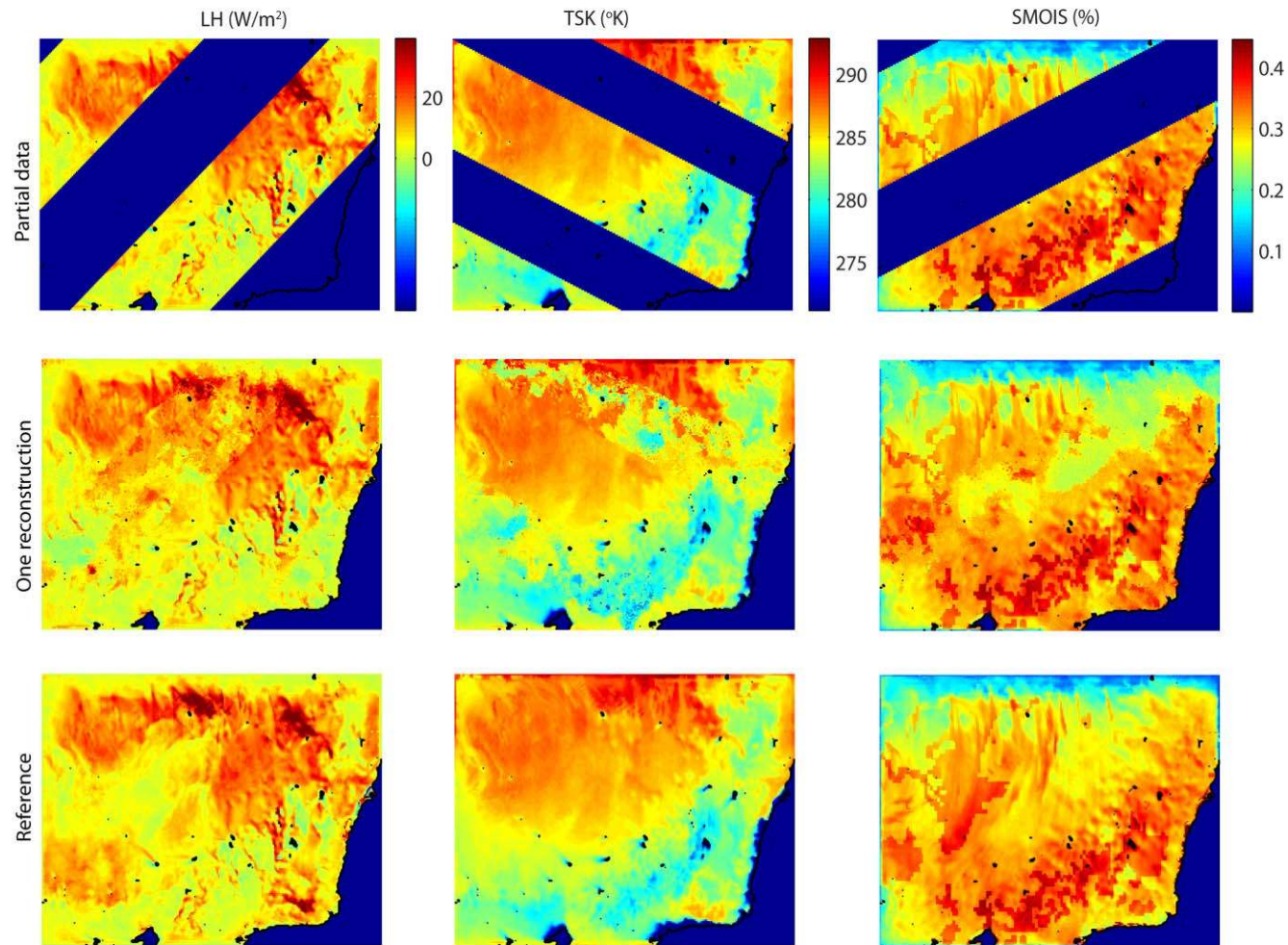
Alice Springs



