

Hybrid prognostics using simulation codes and statistical models – Application to the study of steam generators clogging

Edgar JABER - PhD Defense - 09/02/2026

Examinators:

J. GARNIER
P. M. CONGEDO
C. SCHILLINGS
E. VAZQUEZ

Reviewers:

E. CHATZI
B. LAURENT-BONNEAU

Supervisors:

M. MOUGEOT
D. LUCOR
E. REMY
V. CHABRIDON



Outline

1. Introduction
2. The physical clogging simulation model
 - 2.1 The physical model
 - 2.2 THYC-Puffer-DEPO computational model
 - 2.3 Expert-informed UQ on TPD
3. Non-intrusive surrogate modeling
 - 3.1 General idea
 - 3.2 Gaussian processes
 - 3.3 GP validation with conformal prediction
4. Bayesian fusion of heterogeneous data
 - 4.1 Offline data assimilation
 - 4.2 The BMU algorithm
 - 4.3 Ensemble Kalman smoothing
5. Conclusion
6. Appendix

Clogging of steam generators (SGs)

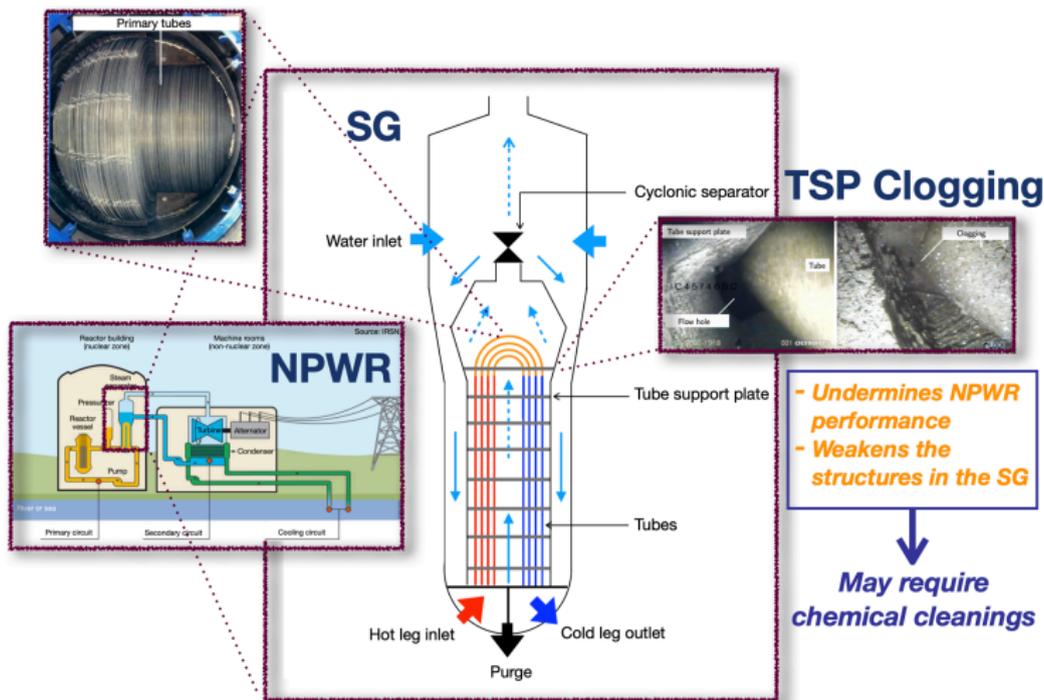


Figure 1: Nuclear pressurized water reactor (NPWR) scheme, an SG and example of video examination (TVE) of a tube support plate (TSP) during a NPWR outage (© IRSN, EDF)

Notions of prognostics

- ▶ EDF Engineers may have to schedule SG chemical cleanings
- ▶ Degradation level of the system ($t \mapsto \Delta(t)$) → predict the remaining useful life (RUL) [Vachtsevanos et al., 2006] for a fixed threshold $\Delta_* \in \mathbb{R}_+$:

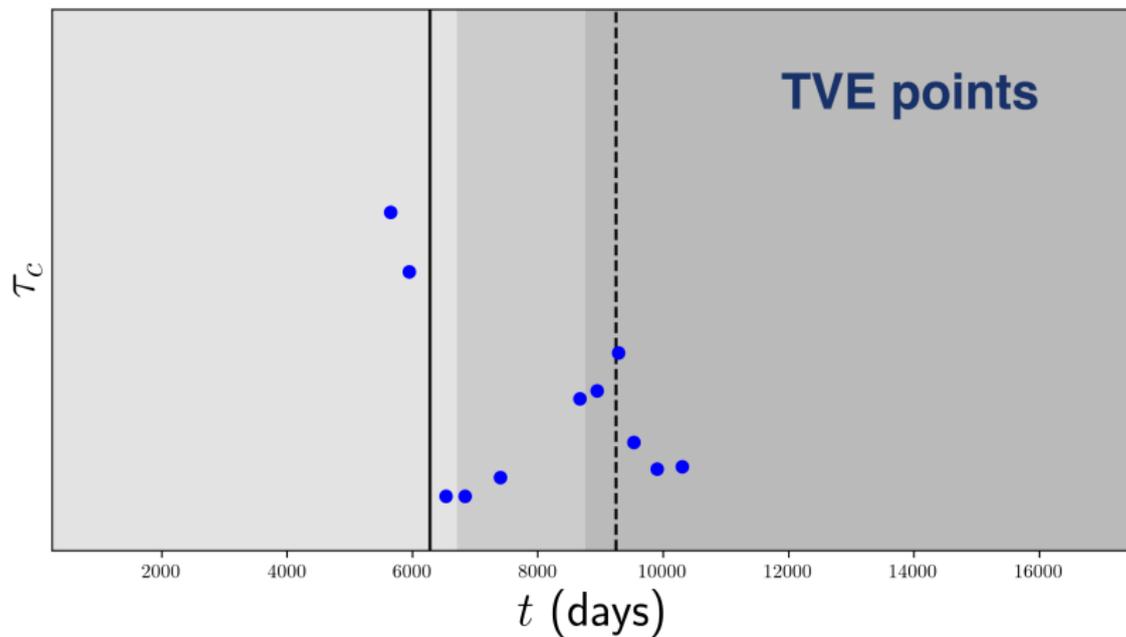
$$\text{RUL}(\Delta_*) = \underset{t > t_P}{\operatorname{argmin}} \{ \Delta(t) \geq \Delta_* \} \quad (1)$$

where t_P is the present time.

- Relies usually on physics-based simulation codes, and/or data driven methods
- ▶ RUL prediction with each individual tool could lack robustness, especially for complex operating systems
- ▶ For clogging, $\Delta = \tau_c$

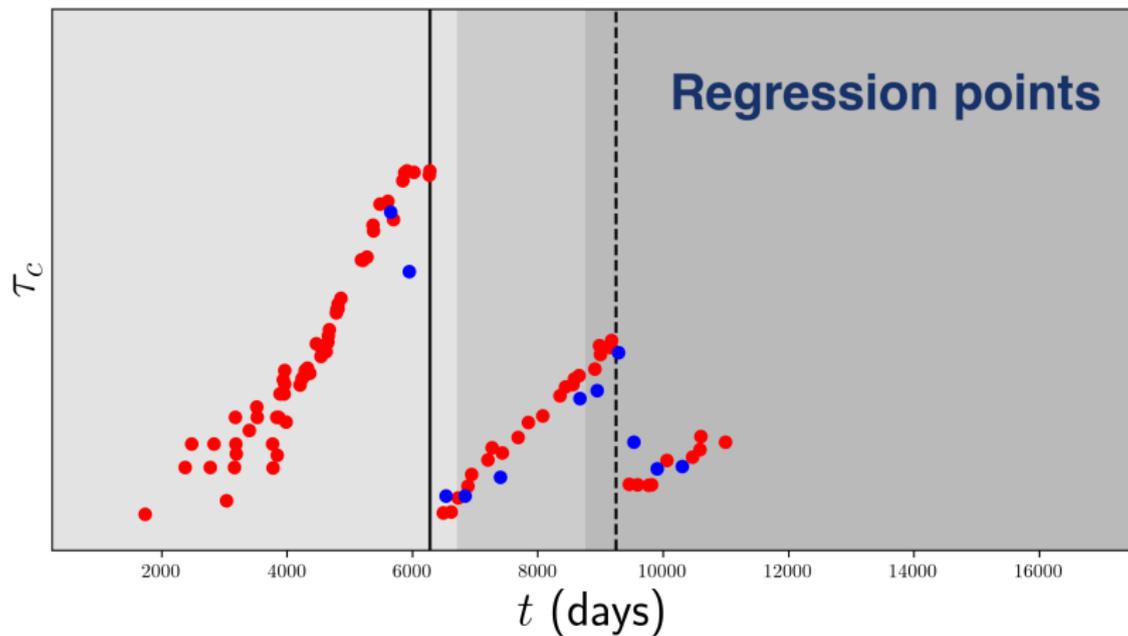
Idea: *build a hybrid framework for evaluating the SG clogging RUL using the physics-based model and the data-driven models*

Available clogging-related informations



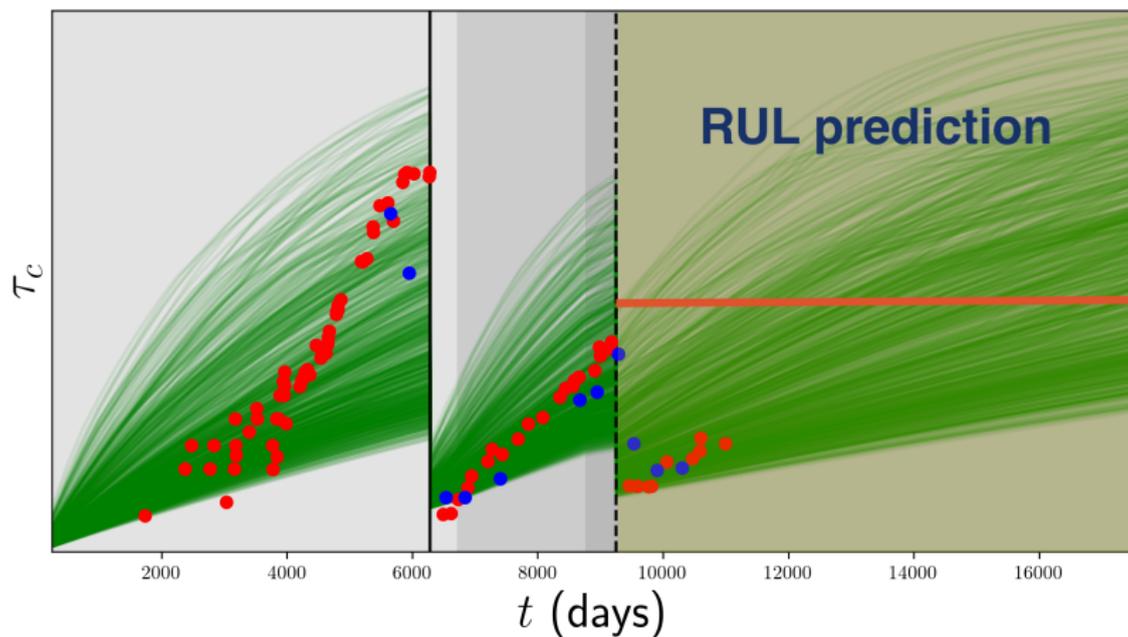
- *Scarce* televised video examinations (TVEs) monitoring data

Available clogging-related informations



- ▶ Data-driven regression algorithms [Pincirolì et al., 2021] based on the TVEs

Available clogging-related informations



How to make use of the available knowledge for achieving reliable RUL predictions?

Digital twin for a SG

- ▶ Use of a digital twin (DT) [AIAA, 2021] platform for predicting the state of health of a steam generator
- ▶ Enabling informed decision-making using all available asset information

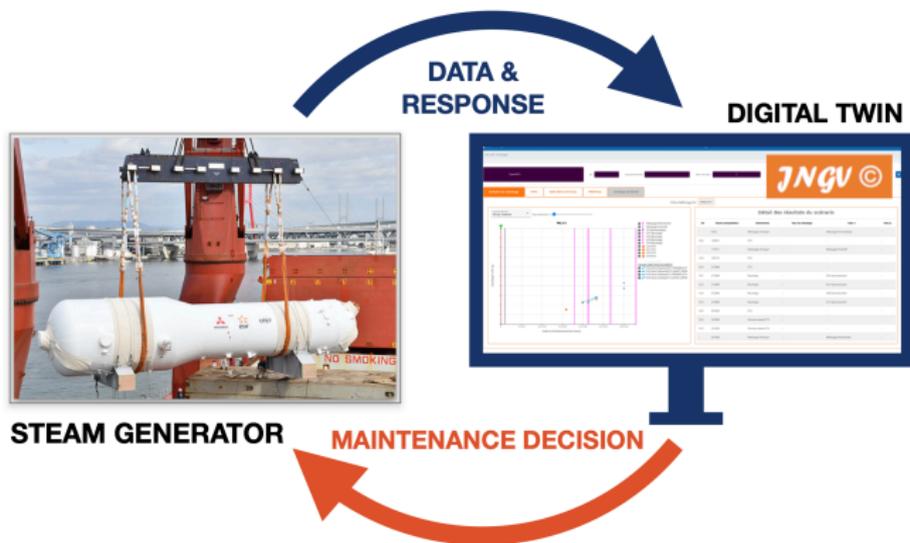


Figure 2: DT architecture for a SG [Deri et al., 2021] ("Jumeau Numérique Générateur de Vapeur" - JNGV)

