

## KERNEL-BASED PARAMETER SCREENING FOR CONDITIONAL BAYESIAN CALIBRATION OF CHAINED NUMERICAL MODELS: APPLICATION TO FUEL PERFORMANCE SIMULATION OF PRESSURIZED WATER REACTORS

OUMAR BALDÉ<sup>1,\*</sup>, GABRIEL SARAZIN<sup>2</sup>, AMANDINE MARREL<sup>1,3</sup>,  
GUILLAUME DAMBLIN<sup>2</sup> AND ANTOINE BOULORÉ<sup>4</sup>

**Abstract.** Numerical simulation is widely used in many fields of engineering to study complex physical systems. The numerical models, designed to faithfully represent the underlying physical phenomena, are subject to uncertainties of different natures (either numerical, stochastic or epistemic) that degrade the accuracy of the simulated outputs. Part of epistemic uncertainty arises from limited knowledge regarding some input model parameters  $\theta$ . This component can be reduced through Bayesian calibration of the model against experimental data. Before calibration itself, sensitivity analysis can be used to better understand how parameter uncertainties impact the model output, and this may help confine calibration to the most impactful parameters. In this work, we show that kernel methods, especially those based on the Hilbert–Schmidt independence criterion (HSIC), are effective tools in support of Bayesian calibration, both for a single model and for two chained models. In the latter case, our main contribution is a screening methodology for the parameters  $\theta$  of the downstream model, which accounts for the posterior distribution of the upstream model parameters  $\lambda$ . By taking the expectation of the HSIC over  $\lambda$ , we define a new sensitivity measure that is able to incorporate the residual uncertainty due to the upstream model calibration. We show that the resulting sensitivity indices can be estimated from the same data used for conditional Bayesian calibration. We further demonstrate that the corresponding estimators are consistent and achieve convergence rates comparable to those of classical Monte Carlo estimators. Importantly, we construct two test procedures that enable rigorous decisions on which parameters among  $\theta$  should be selected. Finally, we apply the proposed approach to nuclear fuel simulation to screen the calibration parameters  $\theta$  of a fission gas behavior model which follows an upstream thermal model whose thermal conductivity  $\lambda$  was calibrated in previous work.

**Mathematics Subject Classification.** 62H20, 62F15, 62G05, 62G10, 62F12, 62P30.

Received September 24, 2025. Accepted February 26, 2026.

---

*Keywords and phrases:* Bayesian calibration, global sensitivity analysis, chained numerical models, dependence measures, Hilbert–Schmidt independence criterion, independence testing.

<sup>1</sup> CEA, DES, IRESNE, DER, SESI, Cadarache, 13108 Saint-Paul-Lez-Durance, France.

<sup>2</sup> Université Paris-Saclay, CEA, DES, ISAS, DM2S, SGLS, 91191 Gif-sur-Yvette, France.

<sup>3</sup> Avignon Université, LMA UPR 2151, 84140 Avignon, France.

<sup>4</sup> CEA, DES, IRESNE, DEC, SESC, Cadarache, 13108 Saint-Paul-Lez-Durance, France.

\* Corresponding author: [oumar.balde@cea.fr](mailto:oumar.balde@cea.fr)

## 1. INTRODUCTION

Like the power generation industry, many high-impact industrial sectors rely on numerical simulations to investigate complex physical phenomena. The numerical models used to mimic the behavior of systems or facilities often involve a large number of control variables, such as experimental conditions or geometric design specifications [1]. Additionally, simulations often involve physical parameters that cannot be measured experimentally, thus carry epistemic uncertainty, and therefore require calibration [2, 3]. Mathematically, a numerical model is a function  $y : (\mathbf{x}, \boldsymbol{\theta}) \in \mathcal{D}_{\mathbf{X}} \times \mathcal{D}_{\boldsymbol{\theta}} \subseteq \mathbb{R}^d \times \mathbb{R}^p \rightarrow \mathcal{D}_Y \subseteq \mathbb{R}$  where  $\mathbf{x}$  denotes the control variables and  $\boldsymbol{\theta}$  the calibration parameters.

In the context of uncertainty quantification [4], model calibration is a crucial step to improve the accuracy and robustness of numerical simulations [2, 5]. The calibration step aims to set the optimal parameter value(s) for  $\boldsymbol{\theta}$  so that the simulator  $y$  sticks as closely as possible to the underlying physical reality. To do so, some experimental data of the system must be available. There are basically two types of calibration: deterministic calibration and Bayesian calibration. The first one tries to find the optimal value(s) of  $\boldsymbol{\theta}$  by minimizing a loss function (most often based on the least-squares criterion) which quantifies the discrepancy between the available experimental data and the corresponding simulation data [6, 7]. In contrast, the second approach treats uncertainty in the calibration parameters  $\boldsymbol{\theta}$  within a fully probabilistic framework. Bayesian inference then leverages both experimental data and expert judgment to extract maximal information [8, 9].

Bayesian calibration improves the accuracy of numerical models by estimating the parameters  $\boldsymbol{\theta}$  from prior knowledge, available experimental data  $\mathbf{z} = \{\mathbf{z}_i\}_{i=1}^{n_{\text{exp}}}$  and a probabilistic model. The prior distribution  $\pi(\boldsymbol{\theta})$  encodes prior knowledge about  $\boldsymbol{\theta}$  typically obtained through expert elicitation or statistical considerations [10]. The experimental data  $\mathbf{z}$ , reflecting the physical system with inherent uncertainty, are then used to update this prior with Bayes' rule, which leads to the posterior distribution  $\pi(\boldsymbol{\theta} | \mathbf{z})$ . In practice, the maximum *a posteriori* (MAP) estimate of  $\boldsymbol{\theta}$  (*i.e.*, the mode of the posterior density) is most commonly used to run the calibrated model.

For complex models, the posterior distribution  $\pi(\boldsymbol{\theta} | \mathbf{z})$  is almost always intractable analytically, but it can be approximated via simulation. In particular, Markov Chain Monte Carlo (MCMC) methods, such as the Metropolis-Hastings (MH) algorithm [11], allow simulating samples from  $\pi(\boldsymbol{\theta} | \mathbf{z})$  by constructing a Markov chain whose stationary distribution is exactly  $\pi(\boldsymbol{\theta} | \mathbf{z})$ . However, due to the burn-in period required to reach stationarity, MCMC typically demands thousands of model evaluations, which is prohibitively expensive in many applications. In practice, to mitigate this cost, the model  $y$  is replaced by an inexpensive metamodel<sup>1</sup> (also called a surrogate model)  $\hat{y}$  [12, 13]. When the number  $p$  of parameters to be calibrated increases, additional difficulties arise. The most significant of these is the slow convergence of the MH algorithm. Adaptive MH variants [14] overcome this limitation by iteratively updating the proposal distribution, which improves convergence. However, such adaptive MH algorithms are difficult to tune efficiently. For this reason, in many recent works dedicated to high-dimensional Bayesian calibration, a preliminary step of dimension reduction is completed before using the MH algorithm. Bayesian calibration is thus limited to a subset of parameters which are selected by means of a preliminary sensitivity analysis [15–18].