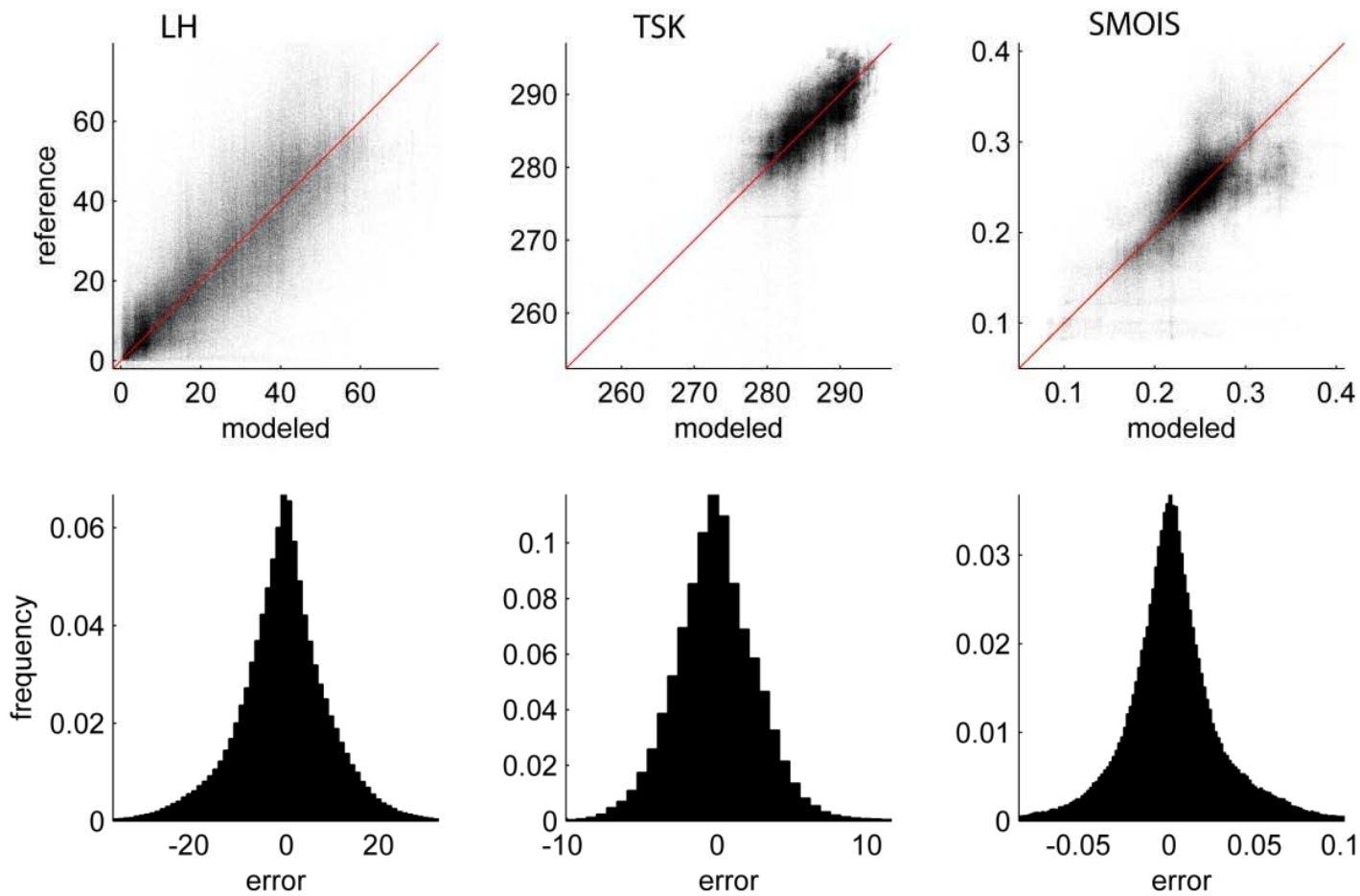
Darwin



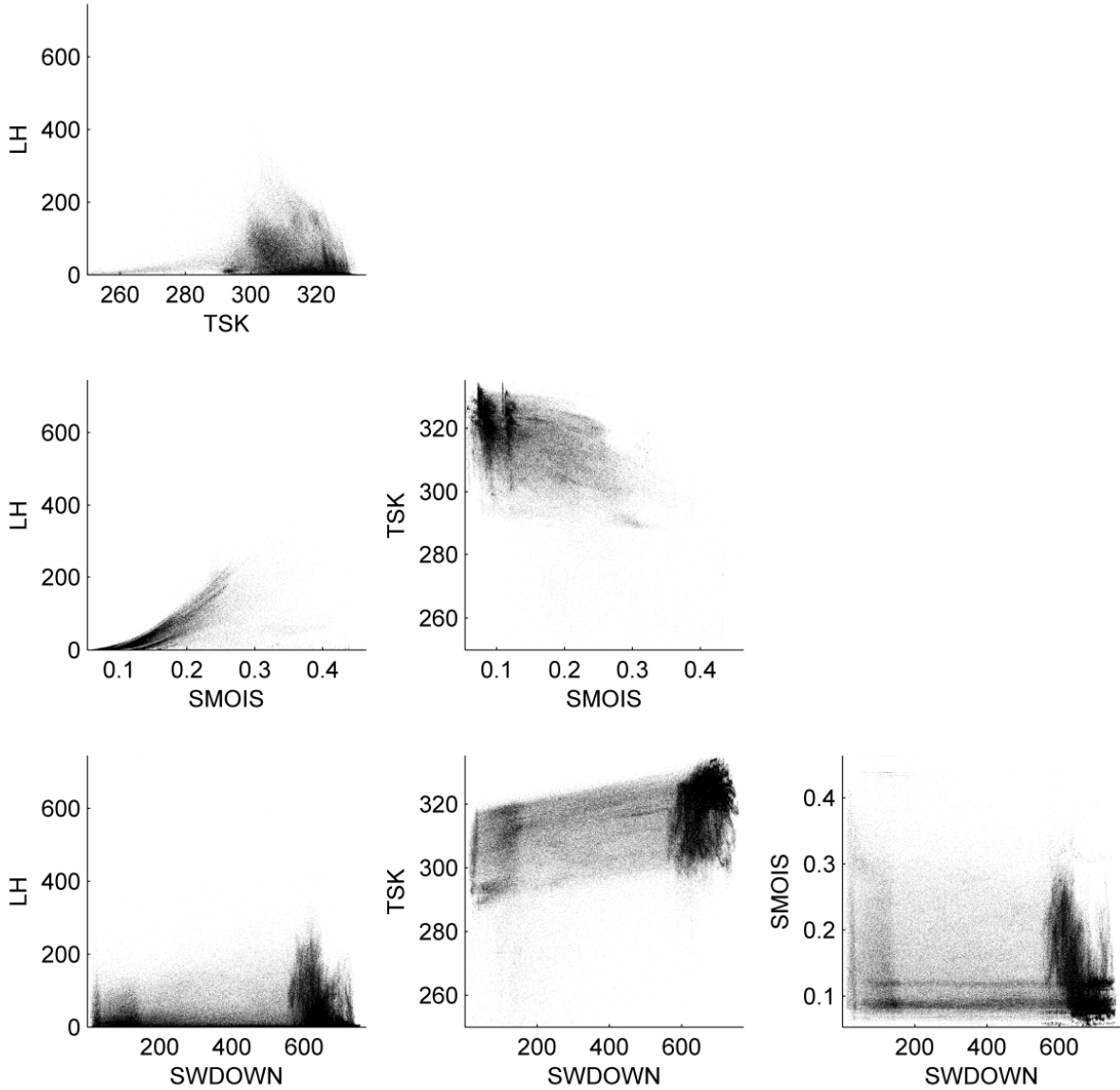
Reconstruction of gaps caused by orbital passages on remote sensing images