Hybrid framework for clogging prognostics

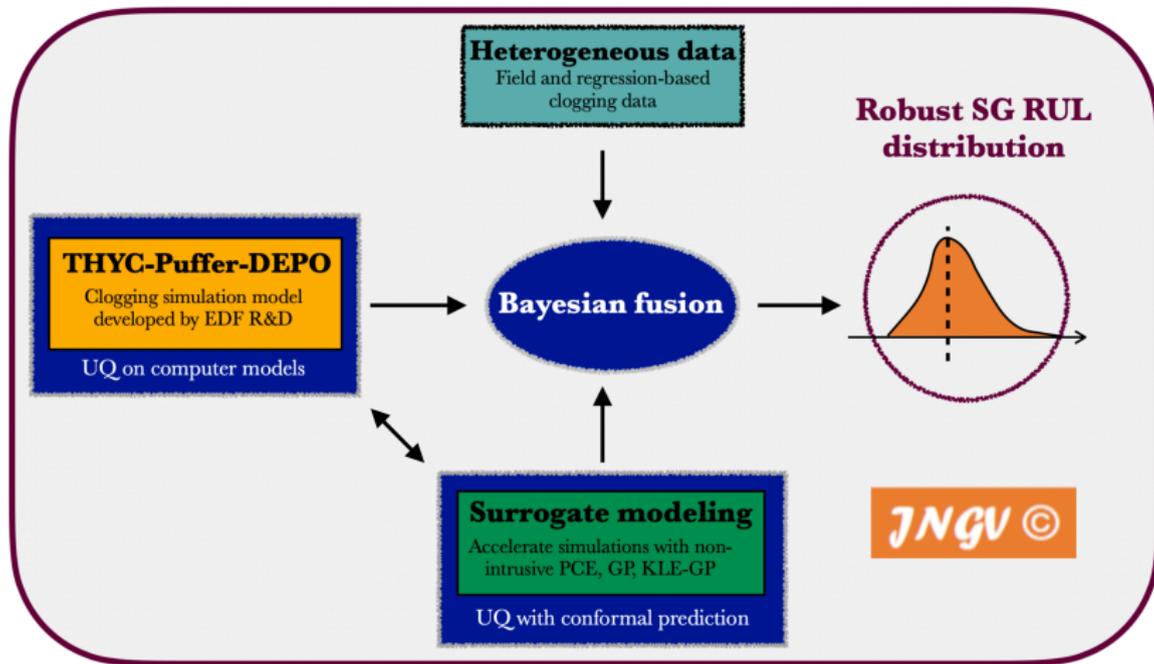


Figure 3: Hybrid framework for a specific SG clogging prognostics

DT integration for clogging prognostics

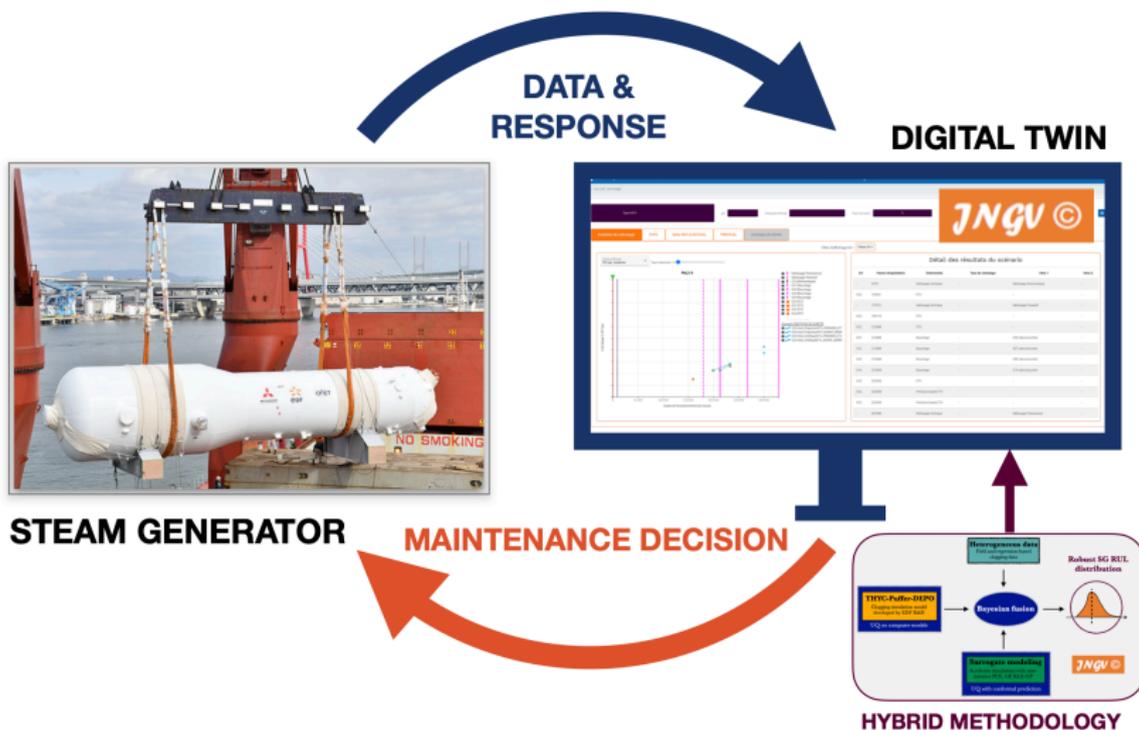


Figure 4: Integration of the hybrid framework in the JNGV

Summary of the introduction

- ▶ Available tools: clogging simulation code, field data, regression data, EDF expert and field knowledge
- ▶ Objective: estimate as accurately as possible the probabilistic RUL for SG clogging by developing a hybrid framework involving all the available tools
- ▶ Challenges: complex physics, heavy simulation code, inherent uncertainty, operational variability, field data scarce both in time and space, different degrees of data fidelity, costly simulation code

*Proposed solution → **Bayesian fusion methodology based on the clogging simulation code, its surrogate models and the heterogeneous data***

Outline

1. Introduction
2. The physical clogging simulation model
 - 2.1 The physical model
 - 2.2 THYC-Puffer-DEPO computational model
 - 2.3 Expert-informed UQ on TPD
3. Non-intrusive surrogate modeling
 - 3.1 General idea
 - 3.2 Gaussian processes
 - 3.3 GP validation with conformal prediction
4. Bayesian fusion of heterogeneous data
 - 4.1 Offline data assimilation
 - 4.2 The BMU algorithm
 - 4.3 Ensemble Kalman smoothing
5. Conclusion
6. Appendix

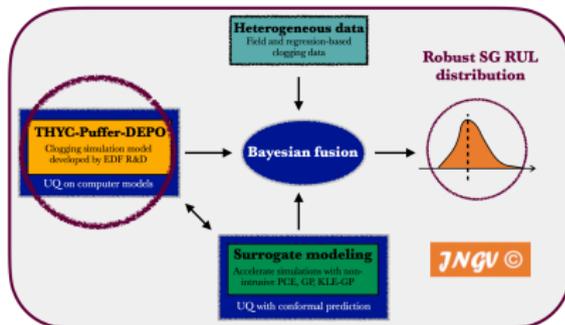
Hybrid framework for the JNGV

Objectives:

- ▶ Evaluate uncertainty on EDF clogging simulation code
- ▶ Extend the preliminary work of [Lefebvre et al., 2023]
- ▶ Obtain a first RUL distribution based on the physical tool

Methodology:

- ▶ *State of the art:*
 - ➔ Uncertainty analysis (UA) of industrial simulation codes [De Rocquigny et al., 2008]
 - ➔ High-performance computing
 - ➔ Surrogate models, sensitivity analysis
- ▶ *Strategy:* Couple the EDF code with the OpenTURNS UQ Python module and compute different sensitivity analysis indices



Clogging physical model

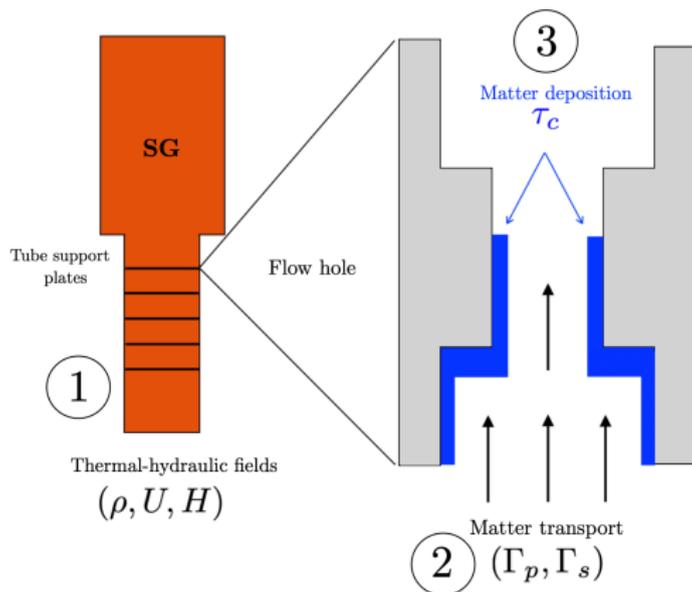


Figure 5: Clogging physical model

- ▶ Clogging results from two main mechanisms → *vena contracta* and *flashing* [Prusek et al., 2013]
- ▶ Long-term clogging model [Feng et al., 2023] → must change stationary thermohydraulics + compute chemical conditioning

Clogging computational model

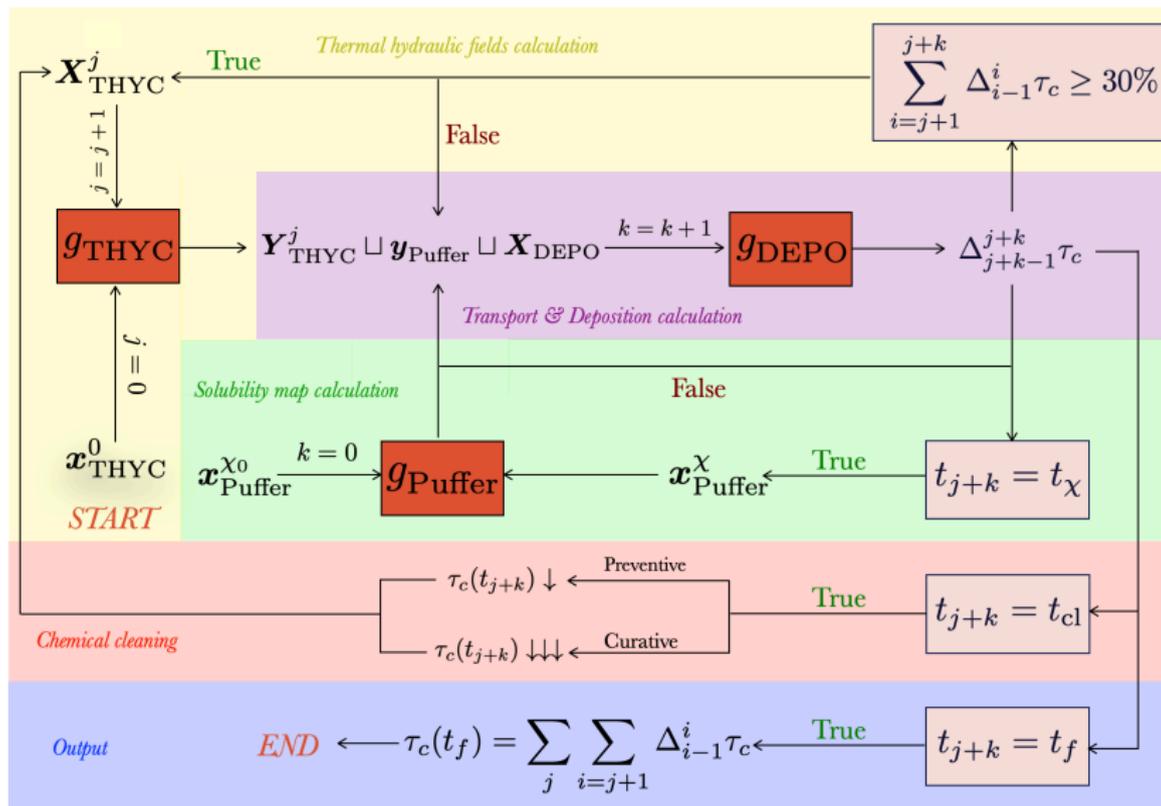
- ▶ THYC-Puffer-DEPO (TPD) [Feng et al., 2023] is the chaining of 3 codes:
 - THYC [David, 1999] is based on a finite-volume numerical scheme for the two-phase conservation equations
 - Puffer is an in-house chemical code allowing to compute the solubility of iron oxides as a function of pH
 - DEPO [Prusek et al., 2013] is the deposit module, solving the transport and clogging equations with iterative finite-differences schemes methods

- ▶ Allows to simulate SG clogging on entire lifespan of the asset integrating past chemical cleanings and predicting future τ_c states → Unitary call is $\sim 5h$ on HPC infrastructure