Sensitivity analysis (SA) may be used to rank or screen model inputs according to their influence on the output(s) [19, 20]. Unlike local methods that measure sensitivity in the neighborhood of a nominal point, global sensitivity analysis (GSA) endeavors to account for the entire joint distribution of all sources of uncertainty [21]. In the 40 years that SA has been around, very different approaches have emerged, such as those based on the decomposition of the output variance [22, 23], the output gradient [24], linear regression [25], cooperative game theory [26, 27] or kernel methods [28, 29]. For a given problem, the most suitable SA method depends on many factors, including the dimensionality of the input space (*i.e.*, the number of input variables), the simulation budget (*i.e.*, the maximum number of model evaluations), correlations between the inputs, the mathematical nature of the output object, and specific data-related constraints. In particular, the given-data context arises when GSA must be performed using a provided set of input-output data originally generated for another purpose (*e.g.*, building a metamodel). This typically corresponds to a small dataset obtained from a Monte

<sup>1</sup>Mathematical function  $\hat{y}$  that is trained on input-output simulation data to provide a sufficiently accurate approximation of  $y$ .

Carlo or space-filling design. Many GSA methods have been proposed to tackle these challenges. Among them, the Hilbert–Schmidt independence criterion (HSIC) approach [28, 30] handles these difficulties effectively and has become widely used. In this work, the main technical challenge stems from a specific engineering context, involving the chaining of multiple numerical models within a multiphysics application.

The ALCYONE fuel application is composed of interdependent physical models that imitate the mechanical, thermal and chemical behaviors of fuel rods in the core of pressurized water reactors [31]. Here, we focus on the chaining of two models: the thermal model (upstream model) and the fission gas behavior model (downstream model). The thermal model simulates the evolution of the temperature within the fuel rod during the fission reaction and returns the associated temperature field as output. This field is subject to uncertainty because it is computed from an uncertain quantity  $\lambda$  called the thermal conductivity. In previous work, the posterior distribution of  $\lambda$  was inferred from experimental data of thermal component [32]. Then, the temperature field produced by the thermal model is used by the fission gas behavior model to simulate the evolution of the fission products (fuel swelling and fission gas release). This model returns several scalar outputs, among which the released gas fraction (RGF) will be the only output of interest in this work. It also involves parameters  $\theta$  to be calibrated from experimental RGF data  $\mathbf{z}$ . The uncertainty in the RGF results from both the uncertainty in  $\theta$  and that in  $\lambda$ , the latter being propagated through the temperature field. Estimating the posterior distribution of  $\theta$  conditional on  $\lambda$  (which is known as conditional Bayesian calibration) can be difficult when the dimension of  $\theta$  is high (beyond a dozen parameters).

The goal of this paper is to develop a screening methodology to reduce the dimension of  $\theta$ , while accounting for the additional uncertainty brought by  $\lambda$ . To this end, a dedicated sensitivity measure must be defined to handle the bi-level uncertainty framework. Moreover, given the computational cost of the two models, the associated sensitivity indices must be estimable from a reasonable number of evaluations of the chained model, or from the same data used for the conditional calibration phase. For these reasons, we turned to the HSIC approach, as it converges rapidly, often bypasses problem-specific constraints, and relies on independence testing to control the risk of selecting weakly influential inputs. We further tailored the HSIC methodology to our problem, and we demonstrated that the customization retains its key advantages, particularly rapid convergence and reliable decision-making.

The outline of the paper is as follows. Section 2 provides the unified background material necessary for this work. First, Section 2.1 recalls the main principles of Bayesian calibration and highlights why parameter screening is crucial in high-dimensional settings. Then, Section 2.2 discusses the specific challenges of parameter screening in the context of Bayesian calibration, describes the limitations of classical methods, and shows that the HSIC methodology fulfills all requirements. In Section 3, to illustrate the previous section in the simplest possible context, the HSIC approach is applied to the fission gas behavior model of the ALCYONE application, with the conductivity  $\lambda$  fixed at its nominal value. Parameter screening is subsequently performed in three different ways, still using HSIC indices, but examining different variants of the output of interest. Section 4, where it is shown how to adapt the HSIC approach to the particular context of chained numerical models, is the heart of this work. An integrated version of the HSIC is notably introduced in Section 4.1 to incorporate the uncertainty coming from the upstream numerical model. It is demonstrated that the resulting sensitivity measure properly characterizes independence in this bi-level uncertainty framework, which underlines its usefulness from a screening perspective. Then, four different estimators are proposed and their statistical properties are deeply investigated. Since the values taken by such kernel-based indices are always difficult to interpret, the final decision (regarding which parameters to select) is only trustworthy if it is supported by independence tests. For this reason, two test procedures, which are much inspired from what is done for standard HSIC indices, are developed in Section 4.2. Using these procedures, one can detect statistical dependence between the output and a given downstream model parameter, even when it only appears for certain values of the upstream model parameters. In Section 5, the proposed methodology is implemented on the fission gas behavior model of the ALCYONE application. Unlike what is done in Section 3, the posterior distribution of the thermal conductivity is this time considered in the selection process. Section 6 concludes the paper.

## 2. PARAMETER SCREENING FOR BAYESIAN CALIBRATION

### 2.1. Bayesian calibration for black-box numerical models

The physical system  $r : \mathcal{D}_{\mathbf{X}} \rightarrow \mathcal{D}_Y$  is represented as a function that maps inputs  $\mathbf{x} \in \mathcal{D}_{\mathbf{X}} \subseteq \mathbb{R}^d$  to an output  $y = r(\mathbf{x}) \in \mathcal{D}_Y \subseteq \mathbb{R}$ . The probabilistic equation linking the experimental data  $\mathbf{z}$  to the model outputs  $\{y_{\theta}(\mathbf{x}_i)\}_{i=1}^{n_{\text{exp}}}$  for a set of experimental input configurations  $\{\mathbf{x}_i\}_{i=1}^{n_{\text{exp}}}$  is given by:

$$\forall 1 \leq i \leq n_{\text{exp}}, \quad z_i = y_{\theta}(\mathbf{x}_i) + b(\mathbf{x}_i) + \epsilon_i. \quad (2.1)$$

For a given input configuration  $\mathbf{x}_i$ , the term  $\epsilon_i$  represents the experimental uncertainty, mainly due to measurement error. The vector  $\boldsymbol{\epsilon} = \{\epsilon_i\}_{i=1}^{n_{\text{exp}}}$  is typically assumed to follow a multivariate Gaussian distribution with covariance matrix  $\boldsymbol{\Sigma}_{\boldsymbol{\epsilon}}$ . The function  $b(\mathbf{x})$ , referred to as the model discrepancy in [8], quantifies the systematic mismatch between  $y_{\theta}(\mathbf{x})$  and  $r(\mathbf{x})$  when the model is operated at the optimal but unknown calibration parameter value  $\boldsymbol{\theta} \in \mathcal{D}_{\boldsymbol{\theta}} \subseteq \mathbb{R}^p$ . The discrepancy  $b(\mathbf{x})$  is generally modeled by a Gaussian process and its hyperparameters are estimated jointly with the calibration parameters. In practice, this component may be neglected and merged with experimental uncertainty [33]. We will make this assumption in this work. The probabilistic model then simplifies to:

$$\forall 1 \leq i \leq n_{\text{exp}}, \quad z_i = y_{\theta}(\mathbf{x}_i) + \epsilon_i. \quad (2.2)$$

Then, the posterior distribution  $\pi(\boldsymbol{\theta} \mid \mathbf{z})$  is computed from Bayes' formula:

$$\pi(\boldsymbol{\theta} \mid \mathbf{z}) \propto \mathcal{L}(\mathbf{z} \mid \boldsymbol{\theta}) \pi(\boldsymbol{\theta}), \quad (2.3)$$

where  $\pi(\boldsymbol{\theta})$  is the prior distribution quantifying the uncertainty of the calibration parameters before collecting  $\mathbf{z}$ . The likelihood  $\mathcal{L}(\mathbf{z} \mid \boldsymbol{\theta})$  is a measure of the agreement between the observed data  $\mathbf{z}$  and the outputs of the model  $y_{\theta}(\mathbf{x})$  for a given value of the parameters  $\boldsymbol{\theta}$ .

