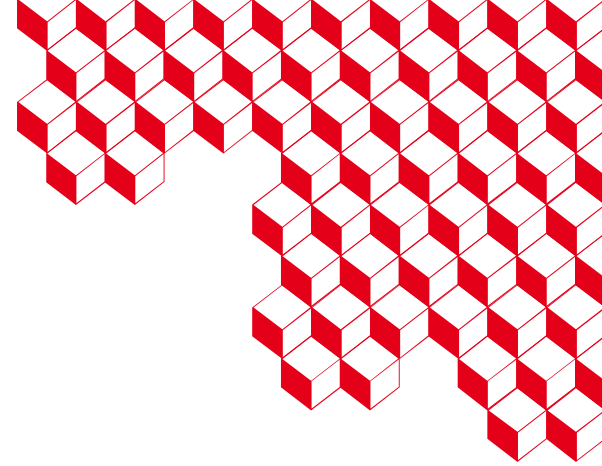




isds



# **A new exact and scalable model for multitask gaussian process regression : application to the reconstruction of nuclear data in neutronics codes**

Olivier TRUFFINET : CEA/DES/ISAS/DM2S/SERMA/LPEC

Supervisors : Karim AMMAR, Bertrand BOURIQUET, Jean-Philippe ARGAUD, Nicolas GERARD CASTAING

# A new exact and scalable model for multitask gaussian process regression

## 1. Multitask gaussian processes

- Reminders on gaussian processes
- The LMC framework
- Common simplifications

## 2. The projected LMC

- Decoupled expressions
- Resolution : noise model and loss
- Experimental results

## 3. Application : nuclear data in deterministic neutronics

- Description of the problem
- Results
- Corollary : adaptive sampling

## 4. Conclusion & Prospects



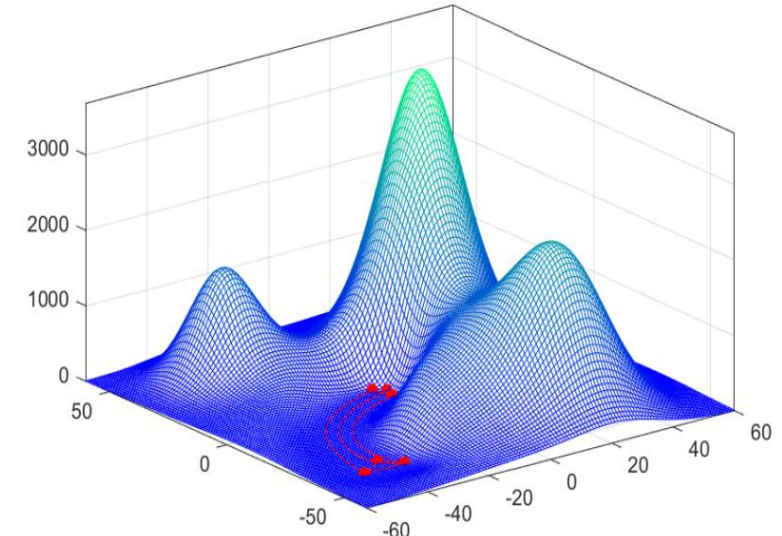


# 1 ■ Multitask GP

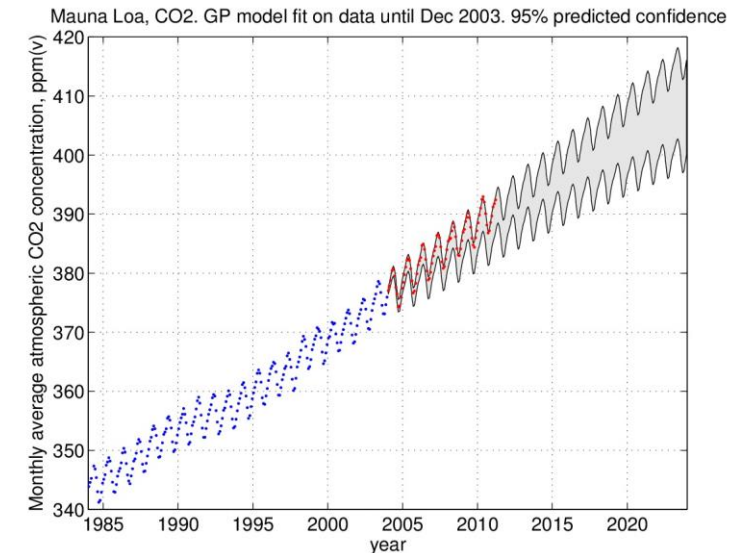
Elements on single-output and multi-output gaussian processes

# Introduction to gaussian processes

- Very popular methods for regression, classification, optimization...
- Apparented methods : Radial Basis Interpolation, Kernel regression, kriging (+Support Vector Regression, splines)
- Advantages :
  - ❖ Accuracy
  - ❖ No support restriction, trivial handling of multidimensionality
  - ❖ Flexibility (choice of kernels, mean functions...)
  - ❖ Safe optimization with maximum likelihood estimation
  - ❖ Uncertainty metrics thanks to the Bayesian framework
  - ❖ « Non-parametric »
- Flaws : Scalability? **Not so sure...**



Source: T.X.Lin & al, *A Distributed Scalar Field Mapping Strategy for Mobile Robots*



Source: C.E.Rasmussen, personal website

# Definition

$\theta$  : parameters of mean, kernel and noise

- Definition : a gaussian process is a **random variables field** over  $\mathbb{R}^d$ , such that each of its subsets follows a joint Gaussian distribution. It is entirely defined by its mean function  $m(x)$  and its covariance function (*kernel*)  $k(x, x')$  :

$$\begin{cases} \mathbb{E}[f(x)] = m(x) \\ \text{Cov}[f(x), f(x')] = k(x, x') \end{cases}$$

$f(x)$ : function to interpolate

- To estimate the process at an unobserved point  $x_*$ , one expresses the joint probability of known observations  $f(\mathbf{X})$  and unknown value  $f(x_*)$ :

$$\begin{bmatrix} \mathbf{y} \\ f(x_*) \end{bmatrix} = \mathcal{N} \left( \begin{bmatrix} m(\mathbf{X}) \\ m(x_*) \end{bmatrix}, \begin{bmatrix} k(\mathbf{X}, \mathbf{X}) + \sigma^2 I & k(x_*, \mathbf{X}) \\ k(x_*, \mathbf{X}) & k(x_*, x_*) \end{bmatrix} \right)$$

$\mathbf{y} = f(\mathbf{X}) + \text{noise}$

- One can then condition over  $f(\mathbf{X})$  to obtain the conditional distribution of  $f(x_*)$ :

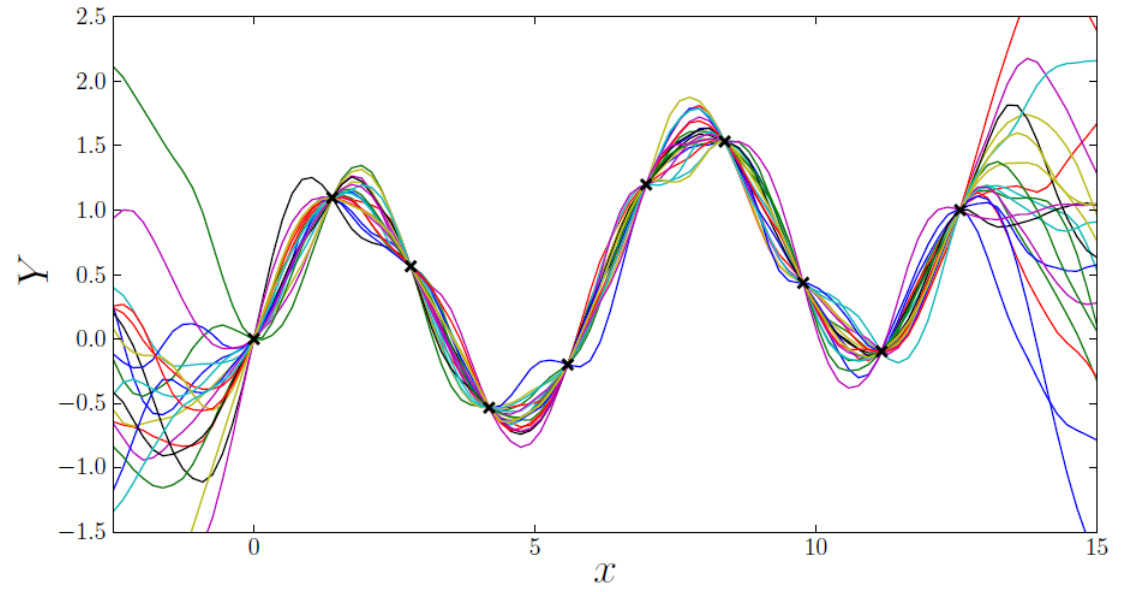
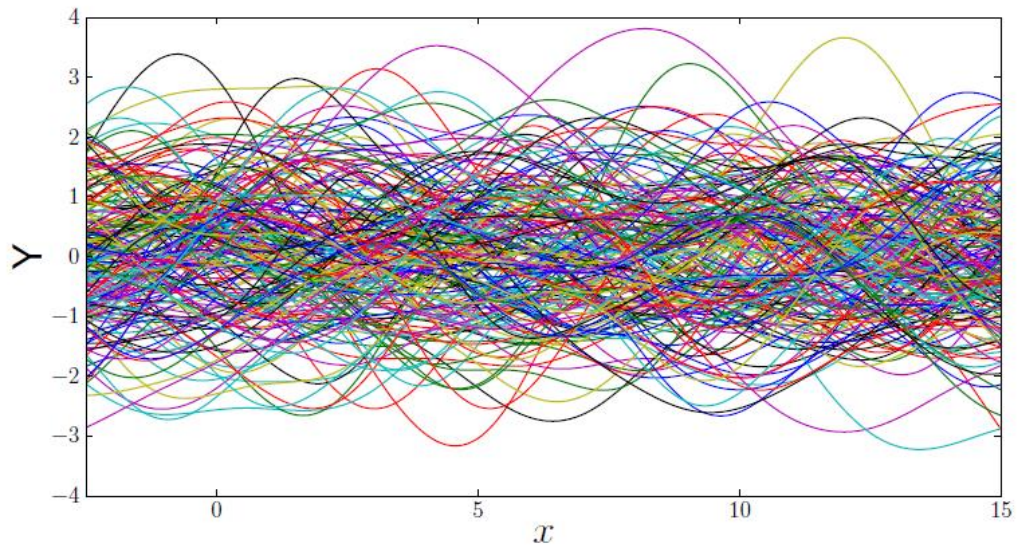
$$f(x_*) | (x_*, \mathbf{X}, \mathbf{y}) \sim \mathcal{N} \left( \underbrace{k(x_*, \mathbf{X})(k(\mathbf{X}, \mathbf{X}) + \sigma^2 I)^{-1} \mathbf{y}}_{\text{Mean}}, \underbrace{k(x_*, x_*) - k(x_*, \mathbf{X})(k(\mathbf{X}, \mathbf{X}) + \sigma^2 I)^{-1} k(x_*, \mathbf{X})}_{\text{Variance}} \right)$$

(here in the case  $m(x) = 0$ )