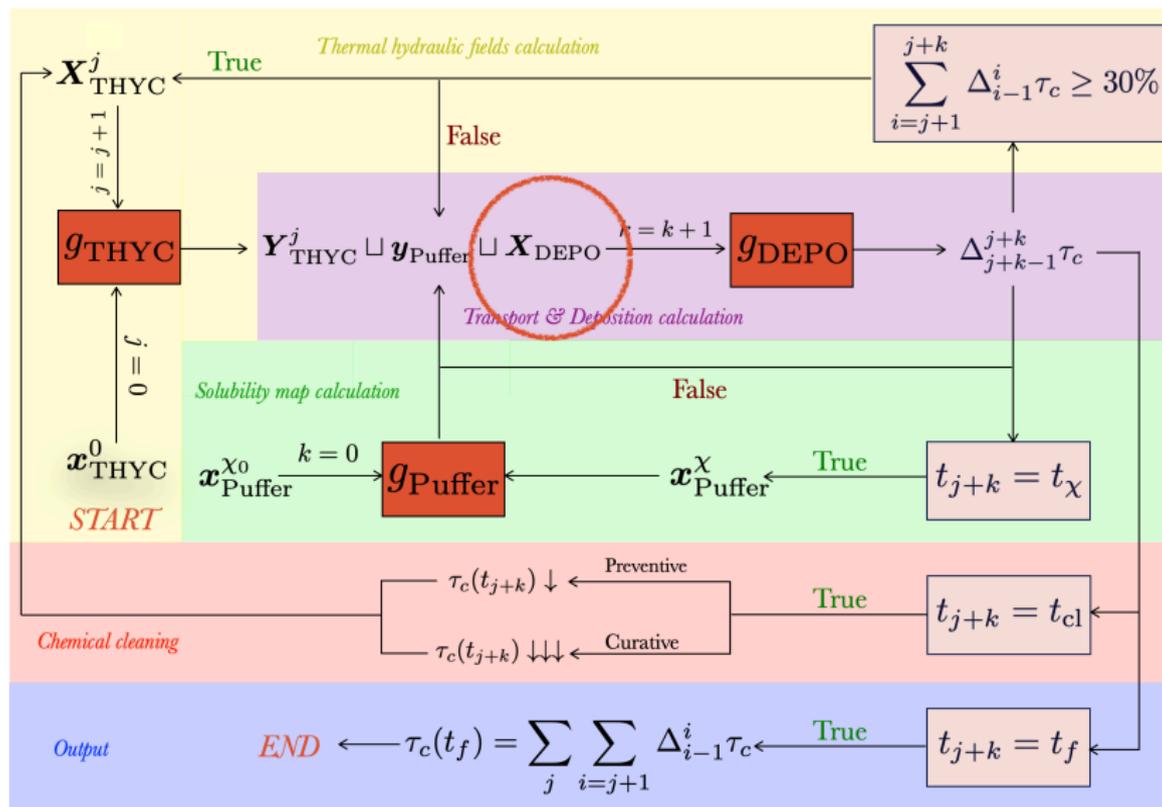
We denote it as a function $g_{\text{TPD}} =: g : \mathcal{X} \subset \mathbb{R}^d \rightarrow \mathbb{R}^N$ with $N \simeq 10^2$, TPD simulates a trajectory of τ_c for each input \mathbf{X} ,

$$g(\mathbf{X}) = (g(t_1, \mathbf{X}), \dots, g(t_N, \mathbf{X})) \in \mathbb{R}^N \quad (2)$$

TPD simulation chain



TPD simulation chain



Initial Design of Experiments (DoE)

- ▶ $d = 7$ input variables of the deposit module g_{DEPO} with uncertainty:

$$\mathbf{X}_{\text{DEPO}} =: \mathbf{X} = (\alpha, \beta, \epsilon_e, \epsilon_c, d_p, \Gamma_p(0), a_v) \sim p_{\mathbf{X}} = \otimes_{i=1}^d p_{X_i},$$

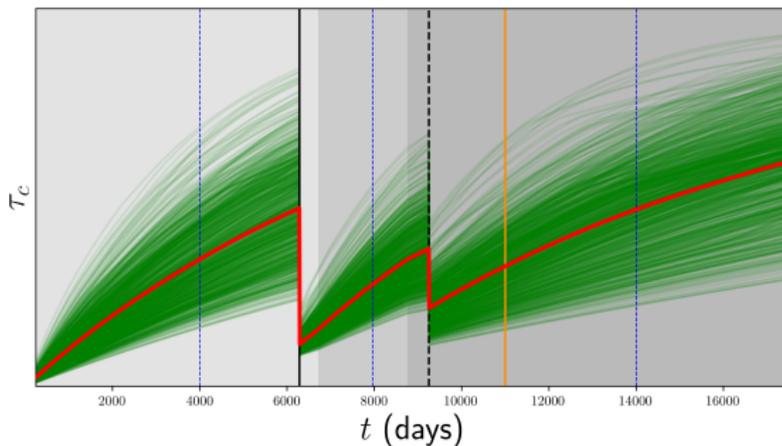
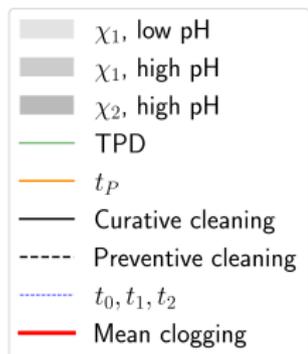
with supports and distributions provided and validated [Lefebvre et al., 2023]

- ▶ $n = 1000$ i.i.d. Monte Carlo samples drawn according to the distributions below, $\text{DoE}_g^{p_{\mathbf{X}}} = \{(\mathbf{X}^{(i)}, g(\mathbf{X}^{(i)}))\}_{i=1}^n$

Variable	Physical meaning	Distribution
α	First empirical correlation parameter	$\mathcal{N}(101.6, 4.0)$
β	Second empirical correlation parameter	$\mathcal{N}(0.0233, 0.0005)$
ϵ_e	Porosity of fouling deposits	$\mathcal{T}(0.2, 0.3, 0.5)$
ϵ_c	Porosity of clogging deposits	$\mathcal{T}(0.01, 0.05, 0.3)$
d_p	Diameter of iron oxide particles (m)	$\mathcal{T}(0.5, 5.0, 10.0) \times 10^{-6}$
$\Gamma_p(0)$	Initial solid mass fraction	$\mathcal{T}(1.0, 4.5, 8.0) \times 10^{-9}$
a_v	Calibration parameter	$\mathcal{T}(0.1, 7.8, 12) \times 10^{-4}$

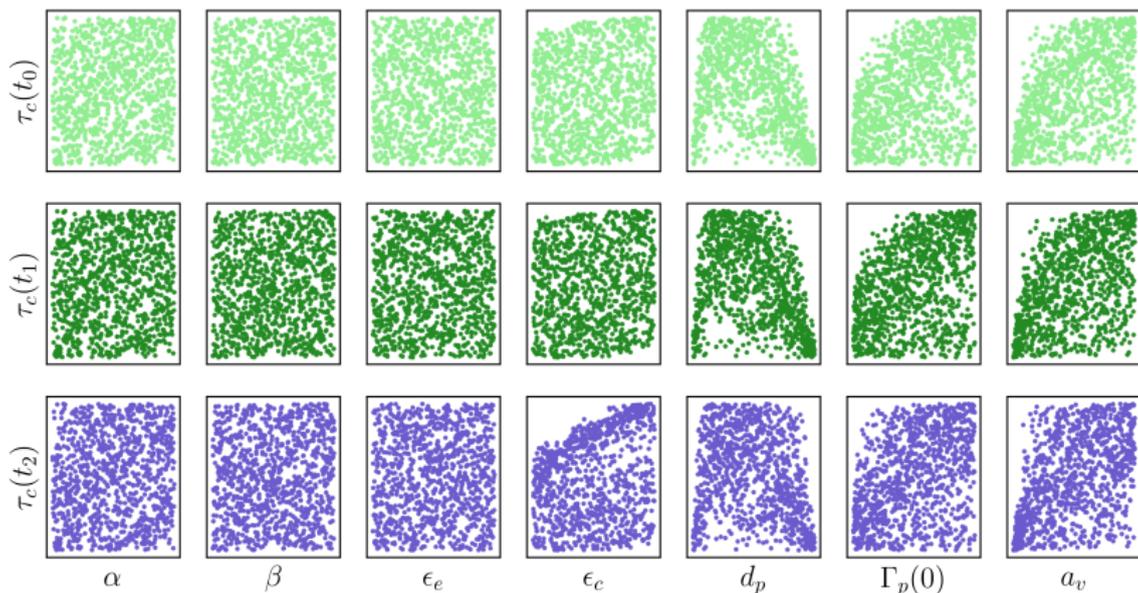
Table 1: Uncertain input variables, physical meaning and distribution

Output results of the uncertainty propagation



- High-level of uncertainty observed throughout the simulation times + very high uncertainty on RUL window

Scatter-plots in rank space



- ▶ Linear correlations of $\Gamma_p(0)$ and a_v during all simulation times
- ▶ Non-linear correlations of d_p → confirms findings in [Lefebvre et al., 2023], novel finding: non-linear correlation of ϵ_c that evolves in time

Given-data sensitivity analysis with HSIC

- ▶ Hilbert-Schmidt Independence Criterion (HSIC) [Gretton et al., 2005], kernel method
 - ➔ Evaluates sensitivity of a single-input in a given-data context
- ▶ Theoretical result for all $i \in \{1, \dots, d\}$, $k \in \{1, \dots, N\}$:

$$\text{HSIC}(X_i, g(t_k, \mathbf{X})) = 0 \iff X_i \text{ and } g(t_k, \mathbf{X}) \text{ are independent} \quad (3)$$

- ▶ Index disposes of U-stat and V-stat estimators [Da Veiga, 2015]
 - + hypothesis testing with corresponding p -value
- ▶ Normalized R_{HSIC}^2 index is better suited for interpretation:

$$R_{\text{HSIC}}^2(X_i, g(t_k, \mathbf{X})) = \frac{\text{HSIC}(X_i, g(t_k, \mathbf{X}))}{\sqrt{(\text{HSIC}(X_i, X_i)\text{HSIC}(g(t_k, \mathbf{X}), g(t_k, \mathbf{X})))}} \in [0, 1]$$

Results

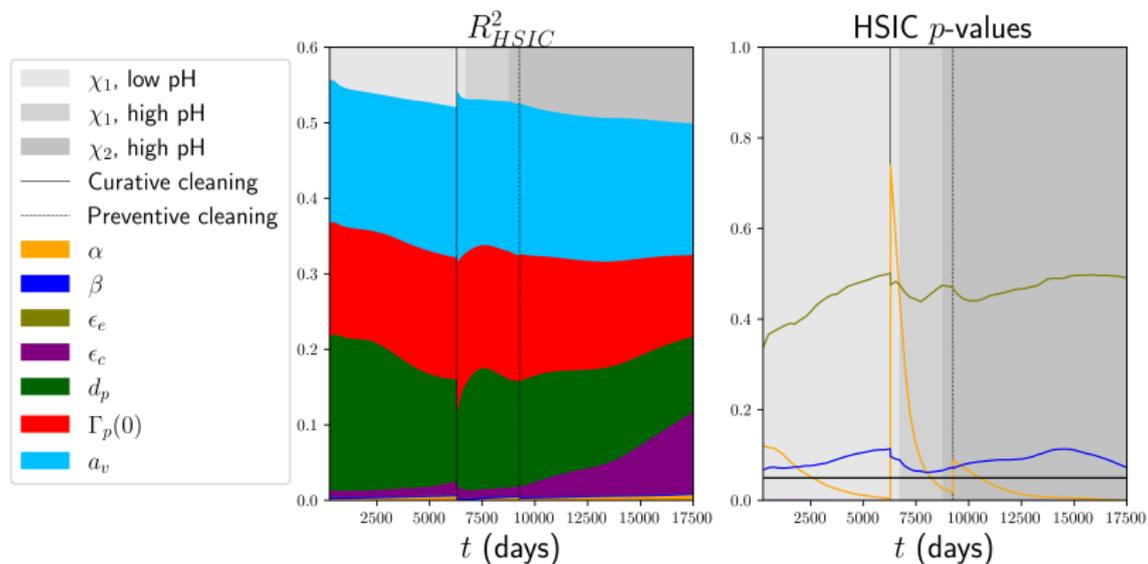


Figure 6: Time variation of the estimated R^2_{HSIC} indices and associated p -values

Summary of part 2.

Take-home messages:

- Uncertainty analysis on TPD with expert-informed p_X reveals a high level of uncertainty in prognostics
- Sensitivity analysis results (both HSIC and Sobol') are robust with the preliminary work [Lefebvre et al., 2023]
- New findings related to the influence of clogging porosity ϵ_c in high pH paves the way for future clogging modeling

Contributions:

- ▶ Paper [Jaber et al., 2025b] published in the **International Journal of Uncertainty Quantification**
- ▶ **Github repository** with scripts based on the **OpenTURNS** Python module [Baudin et al., 2017]
- ▶ Poster presentation at **MASCOT-NUM 2023**

Outline

1. Introduction
2. The physical clogging simulation model
 - 2.1 The physical model
 - 2.2 THYC-Puffer-DEPO computational model
 - 2.3 Expert-informed UQ on TPD
- 3. Non-intrusive surrogate modeling**
 - 3.1 General idea**
 - 3.2 Gaussian processes**
 - 3.3 GP validation with conformal prediction**
4. Bayesian fusion of heterogeneous data
 - 4.1 Offline data assimilation
 - 4.2 The BMU algorithm
 - 4.3 Ensemble Kalman smoothing
5. Conclusion
6. Appendix

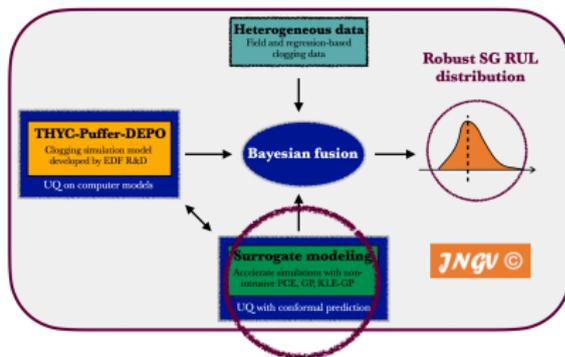
Hybrid framework for the JNGV

Objectives:

- ▶ Obtain fast evaluations of TPD trajectories for downstream tasks
- ▶ Ensure the validity of these evaluations with proven methodologies

Methodology:

- ▶ *State of the art:*
 - ➔ Vector polynomial chaos expansions, Gaussian processes, Reduced order models
 - ➔ Predictivity coefficients, Bayesian credibility intervals [Acharki et al., 2023]
- ▶ *Strategy:* Build a variety of different surrogate models, develop a validation methodology for Gaussian processes based on conformal prediction