The posterior distribution  $\pi(\boldsymbol{\theta} \mid \mathbf{z})$  is often estimated using MCMC algorithms [34], particularly the MH algorithm [11]. Once its Markov chain has reached stationarity, the algorithm generates samples from  $\pi(\boldsymbol{\theta} \mid \mathbf{z})$ . Achieving stationarity in the Markov chain can demand thousands of model evaluations, which makes this approach impractical when the model is computationally expensive. To overcome this limitation, the model  $y_{\theta}(\mathbf{x})$  is often replaced by a metamodel  $\hat{y}(\mathbf{x}, \boldsymbol{\theta})$  [12]. After being trained on simulation data, this metamodel provides predictions of the output  $y_{\theta}(\mathbf{x})$  for any input pair  $(\mathbf{x}, \boldsymbol{\theta})$ . Sometimes, a metamodel  $\hat{y}_i(\boldsymbol{\theta})$  is built for each function  $y_i(\boldsymbol{\theta}) = y_{\theta}(\mathbf{x}_i)$  obtained by restricting  $y$  to a specific experimental configuration  $\mathbf{x}_i$ . This metamodeling strategy is particularly relevant when the experimental configurations are significantly different from each other, such as for example in nuclear engineering [35]. However, a drawback of this approach is that it requires training a separate metamodel for each experimental configuration. For Bayesian calibration, Gaussian process (GP) regression is by far the most widely employed technique [33, 36], while polynomial chaos expansion (PCE) is also frequently advocated in the literature [15, 17].

Unfortunately, everything becomes much harder when the number of calibration parameters becomes large. This difficulty is encountered in many application fields, as illustrated in [15, 37, 38] where  $p \in \{36, 48, 11\}$  respectively. Indeed, increasing the number of calibration parameters raises two majors issues:

- **Poor metamodel accuracy.** Building an accurate metamodel becomes increasingly difficult as the number of hyperparameters to optimize grows. For GP metamodels, recent strategies to address this problem include advanced screening techniques [39], multi-objective optimization [40], or active subspace construction [41], among many others. A comprehensive review dedicated to GP surrogate modeling in high-dimensional input spaces can be found in [42].
- **Poor MCMC mixing.** Working in a high-dimensional parameter space often leads to dysfunction of the MH algorithm. Suboptimal tuning of the proposal distribution can severely impair the mixing of the

Markov chain, resulting in very slow convergence to the stationary distribution [43]. This may create identifiability issues and yield uninformative posterior distributions with highly correlated marginals [44]. Adaptive MH variants have been proposed to mitigate this phenomenon [14].

Both problems are classic manifestations of the well-known curse of dimensionality. The simplest solution is to restrict calibration to the most influential parameters, while setting all other parameters to their nominal values. This reduces the dimensionality without losing significant information about the output distribution. GSA provides a rigorous framework for making such decisions.

### 2.2. HSIC-based screening for efficient Bayesian calibration

Whether used to facilitate surrogate modeling or to improve Markov chain mixing, GSA is a fairly standard technique in model calibration [37, 38, 45–47]. When performed as a preliminary step to calibration, GSA has some specific features.

- **GSA can be limited to screening.** There is no need to rank the inputs according to their influence on the output, since the goal is only dimension reduction. This is actually an advantage, as screening methods [20, 24, 28] are computationally cheaper than ranking methods [23, 48].
- **GSA must be performed from given data.** Most often, the design of experiments is built to facilitate the (forthcoming) construction of a metamodel. No specific computational budget is allocated to GSA, which must therefore work with the available simulation data. These data typically take the form of a training set of input-output observations. Here, they will be assumed to come from a Monte Carlo design, even though in practice they are more likely to originate from a space-filling design.
- **GSA may target different output quantities.** Depending on the objective, it is not always the same output object that is analyzed. Three common situations can be distinguished. Let  $\Theta$  denote the random vector modeling the prior uncertainty on the calibration parameters  $\theta$ .
  - **Approach A.** The output of interest is the model output for one specific input configuration  $\mathbf{x}_i$ :

$$Y_i(\Theta) = y(\mathbf{x}_i, \Theta). \tag{2.4}$$

For an exhaustive study over all configurations, the GSA procedure must be repeated  $n_{\text{exp}}$  times, resulting in  $n_{\text{exp}}$  sets of sensitivity indices (one set for each configuration).

- **Approach B.** The output of interest is the full output vector:

$$\mathbf{Y}(\Theta) = \left[ Y_i(\Theta) \right]_{1 \leq i \leq n_{\text{exp}}}. \tag{2.5}$$

Here, the output object is a random vector with  $n_{\text{exp}}$  components. This constitutes a genuine methodological constraint, as many GSA methods do not naturally handle vector outputs.

- **Approach C.** The output of interest is the (generalized) least-squares criterion:

$$L(\Theta) = \left( \mathbf{z} - \mathbf{Y}(\Theta) \right)^t \Sigma_{\epsilon}^{-1} \left( \mathbf{z} - \mathbf{Y}(\Theta) \right). \tag{2.6}$$

These three approaches provide complementary insights, as each has its own scope of application. Approaches A and B are particularly suitable when a metamodel needs to be constructed. Screening is carried out with the goal of selecting the parameters that offer the best predictive performance. By contrast, in scenarios where no metamodel is required, Approach C should be preferred because the analyzed output is more directly aligned with the final objective of Bayesian calibration. Indeed, under the assumption of Gaussian noise, minimizing the least-squares criterion is equivalent to maximizing the likelihood which is a central component of Bayesian inference. Consequently, it will be interesting to compare the results of Approaches A and B against those of Approach C.

Only few GSA methods are able to accommodate all of these constraints. Derivate-based techniques [20, 24] are not suitable, since they require specific experimental designs to estimate model derivatives. Linear methods [25] are compatible with a prescribed Monte Carlo design, but they fail to detect nonlinear effects and interactions, and therefore cannot be regarded as reliable screening alternatives.

The Sobol' approach [23], also known as ANOVA (ANalysis Of VAriance) deserves a more thorough discussion. For a given input variable, the total-order Sobol' index equals zero if and only if the model output does not functionally depend on that input [49]. This property makes the total-order Sobol' indices particularly well-suited for screening. Moreover, extensions to vector and functional outputs have been proposed [50, 51], thereby broadening their scope. However, their estimation is not straightforward in the given-data context. In fact, the most natural way of estimating Sobol' indices is to generate a pick-freeze design [52] and this often results in a prohibitive simulation budget.

From a given Monte Carlo sample, two strategies can be considered to estimate the total-order Sobol' indices. One option is to train, validate, and emulate a metamodel, for instance a GP model as in [53]. Of course, the accuracy of the resulting estimates strongly depends on the predictive performance of the metamodel, which can be challenging to achieve in high dimensions, as previously discussed. Other option is to employ the kernel-based estimator [54], obtained by plugging the Nadaraya–Watson regression estimator into the naive Sobol' estimator. Again, this approach is difficult to apply in high dimensions, this time because the assumptions made on the smoothing kernels become increasingly restrictive and often unrealistic. In the Bayesian calibration literature, parameter screening using Sobol' indices is typically restricted to applications where either the computational cost or the dimensionality of the parameter space is low.

Dependence measures form another important family of GSA methods [28, 55, 56]. For each input-output pair, dependence is measured by quantifying the discrepancy between the joint distribution (capturing the true dependence) and the product of the marginal distributions (which would arise under independence). This framework brings together sensitivity measures that may target fundamentally different aspects of the input-output dependence. Rank-based measures such as Spearman's  $\rho$  or Kendall's  $\tau$  are easy to estimate but do not characterize independence, whereas divergence-based measures can detect any form of dependence but rely on density estimation and are poorly suited to multidimensional outputs.