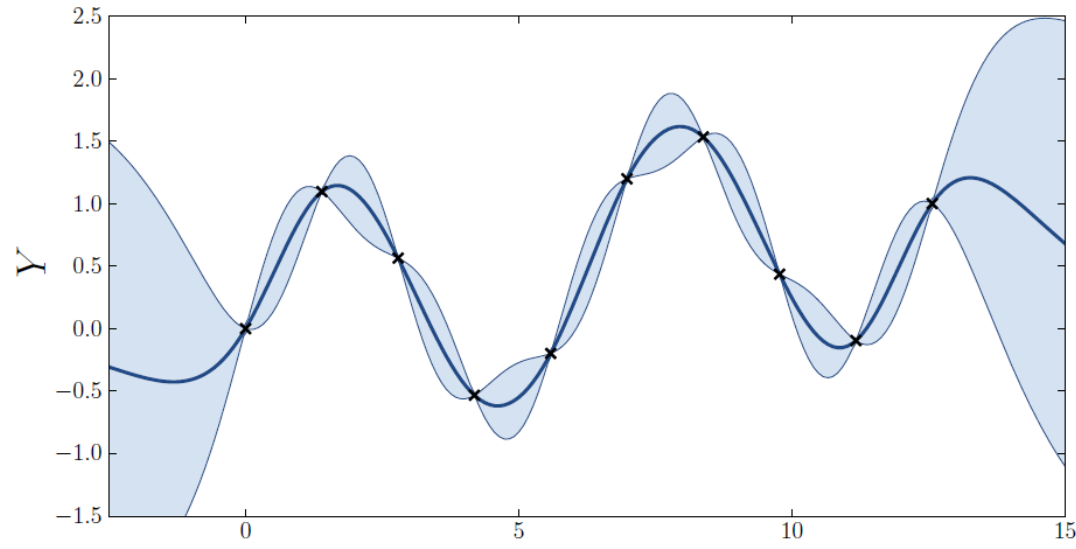
- Marginal log-likelihood (loss function) :  $\log p(\mathbf{y} | \mathbf{X}, \theta) = -\frac{1}{2} \mathbf{y}^T k(\mathbf{X}, \mathbf{X})^{-1} \mathbf{y} - \frac{1}{2} \log |k(\mathbf{X}, \mathbf{X})| + \text{cste}$

Since  $f()$  is unknown, we make the general assumption that it is the sample path of a Gaussian process  $Y \sim \mathcal{N}(\mu(), k(,))$ :

If we remove all the samples that do not interpolate the observations we obtain:



It can be summarized by a mean function and 95% confidence intervals.



Each curve represents a realization of the GP.

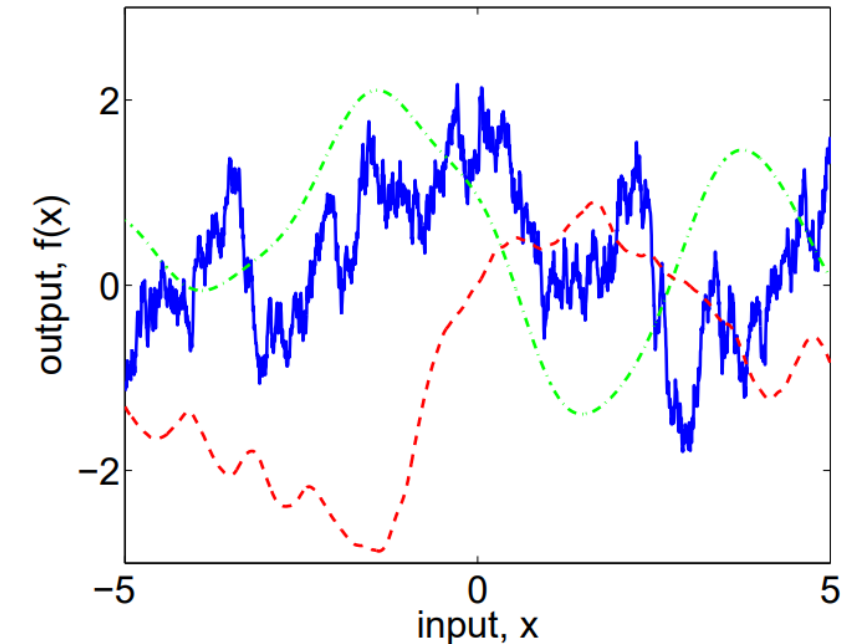
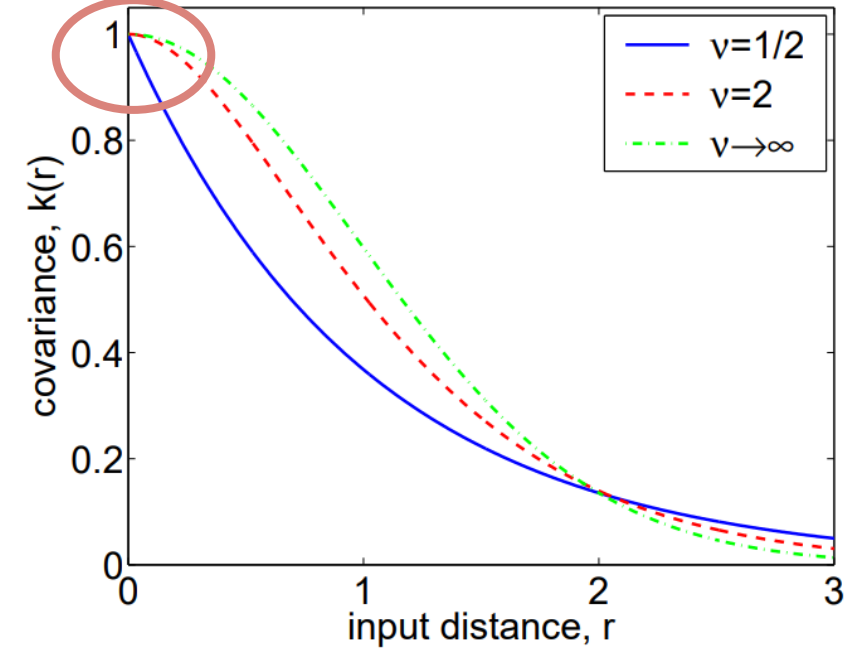
The modeled function is supposed to be such a realization.

Source : Le Riche & Durrande, *An overview of kriging for researchers*



# Kernel cookbook

- Stationary kernels :  $k(\mathbf{x}, \mathbf{x}') = \psi(\|\mathbf{x} - \mathbf{x}'\|) = \psi(r)$ . Highly recommended unless the data is strongly anisotropic
- Standard stationary kernels :  $\psi(r) = e^{-r^2}$  (*squared exponential*),  
 $\psi_\alpha(r) = (1 + \frac{r^2}{2\alpha})^{-1}$  (*rational quadratic* with parameter  $\alpha$ ),  
 $\psi_\nu(r) = 2^{1-\nu}(\sqrt{2\nu}r)^\nu K_\nu(\sqrt{2\nu}r)$  (*Matérn kernel* with parameter  $\nu$ )...
- The regularity of the estimator is controlled by the derivatives of  $\psi$  at  $r = 0$  (see figures opposite for processes obtained for Matérn kernels with different values of  $\nu$ )
- The distance  $r$  between two points can be customized, typically by parametrizing a different characteristic length  $l_i$  for each variable :  
$$r = \|\mathbf{x} - \mathbf{x}'\|_{\text{custom}} = \sqrt{\sum_{i=1}^d \left(\frac{x_i - x'_i}{l_i}\right)^2}$$
- A sum or product of kernels remains a valid kernel: we can for example build a tensorized kernel,  $k(\mathbf{x}, \mathbf{x}') = \prod_{i=1}^d k(x_i, x'_i)$



# Multi-output gaussian processes



- Reminder : equations of a single-output GP:

$$\begin{cases} f(x_*) | (x_*, X, y) \sim \mathcal{N} \left( k_*^T (K + \sigma^2 I)^{-1} y, \quad k_{**} - k_*^T (K + \sigma^2 I)^{-1} k_* \right) \\ -2 \log p(y | X, \theta) = y^T K^{-1} y - \log |K| + n\pi \end{cases}$$

$y$  : training output values  
 $X$  : training input values  
 $\sigma$  : noise  
 $\theta$  : hyperparameters

- General framework of multi-output GPs:

$$\begin{cases} f(x_*) | (x_*, X, Y) \sim \mathcal{N} \left( k_*^{M,T} (K^M + \Sigma \otimes I)^{-1} Y_v, \quad k_{**}^M - k_*^{M,T} (K^M + \Sigma \otimes I)^{-1} k_*^M \right) \\ -2 \log p(Y | X, \theta) = Y_v^T (K^M + \Sigma \otimes I)^{-1} Y_v - \log |K^M + \Sigma \otimes I| + np\pi \end{cases}$$

With  $k(\cdot, \cdot)^M = \begin{pmatrix} k_{1,1}(\cdot, \cdot) & \cdots & k_{1,p}(\cdot, \cdot) \\ \vdots & \ddots & \vdots \\ k_{1,p}(\cdot, \cdot) & \cdots & k_{p,p}(\cdot, \cdot) \end{pmatrix}$  a matricial kernel detailing cross-tasks covariances,  $Y_v = \text{vec}(Y)$

the flattened matrix  $Y$ ,  $\Sigma = (\sigma_{i,j})_{\substack{1 \leq j \leq p \\ 1 \leq i \leq p}}$  a cross-tasks noise matrix and  $\otimes$  the Kronecker product.

Kronecker product :  $A \otimes B = \begin{pmatrix} a_{1,1}B & \cdots & a_{1,m}B \\ \vdots & \ddots & \vdots \\ a_{n,1}B & \cdots & a_{n,m}B \end{pmatrix}$ . Its shape is the product of shapes of  $A$  and  $B$ .



# Linear model of co-regionalization

- We model each of the  $p$  functions to be reconstructed as a linear combination of  $q$  independent latent (unobserved) Gaussian processes  $u_j$  :

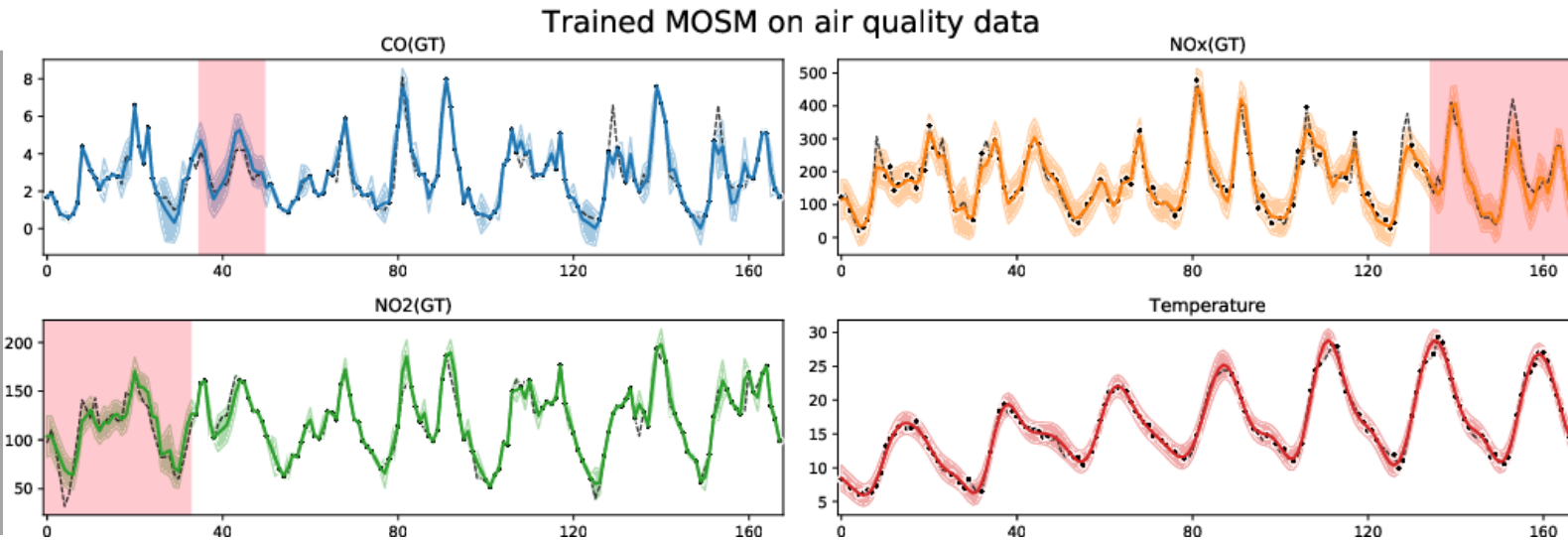
$$\widehat{y}_m(\mathbf{x}) = \sum_{j=1}^q H_{m,j} g_j(\mathbf{x}), \quad \forall m \in \llbracket 1; p \rrbracket$$