Surrogate modeling

- ▶ General idea is to create a model $\hat{g} \approx g$ that enables faster evaluations of the computer simulation model [Sullivan, 2015]
 - ➔ Purpose is to sample rapidly from the *solution manifold*:

$$\mathcal{M}_{\mathbf{X}} = \{g(\mathbf{X}), \mathbf{X} \sim p_{\mathbf{X}} \in \mathcal{P}(\mathbf{X})\} \quad (4)$$

- ▶ Depending on the nature of the computer code g and its outputs, one can adopt an:
 - ➔ *Intrusive* approach: modifying the source-code g (e.g. Reduced basis methods) [Le Maître and Knio, 2010] ➔ approximation guarantees for many families of problems
 - ➔ *Non-intrusive* approach [De Rocquigny et al., 2008]: using methods similar to supervised learning techniques ➔ more useful in industry where there are legacy codes
- ▶ Problem specific and task-oriented ➔ method should require validity in well chosen scenarios, specific validation methodologies have to be developed

Stochastic time-collocation DoE

- ▶ Augment N to $D \gg N$ post-simulations with the help of linear interpolators $\mathcal{G} : \mathbb{R}^N \rightarrow \mathbb{R}^D$ (clogging is strictly increasing) to get:

$$\mathcal{G}(g(\mathbf{X})) = \left(\sum_{k=1}^N g(t_k, \mathbf{X}) \varphi_k(t_1), \dots, \sum_{k=1}^N g(t_k, \mathbf{X}) \varphi_k(t_D) \right) =: Z(\mathbf{X}) \in \mathbb{R}^D, \quad (5)$$

where φ_k are linear interpolators

- ▶ Choose a uniform probability density on the time-indices and draw independent samples:

$$(t, \mathbf{X}) \sim \mathcal{U}\{1, \dots, D\} \otimes p_{\mathbf{X}}$$

- ▶ Apply to all $i \in \{1, \dots, n\}$, the following strategy:

$$(t^{(i)}, \mathbf{X}^{(i)}) \xrightarrow[\text{interpolation}]{\mathcal{L}} (\mathbf{X}^{(i)}, \mathcal{G}(g(\mathbf{X}^{(i)}))) \xrightarrow[\text{projection}]{\text{pr}_i} ((t^{(i)}, \mathbf{X}^{(i)}), \text{pr}_i \circ Z(\mathbf{X}^{(i)}))$$

Gaussian Process (GP) surrogate

- ▶ DoE with a *scalar output* $\mathcal{D}_n = \{((t^{(i)}, \mathbf{X}^{(i)}), \text{pr}_i \circ Z(\mathbf{X}^{(i)}))\}_{i=1}^n$, construct a regular scalar GP [Rasmussen and Williams, 2006]
 → split $\mathcal{D}_{n_{\text{train}}} \cup \mathcal{D}_{n_{\text{test}}}$

- ▶ We choose a prior $\mathcal{G} \sim \mathcal{GP}(m_\beta, \gamma_\varphi)$ and compute and optimize the conditioned $\tilde{\mathcal{G}} = \mathcal{G} \mid \mathcal{D}_{n_{\text{train}}} \sim \mathcal{GP}(m_{\beta^*}, \gamma_{\varphi^*})$ using maximum likelihood estimator for the prior hyperparameters $\theta = (\beta, \varphi)$
- ▶ GP metamodel of TPD is then estimated by parametrizing the first marginal of the GP of Z , to the time indices of TPD, in other words:

$$\hat{g}(\mathbf{X}) := (m_{\beta^*}(t_1, \mathbf{X}), \dots, m_{\beta^*}(t_N, \mathbf{X})) \in \mathbb{R}^N \quad (6)$$

- ▶ Validation method: compute the predictivity coefficient on $\mathcal{D}_{\text{test}}$ at each time step $Q^2(t_k)$, and then average it:

$$Q^2(t_k) = 1 - \sum_{j=1}^{n_{\text{test}}} \frac{|g(t_k, \mathbf{X}^{(j)}) - m_{\beta^*}(t_k, \mathbf{X}^{(j)})|^2}{\widehat{\text{Var}}(g(t_k, \cdot))}, \quad \overline{Q^2} = \frac{1}{N} \sum_{k=1}^N Q^2(t_k) \quad (7)$$

Results for different priors

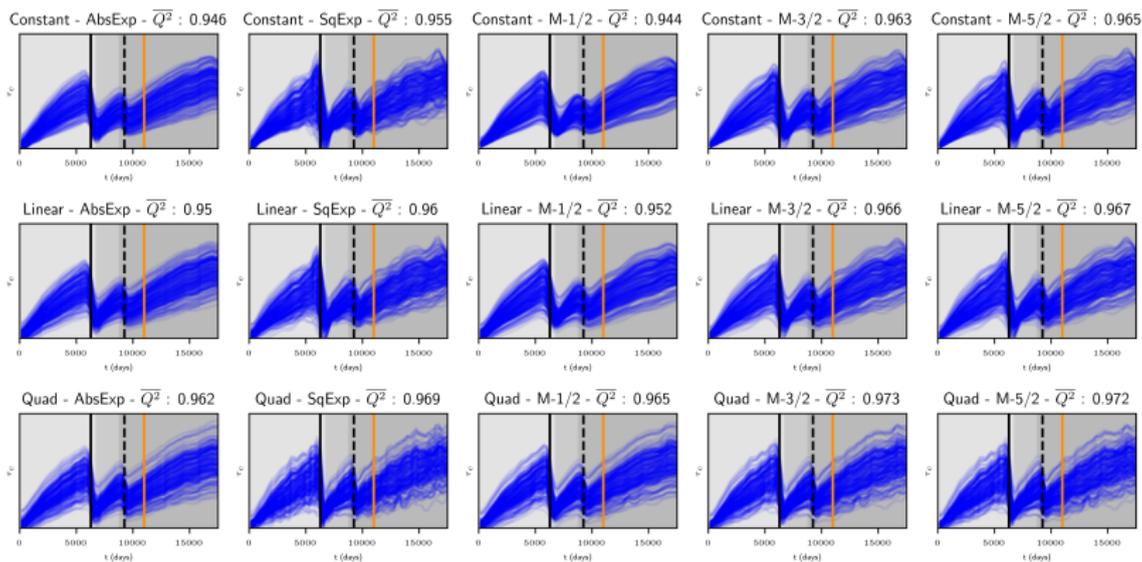


Figure 7: GP surrogate models of TPD with time stochastic collocation method and with different sets of priors

- All surrogate models have a very competitive predictive coefficient $\overline{Q^2} \sim 0.9$

→ *How to differentiate between them, which one is more robust?
For which application?*

Conformal Prediction (CP)

- ▶ CP paradigm [Vovk et al., 2005] is a generic, model-agnostic theory allowing to build *prediction sets* for machine learning models like \hat{g} with frequentist coverage guarantees

Definition

[Vovk et al., 2005] Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space and $\mathcal{Z} := \mathcal{X} \times \mathcal{Y}$. Let $n \in \mathbb{N}$ and $\mathcal{D}_n = \{\mathbf{Z}^{(1)}, \dots, \mathbf{Z}^{(n)}\} \in \mathcal{Z}^n$ a training sample. For $\alpha \in (0, 1)$, a conformal predictor of coverage $1 - \alpha$ is any measurable function of the form:

$$\begin{aligned} \mathcal{C}_\alpha: \mathcal{Z}^n \times \mathcal{X} &\rightarrow 2^{\mathcal{Y}} \\ (\mathcal{D}_n, \mathbf{X}) &\mapsto \mathcal{C}_{n,\alpha}(\mathbf{X}), \end{aligned} \tag{8}$$

s.t. for any new couple of points $\mathbf{Z}^{(n+1)} = (\mathbf{X}^{(n+1)}, Y^{(n+1)}) \in \mathcal{Z}$ (marginal coverage property):

$$\mathbb{P} \left(Y^{(n+1)} \in \mathcal{C}_{n,\alpha}(\mathbf{X}^{(n+1)}) \right) \geq 1 - \alpha \tag{9}$$

Jackknife+/minmax interval estimators

- ▶ There are three types of conformal estimators [Angelopoulos and Bates, 2023] → full, split and cross → last one is adapted for *small dataset*, i.e $n \propto 10^2$
- ▶ \mathcal{D}_n an i.i.d. design of experiments, with scalar outputs
- ▶ \hat{g} a surrogate model trained on \mathcal{D}_n , \hat{g}_{-i} Leave-One-Out (LOO) surrogate model trained on $\mathcal{D}_n \setminus \{i\text{-th data-point}\}$
- ▶ With empirical quantiles of LOO residuals → can build interval estimators: Jackknife+ $\hat{C}_{n,\alpha}^+$ and Jackknife-minmax $\hat{C}_{n,\alpha}^{\text{J-minmax}}$ [Barber et al., 2021]:

Estimators	$\hat{C}_{n,\alpha}^+$	$\hat{C}_{n,\alpha}^{\text{J-minmax}}$
Marginal coverage	$\alpha \in (0, 1/2)$	$\alpha \in (0, 1)$
+	Cross-validation method Fast to compute	Coverage property
-	Constant width	Width is too large

Table 2: Pros and cons of main cross-conformal estimators

GP credibility intervals

- ▶ Credibility intervals of the posterior $\tilde{\mathcal{G}}$, for any new point $\mathbf{X}^{(n+1)} \in \mathcal{X}$, $\alpha \in (0, 1)$ and all $k \in \{1, \dots, N\}$:

$$\mathcal{CR}_\alpha(t_k, \mathbf{X}^{(n+1)}) = \left[m_{\beta^*}(t_k, \mathbf{X}^{(n+1)}) \pm u_{1-\alpha/2} \gamma_{\varphi^*}(t_k, \mathbf{X}^{(n+1)}) \right]$$

- ▶ Has the training-conditional coverage property (stronger than marginal):

$$\mathbb{P} \left(g(t_k, \mathbf{X}^{(n+1)}) \in \mathcal{CR}_\alpha(t_k, \mathbf{X}^{(n+1)}) \mid \mathcal{D}_{n_{\text{train}}} \right) = 1 - \alpha \quad (10)$$

- ▶ However, the above equality relies on **two** hypotheses:
 1. The output is modeled by a GP
 2. The prior mean and covariance functions m_β, γ_φ are suitable
- ▶ No generic way to test **misspecification** when metamodeling black-box computer codes → *a real challenge for industrial application of UQ methodology*

Conformalized GPs

- ▶ Fix for simplicity a certain $k \in \{1, \dots, n\}$ and denote $\tilde{g}(\cdot) = m_{\beta^*}(t_k, \cdot)$ and similarly for $\tilde{\gamma}$
- ▶ Define the Leave-One-Out-Gaussian (LOO $_{\gamma}$) error (with $\varepsilon > 0$ a small constant):

$$R_i^{\text{LOO}_{\gamma}} := \frac{|g(\mathbf{X}^{(i)}) - \tilde{g}_{-i}(\mathbf{X}^{(i)})|}{\max(\varepsilon, \tilde{\gamma}_{-i}(\mathbf{X}^{(i)}))}, \quad \forall i \in \{1, \dots, n\} \quad (11)$$

Main result and consequences [Jaber et al., 2025a]:

$$\widehat{\mathcal{C}}_{n,\alpha}^{\text{J+GP}}(\mathbf{X}^{(n+1)}) = \left[\widehat{q}_{n,\alpha}^{\pm} \left\{ \tilde{g}_{-i}(\mathbf{X}^{(n+1)}) \pm R_i^{\text{LOO}_{\gamma}} \times \max(\varepsilon, \tilde{\gamma}_{-i}(\mathbf{X}^{(n+1)})) \right\} \right]$$

- ▶ Coverage property **still verified** for $\alpha \in (0, 1/2)$
- ▶ Intervals have adaptive width → **more informative**
- ▶ **No hypotheses** for interpreting the interval!
- ▶ The J+GP-minmax variant **has the same properties**

Example on a GP

Conformalized stochastic collocation GP, $1 - \alpha = 0.9$

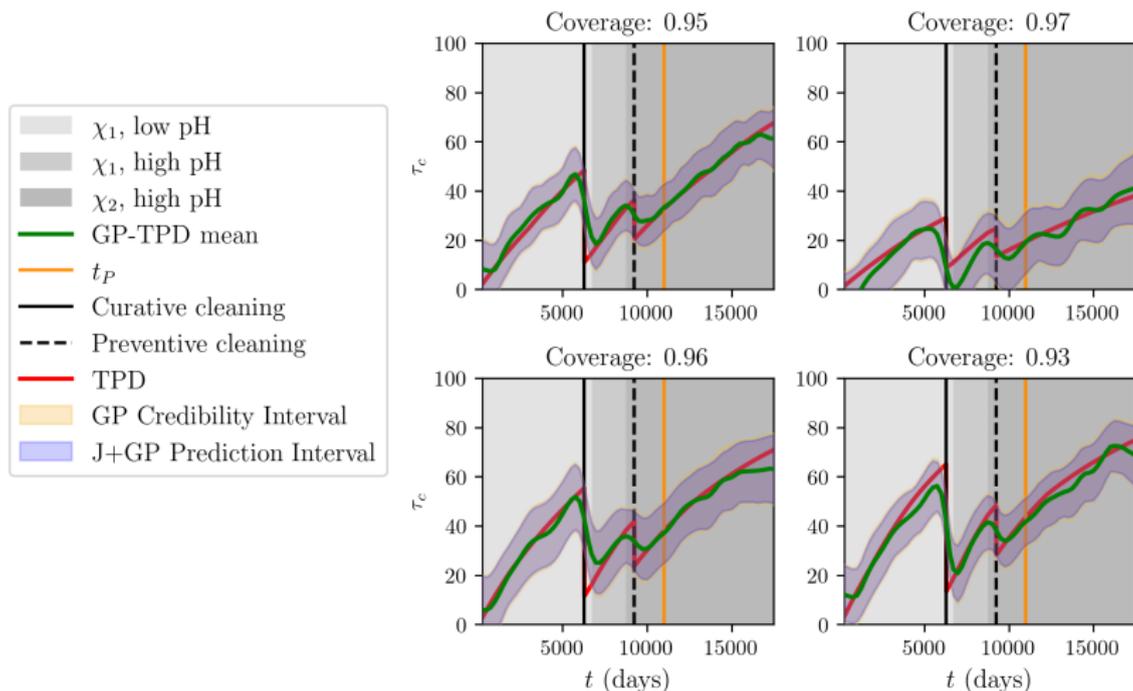


Figure 8: Stochastic GP metamodel of TPD with a linear trend and a Matérn-1/2 kernel