These limitations are overcome by the Hilbert-Schmidt independence criterion (HSIC) [30], a kernel-based dependence measure that transports the distributions into a reproducing kernel Hilbert space (RKHS) and computes the squared norm between their image functions. The use of HSIC indices for screening purposes was popularized by [28]. The method was subsequently extended in [29, 57, 58] and has now become a widely adopted standard practice [59–61]. Several HSIC-based methods have recently been developed for a variety of contexts, including fuzzy systems, knowledge-based systems, and chemometrics [62–64]. However, to the best of our knowledge, HSIC indices have never been applied to parameter screening in Bayesian calibration. This is rather surprising given that they meet all the necessary requirements.

- Under some mild conditions on the kernels, the HSIC is equal to zero if and only if the two random objects under study are independent of each other. However, it is not straightforward to determine which inputs have low influence on the output because each HSIC index follows its own maximum mean discrepancy (MMD) metric, which depends on the two kernels and the two marginal distributions. Consequently, for each input-output pair, a test of independence must be run to decide whether the HSIC value is sufficiently close to zero for accepting input-output independence. The test returns a  $p$ -value that allows controlling Type-I error (*i.e.*, the probability of incorrectly selecting an input whose influence should be neglected).
- The HSIC can be accurately estimated from a small amount of input-output data.
- The HSIC is well-defined for any type<sup>2</sup> of output objects, provided that an appropriate kernel is available. Moreover, the estimation procedure remains unchanged.

Further details on the HSIC methodology are provided in Appendix C. The basic concepts are presented in Appendix C.1. This includes the definition via the MMD, the alternative expression in terms of kernel-based

<sup>2</sup>In the seminal paper [65], it is only assumed that the two random objects lie in (possibly different) separable spaces.

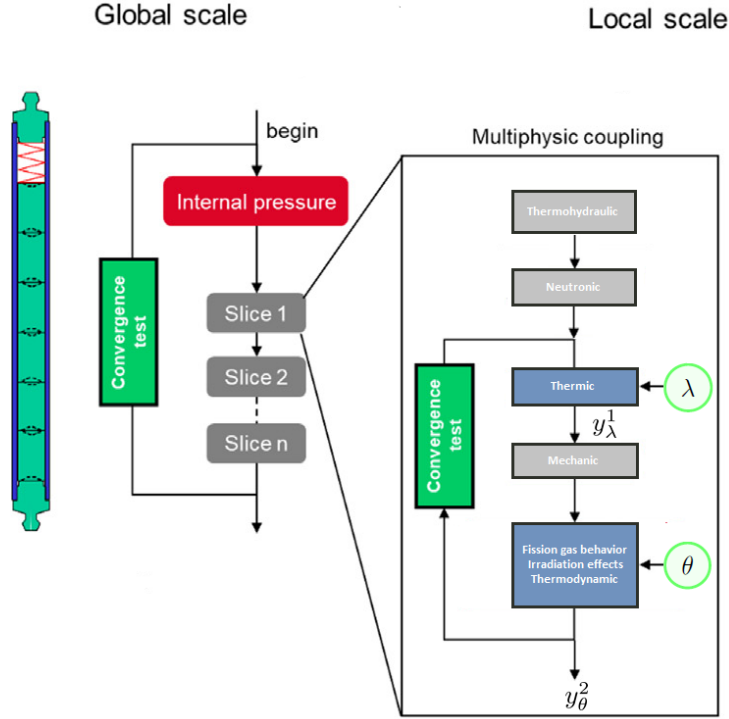


FIGURE 1. Simplified representation of the multiphysics coupling at work in the ALCYONE application [31]. The thermal model  $y_\lambda^1$  feeds into the fission gas behavior model  $y_\theta^2$ . Both are depicted as blue boxes.

moments, and estimation using either a U- or V-statistic estimator. Then, HSIC-based independence testing is outlined in Appendix C.3, with emphasis on the three possible test procedures. Finally, Appendix C.5 discusses implementation aspects in the specific context of Bayesian calibration. In particular, guidelines are provided on how to choose kernels depending on the selected approach.

### 3. APPLICATION TO THE FISSION GAS BEHAVIOR MODEL

In this section, we apply the HSIC methodology to a numerical model from the nuclear industry whose parameters need to be calibrated. More specifically, the goal is to perform parameter screening for the fission gas behavior model of the ALCYONE application [31].

#### 3.1. Presentation of the ALCYONE application

ALCYONE is a multidimensional fuel performance simulation code developed at CEA within the PLEIADES software environment [66]. It is specifically designed to simulate the thermo-mechanical behavior of nuclear fuel rods under irradiation in pressurized water reactors (PWRs). A simplified diagram is shown in Figure 1. Although the application integrates numerical models from five physical disciplines, we focus here only on the chaining between two models (represented in the diagram as blue rectangular boxes): the thermal model (upstream model) and the fission gas behavior model (downstream model). The thermal model simulates how the temperature evolves inside the fuel rod during fission and outputs the corresponding temperature field throughout the rod. This temperature field feeds the fission gas behavior model which predicts the release and redistribution of fission gases within the fuel rod.

Fission gases produced during irradiation are distributed in three forms: dissolved in the fuel grains, trapped in intragranular bubbles, and trapped in intergranular bubbles. When the intergranular bubbles at the grain boundaries of uranium dioxide connect to form continuous pathways, a phenomenon known as fission gas release occurs into the free volume of the fuel rod. Under high burnup conditions, this process may induce fuel restructuring (finer grain structure, increased porosity), which strongly affects gas behavior. The fission gas behavior model simulates these phenomena and provides a scalar output of interest, the released gas fraction (RGF), which quantifies the proportion of gas released from the fuel.

The thermal model will be denoted as  $y_\lambda^1$  where  $\lambda \in \mathcal{D}_\Lambda \subseteq \mathbb{R}$  is the thermal conductivity. Experimental data were used to calibrate this parameter in previous work [32]. The fission gas behavior model will be denoted by  $y_\theta^2$  where  $\theta \in \mathcal{D}_\Theta \subseteq \mathbb{R}^p$  are the calibration parameters. Accurate calibration of these  $p = 11$  parameters, while rigorously accounting for the posterior uncertainty on  $\lambda$ , is a critical step to guarantee the reliability of ALCYONE simulations.

For fixed  $\lambda$ , the model  $y_\lambda^1$  takes a control variable  $\mathbf{x} \in \mathcal{D}_\mathbf{X}$  describing the fuel rod configuration (linear heat rate, burnup rate, geometry) and returns the temperature field throughout the rod. For fixed  $\theta$ , the model  $y_\theta^2$  takes this temperature field as input and returns the RGF. The chaining between the two models can be summarized by the following deterministic function:

$$\begin{aligned} y : \mathcal{D}_\mathbf{X} \times \mathcal{D}_\Lambda \times \mathcal{D}_\Theta &\longrightarrow \mathcal{D}_Y \\ (\mathbf{x}, \lambda, \theta) &\longmapsto y(\mathbf{x}, \lambda, \theta) = [y_\theta^2 \circ y_\lambda^1](\mathbf{x}). \end{aligned} \quad (3.1)$$

At this step, the epistemic uncertainty affecting  $\lambda$  is voluntarily ignored. The conductivity is fixed at its nominal value  $\lambda_{\text{nom}} = 1$ . This reduces the chained model to:

$$\begin{aligned} y_s : \mathcal{D}_\mathbf{X} \times \mathcal{D}_\Theta &\longrightarrow \mathcal{D}_Y \\ (\mathbf{x}, \theta) &\longmapsto y_s(\mathbf{x}, \theta) = [y_\theta^2 \circ y_{\lambda_{\text{nom}}}^1](\mathbf{x}). \end{aligned} \quad (3.2)$$

The uncertainty on  $\lambda$  will be reintroduced in the study from Section 3.4 onwards. For now, the goal is to identify the components of  $\theta$  that have the greatest influence on the RGF  $y_s(\mathbf{x}, \theta)$ .