- This is equivalent to the previous formulation of a « standard » GP with matrix kernel:

$$k^M(\mathbf{x}, \mathbf{x}') = \sum_{j=1}^q \mathbf{B}_j k_j(\mathbf{x}, \mathbf{x}'), \quad \text{where } k_j \text{ is the kernel of process } g_j \text{ and } \mathbf{B}_m = \left( H_{i,m} H_{j,m} \right)_{\substack{1 \leq j \leq q \\ 1 \leq i \leq p}}.$$

The cross-task kernels are then :  $k_{i,l}(\cdot, \cdot) = \sum_{j=1}^q H_{i,j} H_{l,j} k_j(\cdot, \cdot)$ .

Problem : treating this exactly has time complexity  $O((np)^3)$ !



Source : T.Wolff & al, *MOGPTK: The Multi-Output Gaussian Process Toolkit*

# Intrinsic model of coregionalization (IMC)

- Simplification of the LMC : all latent processes have the same kernel  $k_x$ .  $K_x = k_x(\mathbf{X}, \mathbf{X})$

- $\Rightarrow K^M = \sum_{j=1}^q (\mathbf{B}_j \otimes K_x) = \left( \sum_{j=1}^q \mathbf{B}_j \right) \otimes K_x = \mathbf{H}\mathbf{H}^T \otimes K_x = \mathbf{K}_t \otimes K_x$

- Easy computation thanks to the trick :

$$\mathbf{K}_t \otimes K_x + \sigma \mathbf{I} = (\mathbf{U}_t \otimes \mathbf{U}_x) (\mathbf{D}_t \otimes \mathbf{D}_x + \sigma \mathbf{I}) (\mathbf{U}_t^T \otimes \mathbf{U}_x^T), \text{ with } \mathbf{U}_i \mathbf{D}_i \mathbf{U}_i^T \text{ the eigendecomposition of } \mathbf{K}_i$$

→ Runtime complexity in  $\mathcal{O}(n^3 + p^3)$  instead of  $\mathcal{O}((np)^3)$

- Other similar tricks whenever the covariance matrix is a Kronecker product of kernels
- BUT : much lesser expressivity (what if the data has several characteristic lengthscales ?)
- By comparison, for the LMC we have :  $K^M = (\mathbf{H} \otimes \mathbf{I}) \text{Diag}(\mathbf{K}_i) (\mathbf{H}^T \otimes \mathbf{I})$ , with  $\mathbf{H} \in \mathbb{R}^{p \times q}$  ( $p$  number of tasks,  $q$  number of latent processes,  $n$  number of points)

# Inducing points

- Most common GP approximation
- Principle: few "pseudo-inputs" points with inferred output values replace the much more numerous true datapoints
- Pseudo-points locations are optimized parameters
- Rationale: the set of training points is often redundant. Replacing it with a smaller number of optimized pseudo-inputs greatly speeds up matrix inversions, without degrading performance too much.

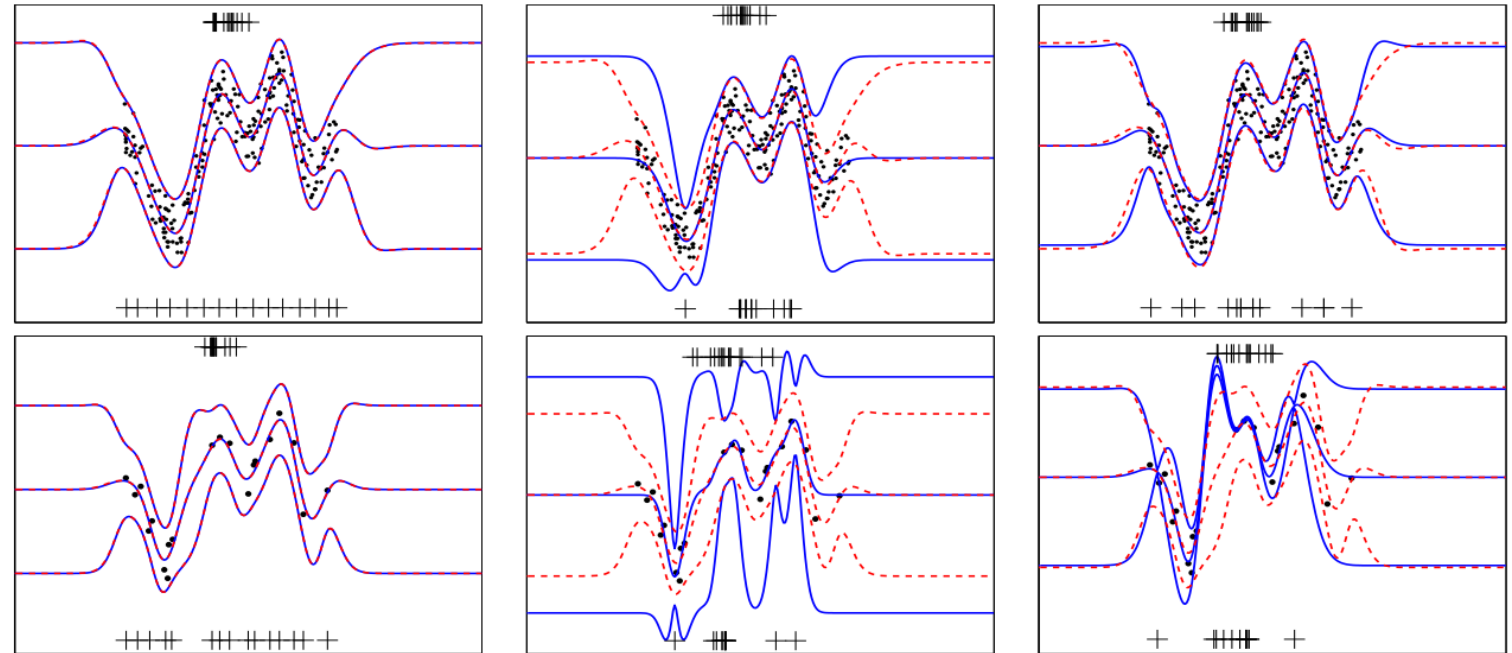


Figure 1: The first row corresponds to 200 training points and the second row to 20 training points. The first column shows the prediction (blue solid lines) obtained by maximizing  $F_V$  over the 15 pseudo-inputs and the hyperparameters. The full GP prediction is shown with red dashed lines. Initial locations of the pseudo-inputs are shown on the top as crosses, while final positions are given on the bottom as crosses. The second column shows the predictive distributions found by PP and similarly the third column for SPGP.

# Variational approach

- Built upon the inducing points approximation
- $\log p(\mathbf{y}) = \log \langle p(\mathbf{y}|\mathbf{f}) \rangle_{p(\mathbf{f}|\mathbf{u})} \geq \langle \log p(\mathbf{y}|\mathbf{f}) \rangle_{p(\mathbf{f}|\mathbf{u})} \equiv \mathcal{L}$  by convexity, with  $\mathbf{u}$  the fictive values of  $f$  at inducing points locations, and  $\mathbf{f}$  its possible (non-noisy) values at training locations.
- **The lower bound  $\mathcal{L}$  is optimized instead of the true marginal likelihood.** To compute it efficiently,  $p(\mathbf{f}, \mathbf{u} | \mathbf{y})$  is replaced by a variational distribution  $q(\mathbf{u})$  which parameters are optimized
- Predictions :  $p(f(x_*) | \mathbf{y}) \approx \int \underbrace{p(f(x_*)|\mathbf{u})}_{\text{Standard GP expression}} \underbrace{q(\mathbf{u})}_{\text{Gaussian}} d\mathbf{u}$
- Numerous variations are possible within this general approach (modifications to  $\mathcal{L}$ , choice of  $q$ ), which can in particular enable decoupling of inputs and outputs, or minibatch training (=stochastic optimization with subsets of the dataset).
- **A variational LMC can be built this way** : the variational distribution becomes a product of independent latent GPs,  $q(\mathbf{u}) = \prod_{j=1}^q q_j(\mathbf{u}) \rightarrow$  Runtime complexity in  $\mathcal{O}(n^3 + p^3)$



# 2 ■ The projected LMC

A nice finding on the model structure !

# Exact posteriors of latent processes

- Question (initially for leave-one-out error computation) : can we « open the box » of the LMC and compute the values taken by latent processes ?

- It turns out that we can : 
$$\begin{cases} \mathbb{E}[\mathbf{y}_*|\mathbf{Y}] = \text{Diag}(\mathbf{k}_{i_*}^T)(\text{Diag}(\mathbf{K}_i) + \boldsymbol{\Sigma}_P \otimes \mathbf{I})^{-1} \text{vec}(\mathbf{T}\mathbf{Y}) \\ \mathbb{V}[\mathbf{y}_*|\mathbf{Y}] = \text{Diag}(k_{i_{**}}) - \text{Diag}(\mathbf{k}_{i_*}^T)(\text{Diag}(\mathbf{K}_i) + \boldsymbol{\Sigma}_P \otimes \mathbf{I})^{-1} \text{Diag}(\mathbf{k}_{i_*}) \end{cases}$$

With :  $\mathbf{k}_{i_*} = k_i(\mathbf{x}_*, \mathbf{X})$ ,  $k_{i_{**}} = k_i(\mathbf{x}_*, \mathbf{x}_*)$ ,  $\mathbf{K}_i = k_i(\mathbf{X}, \mathbf{X})$ ,

And:  $\boldsymbol{\Sigma}_P = (\mathbf{H}^T \boldsymbol{\Sigma} \mathbf{H})^{-1}$ ,  $\mathbf{T} = \boldsymbol{\Sigma}_P \mathbf{H}^T \boldsymbol{\Sigma}^{-1}$ .