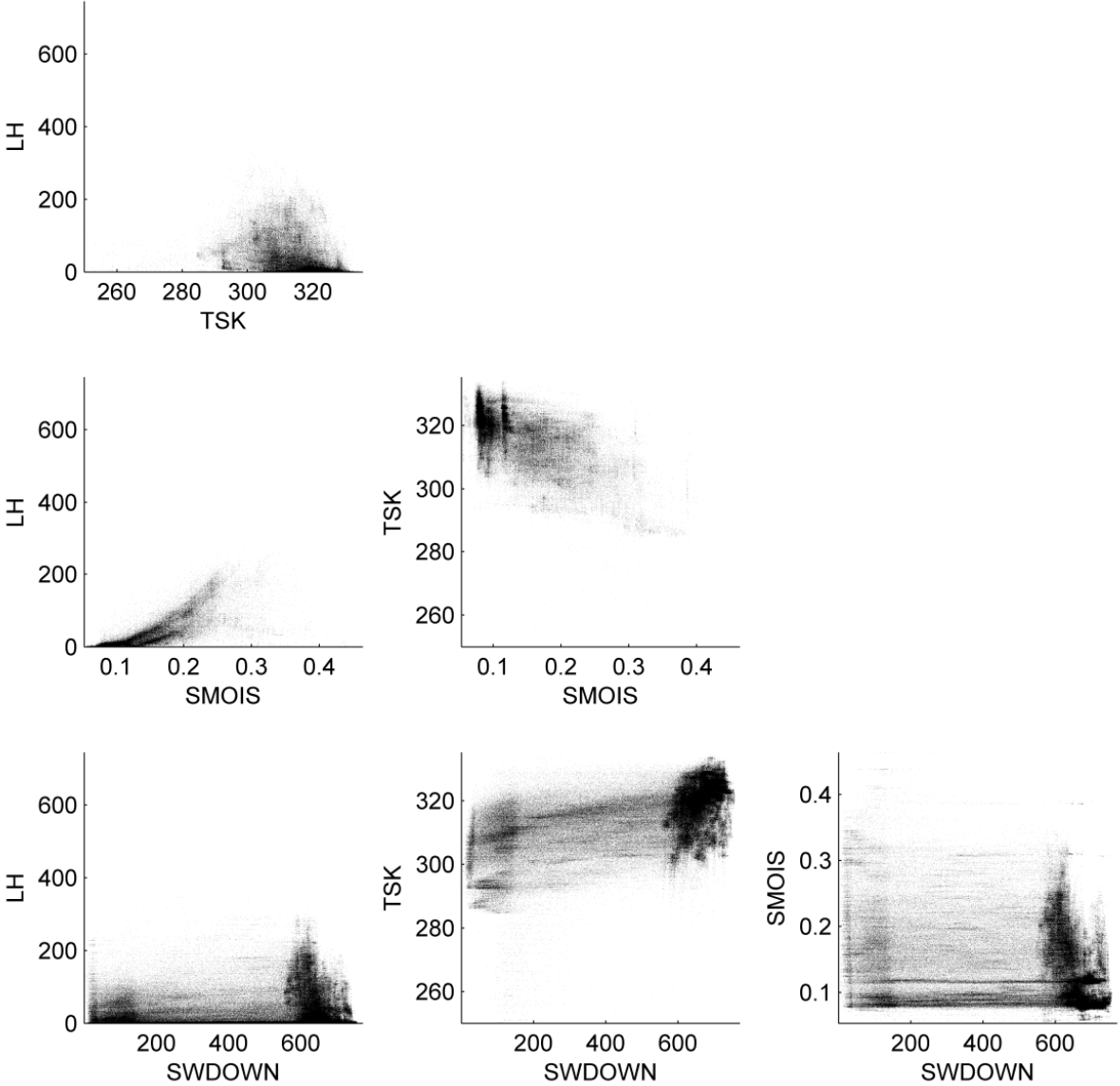
Error analysis



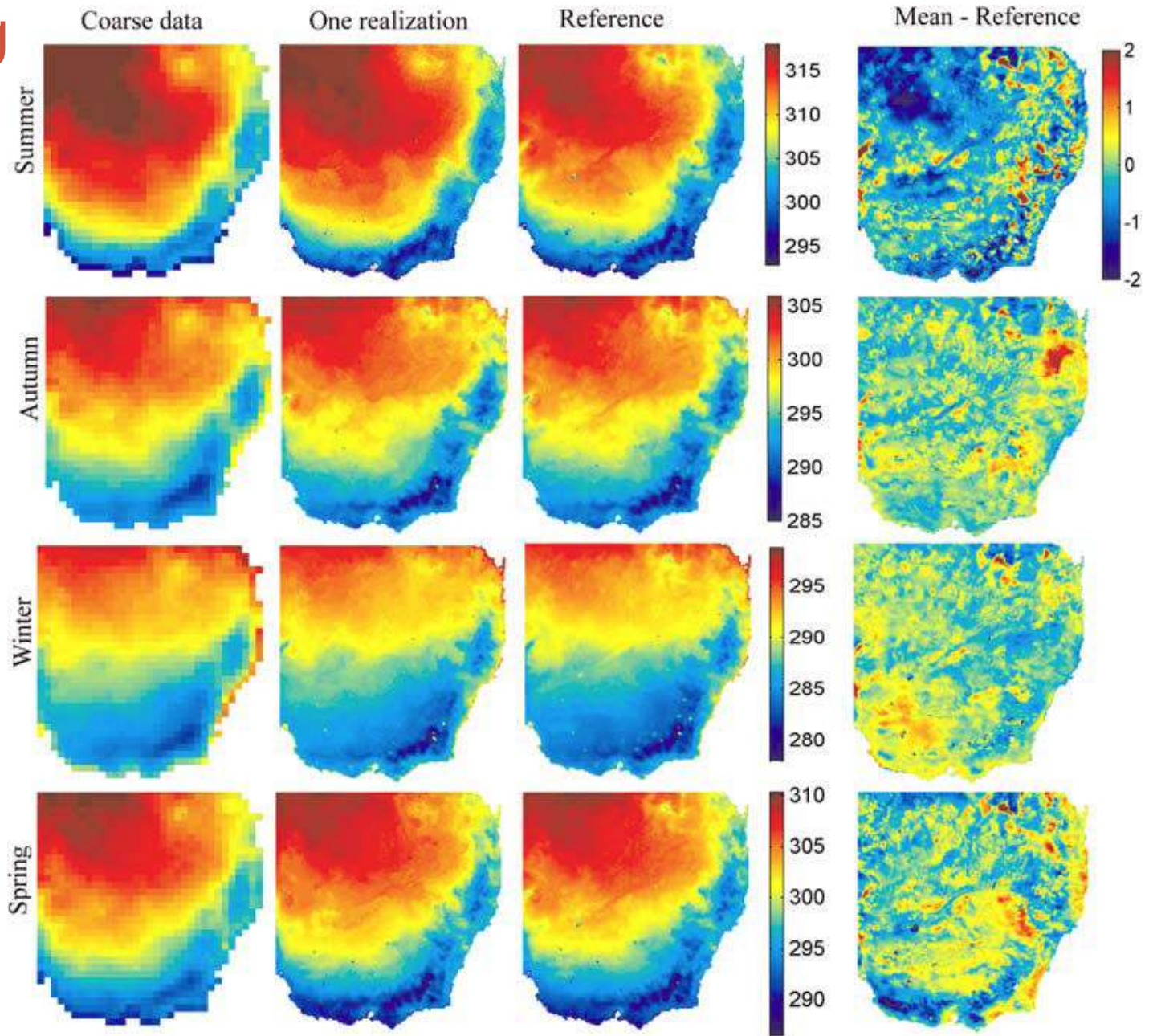
Dependence between reference values