Results on different GPs

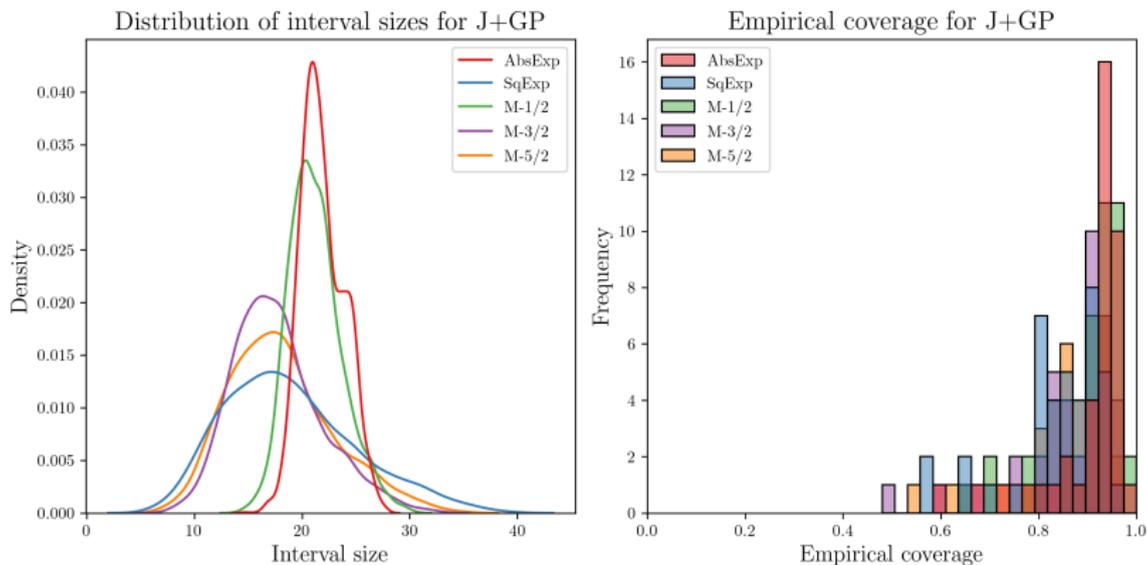


Figure 9: Distribution of interval sizes of the J+GP conformal predictor and empirical coverage distributions for a stochastic collocation GP with a linear trend and different kernels

Summary of 3

Take-home messages:

- TPD can be approximated using a variety of non-intrusive surrogate modeling techniques
- Stochastic time-collocation GP method can be applied for time-interpolating surrogates if monotonicity properties of the code are satisfied
- Cross-CP can be applied to go beyond the predictivity coefficient for qualifying non-intrusive GP surrogate models

Contributions:

- ▶ Work on CP+GP with **Dr. Vincent BLOT** at the **CEMRACS 2023** summer school on Scientific Machine Learning; Oral communication at **SIAM-UQ 2024**
- ▶ Paper [**Jaber et al., 2025a**] co-authored published in the **Journal of Machine Learning for Modeling and Computing**
- ▶ **Github repository** with plug-and-play scripts based on the **OpenTURNS** and **MAPIE [Cordier et al., 2023]** Python libraries

Outline

1. Introduction
2. The physical clogging simulation model
 - 2.1 The physical model
 - 2.2 THYC-Puffer-DEPO computational model
 - 2.3 Expert-informed UQ on TPD
3. Non-intrusive surrogate modeling
 - 3.1 General idea
 - 3.2 Gaussian processes
 - 3.3 GP validation with conformal prediction
- 4. Bayesian fusion of heterogeneous data**
 - 4.1 Offline data assimilation**
 - 4.2 The BMU algorithm**
 - 4.3 Ensemble Kalman smoothing**
5. Conclusion
6. Appendix

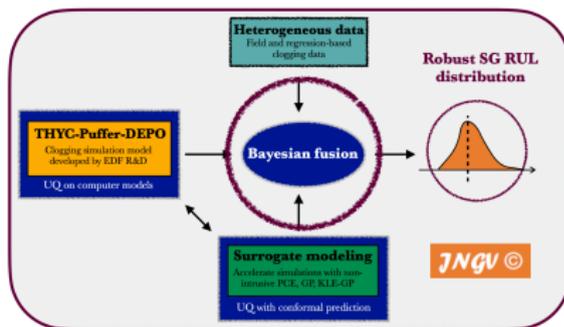
Hybrid framework for the JNGV

Objectives:

- ▶ Reduce the parametric uncertainty of TPD using a Bayesian methodology to get a more informed RUL distribution
- ▶ Use the heterogeneous data (field and regression-based), the information from the UA on TPD; the surrogates

Methodology:

- ▶ *State of the art:*
 - ➔ Data assimilation
 - ➔ Bayesian calibration
 - ➔ MCMC algorithms, ensemble methods
- ▶ *Strategy:* Generalize the Bayesian modeling in [Keller et al., 2022] to heterogeneous groups of data, include it in an offline data assimilation methodology



Available tools

- *Physics-based computer simulation model* $g : \mathcal{X} \subset \mathbb{R}^d \rightarrow \mathbb{R}^N$ with prior uncertainty on input variables \mathbf{X} , one input \mathbf{x}_0 gives a full degradation trajectory:

$$g(\mathbf{x}_0) = (g(t_1, \mathbf{x}_0), \dots, g(t_N, \mathbf{x}_0)) \quad (12)$$

and $\text{pr}_\ell \circ g(\mathbf{X}) := g(t_\ell, \mathbf{X})$ models $\Delta(t_\ell) \rightarrow$ grey-box, physics known with non-intrusive surrogate modeling strategy \hat{g}

- *q heterogeneous degradation data groups* (from different sensors, statistical models,...) $\mathcal{D} = (\mathbf{y}^1, \dots, \mathbf{y}^q)$ with different sizes $\mathbf{y}^i \in \mathbb{R}^{m_i} \rightarrow$ corresponding to different time indices in \mathcal{J}_i so that $\mathcal{J} = \cup_{i=1}^q \mathcal{J}_i$ and $|\mathcal{J}| = m_1 + \dots + m_q$, such that:

$$\mathbf{y}^i(t_\ell) = \Delta(t_\ell) + \eta_\ell^i \quad (13)$$

with $\eta_\ell^i \sim \mathcal{N}(0, R^i) \rightarrow$ homoskedastic noise for each data group

- ***How to fuse these tools for hybrid RUL estimation of the industrial system?***

Offline data assimilation

- ▶ Perform offline sequential data assimilation over time windows [Geir Evensen, 2022] using ensemble methods
- ▶ Let $\mathbf{Z} = (\mathbf{X}, \mathbf{Y})$ where $\mathbf{Y} = g(\mathbf{X}) + \epsilon = \widehat{g}(\mathbf{X})$ are the computer model emulator outputs on a prescribed time window and \mathbf{X} are uncertain parameters that do not vary over time
- ▶ The goal of assimilation is to estimate the distribution $p(\mathbf{Z} | \mathcal{D})$
- ▶ Use *modular* approach on each time window:
 1. Tailored Bayesian model updating (BMU) to obtain a posterior distribution $p_{\mathbf{X}|\mathcal{D}}$ for the input parameters: $(\mathbf{X}|\mathcal{D}) \sim p_{\mathbf{X}|\mathcal{D}}$.
 2. Apply smoothing to estimate $p(\mathbf{Y}|\mathbf{X} \sim p_{\mathbf{X}|\mathcal{D}}, \mathcal{D})$, yielding the assimilated posterior approximation:

$$p(\mathbf{Z}|\mathcal{D}) \approx p(\mathbf{Y}|\mathbf{X} \sim p_{\mathbf{X}|\mathcal{D}}, \mathcal{D})p_{\mathbf{X}|\mathcal{D}}. \quad (14)$$

Illustration of the methodology

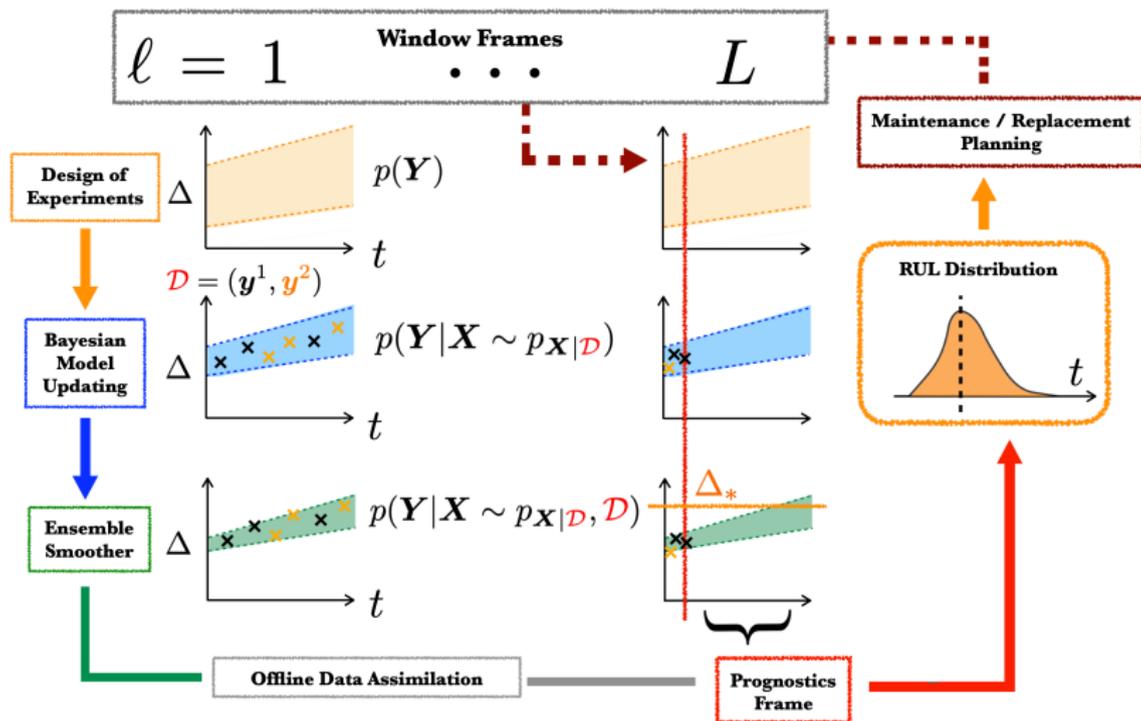


Figure 10: A sketch of the offline data assimilation methodology

On the prognostics window

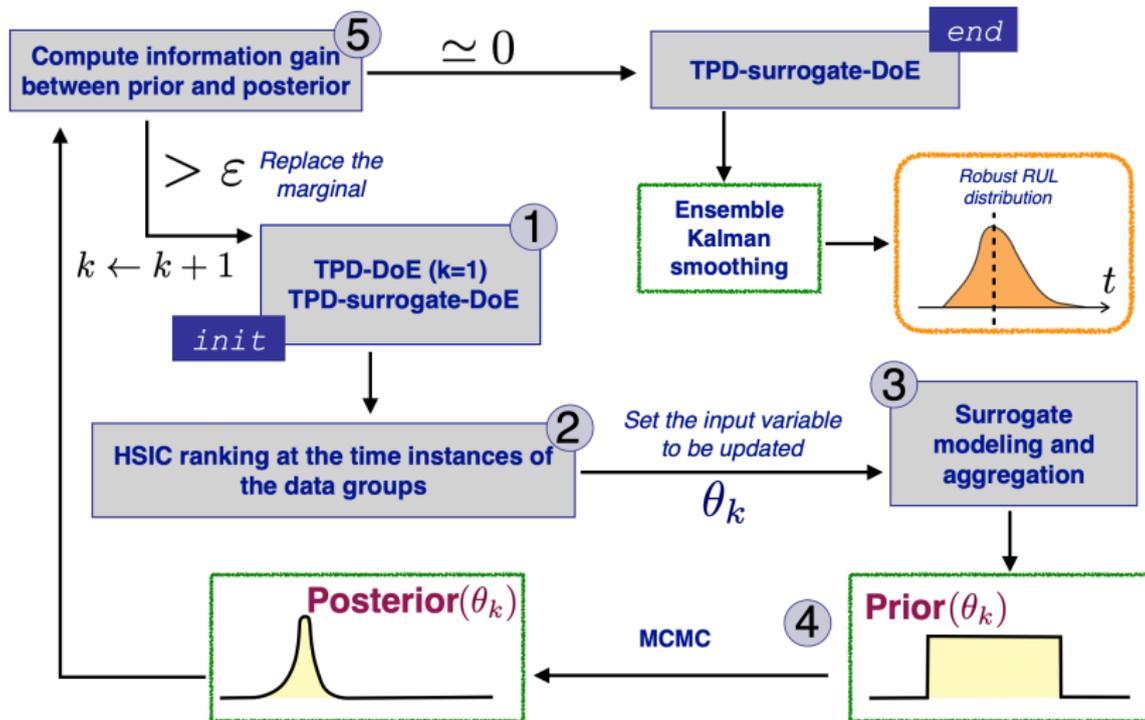
- Assimilation is performed offline until the prognostics window $\ell = L$ is reached, then the RUL distribution can be computed:

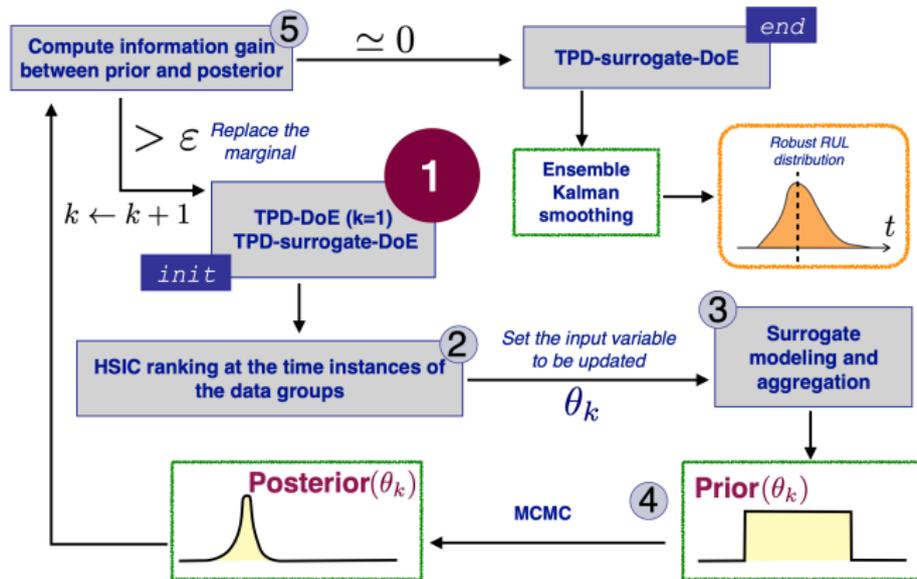
$$\mathbb{P}(\text{RUL}(\Delta_*) \leq t_j \mid \mathcal{D}) = \int_{\mathbb{R}} \mathbf{1}\{\text{pr}_{j+1}(\mathbf{y}) \geq \Delta_*\} p(\mathbf{y} \mid \mathcal{D}, \mathbf{X} \sim p_{\mathbf{X} \mid \mathcal{D}}) d\mathbf{y}, \quad (15)$$

- In practice this probability is estimated using a Monte Carlo ensemble $\{(\mathbf{X}^{(i)}, \hat{\mathbf{g}}(\mathbf{X}^{(i)}))\}_{i=1}^n \sim p_{\mathbf{X} \mid \mathcal{D}} \otimes \hat{\mathbf{g}} \# p_{\mathbf{X} \mid \mathcal{D}}$:

$$\mathbb{P}(\text{RUL}(\Delta_*) \leq t_j \mid \mathcal{D}) \approx \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{\text{pr}_{j+1} \circ \hat{\mathbf{g}}(\mathbf{X}^{(i)}) \geq \Delta_*\} \quad (16)$$

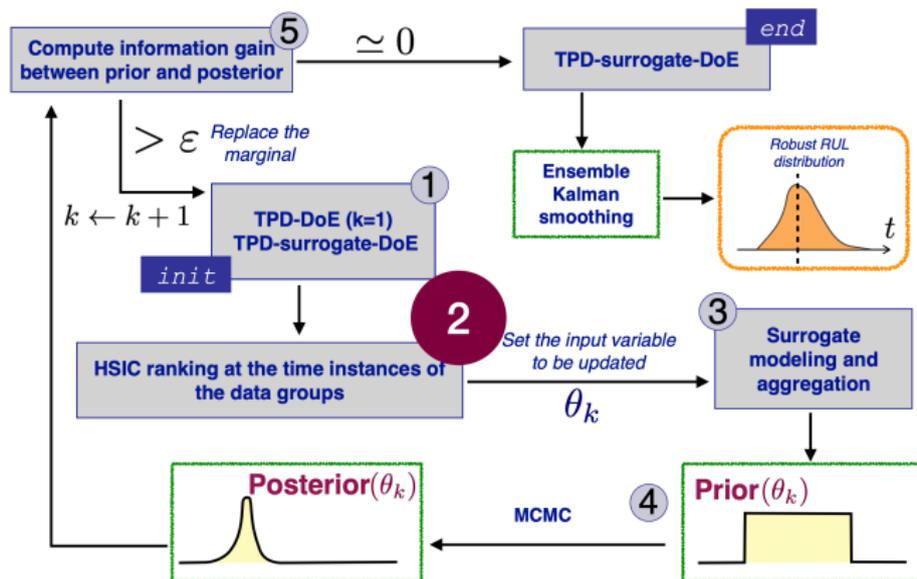
The BMU algorithm



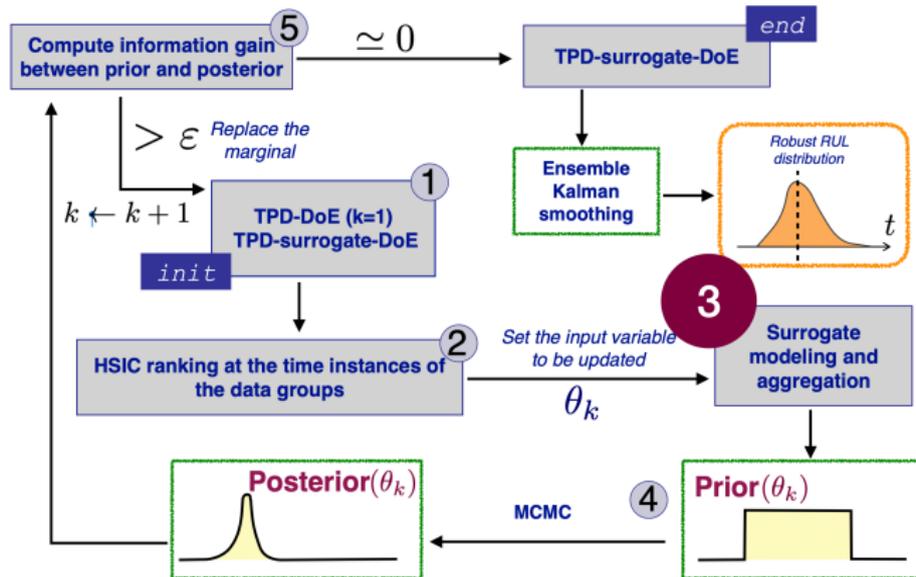


Perform k iterations where $1 \leq k \leq d$:

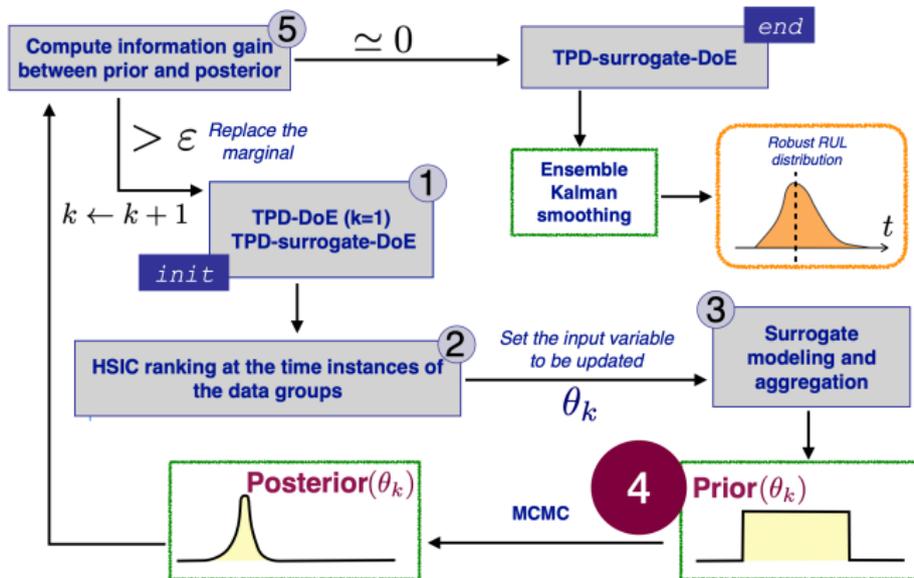
1. If $k = 0$, assume uniform independent priors $\mu_{\mathbf{X},k} \simeq \mathcal{U}[-1, 1]^{\otimes d} \rightarrow$ generate a design of experiments $\text{DoE}_g^{\mu_{\mathbf{X},k}} = \{(\mathbf{X}^{(j)}, \hat{g}(\mathbf{X}^{(j)}))\}_{1 \leq j \leq n}$



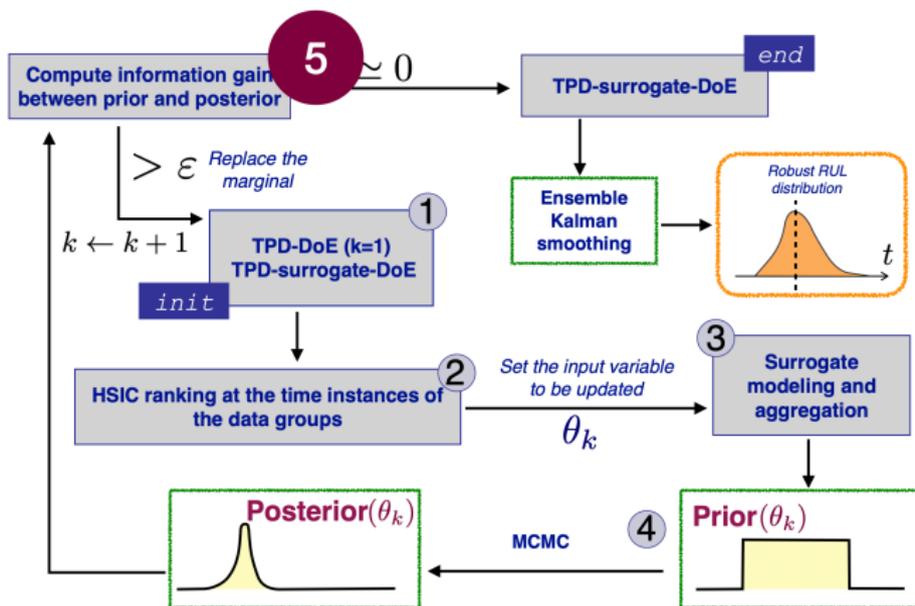
2. Compute **HSIC indices** [Gretton et al., 2005] between input variables and outputs at data time instances \rightarrow evaluate local sensitivity



3. If g is time-costly, build and validate p metamodels $\hat{\mathbf{g}} = (\hat{g}^{(1)}, \dots, \hat{g}^{(p)})$ with chosen strategy \rightarrow avoid metamodeling bias with *convex aggregation* on the unit-simplex choosing $\mathbf{w} \in \Delta^{p-1} := \{\mathbf{w} \in [0, 1]^p, \|\mathbf{w}\|_1 = 1\}$, fix nominal value of $\mathbf{U}_{0,k} = \mathbf{u}_{0,k}$ by taking the mean



4. Estimate the posterior distribution $\hat{p}(\theta_k | \mathcal{D}, \mathbf{u}_{0,k})$ with an Monte Carlo Markov Chain (MCMC) sampling procedure



5. Compute the Kullback-Leibler divergence d_{KL} between prior distribution $\mathcal{U}(\theta_k)$ and the estimated density:

- ▶ If $d_{KL} > \epsilon$, update the prior $\mu_{X,k}$ by replacing marginal $\mathcal{U}(\theta_k)$ with $\hat{p}(\theta_k | \mathcal{D}, \mathbf{u}_{0,k})$ and continue $k \leftarrow k + 1$
- ▶ Otherwise, stop and obtain an *updated* RUL prediction by computing $g \# \mu_{X,k^*}$

Bayesian updating step

$$p(\theta_k | \mathcal{D}, \mathbf{u}_{0,k}) \propto \frac{1}{M} \sum_{r=1}^M \prod_{i=1}^q \| \mathbf{y}^i - \langle \mathbf{w}^{(r)}, \hat{\mathbf{g}}(\mathbf{u}_{0,k}, \theta_k) \rangle \|^2 \quad (17)$$

- ▶ Use Random Walk Metropolis-Hastings (RWMH) MCMC algorithm [Sullivan, 2015] to sample from (17)
- ▶ Monte-Carlo integration using sample $\{\mathbf{w}^{(r)}\}_{r=1}^M$ from the Dirichlet- $\mathbf{1}_p$ distribution on the simplex → integrate hyperparameter
- ▶ Updated densities are conditioned on nominal values $\mathbf{u}_{0,k}$ of other $d - 1$ input variables → future work on how to integrate uncertainty

Ensemble Kalman smoothing step

- ▶ We don't need sequential updating at each time-step $p(g(t_k, \mathbf{X})|\mathcal{D})$ (which is the goal of filtering), but to assimilate the data on the entire window frame
- ▶ Indeed the data is *unavailable* at each time-step, so we need to account of all the information on a window
- ▶ The relevant paradigm for this is smoothing, i.e getting estimates from the full-posterior distribution $p(g(t_1, \mathbf{X}), \dots, g(t_k, \mathbf{X})|\mathcal{D}) \rightarrow$ Ensemble Kalman smoothing (EnKS) [Evensen and Van Leeuwen, 2000]