### 3.2. Strategies considered for parameter screening

In the existing literature, the HSIC approach is mostly used to reduce dimensionality before building a metamodel [39, 58, 67]. Here, the presence of control variables  $\mathbf{x}$  encoding multiple experimental configurations introduces an additional layer of complexity, which motivates the use of Approaches A, B and C presented in Section 2.2.

As the calibration parameters  $\theta$  of the fission gas behavior model are uncertain, they are modeled by a random vector  $\Theta$  following a distribution  $\mathbb{P}_\Theta$ , with support  $\mathcal{D}_\Theta \subseteq \mathbb{R}^p$  and density  $\pi(\theta)$ . The components  $\Theta_j$  of  $\Theta$  are assumed to be independent. As parameter screening precedes model calibration,  $\mathbb{P}_\Theta$  corresponds to the prior distribution of  $\theta$ , informed by expert knowledge.

One key advantage of the HSIC approach is that the related indices can be estimated using a Monte Carlo design of (relatively small) size  $n$ . In the parameter space, the input data can be written as:

$$D_{\Theta, n} := \{\Theta^{(l)}\}_{1 \leq l \leq n_{\text{exp}}} \in \mathbb{R}^{n \times p} \quad \text{with} \quad \Theta^{(l)} = [\Theta_1^{(l)}, \dots, \Theta_p^{(l)}]^t \in \mathbb{R}^p. \quad (3.3)$$

The fission gas behavior model is then evaluated for all pairs  $(\mathbf{x}_i, \Theta^{(l)})$  where  $\mathbf{x}_i$  is the  $i$ -th experimental configuration and  $\Theta^{(l)}$  is the  $l$ -th parameter value. For a fixed configuration  $\mathbf{x}_i$ , the  $n$ -sample obtained by evaluating  $y_s(\mathbf{x}_i, \cdot)$  on all parameter values  $\Theta^{(l)}$  yields:

$$D_{Y_i, n} := \{Y_i^{(l)}\}_{1 \leq l \leq n} \quad \text{with} \quad Y_i^{(l)} = y_s(\mathbf{x}_i, \Theta^{(l)}). \quad (3.4)$$

Combining all such samples  $D_{Y_i, n}$  forms the complete output dataset:

$$\begin{aligned} D_{\mathbf{Y}, n} &:= \{D_{Y_i, n}\}_{1 \leq i \leq n_{\text{exp}}} \\ &= \left\{ \mathbf{Y}^{(l)} \right\}_{1 \leq l \leq n} \quad \text{with} \quad \mathbf{Y}^{(l)} = \left[ Y_i^{(l)} \right]_{1 \leq i \leq n_{\text{exp}}} \end{aligned} \quad (3.5)$$

The total simulation budget amounts to  $n_{\text{exp}} \times n$  evaluations of  $y_s$ . At this stage, all the data required to implement the three screening approaches are available. The next steps depend on the selected screening strategy.

- **Approach A.** GSA is performed on the RGF of each fuel rod, that is on the output of  $y_s$  for a given experimental configuration  $\mathbf{x}_i$ :

$$Y_i = Y_i(\Theta) = y(\mathbf{x}_i, \lambda_{\text{nom}}, \Theta) =: f_i^A(\Theta) =: \tilde{f}_i^A(\Theta, \lambda_{\text{nom}}). \quad (3.6)$$

The function  $f_i^A$ , which encapsulates the chained model  $y$ , is introduced to establish a direct link between the analyzed output and the parameters  $\Theta$ . The function  $\tilde{f}_i^A$  is simply a rewriting of  $f_i^A$  that makes the role of  $\lambda$  explicit. To estimate HSIC indices, the following data are required:

$$D_{i, n}^A := \{D_{\Theta, n}; D_{Y_i, n}\}. \quad (3.7)$$

The resulting HSIC estimates are denoted by  $\hat{H}_{ij}$  for a fixed  $i$  and  $1 \leq j \leq p$ . These quantities are then used in independence tests. For each parameter  $\Theta_j$ , the null hypothesis ( $H_0$ ):  $\Theta_j \perp\!\!\!\perp Y_i$  is tested. For this, a test procedure (chosen depending on the sample size) is run and provides a  $p$ -value  $\hat{p}_{ij}$ . The significance level (controlling Type-I error) is set to  $\alpha = 5\%$ . If  $\hat{p}_{ij} < \alpha$ , the null hypothesis is rejected and  $\Theta_j$  is considered influential on  $Y_i$ . Repeating the procedure across all rods of the fuel assembly produces  $n_{\text{exp}}$  sets of  $p$ -values, which may lead to as many distinct conclusions. To make a global decision, one may compute the detection rate of each parameter  $\Theta_j$ , defined as the percentage of rods for which that parameter is detected as influential:

$$\rho_j := \sum_{i=1}^{n_{\text{exp}}} \mathbf{1}_{\{\hat{p}_{ij} < \alpha\}}. \quad (3.8)$$

A parameter  $\Theta_j$  is then selected if its detection rate  $\rho_j$  exceeds a (user-chosen) threshold, for example  $\beta = 80\%$ . In other words, a parameter is considered influential on the fuel assembly if it is influential on at least a fraction  $\beta$  of the fuel rods.

- **Approach B.** GSA is performed on the vector of all RGFs, that is on the vector:

$$\mathbf{Y} = \mathbf{Y}(\Theta) = [Y_i]_{1 \leq i \leq n_{\text{exp}}} =: f^B(\Theta) =: \tilde{f}^B(\Theta, \lambda_{\text{nom}}). \quad (3.9)$$

To estimate HSIC indices, the following data are required:

$$D_n^B := \{D_{\Theta, n}; D_{\mathbf{Y}, n}\}. \quad (3.10)$$

The null hypothesis becomes ( $H_0$ ):  $\Theta_j \perp\!\!\!\perp \mathbf{Y}$ . The HSIC methodology naturally extends to the case of a vector output, provided that a characteristic kernel on  $\mathbb{R}^{n_{\text{exp}}}$  is used to handle  $\mathbf{Y}$  [68, 69]. A first option is to take the multivariate generalization of standard kernels (such as Gaussian and Matérn kernels). A second option is to equip each component  $Y_i$  with a characteristic kernel on  $\mathbb{R}$  (for instance a Gaussian kernel) and then take the tensor product of these kernels. A third option is to consider kernel functions specifically designed for high-dimensional data, such as PCA kernels [58]. We adopted the second option,

as tensor-product kernels are characteristic in most cases (see Appendix A.3) while being very flexible (with a separate bandwidth parameter for each component  $Y_i$ ). A more detailed discussion of the pros and cons of the different kernels is provided in Appendix C.5. Once an appropriate kernel is chosen for  $\mathbf{Y}$ , the standard HSIC procedure can be applied (see Appendix C). This produces HSIC estimates  $\widehat{H}_j$  and corresponding  $p$ -values  $\widehat{p}_j$ , from which a global decision can ultimately be made.