- Many block-diagonal matrices here; if  $\boldsymbol{\Sigma}_P$  is also diagonal, both expression decouple !! Then :

$$\begin{cases} \mathbb{E}[y_{j_*}|\mathbf{Y}] = \mathbf{k}_{j_*}^T (\mathbf{K}_j + \sigma_j \mathbf{I})^{-1} \mathbf{T}_i \mathbf{Y} \\ \mathbb{V}[y_{j_*}|\mathbf{Y}] = k_{j_{**}} - \mathbf{k}_{j_*}^T (\mathbf{K}_j + \sigma_j \mathbf{I})^{-1} \mathbf{k}_{j_*} \end{cases} \rightarrow \text{expressions of standard single-output GPs} \quad (\boldsymbol{\Sigma}_P = \text{Diag}(\sigma_i))$$

- Same is true for the likelihood :  $\prod_{l=1}^n \frac{\mathcal{N}(\mathbf{Y}_j | \mathbf{0}, \boldsymbol{\Sigma})}{\mathcal{N}(\mathbf{T}\mathbf{Y}_j | \mathbf{0}, \boldsymbol{\Sigma}_P)} \times \int p(\mathbf{G}) \prod_{l=1}^n \mathcal{N}(\mathbf{T}\mathbf{Y}_j | \mathbf{G}_j, \boldsymbol{\Sigma}_P) d\mathbf{G}$   
( $\mathbf{G}$  values of latent processes at observed locations)



# Sketch of the proof(s)

- Block-diagonal matrix algebra :

- Start from  $f(x_*)|x_*, X, Y \sim \mathcal{N} \left( k_*^{M,T} (K^M + \Sigma \otimes I)^{-1} Y_v, k_{**}^M - k_*^{M,T} (K^M + \Sigma \otimes I)^{-1} k_*^M \right)$ .

- In each expression, expand inverses with Woodburry's formula :

$$(A + UCV)^{-1} = A^{-1} - A^{-1}U(C^{-1} + VA^{-1}U)^{-1}VA^{-1}$$

- Rearrange using properties of the Kronecker product :

$$(A \otimes B)(C \otimes D) = (AC \otimes BD), \quad (C^T \otimes A)vec(B) = vec(ABC), \quad (A \otimes B)^{-1} = A^{-1} \otimes B^{-1}$$

- Use backward Woodburry's formula to complete the proof.

- Ab initio bayesian computation :

- Start by computing the posterior of latent processes at training points,  $p(G|Y)$ . By Baye's rule,  $p(G|Y) \propto p(G|Y)p(G)$  ; by properties of gaussians,  $p(G|Y)$  is a gaussian in the variables  $G$ . One can therefore expand  $p(G|Y)$  and  $p(G)$  which are given by model definition, and identify terms in  $G$ .

- From there, deduce the posterior of latent processes at test point,  $p(g_*|Y)$ . For this, use the laws of total expectation and total variance, conditioning on  $G$  :

$$\mathbb{E}[Z] = \mathbb{E}[\mathbb{E}[Z|Z']] ; \mathbb{V}[Z] = \mathbb{E}[\mathbb{V}[Z|Z']] + \mathbb{V}[\mathbb{E}[Z|Z']] , \text{ and GP identities yielding } \mathbb{E}[g_*|G].$$

# Projected data and projected noise

- What is the meaning of the recurring terms  $\Sigma_P = (H^T \Sigma H)^{-1}$  and  $T = \Sigma_P H^T \Sigma^{-1}$  ?
- We have the following set of properties :
  - $T$  is a generalized inverse of  $H$  :  $TH = I_q$ . Therefore,  $HT$  is a projection.
  - $TY$  is a maximum likelihood estimator of  $G$  :  $TY = \underset{G}{\operatorname{argmax}} p(Y|G)$ .
  - $TY$  is a sufficient statistic of the model for the variable  $G$  :  $p(Y|TY, G) = p(Y|TY)$
  - $T\Sigma T^T = \Sigma_P$ , so that  $TY|G \sim \mathcal{N}(\operatorname{vec}(G), \Sigma_P)$
  - Latent processes are conditionally independent  $\Leftrightarrow \Sigma_P$  is diagonal
- Interpretation :
  - $TY$  is the data « seen » by the latent processes. In some sense, it is a projection of the full data  $Y$ .
  - It contains all useful information regarding latent processes, and is the best estimate of their observed values.
  - $\Sigma_P$  is the covariance matrix of this projected data ; if it is diagonal, latent processes are independent.



# A compatible parametrization

- **Diagonal Projected Noise (DPN) hypothesis** : enforcing  $\Sigma_p$  to be diagonal decouples latent processes and yields runtime complexity  $\mathcal{O}(n^3 + p^3)$
- **Problem** : this condition couples the mixing matrix  $H$  and noise matrix  $\Sigma \rightarrow$  need for an adapted parametrization.
- Final retained parametrization : set the  $QR$  decomposition of  $H$ , define  $Q_\perp$  an orthogonal supplement of  $Q$ , decompose  $\Sigma^{-1} = QAQ^T + Q_\perp BQ_\perp^T + QCQ_\perp^T + Q_\perp C^T Q^T$ , set  $D = R^T A R$  and  $M = R^T C$ . Then the DPN condition is equivalent to ( $D$  diagonal), and :

$$\Sigma = Q_+ R_+ D_+ R_+^T Q_+^T, \quad \text{with } Q_+ = (Q | Q_\perp), \quad R_+ = \begin{pmatrix} R & \mathbf{0} \\ \mathbf{0} & I_{p-q} \end{pmatrix}, \quad D_+^{-1} = \begin{pmatrix} D & M \\ M^T & B \end{pmatrix}; \quad D_+ \equiv \begin{pmatrix} \tilde{D} & \tilde{M} \\ \tilde{M}^T & \tilde{B} \end{pmatrix}$$

- This yields a convenient form for the marginal log-likelihood :

$$-2 \log p(Y) = n(p - q)\pi + \underbrace{2n \log |R| + n \log \tilde{B}}_{\text{Discarded noise}} + \underbrace{\text{Tr}(Y^T Q_\perp \tilde{B}^{-1} Q_\perp^T Y)}_{\text{Discarded data}} + \sum_{j=1}^q \underbrace{\log \mathcal{N}(Y T_j | \mathbf{0}, K_j + \sigma_j I_n)}_{\text{Latent processes}}$$

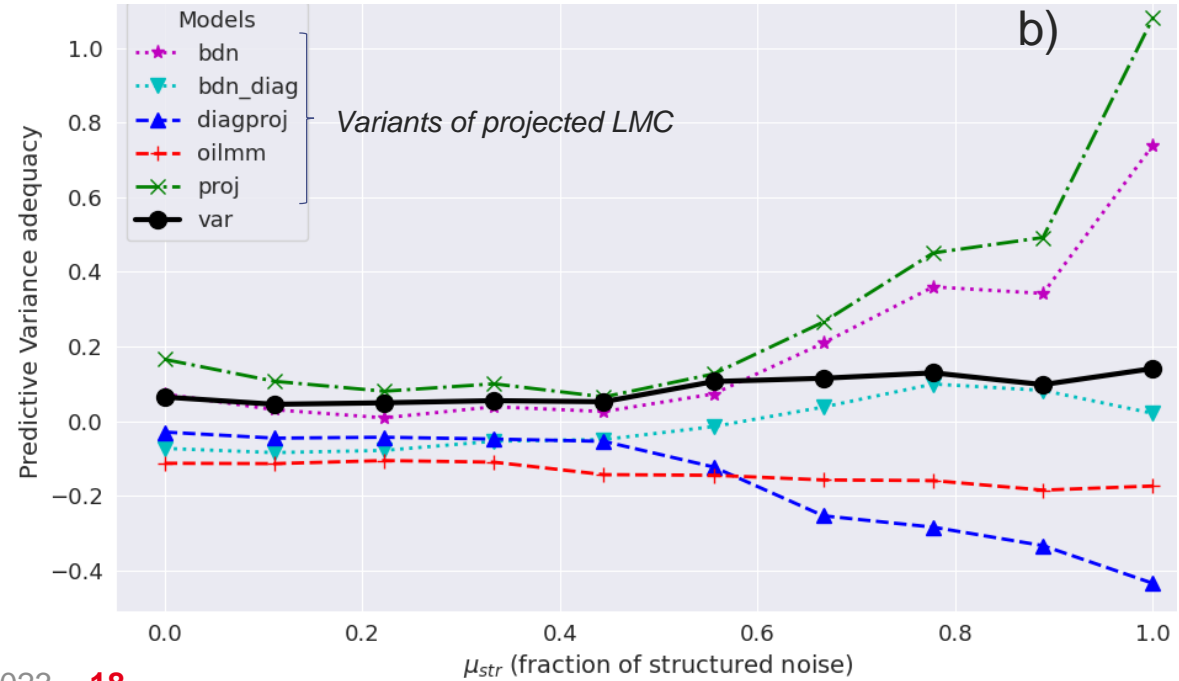
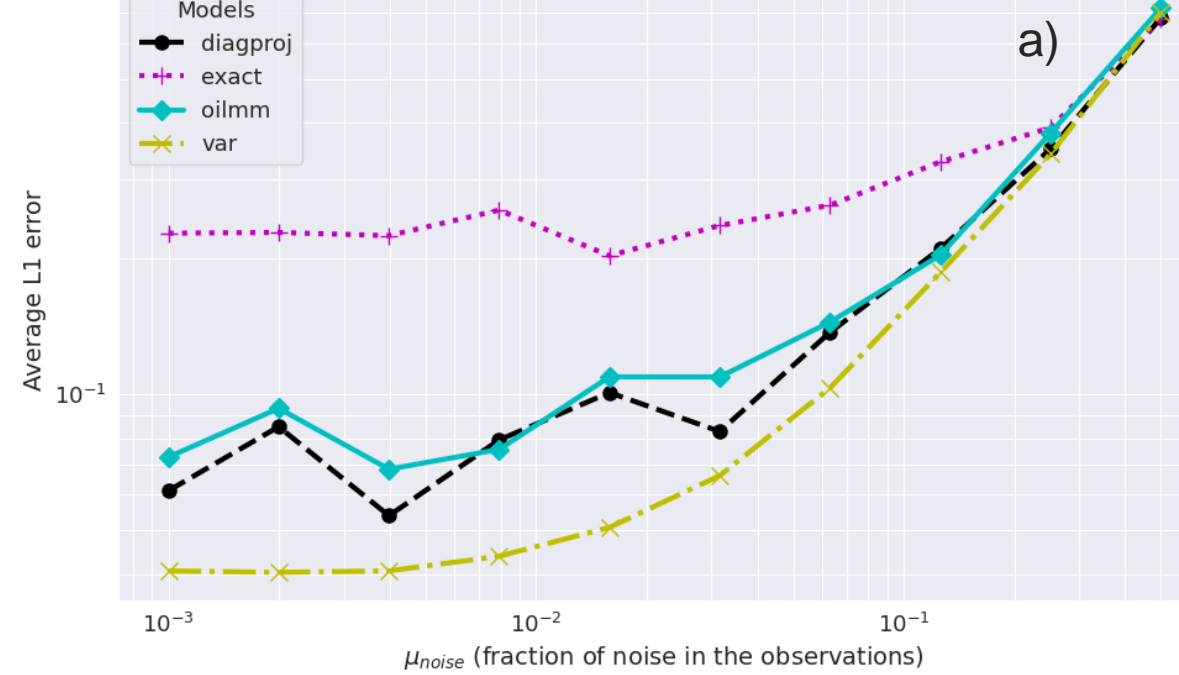
# Results : impact of the noise hypothesis