Results: posterior distributions

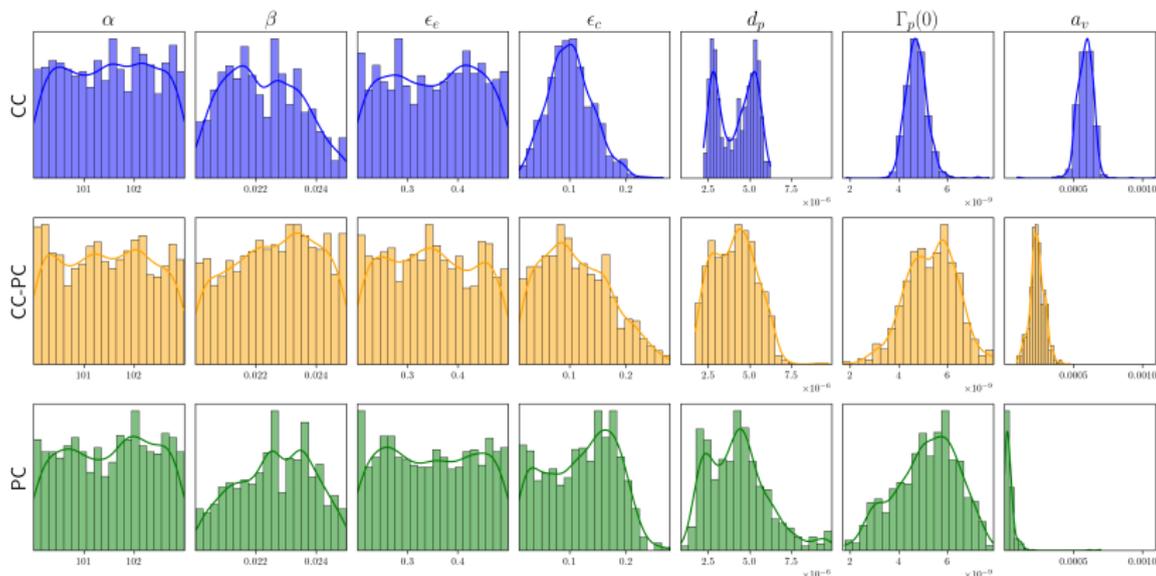


Figure 11: Posterior distributions of TPD clogging simulation code

- ▶ $L = 3$ time windows corresponding to cleanings CC/CC-PC/PC
- ▶ 5/7 distributions are informed by the data, distinct modes for a_v

Results: posterior trajectories

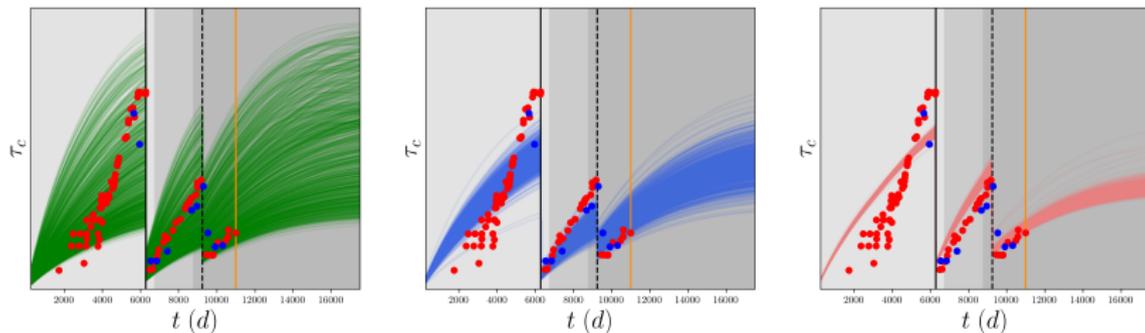
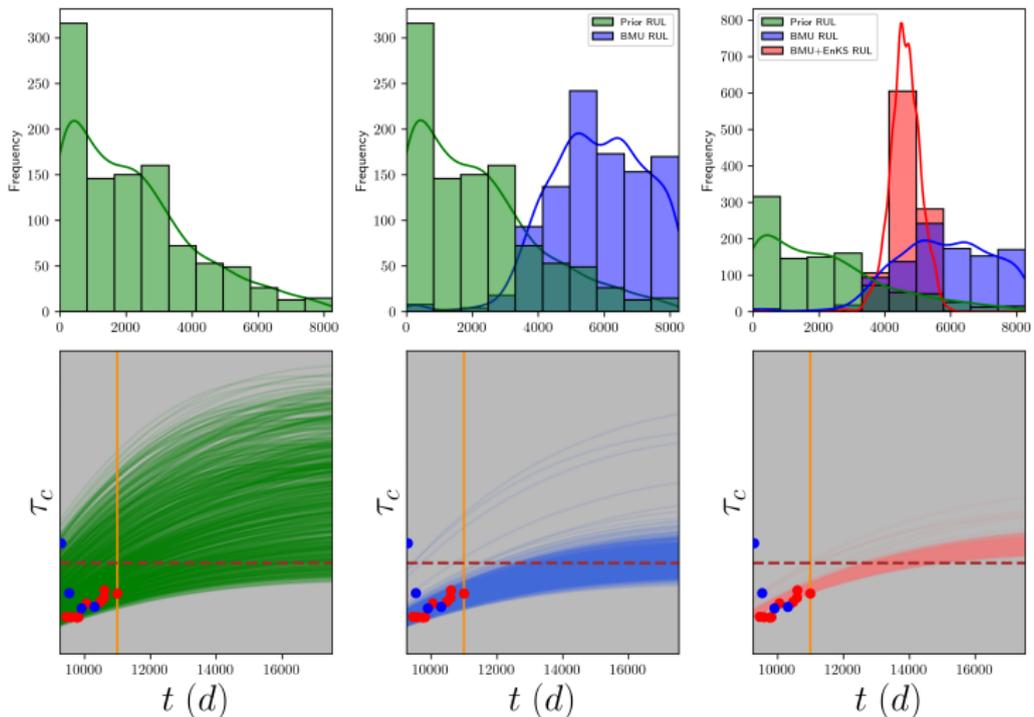


Figure 12: Prior/posterior and smoothed TPD emulations with Karhunen-Loève expansion as main surrogate

Results: posterior RUL



- RUL prediction uncertainty substantially reduced and mean of the distribution shifted compared to the prior

➔ **More informed maintenance planning!**

Summary of 4

Take-home messages:

- Offline data assimilation generic modular approach involving iterative BMU, TPD emulators, heterogeneous data groups and ensemble Kalman smoothing
- Reduces parametric uncertainty of the TPD - RUL prediction, more informed decision making

Contributions:

- ▶ Paper [[Jaber et al., 2025c](#)] under review in **Journal of Reliability Engineering & System Safety**
- ▶ Reproducible code in [GitHub repository](#)
- ▶ Communications: [INI Workshop on calibrating prediction uncertainty](#) - Cambridge, U.K, [Mathematical foundations of digital twins](#) - CIRM, Marseille

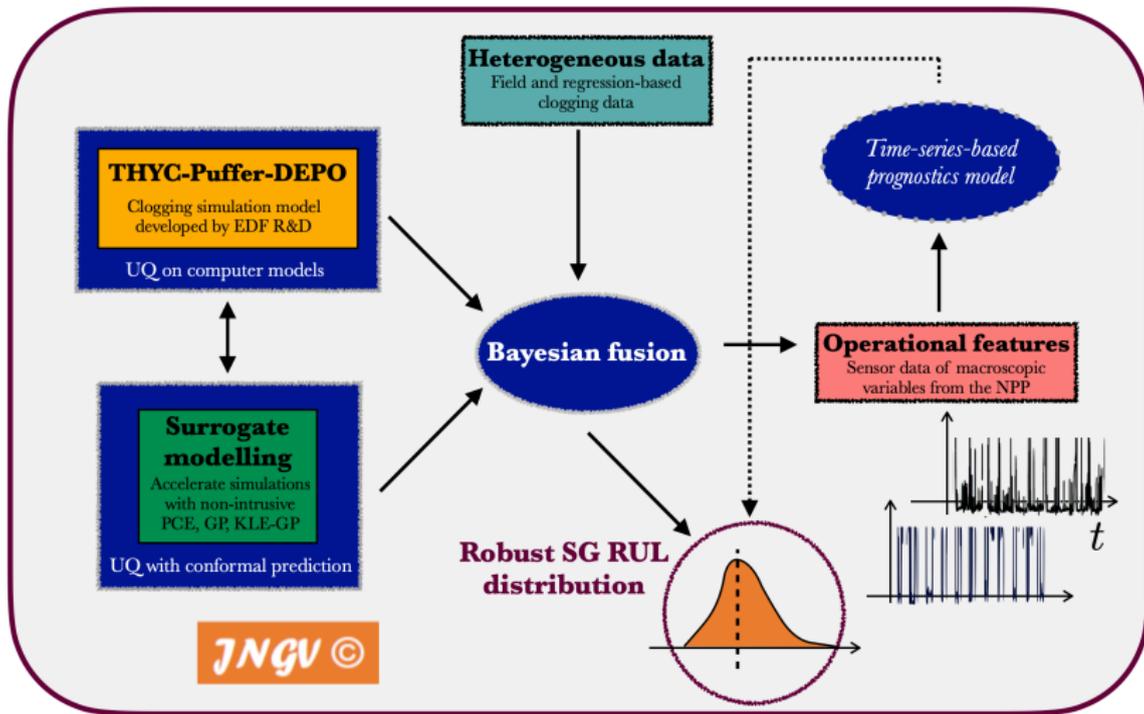
Outline

1. Introduction
2. The physical clogging simulation model
 - 2.1 The physical model
 - 2.2 THYC-Puffer-DEPO computational model
 - 2.3 Expert-informed UQ on TPD
3. Non-intrusive surrogate modeling
 - 3.1 General idea
 - 3.2 Gaussian processes
 - 3.3 GP validation with conformal prediction
4. Bayesian fusion of heterogeneous data
 - 4.1 Offline data assimilation
 - 4.2 The BMU algorithm
 - 4.3 Ensemble Kalman smoothing
- 5. Conclusion**
6. Appendix

Key take-aways

- ▶ Hybrid framework methodology for reliable RUL distribution involving:
 - Clogging simulation code TPD with parametric uncertainty
 - Different non-intrusive surrogate modeling techniques of TPD and validation of GPs enhanced with CP
 - Bayesian fusion methodology based on an offline data assimilation framework
- ▶ Three papers on each subject [[Jaber et al., 2025b,a,c](#)] and numerous communications in international conferences and workshops
- ▶ Exploration of time-series prognostic models to combine NPWR monitoring data with simulation outputs (in manuscript)

Full hybrid framework



Some scientific perspectives

- ▶ Enable online (real-time) data assimilation by developing non-destructive measurements to update TPD model continuously and time-series prognostics models
- ▶ Create dedicated experimental validation datasets and strengthen collaboration between experimentalists and modelers to better calibrate and benchmark physical formulations and closure laws
- ▶ Benchmark clogging models at different fidelity levels, perform higher-order sensitivity analyses of THYC closure laws, and propagate parametric uncertainty across THYC, Puffer, and DEPO (e.g., via stochastic field representations)

Industrial perspectives

- ▶ Embed the UA and hybrid modules into EDF's SG DT (JNGV) platform, provide clear dashboards, secure data pipelines and user workflows, and deploy the method across the SG fleet
- ▶ The offline data assimilation approach can be adapted to other SG failure modes (if a physical model is possible) and to DTs in different industries
- ▶ Develop Verification and Validation procedures, regulatory-aligned validation and explainability tools to build operator and regulator trust in nuclear DTs

Thank you for your attention!



[Jaber et al., 2025b] **Jaber, E.**, Chabridon, V., Remy, E., Baudin, M., Lucor, D., Mougeot, M., Iooss, B., *Sensitivity Analyses of a Multi-Physics Long-Term Clogging Model For Steam Generators*, **International Journal of Uncertainty Quantification**, 2025, DOI: [10.1615/Int.J.UncertaintyQuantification.2024051489](https://doi.org/10.1615/Int.J.UncertaintyQuantification.2024051489)

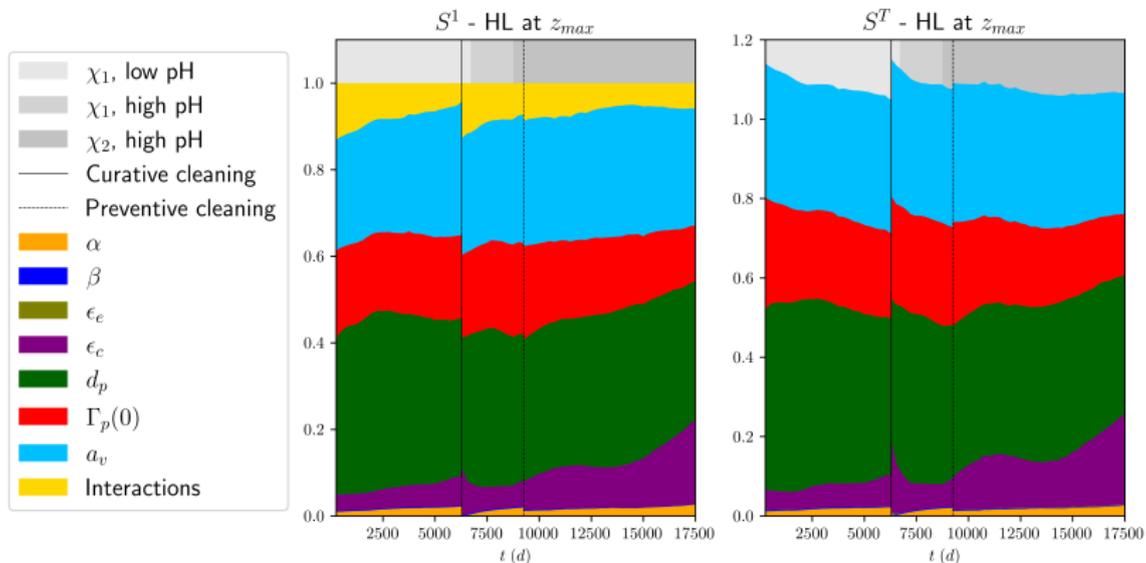
[Jaber et al., 2025a] **Jaber, E.**, Blot, V., Brunel, N., Chabridon, V., Remy, E., Iooss, B., Lucor, D., Mougeot, M., Leite, A., *Conformal Approach to Gaussian Process Surrogate Evaluation with Marginal Coverage Guarantees*, **Journal of Machine Learning for Modeling and Computing**, 2025, DOI: [10.1615/JMachLearnModelComput.2025054687](https://doi.org/10.1615/JMachLearnModelComput.2025054687)