- **Approach C.** GSA is performed on the least-squares criterion (measuring the discrepancy between the observed and simulated RGFs):

$$L = L(\Theta) = \left( \mathbf{z} - \mathbf{Y} \right)^t \Sigma_\epsilon^{-1} \left( \mathbf{z} - \mathbf{Y} \right) =: f^C(\Theta) =: \widetilde{f}^C(\Theta, \lambda_{\text{nom}}). \quad (3.11)$$

To estimate HSIC indices, the following data are required:

$$D_n^C := \{D_{\Theta,n}; D_{L,n}\} \quad \text{with} \quad D_{L,n} = \{L^{(l)}\}_{1 \leq l \leq n} \quad \text{and} \quad L^{(l)} = L(\Theta^{(l)}). \quad (3.12)$$

Since  $L$  is scalar, applying the HSIC methodology is straightforward. This time, the null hypothesis is  $(H_0) : \Theta_j \perp\!\!\!\perp L$ . As in Approach B, one obtains HSIC estimates  $\widehat{H}_j$ , and then  $p$ -values  $\widehat{p}_j$  from which the influential parameters are selected.

**Remark 3.1.** All three approaches can be carried out without much extra computational effort. Indeed, once  $D_{\mathbf{Y},n}$  is obtained, at the cost of  $n_{\text{exp}} \times n$  evaluations of  $y_s$ , everything that follows may be viewed as inexpensive post-processing. In particular, the datasets built at equations (3.7), (3.10) and (3.12) are directly derived from  $D_{\Theta,n}$  and  $D_{\mathbf{Y},n}$ .

**Remark 3.2.** Rank-based dependence measures also offer a robust and distribution-free alternative to kernel-based dependence measures. However, two important limitations motivated our decision not to use them. First, unlike the HSIC built on two characteristic kernels, rank correlations (such as Spearman's  $\rho$  or Kendall's  $\tau$ ) do not characterize independence, as they mainly capture monotonic forms of dependence. Second, extending rank-based measures to vectors or non-Euclidean objects is non-trivial, whereas this is immediate for the HSIC, provided that appropriate kernels are available. For these two reasons, the HSIC appears better suited for parameter screening, although rank-based measures could in principle be used as complementary diagnostics.

**Remark 3.3.** Due to the strong similarities between the three approaches, developments will be conducted within a unified framework where the output of interest is denoted as:

$$G_{\text{nom}} := f(\Theta) = \widetilde{f}(\Theta, \lambda_{\text{nom}}).$$

The superscripts  $A$ ,  $B$  and  $C$  will be omitted, except when needed to avoid ambiguity.

### 3.3. Results

In this section, we apply the three proposed approaches to the fission gas behavior model  $y_s$  and compare their results. The experimental input configurations  $\{\mathbf{x}_i\}_{i=1}^{n_{\text{exp}}}$  correspond to the specifications of the  $n_{\text{exp}} = 40$  fuel rods within the fuel assembly.

Due to computational constraints, the size of the Monte Carlo design  $D_{\Theta,n}$  is restricted to  $n = 200$ . Although this sample size is already large, it is not yet sufficient to use the asymptotic test [30]. Instead, the permutation-based test [29] and the non-asymptotic Gamma test [58] are employed. The significance level is set to  $\alpha = 5\%$ .

The results of all independence tests performed in Approach A are shown in Figure 2. The parameters  $\Theta_j$  are plotted on the  $y$ -axis, the outputs  $Y_i$  on the  $x$ -axis, sorted by decreasing burnup rates. Green points indicate that  $(H_0) : \Theta_j \perp\!\!\!\perp Y_i$  was rejected (statistical dependence detected), whereas gray points indicate that  $(H_0)$  was accepted (no influence). We specifically distinguish results obtained with U-statistics from those obtained with V-statistics. Decisions are made from the  $p$ -values returned by the permutation test. The non-asymptotic Gamma test yields identical decisions.

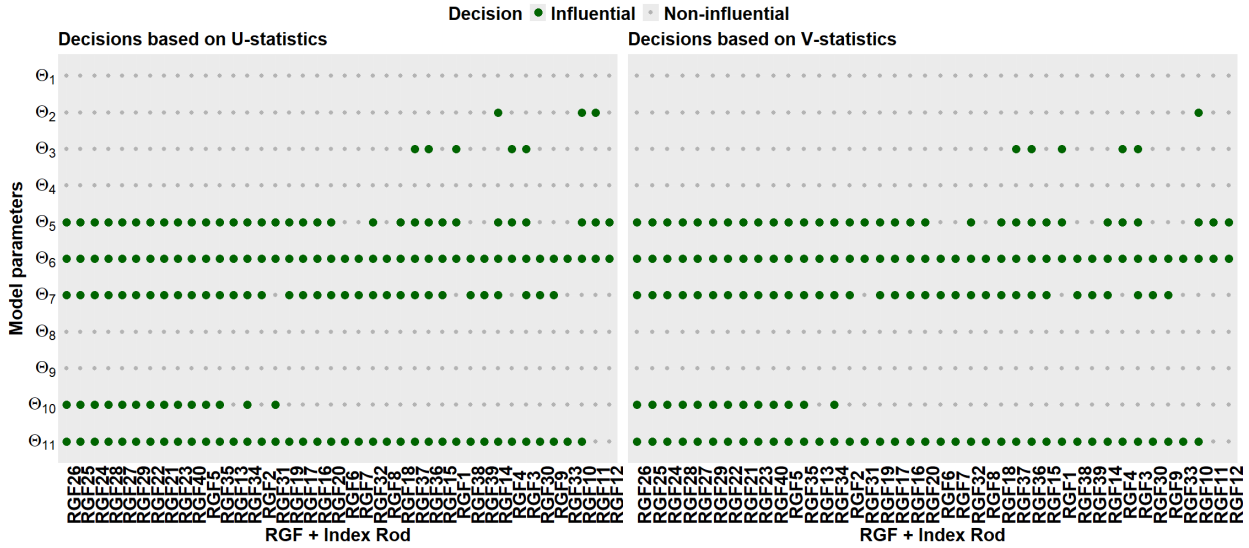


FIGURE 2. Results obtained for Approach A. For each fuel rod (characterized by the experimental configuration  $x_i$ ), the influence (weak or strong) of each parameter  $\Theta_j$  is assessed using an HSIC test of independence. The test statistic is either the U-statistic estimator (left) or the V-statistic estimator (right). In both cases, a permutation-based test procedure is used. The significance level is set to  $\alpha = 5\%$ . The detection rates (percentages of rods where the parameters are detected as influential) are reported in Table 1.