- Tests on synthetic data with controlled noise magnitude and share of structured (= mixed cross-tasks) noise // per-task noise
- The model doesn't seem to underperform even for large and heavily-structured noises !

$$PVA = \log \left( \frac{1}{np} \sum_{i=1}^n \sum_{j=1}^p \frac{(y_{ij} - \widehat{y}_{ij})^2}{\widehat{v}_{ij}} \right)$$

**Fig.1:**

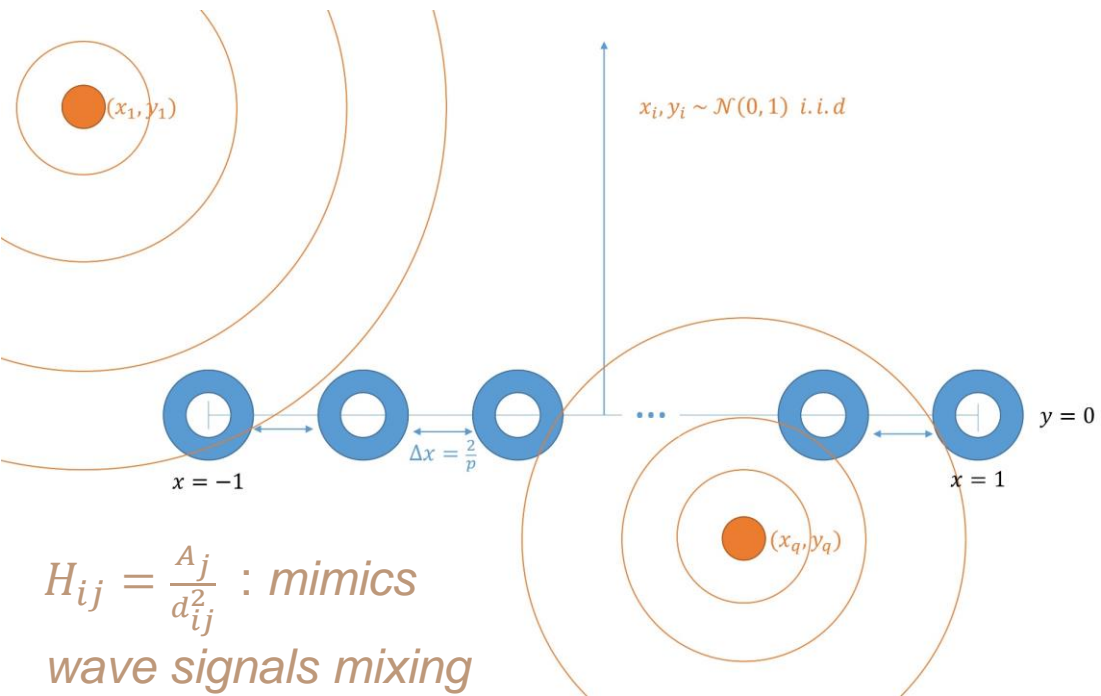
- a) Average test L1 error of the projected model ('diagproj') compared with reference LMC approaches, in function of the proportion of noise
- b) Average PVA errors of reference and projected models, in function of the proportion of structured-to-unstructured noise



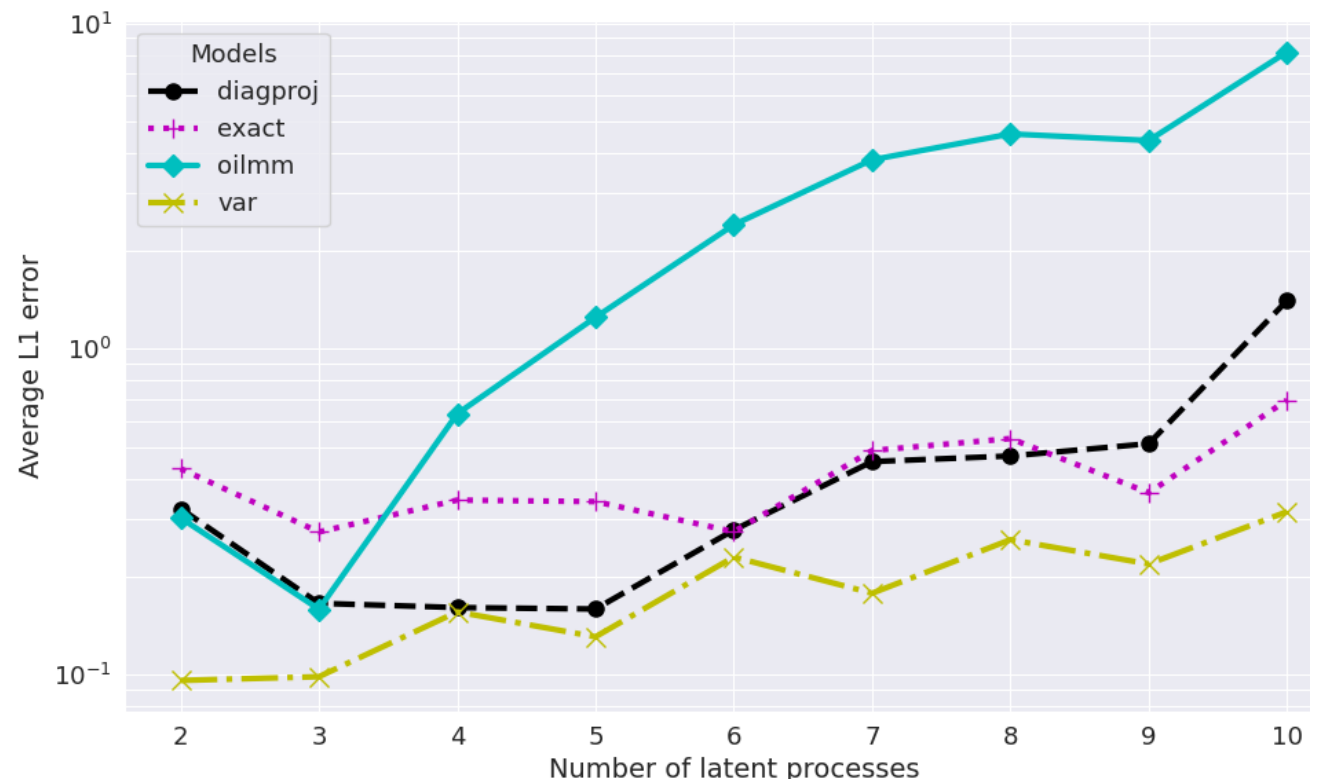
# Results : benefits of a general mixing matrix



- This work is based on a paper on OILMM, a model which derived the same results in the special case of an orthogonal mixing matrix  $H$  (see previous graphs)
- This restriction seems to cope very well with the previous case of i.i.d normally-distributed  $H$  coefficients
- Can we find realistic cases where this restriction is more of a problem ? → Yes we can :



**Fig.2:** set up of the « physical » toy data leading to non-orthogonal  $H$



**Fig.3:** Average test errors of several models in function of  $N_{latents}$



# 3 ■ Application

Homogenized cross-sections in deterministic neutronics codes

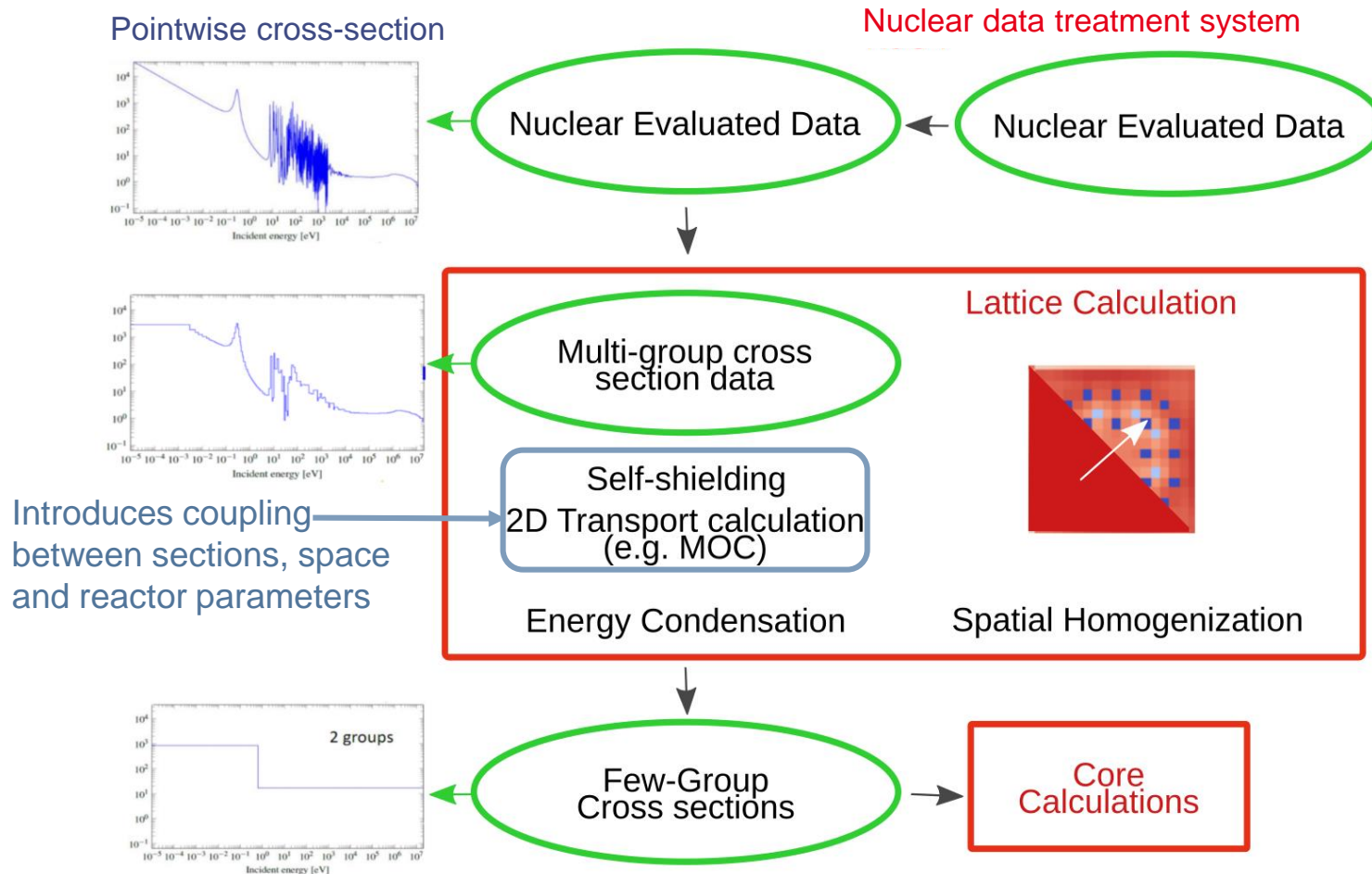
# Context

- **Neutron modeling:** determination of important core parameters (neutron flux, local power, etc.) by solving the Boltzmann equation: (+ Bateman equations for evolution of isotopic concentrations)

$$\begin{aligned}
 \underbrace{\frac{1}{v} \frac{\partial \phi(\mathbf{r}, v, \boldsymbol{\Omega}, t)}{\partial t}}_{\text{Variation of neutron population}} &= + \chi(v) \int_0^{+\infty} dv' \frac{1}{4\pi} \int_{4\pi} d^2\Omega' \bar{\nu} \Sigma_f(\mathbf{r}, v') \phi(\mathbf{r}, v', \boldsymbol{\Omega}', t) && \text{Fissions} \\
 &+ \int_{4\pi} d^2\Omega' \int_0^{+\infty} dv' [\Sigma_s(\mathbf{r}, v' \rightarrow v, \boldsymbol{\Omega}' \rightarrow \boldsymbol{\Omega}, t) \underbrace{\phi(\mathbf{r}, v', \boldsymbol{\Omega}', t)}_{\text{Flux = neutron population}}] && \text{Scattering} \\
 &+ S_a(\mathbf{r}, v, \boldsymbol{\Omega}, t) && \text{External source} \\
 &- \underbrace{\Sigma_t(\mathbf{r}, v, \boldsymbol{\Omega}, t)}_{\text{Cross-sections}} \phi(\mathbf{r}, v, \boldsymbol{\Omega}, t) && \text{Absorption} \\
 &- \text{div} [\boldsymbol{\Omega} \phi(\mathbf{r}, v, \boldsymbol{\Omega}, t)] && \text{Leaks}
 \end{aligned}$$