Dependence between simulated values

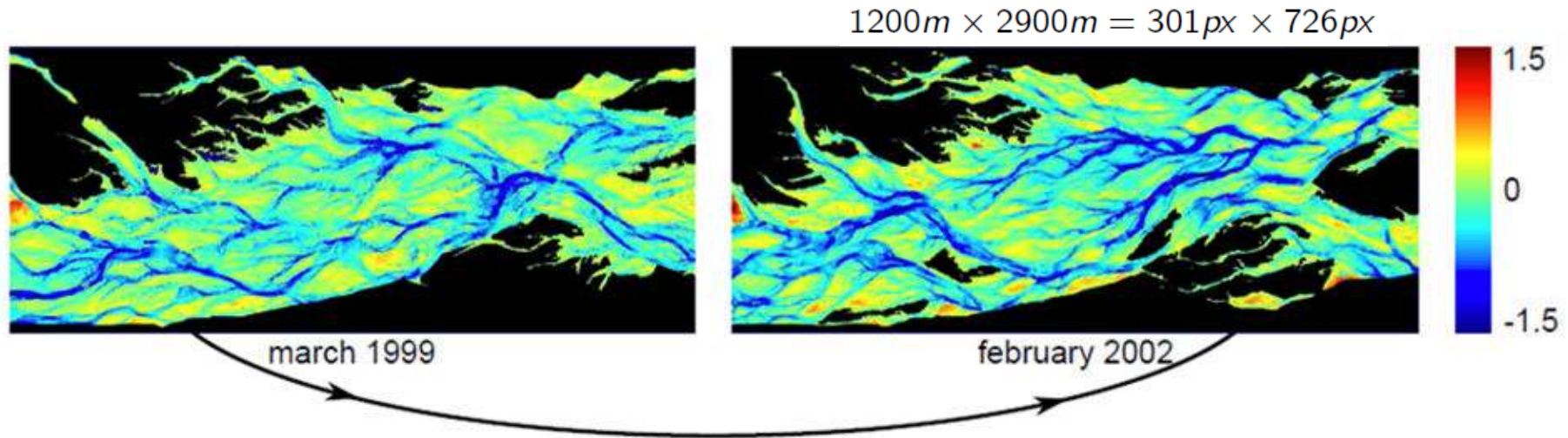


Downscaling

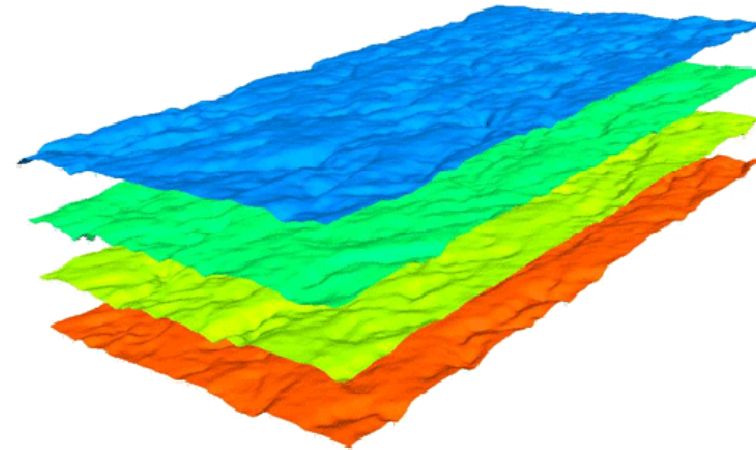