TABLE 1. Comparison of the screening results obtained for the three considered approaches. Only the results associated with U-statistics are shown.

| Parameters    | Results  | Detection rates |       | $p$ -values |
|---------------|----------|-----------------|-------|-------------|
|               | Approach | A               | B     | C           |
| $\Theta_1$    |          | 0%              | 0.972 | 0.538       |
| $\Theta_2$    |          | 7.5%            | 0.25  | 0.741       |
| $\Theta_3$    |          | 12.5%           | 0.372 | 0.127       |
| $\Theta_4$    |          | 0%              | 0.909 | 0.801       |
| $\Theta_5$    |          | 80%             | 0     | 0           |
| $\Theta_6$    |          | 100%            | 0     | 0.006       |
| $\Theta_7$    |          | 82.5%           | 0     | 0           |
| $\Theta_8$    |          | 0%              | 0.869 | 0.757       |
| $\Theta_9$    |          | 0%              | 0.856 | 0.729       |
| $\Theta_{10}$ |          | 35%             | 0.022 | 0.081       |
| $\Theta_{11}$ |          | 95%             | 0     | 0           |

The parameters  $\Theta_5$ ,  $\Theta_6$ ,  $\Theta_7$ , and  $\Theta_{11}$  are found to be influential on the RGF of all fuel rods. By contrast, the parameter  $\Theta_{10}$  seems to be influential on a small number of RGFs, those corresponding to the highest burnup rates. The remaining parameters appear to have no influence on the RGFs, regardless of the fuel rod considered. To summarize these results into a single set of indices, the detection rates defined in equation (3.8)

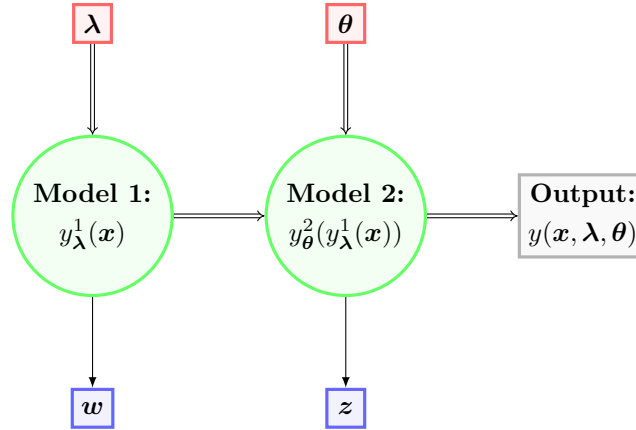


FIGURE 3. Diagram of the chaining between two numerical model: the upstream model  $y_{\lambda}^1$  and the downstream model  $y_{\theta}^2$ . The calibration parameters are put above the models, and the available experimental data below. Double arrows indicate uncertainty propagation, while single arrows indicate data assimilation.

are computed and reported in Column A of Table 1. For the four most influential parameters, the detection rates exceed  $\beta = 80\%$ , whereas for  $\Theta_{10}$ , it is only  $\rho_{10} = 35\%$ .

Table 1 also provides the  $p$ -values obtained for Approaches B and C. One can see that Approach B selects the five expected parameters, including  $\Theta_{10}$ . Approach C selects the four predominant parameters, but rejects  $\Theta_{10}$ . However, it should be noted that the related  $p$ -value is quite close to  $\alpha = 5\%$ , which confirms that  $\Theta_{10}$  exerts only a minor influence. Overall, it is reassuring to see that the three screening strategies lead to the consistent conclusions. Furthermore, the results are largely insensitive to the HSIC estimator (U- vs. V-statistic) and the test procedure (permutations vs. Gamma approximation).

### 3.4. Accounting for the upstream model uncertainty

For a rigorous uncertainty treatment in the simulation chain, the thermal conductivity  $\lambda$  cannot be regarded as a constant parameter. Indeed, it is subject to an epistemic uncertainty resulting from the calibration of the thermal model. Naturally, this uncertainty must be incorporated into the study of the fission gas behavior model by adapting both the Bayesian calibration and GSA methodologies. This bi-level uncertainty framework is illustrated in Figure 3. For the sake of generality, the thermal and fission gas behavior models are labeled as Models 1 and 2. In addition, the calibration parameter(s) of Model 1 are denoted by  $\lambda$  to allow for a more general setting, where multiple parameters may be involved.

It is assumed that Model 1 has already been calibrated with some experimental data  $w$ . For ALCYONE, this corresponds to what was done for the thermal model in [32]. To account for the posterior uncertainty on  $\lambda$ , we introduce a random vector  $\Lambda$  with distribution  $\mathbb{P}_{\Lambda}$ . The challenge in calibrating Model 2 is to make optimal use of the posterior knowledge on  $\lambda$ , the prior knowledge on  $\theta$ , and the experimental data  $z$ . The main objective is to estimate the conditional density:

$$\pi(\theta \mid \lambda, z) \propto \mathcal{L}(z \mid \theta, \lambda) \pi(\theta). \quad (3.13)$$

This inference task is called *conditional Bayesian calibration*.

**Remark 3.4.** Before conditional calibration, the random vectors  $\Lambda$  (posterior uncertainty on the parameters of Model 1) and  $\Theta$  (prior uncertainty on the parameters of Model 2) may reasonably be assumed independent.

After conditional calibration, the joint posterior density  $\pi(\boldsymbol{\theta}, \boldsymbol{\lambda} \mid \mathbf{z})$  generally does not factorize, meaning that  $\boldsymbol{\Lambda}$  and  $\boldsymbol{\Theta}$  are no longer independent.

A conditional Bayesian calibration algorithm was developed in [70] to estimate the conditional posterior density for each value  $\boldsymbol{\lambda}$  drawn from the posterior distribution  $\mathbb{P}_{\boldsymbol{\Lambda}}$ . This algorithm relies on experimental data  $\mathbf{z}$  together with simulated data generated according to a nested Monte Carlo design. At the outer level, an  $m$ -sample of  $\boldsymbol{\Lambda} \sim \mathbb{P}_{\boldsymbol{\Lambda}}$  is generated:

$$D_{\boldsymbol{\Lambda},m} := \left\{ \boldsymbol{\Lambda}^{(k)} \right\}_{1 \leq k \leq m} . \quad (3.14)$$

At the inner level, two sampling strategies can be considered.

- (S1) An  $n$ -sample of  $\boldsymbol{\Theta} \sim \mathbb{P}_{\boldsymbol{\Theta}}$  is generated for each value  $\boldsymbol{\Lambda}^{(k)}$  in  $D_{\boldsymbol{\Lambda},m}$ . This amounts to introducing  $m$  independent copies  $\left\{ \boldsymbol{\Theta}^{(k)} \right\}_{k=1}^m$  of  $\boldsymbol{\Theta}$  and then generating an  $n$ -sample for each copy:

$$D_{\boldsymbol{\Theta}^{(k)},n} := \left\{ \boldsymbol{\Theta}^{(kl)} \right\}_{1 \leq l \leq n} . \quad (3.15)$$

The designs  $D_{\boldsymbol{\Theta}^{(k)},n}$  are mutually independent. They are also independent of the design  $D_{\boldsymbol{\Lambda},m}$ . The chained model  $y$  is evaluated to compute all relevant model outputs:

$$D_{\mathbf{Y}^{(k)},n} := \left\{ \mathbf{Y}^{(kl)} \right\}_{1 \leq l \leq n} \quad \text{with} \quad \mathbf{Y}^{(kl)} = \left[ y(\mathbf{x}_i, \boldsymbol{\Lambda}^{(k)}, \boldsymbol{\Theta}^{(kl)}) \right]_{1 \leq i \leq n_{\text{exp}}} . \quad (3.16)$$

$D_{\mathbf{Y}^{(k)},n}$  is an  $n$ -sample of the output vector  $\mathbf{Y}^{(k)}$  obtained with  $\boldsymbol{\Lambda}^{(k)}$ :

$$\mathbf{Y}^{(k)} = \left[ Y_i^{(k)} \right]_{1 \leq i \leq n_{\text{exp}}} \quad \text{with} \quad Y_i^{(k)} := y(\mathbf{x}_i, \boldsymbol{\Lambda}^{(k)}, \boldsymbol{\Theta}^{(k)}) . \quad (3.17)$$

For convenience, we also introduce notations for the final nested design:

$$D_{m,n}^{S_1} := \left\{ D_{\boldsymbol{\Theta}^{(k)},n} \right\}_{1 \leq k \leq m} \quad \text{with} \quad D_{\boldsymbol{\Theta}^{(k)},n} := \left\{ D_{\boldsymbol{\Theta}^{(k)},n}; D_{\mathbf{Y}^{(k)},n} \right\} . \quad (3.18)$$