- 6+1 phase space, resonant phenomena → complex discretization

# Deterministic neutronics code



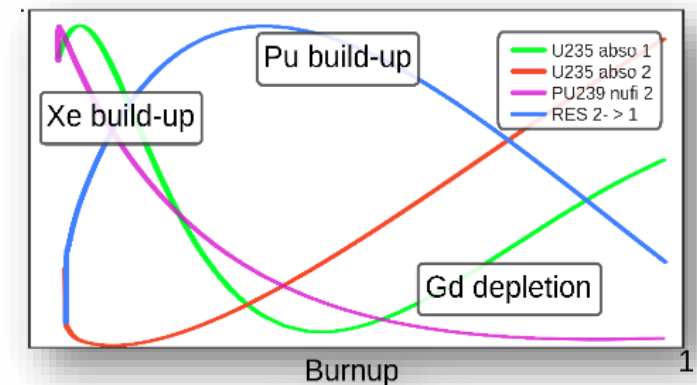
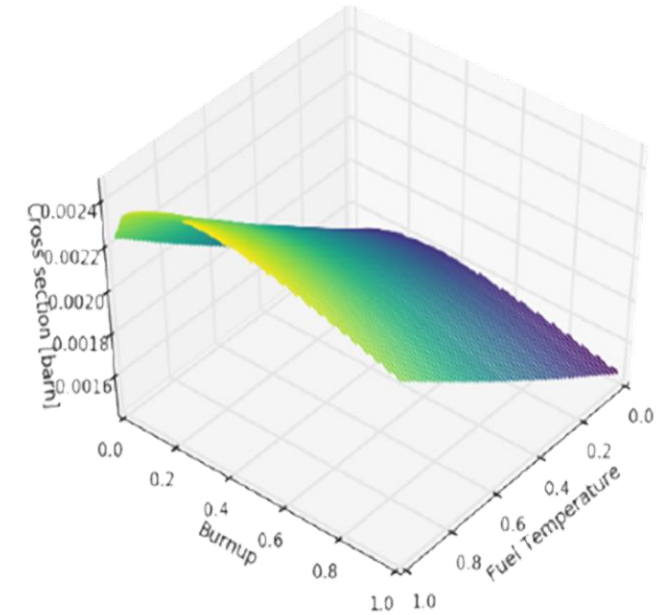
Source: E. Szames, *Few group cross section modeling by machine learning for nuclear reactor*

# Homogenized cross-sections (HXS)

- After a lattice calculation (assembly-wise), point cross sections are: **homogenized in a region  $v$** , **condensed over an energy range  $g$** , for a reaction type  $r$ , and a particularized isotope  $i$ .

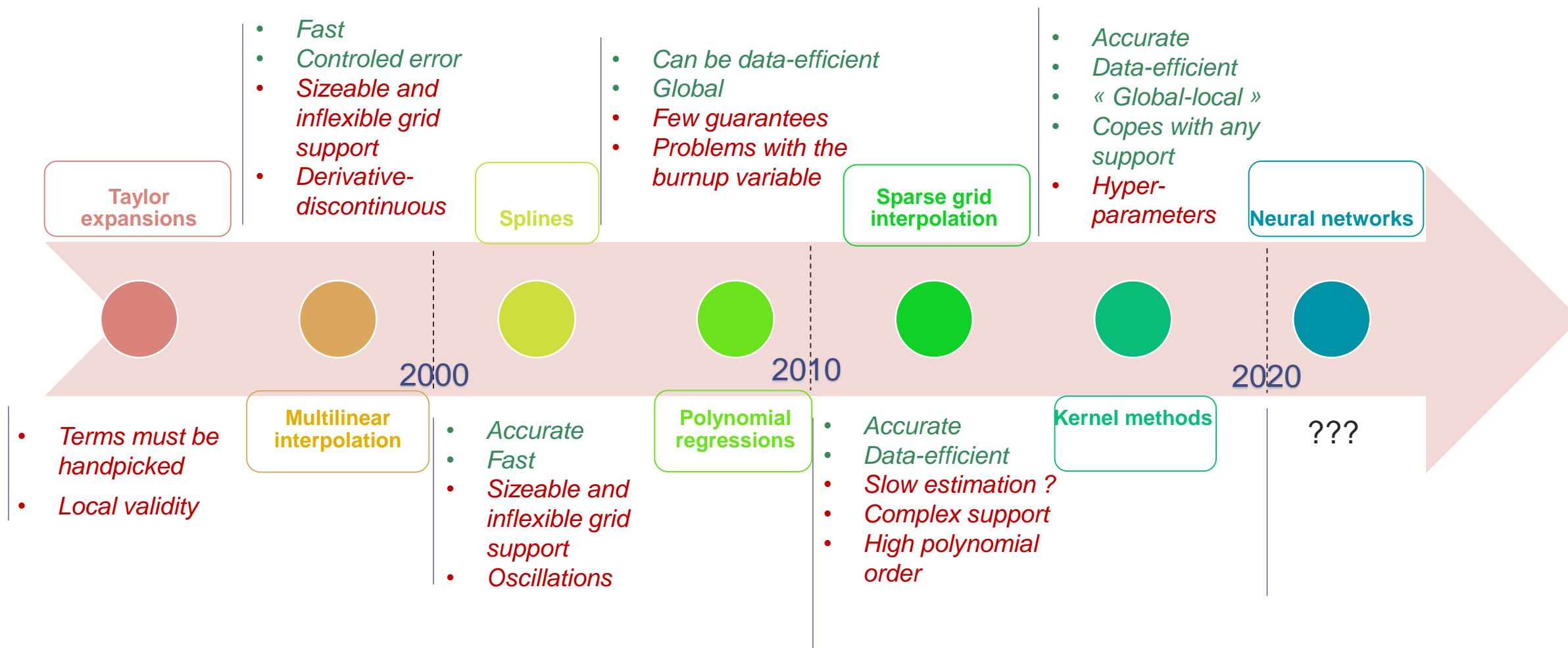
$$\sigma_{v,i,r,g} = \frac{\int_v \int_g (\sigma_{i,r} \phi)(\vec{r}, E) dE d\vec{r}}{\int_v \int_g \phi(\vec{r}, E) dE d\vec{r}}$$

- Problem to be solved: create and store a large number (up to millions) of highly correlated multivariate approximation models, one for each homogenized cross-section.
- Constraints: high accuracy requirement (relative error  $\sim 10^{-3}$ ), low memory footprint, smallest possible support, fast estimation



Source: E.Szames

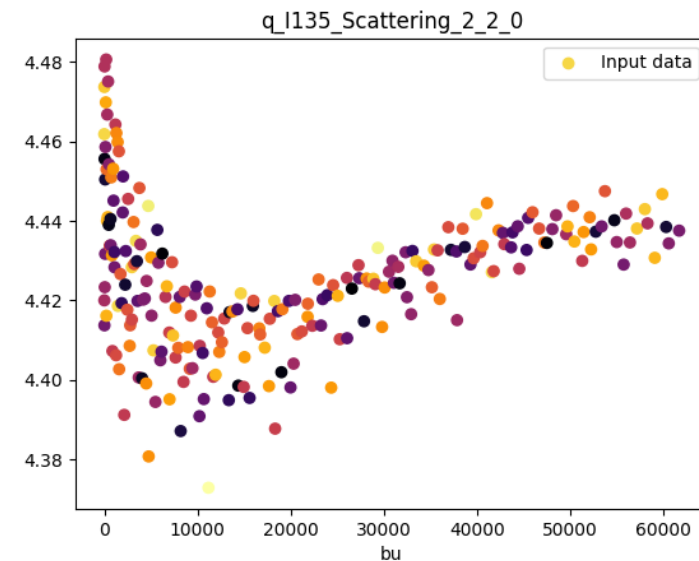
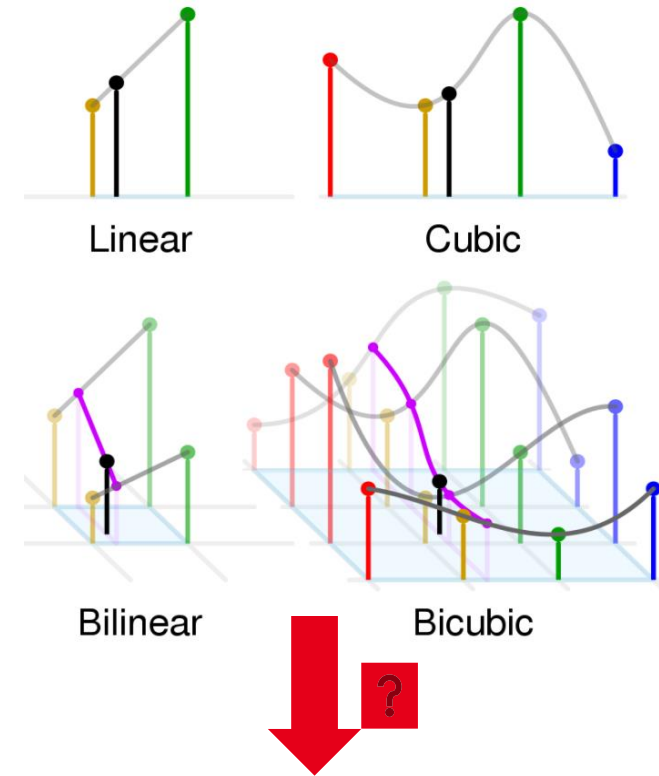
# State of the art