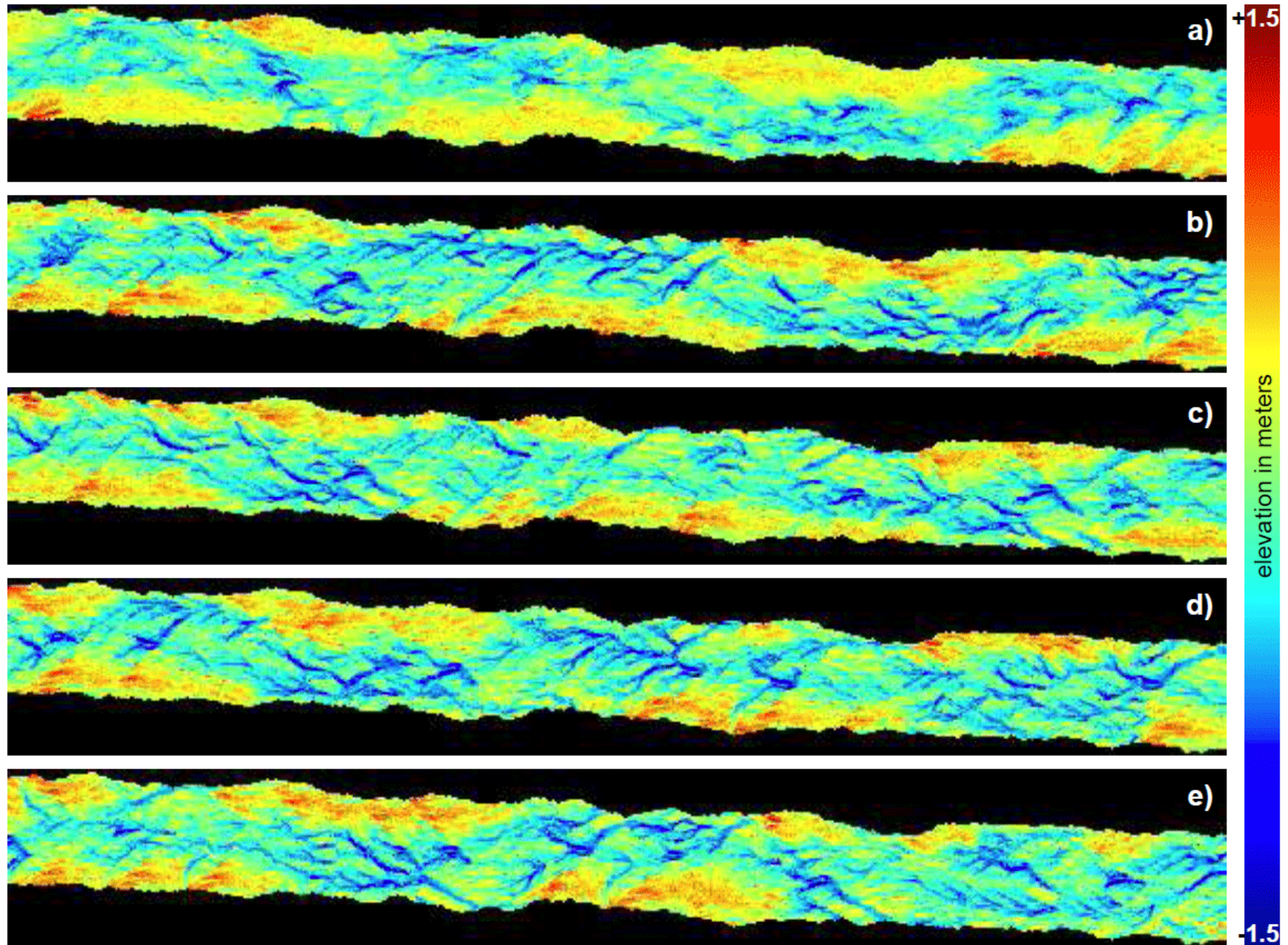
Surface
Temperature

Topography simulations

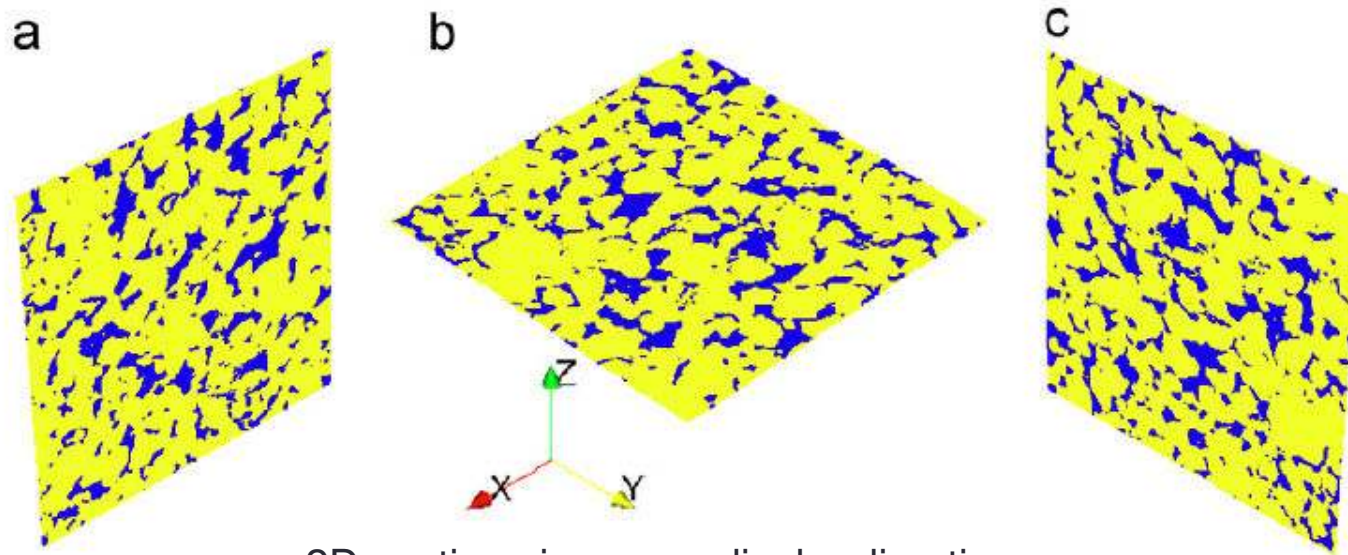


- High resolution DEM
- Multiple-point statistics to model successive topographies
- Stack them
- Fill the volumes with sediments