[Jaber et al., 2025c] **Jaber, E.**, Remy, E., Chabridon, V., Lucor, D., Mougeot, M., *Fusion of heterogeneous data for robust degradation prognostics*, ArXiv [2506.05882](https://arxiv.org/abs/2506.05882), 2025, **Reliability Engineering and System Safety**, under review

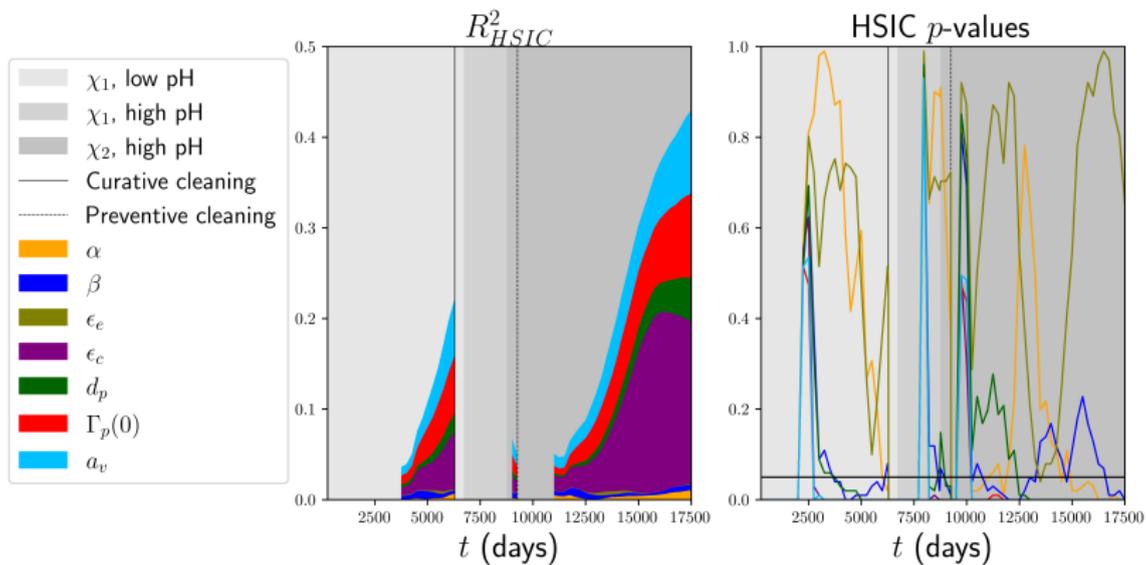
Outline

1. Introduction
2. The physical clogging simulation model
 - 2.1 The physical model
 - 2.2 THYC-Puffer-DEPO computational model
 - 2.3 Expert-informed UQ on TPD
3. Non-intrusive surrogate modeling
 - 3.1 General idea
 - 3.2 Gaussian processes
 - 3.3 GP validation with conformal prediction
4. Bayesian fusion of heterogeneous data
 - 4.1 Offline data assimilation
 - 4.2 The BMU algorithm
 - 4.3 Ensemble Kalman smoothing
5. Conclusion
6. Appendix

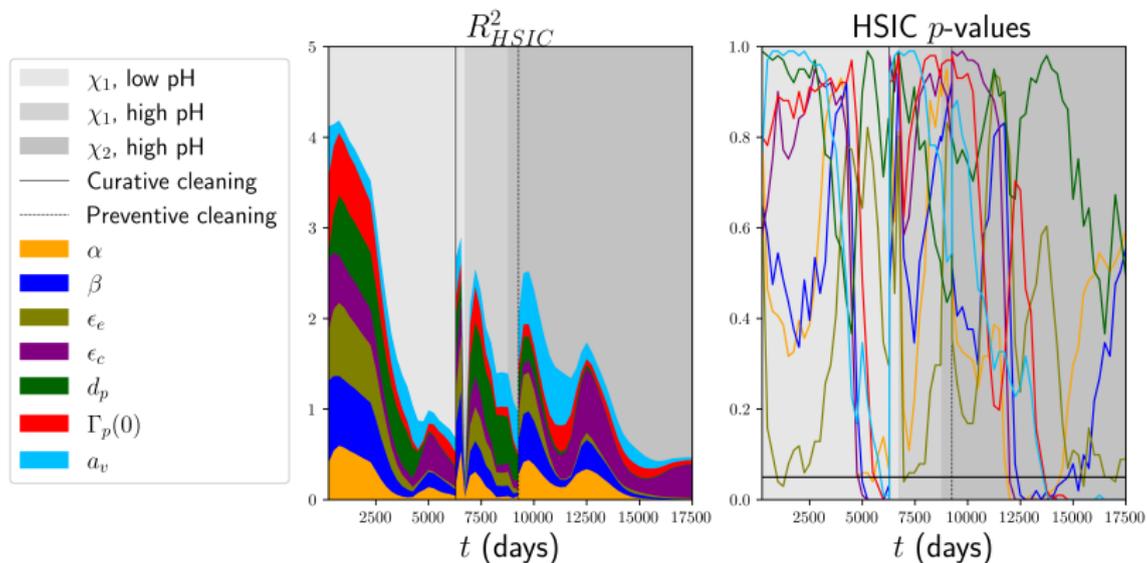
Sobol' indices



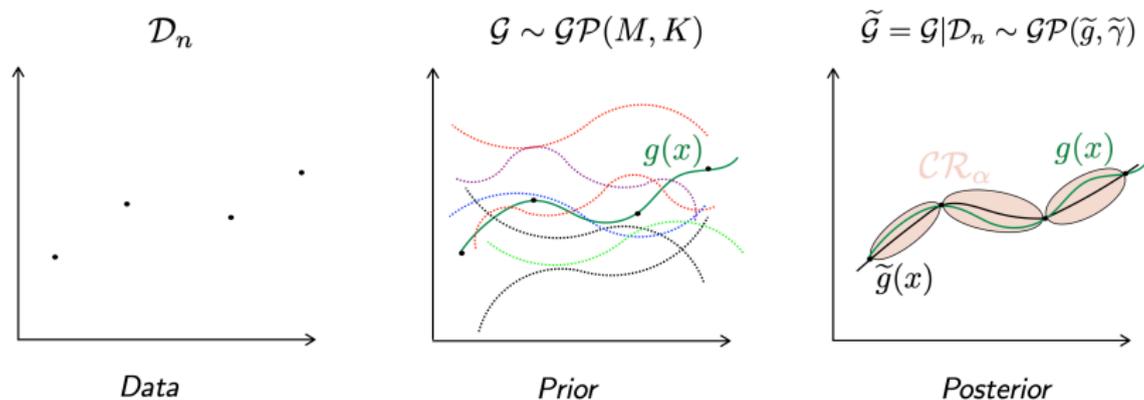
Target HSIC-indices



Conditional HSIC-indices



Gaussian processes



- ▶ GP metamodel prediction is the *posterior mean* $\hat{g} = \tilde{g}$
- ▶ Notion of uncertainty through the posterior covariance $\tilde{\gamma}$ and the Gaussian structure of the metamodel

Bayesian updating step

Proposition

Assume $\lambda := 1/\sigma_\eta^2 \sim \mathcal{G}(\frac{m}{2}, \frac{1}{2}\|\mathbf{y} - f(\theta)\|^2)$ (Gamma distribution), where m is the number of data points in \mathbf{y} ; $\theta \sim \mathcal{U}(\theta)$, and $p(\theta, \lambda) \propto \lambda^{-1}$.

Then:

$$p(\theta|\mathbf{y}) \propto \|\mathbf{y} - f(\theta)\|^{-m} \quad (18)$$

Moreover, if multiple groups of data at different time-instances are considered, $\mathbf{y}^1, \dots, \mathbf{y}^q$, with respective priors on the inverse of their standard deviations $\lambda_i \sim \mathcal{G}(\frac{m_i}{2}, \frac{1}{2}\|\mathbf{y}^i - f(\theta)\|^2)$, then the generalization is:

$$p(\theta|\mathbf{y}^1, \dots, \mathbf{y}^q) \propto \prod_{i=1}^q \|\mathbf{y}^i - f(\theta)\|^{-m_i} \quad (19)$$

Proof. Bayes' theorem and simplifications

Ensemble Kalman smoothing

- ▶ Build an ensemble $\{g(\mathbf{X}^{(k)})\}_{k=1}^n \sim g\#h$, suppose for each data-group $i = 1, \dots, q$:

$$y^i(t_\ell^i) = g(t_\ell^i, \mathbf{X}) + \eta_i, \quad (20)$$

where the variance R^i of the noise *is known*

- ▶ Define the ensemble mean and anomalies:

$$\bar{g}(t) = \frac{1}{n} \sum_{k=1}^n g(t, \mathbf{X}^{(k)}), \quad A^{(k)}(t) = g(t, \mathbf{X}^{(k)}) - \bar{g}(t) \quad (21)$$

- ▶ The cross-covariance between ensemble states at any time t and the observation times t_ℓ^i are approximated empirically:

$$\hat{C}(t, t_\ell^i) = \frac{1}{n-1} \sum_{k=1}^n A^{(k)}(t) A^{(k)}(t_\ell^i) \quad (22)$$

Ensemble Kalman smoothing

- ▶ Define the Kalman gain for each t in the time-window:

$$K(t) = \frac{\widehat{C}(t, t_\ell^i)}{\widehat{C}(t, t_\ell^i) + R^i} \quad (23)$$

- ▶ Defining the innovation term as $d_\ell^{(k)} = y^j(t_\ell^i) - g(t_\ell^i, \mathbf{X}^{(k)})$, we update all the ensembles following:

$$g(\mathbf{X}^{(k)}) \leftarrow g(\mathbf{X}^{(k)}) - K_t \cdot d_\ell^{(k)} \quad (24)$$

- ▶ Process repeated sequentially across all observation times t_ℓ^i , for all data groups i

- Cordier, T., Blot, V., Lacombe, L., Morzadec, T., Capitaine, A., and Brunel, N. (2023). Flexible and Systematic Uncertainty Estimation with Conformal Prediction via the MAPIE library. In *Conformal and Probabilistic Prediction with Applications*.
- Da Veiga, S. (2015). Global sensitivity analysis with dependence measures. *Journal of Statistical Computation and Simulation*, 85:1283–1305.
- David, F. (1999). Three Dimensional Thermal-Hydraulic Simulation In Steam Generators With THYC Exchangers Code - Application To The UTSG Model 73/19. Proceedings of the Ninth International Topical Meeting on Nuclear Reactor Thermal Hydraulics (NURETH-9).
- De Rocquigny, E., Devictor, N., and Tarantola, S., editors (2008). *Uncertainty in industrial practice - A guide to quantitative uncertainty management*. Wiley and Sons.

- Deri, E., Varé, C., and Wintergerst, M. (2021). *Development of Digital Twins of PWR Steam Generators: Description of Two Maintenance-Oriented Use Cases*, volume Volume 1: Operating Plant Challenges, Successes, and Lessons Learned; Nuclear Plant Engineering; Advanced Reactors and Fusion; Small Modular and Micro-Reactors Technologies and Applications of *International Conference on Nuclear Engineering*.
- Evensen, G. and Van Leeuwen, P. J. (2000). An ensemble Kalman smoother for nonlinear dynamics. *Monthly Weather Review*, 128(6):1852–1867.
- Feng, Q., Nebes, J., Bachet, M., Pujet, S., You, D., and Deri, E. (2023). Tube support plates blockage of PWR steam generators: thermalhydraulics and chemical modeling.
- Geir Evensen, Femke C. Vossepoel, P. J. v. L. (2022). *Data assimilation fundamentals*. Springer Cham.

- Gretton, A., Bousquet, O., Smola, A., and Schölkopf, B. (2005). Measuring statistical dependence with Hilbert-Schmidt norms. In *Proceedings Algorithmic Learning Theory*, pages 63–77. Springer-Verlag.
- Jaber, E., Blot, V., Brunel, N., Chabridon, V., Remy, E., looss, B., Lucor, D., Mougeot, M., and Leite, A. (2025a). Conformal approach to Gaussian process surrogate evaluation with marginal coverage guarantees. *Journal of Machine Learning for Modeling and Computing*.
- Jaber, E., Chabridon, V., Remy, E., Baudin, M., Lucor, D., Mougeot, M., and looss, B. (2025b). Sensitivity Analyses of a Multi-Physics Long-Term Clogging Model For Steam Generators. *International Journal for Uncertainty Quantification*, 15:27–45.
- Jaber, E., Remy, E., Chabridon, V., Mougeot, M., and Lucor, D. (2025c). Fusion of heterogeneous data for robust degradation prognostics. *arXiv:2506.05882*.

- Keller, M., Damblin, G., Pasanisi, A., Schumann, M., Barbillon, P., Ruggeri, F., and Parent, E. (2022). Validation of a computer code for the energy consumption of a building, with application to optimal electric bill pricing. *Econometrics*, 10(4).
- Le Maître, O. P. and Knio, O. M. (2010). *Spectral Methods for Uncertainty Quantification: With Applications to Computational Fluid Dynamics*. Scientific Computation. Springer.
- Lefebvre, L., Segond, M., Spaggiari, R., Le Gratiet, L., Deri, E., looss, B., and Damblin, G. (2023). Improving the Predictivity of a Steam Generator Clogging Numerical Model by Global Sensitivity Analysis and Bayesian Calibration Techniques. *Nuclear Science and Engineering*, 197(8):2136–2149.
- Pincirolì, L., Baraldi, P., Shokry, A., Zio, E., Seraoui, R., and Mai, C. (2021). A semi-supervised method for the characterization of degradation of nuclear power plants steam generators. *Progress in Nuclear Energy*, 131:103580.