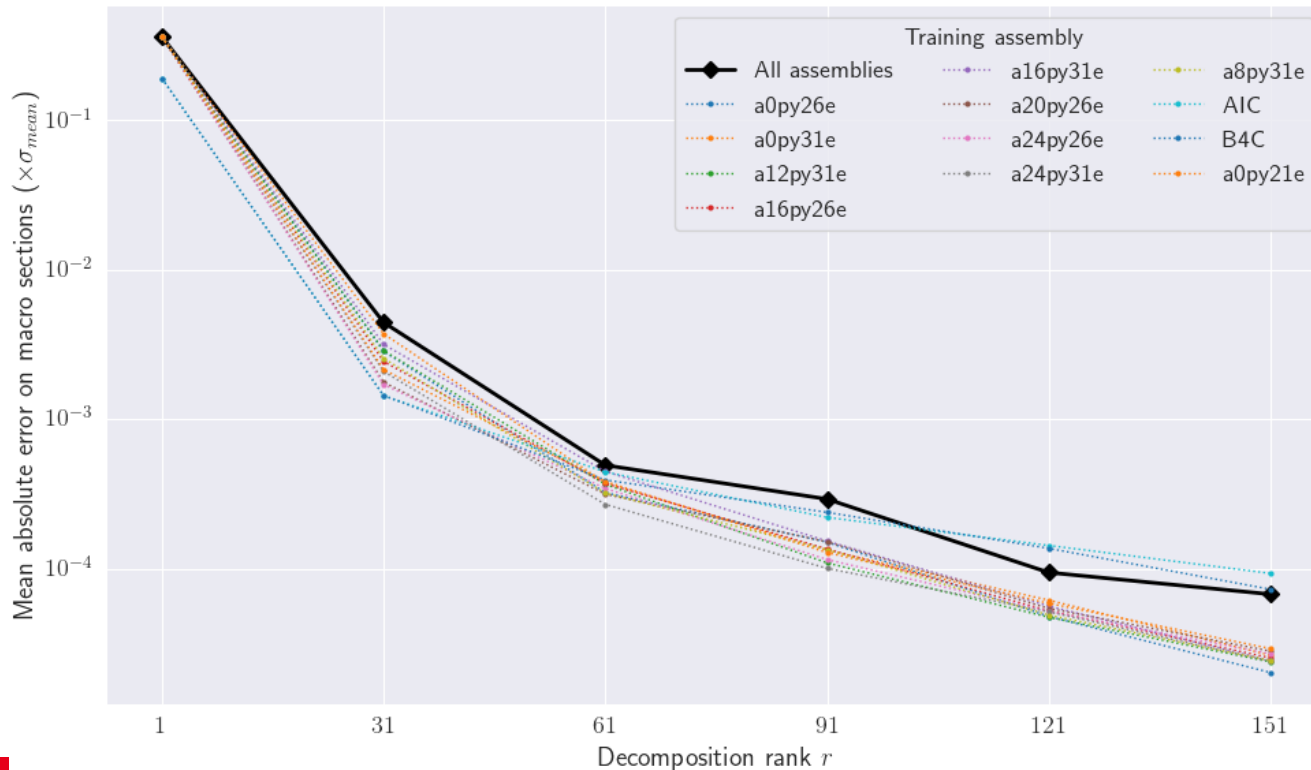
# Current standard

- Most common method today: multilinear interpolation
  - Strengths: simplicity, estimation speed, error control
  - Inconvenients: improvable accuracy, number of coefficients, curse of dimensionality
  - Consequence: cross-sections libraries explode in size (up to several hundred GB)
  - Possible solutions:
    - Get rid of grids
    - Exploit HXS redundancy
- Multitask GPs are natural candidates



# Nuclear data redundancy

- Parallel work: exploiting the redundancy of nuclear data to reduce the size of calculations (typically linear dimension reduction). Reduced Order Models, compression...
- Examples: linear compression of local burnups [Tom1] and HXS [Tom2], PCA of concentration data [Hua], autoencoding of scattering matrices [Whe], multi-output neural network [Sza2]...
- Early thesis work: linear dimension reduction and extrapolation using Empirical Interpolation Method [Tru].

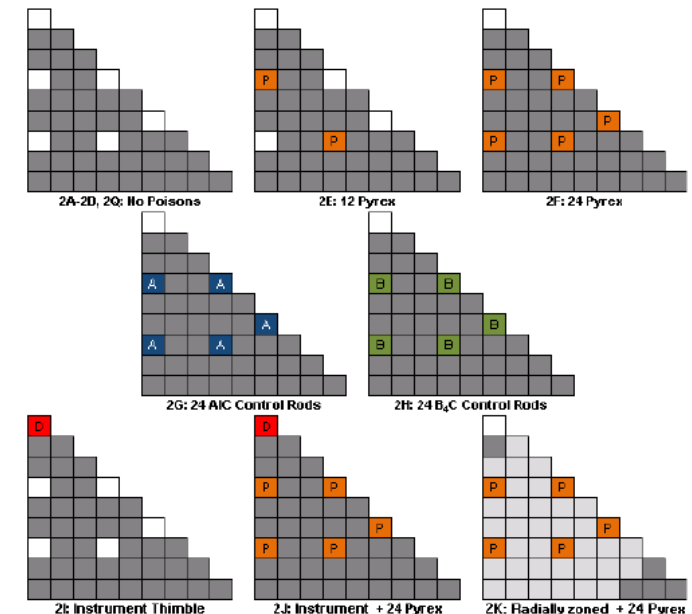
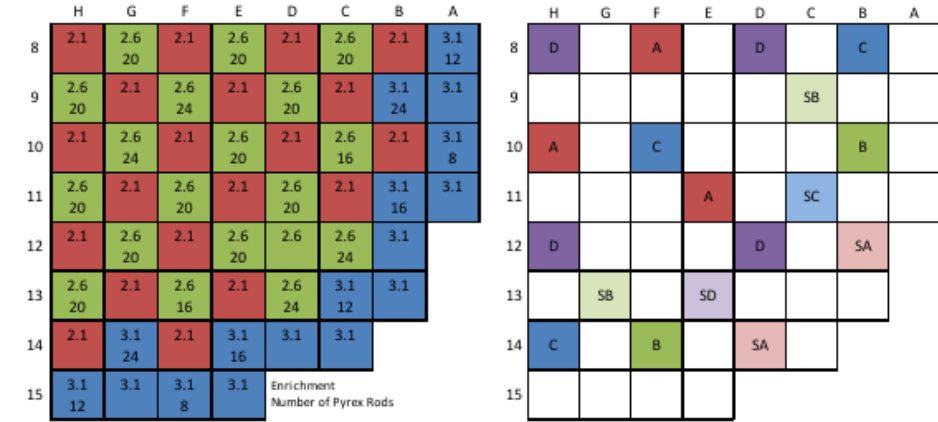


**Source** : An EIM-based compression-extrapolation tool for efficient treatment of homogenized cross-section data, Truffinet, Ammar, Gérard Castaing, Argaud, Bouriquet, Annals of Nuclear Energy, 2023

*Each plot represents the error made when extrapolating - via the EIM algorithm - data from one assembly to all others. We can see that these errors are minor, and that each assembly contains as much useful information as the sum of them (black curve).*

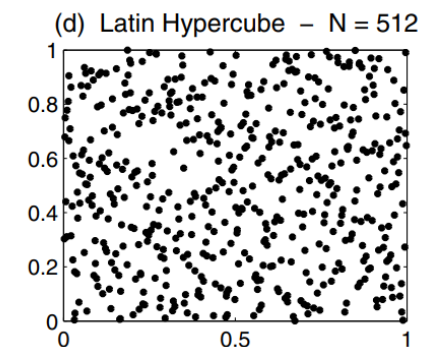
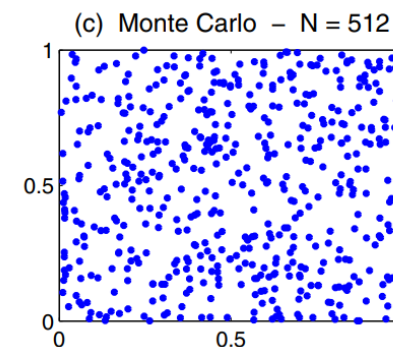
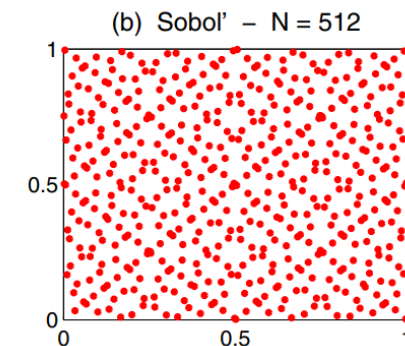
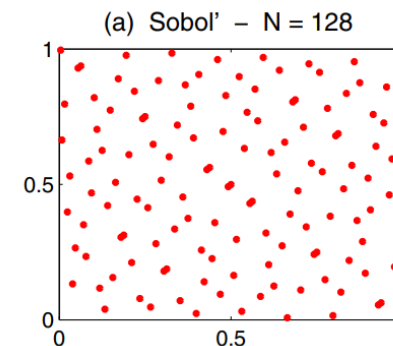
# Experimental test-case

- Data generated with Apollo3® + SIGAR launcher
- 2 assemblies used in this presentation : a «standard» one (2.6% enrichment, no control rods or burnable poisons) and a 3.1% enriched assembly with gadolinium rods
- Variables : burnup  $Bu$ , fuel temperature  $T_f$ , moderator temperature  $T_m$ , boron concentration  $C_b$
- Full grids for comparison with multilinear interpolation
- Several sets generated according to discretization finesse: 2 or 20 energy groups, average per cell or per assembly. From 70 to 100,000 HXS / assembly according to these choices

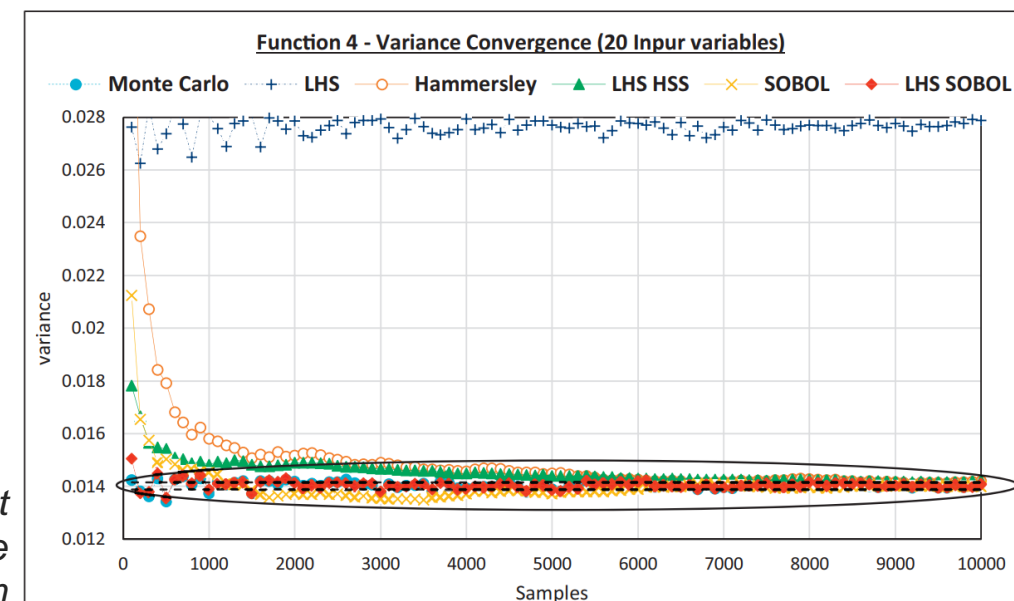


# Low-discrepancy sequences

- Sets of points (deterministic or semi-random) in  $\mathbb{R}^d$  filling space as homogeneously as possible
- Most common methods: Sobol' sequences and Latin Hypercube Sampling (LHS)
- Sobol' sequences generally superior for numerical quadrature and Uncertainty Quantification, tasks close to ours
- Generated point sets:
  - For training, Sobol' sets with 128, 256, 512 and 1024 points (optimal properties when the number of points is a power of 2);
  - For testing, 512-point LHS.