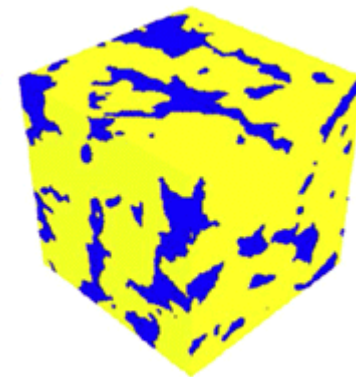


Porous media



2D sections in perpendicular directions
Berea Sandstone

Example of a 3D volume simulated from 2D data sets:

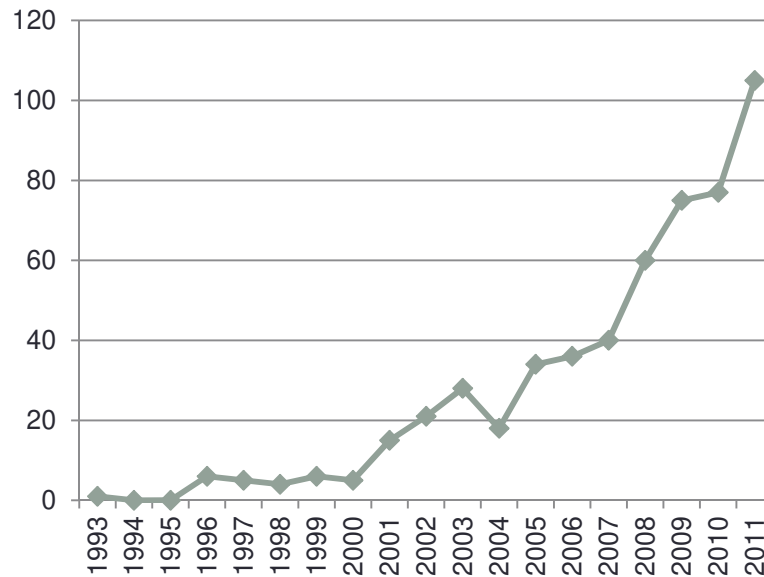


CONCLUSION

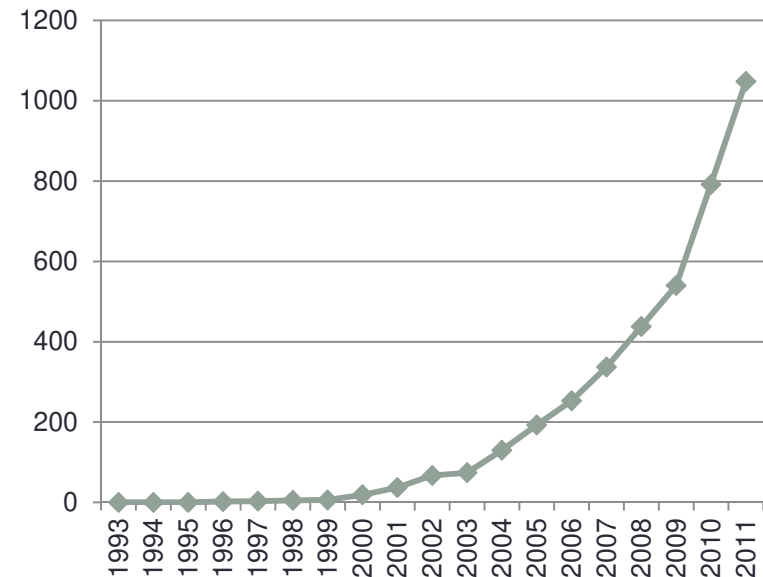
Active field of research

- Guardiano and Srivastava (1993)
- First efficient implementation:
 - Strebelle (2002), probability tree, multi-grid, etc.
- Since then:

of papers/year



of citations/year



MPS Pros / Cons

- Well suited to model complex structures
- General (same code for different structures)
- Easy conditioning
- Integration of secondary data

- Not a well defined Random Function model
- CPU time is longer than other methods
- Where to get the training image?

Current research directions

- Applications / demonstrations
- Braided river systems
- Spatio-temporal fields (Precipitation)

- Algorithmic improvements / acceleration

- Multi-scale
- Inverse problem



The ability to simplify means to eliminate the unnecessary so that the necessary may speak.

Hans Hoffman