Source : G.Blattman & al, *Quasi random numbers in stochastic finite element analysis*



# Performance metrics

- HXS are spread over 10 orders of magnitude (!!), and for most of them vary little around their mean value → important to **center and reduce data to obtain errors compatible with ML standards**
- **Retained metrics :**
  - « Absolute » error (*ML*) measured on standardized data, in % :  $Err_{abs} = |\tilde{y}_{true} - \tilde{y}_{pred}| * 100$
  - Relative error (*neutronics*) measured on raw data, in pcm :  $Err_{rel} = \frac{|y_{true} - y_{pred}|}{y_{true}} * 10^5 = \frac{|\tilde{y}_{true} - \tilde{y}_{pred}|}{\tilde{y}_{true} + \text{mean}(y)/std(y)} * 10^5$
  - "Compression rate": number of coefficients in the multilinear model / number of coefficients in the model
  - $T_{train}$  : duration of model training ;  $N_{points}$  : number of support points
  - $\alpha$ -CI : proportion of predictions within the predicted 5% confidence interval. **Target value : 95%**
  - Predictive Variance Adequacy :  $PVA = \log \left( \frac{1}{np} \sum_{i=1}^n \sum_{j=1}^p \frac{(y_{ij} - \hat{y}_{ij})^2}{\hat{v}_{ij}} \right)$ . **Target value : 0**

# Modeling choices : grid search



- For all discrete or non-optimizable choices, we rely on a grid search (testing all combinations)
- **Kernel type:** *squared exponential* ( $\psi(r) = e^{-r^2}$ ), *rational quadratic* of learned parameter  $\alpha$  ( $\psi_\alpha(r) = (1 + \frac{r^2}{2\alpha})^{-1}$ ), *Matérn* of parameter  $\nu = 2.5$  ( $\psi_\nu(r) = 2^{1-\nu}(\sqrt{2\nu}r)^\nu K_\nu(\sqrt{2\nu}r)$ ) → Little impact on accuracy
- **Variables combination :** → Huge impact on accuracy (factor 3 on errors)
  - $k_1 = \alpha k_{base}(Bu, C_b) + \beta k_{fuel}(Bu, T_f) + \gamma k_{mod}(Bu, T_m)$  ( $\alpha, \beta, \gamma$  learned)
  - $k_2 = \alpha k_{base}(Bu, T_m, C_b) + \beta k_{fuel}(Bu, T_f)$  ( $\alpha, \beta$  learned)
  - $k_3 = \alpha k_{base}(Bu, T_f, T_m, C_b)$  ( $\alpha$  learned)
  - $k_4 = \alpha k_{base}(Bu, T_f, T_m, C_b) + \beta k_{bu}(Bu)$  ( $\alpha, \beta$  learned)
- **Mean function :** fixed to a constant mean (variable means not yet implemented for projected model)
- Training dynamics : test of several optimization algorithms. Selected one : ADAM with exponential learning rate decay between 2 values to be adjusted. Automatic stop in case of stagnation.  
→ Convergent and monotonic learning, little sensitive to the chosen bounds.

# Results on the « light » test case : 70 outputs

	Mean absolute error (%)	Max absolute error (%)	Mean relative error (pcm)	Max relative error (pcm)	alpha-CI (%)	PVA	Compression rate	$T_{train}$ (s)	$N_{points}$
Projected LMC	1,09	28	108	1948	98	-0.86	93	41	128
Variational LMC	0,91	19	92	1635	99	-1,04	114	44	128
IMC	1,09	28	118	2280	99	-2,14	37	19	128
Single-output GPs	1,16	19	129	2040	99	-0,94	35	178	128
Multilinear interpolation	2,81	21	197	1817	-	-	1	-	4704

- Performance of the variational model depends on the number of induction points, set here at 1.5x this of input points
- Accuracy and error estimation can be further improved by adding task-level polynomial mean functions (not available yet for projected LMC...)
- Models have been tested on a heavier dataset with ~100 000 outputs, making for 380 Mo of data ; training is not fully optimized yet, but runs in < 500s and mostly beats multilinear interpolation !



# 4. Conclusion



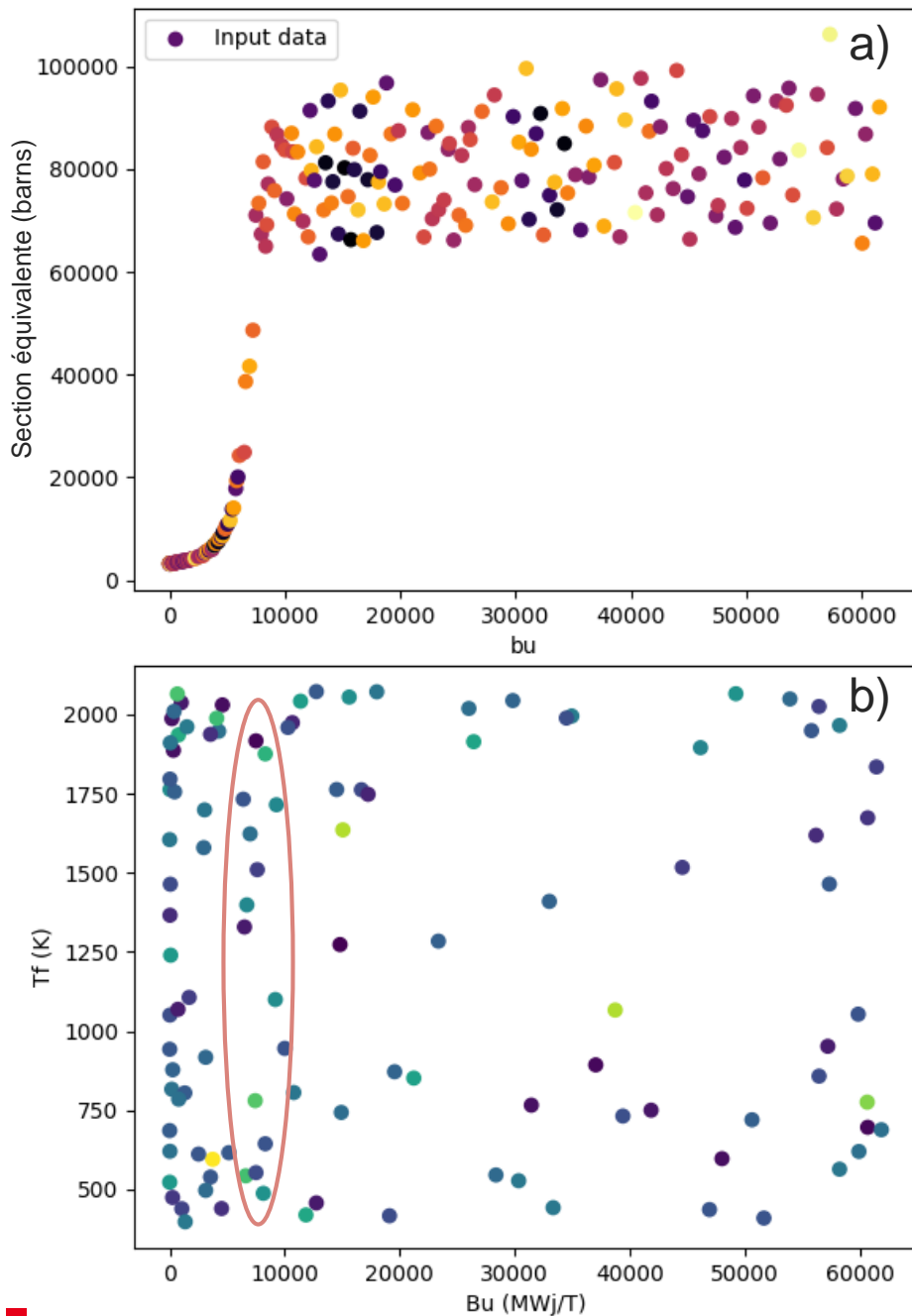
# Conclusion

- **Theoretical aspects :**
  - Better understanding of the structure of the Linear Model of Coregionalization
  - Elaboration of a new multitask GP model, more expressive than the commonly-used IMC, more straightforward and less parametrized than variational ones
  - Authorizes or facilitates some operations on multitask GP, such as fast support updating (online training) or leave-one-out error computation
- **Applicative results :**
  - Improved accuracy with 40x less support → fewer calls to network code
  - Reduction of HXS library size by a factor of **≈100**
  - Reliable estimation of reconstruction uncertainty, easy propagation of uncertainties
  - Accepts any support (useful in case of missing or corrupted data, for flexible support refinement or reduction, etc.)
  - **Possibility of adaptive support construction**





# Adaptive sampling



Section équivalente (barns)

- The GP variance estimator can be used to find areas in space where the model is unsure of its predictions → good candidates for placing new points there.
- Numerous refinements in the literature, e.g. based on LOO-CV error (Le Gratiet & Cannamela, 2012)
- Significant impact on maximum errors ( $\div 4$  for small number of points) Much smaller impact on mean errors
- So far, this has been done with single-outputs GP, summing scores of all outputs with equal weights. **Projected LMC will allow to work at the level of latent functions, with natural weighting → better results ?**

**Fig.4:**

- a) Thermal absorption section of Gd157
- b) Aspec of supports obtained by LOO-CV adaptive sampling on the gadolinium dataset (113 points)



***Thanks for your  
attention !***

**Olivier TRUFFINET**

**olivier.truffinet@cea.fr**

**06 78 07 69 06**

# Références



- [Bon] *Generic Inference in Latent Gaussian Process Models*, Bonilla, Krauth & Dezfouli, Journal of Machine Learning Research, 2019
- [Bot] *Polynomial interpolation of few-group neutron cross sections on sparse grids*, Botes & Bokov, Annals of Nuclear Energy, 2014
- [Bru] *Scalable Exact Inference in Multi-Output Gaussian Processes*, Bruinsma & al, Proceedings of Machine Learning Research, 2020
- [Hua] *PWR pin-homogenized cross-sections analysis using big-data technology*, Hua, Li & Wang, Progress in Nuclear Energy, 2020
- [Ras] *Gaussian Processes for Machine Learning*, Rasmussen & Williams, The MIT Press, 2006
- [Sza1] *Reconstruction of few-group homogenized cross section by kernel method and active learning*, Szames, Tomatis, Ammar & Martinez, M&C 2019
- [Sza2] *Few group cross section modeling by machine learning for nuclear reactor*, Neural and Evolutionary Computing. Université Paris-Saclay, 2020
- [Tom1] *Compression of 3D pin-by-pin burnup data*, Tomatis & Dall'Osso, Annals of Nuclear Energy, 2020
- [Tom2] *A multivariate representation of compressed pin-by-pin cross-sections*, Tomatis, EPJ Nuclear Sci. Technol, 2021
- [Tru] *An EIM-based compression-extrapolation tool for efficient treatment of homogenized cross-section data*, Truffinet, Ammar, Gérard Castaing, Argaud & Bouriquet, Annals of Nuclear Energy, 2023
- [Wat] *Improved cross-section modeling methodology for coupled three-dimensional transient simulations*, Watson, Ivanov, Annals of Nuclear Energy, 2002
- [Whe] *Data Reduction in Deterministic Neutron Transport Calculations Using Machine Learning*, Whewell & McClarell, Annals of Nuclear Energy, 2022
- [Zim] *Building neutron cross-section dependencies for few-group reactor calculations using stepwise regression*, Zimin & Semenov, Annals of Nuclear Energy, 2005