- (S2) A unique  $n$ -sample of  $\boldsymbol{\Theta} \sim \mathbb{P}_{\boldsymbol{\Theta}}$  is generated and used for all values  $\boldsymbol{\Lambda}^{(k)}$ :

$$D_{\boldsymbol{\Theta},n} := \left\{ \boldsymbol{\Theta}^{(l)} \right\}_{1 \leq l \leq n} . \quad (3.19)$$

The designs  $D_{\boldsymbol{\Lambda},m}$  and  $D_{\boldsymbol{\Theta},n}$  remain independent. Evaluating the chained model  $y$  for  $\boldsymbol{\Lambda}^{(k)}$  and all points in  $D_{\boldsymbol{\Theta},n}$  produces the corresponding output data:

$$D_{\mathbf{Y}^{(k)},n} := \left\{ \mathbf{Y}^{(kl)} \right\}_{1 \leq l \leq n} \quad \text{with} \quad \mathbf{Y}^{(kl)} = \left[ y(\mathbf{x}_i, \boldsymbol{\Lambda}^{(k)}, \boldsymbol{\Theta}^{(l)}) \right]_{1 \leq i \leq n_{\text{exp}}} . \quad (3.20)$$

The complete nested design is then given by:

$$D_{m,n}^{S_2} := \left\{ D_{\boldsymbol{\Theta}^{(k)},n} \right\}_{1 \leq k \leq m} \quad \text{with} \quad D_{\boldsymbol{\Theta}^{(k)},n} := \left\{ D_{\boldsymbol{\Theta},n}; D_{\mathbf{Y}^{(k)},n} \right\} . \quad (3.21)$$

Both strategies incur the same computational cost. More precisely, they require  $N_1 = n_{\text{exp}} \times m$  evaluations of  $y_{\boldsymbol{\lambda}}^1$  and  $N_2 = n_{\text{exp}} \times m \times n$  evaluations of  $y_{\boldsymbol{\theta}}^2$ .

Conditional calibration performs well in low dimensions but is strongly affected by the curse of dimensionality. Therefore, the number of parameters must be reduced as much as possible. In this bi-level uncertainty framework, a screening method must satisfy two requirements.

1. The method must fully account for the uncertainty coming from the parameters of Model 1. In particular, a parameter  $\Theta_j$  may be influential for some realizations of  $\mathbf{\Lambda}$  while being non-influential for other ones. Consequently,  $\Theta_j$  should only be discarded if it can be shown that it has a negligible impact on the output of Model 2 for all possible realizations of  $\mathbf{\Lambda}$ .
2. The method must accommodate the nested design used for conditional calibration.

The next section is dedicated to the development of such a methodology.

#### 4. SCREENING METHODOLOGY FOR CONDITIONAL CALIBRATION

This section presents the main contribution of the paper. We develop a screening method that identifies the downstream model parameters that must be included in conditional Bayesian calibration. The section is divided into two parts, highlighting the analogy with the HSIC-based methodology (inference first, statistical testing second). Section 4.1 defines a sensitivity measure tailored to the constraints of our problem. In essence, this involves integrating the HSIC with respect to  $\mathbb{P}_{\mathbf{\Lambda}}$ . We show that the resulting criterion achieves the intended goal. First, it satisfies a mathematical property analogous to the HSIC: it vanishes if and only if the calibration parameter is independent of the output for almost all values of  $\boldsymbol{\lambda}$ . Second, it can be naturally estimated from the available data, with four possible estimators exhibiting attractive properties. However, these estimates alone are insufficient for screening, just as with HSIC indices. This motivates the use of hypothesis testing in Section 4.2. We explain how to formulate the problem and we show that some HSIC-based test procedures can be easily adapted to our bi-level uncertainty context.

##### 4.1. Measuring sensitivity in the bi-level uncertainty framework

In Section 3, parameter screening was performed in a simplified context where  $\boldsymbol{\lambda} = \boldsymbol{\lambda}_{\text{nom}}$ . The HSIC was used to measure the dependence between any parameter  $\Theta_j$  and any output of interest  $G_{\text{nom}} = \tilde{f}(\boldsymbol{\Theta}, \boldsymbol{\lambda}_{\text{nom}})$ . One can write:

$$H_j(\boldsymbol{\lambda}_{\text{nom}}) := \text{HSIC}(\Theta_j, G_{\text{nom}}) . \quad (4.1)$$

Let  $K_{\theta_j} : \mathcal{D}_{\Theta_j} \times \mathcal{D}_{\Theta_j} \rightarrow \mathbb{R}$  and  $K_g : \mathcal{D}_G \times \mathcal{D}_G \rightarrow \mathbb{R}$  denote covariance kernels suitable for handling the variables  $\Theta_j$  and  $G_{\text{nom}}$ . The quantity  $H_j(\boldsymbol{\lambda}_{\text{nom}})$  can be expressed only in terms of kernel-based moments:

$$\begin{aligned} H_j(\boldsymbol{\lambda}_{\text{nom}}) &= \mathbb{E}_{\boldsymbol{\Theta}\boldsymbol{\Theta}' } \left[ K_{\theta_j}(\Theta_j, \Theta'_j) K_g \left( \tilde{f}(\boldsymbol{\Theta}, \boldsymbol{\lambda}_{\text{nom}}), \tilde{f}(\boldsymbol{\Theta}', \boldsymbol{\lambda}_{\text{nom}}) \right) \right] \\ &\quad + \mathbb{E}_{\boldsymbol{\Theta}\boldsymbol{\Theta}'} \left[ K_{\theta_j}(\Theta_j, \Theta'_j) \right] \mathbb{E}_{\boldsymbol{\Theta}\boldsymbol{\Theta}'} \left[ K_g \left( \tilde{f}(\boldsymbol{\Theta}, \boldsymbol{\lambda}_{\text{nom}}), \tilde{f}(\boldsymbol{\Theta}', \boldsymbol{\lambda}_{\text{nom}}) \right) \right] \\ &\quad - 2 \mathbb{E}_{\boldsymbol{\Theta}\boldsymbol{\Theta}'' } \left[ K_{\theta_j}(\Theta_j, \Theta'_j) K_g \left( \tilde{f}(\boldsymbol{\Theta}, \boldsymbol{\lambda}_{\text{nom}}), \tilde{f}(\boldsymbol{\Theta}'', \boldsymbol{\lambda}_{\text{nom}}) \right) \right] . \end{aligned} \quad (4.2)$$

where  $\boldsymbol{\Theta}$ ,  $\boldsymbol{\Theta}'$  and  $\boldsymbol{\Theta}''$  are independent and identically distributed according to  $\mathbb{P}_{\boldsymbol{\Theta}}$ . We refer the reader to Appendix C for further details on the definition, reformulation, and estimation of HSIC indices. The definition of  $H_j(\boldsymbol{\lambda}_{\text{nom}})$  naturally extends to any  $\boldsymbol{\lambda} \in \mathcal{D}_{\mathbf{\Lambda}}$ :

$$\begin{aligned} H_j(\boldsymbol{\lambda}) &= \text{HSIC}(\Theta_j, G_{\boldsymbol{\lambda}}) \\ &= \text{HSIC}(\Theta_j, G \mid \mathbf{\Lambda} = \boldsymbol{\lambda}) \end{aligned} \quad \text{with} \quad \begin{cases} G_{\boldsymbol{\lambda}} := \tilde{f}(\boldsymbol{\Theta}, \boldsymbol{\lambda}) \\ G := \tilde{f}(\boldsymbol{\Theta}, \mathbf{\Lambda}). \end{cases}$$

The notation introduced in the second line means that  $\mathbf{\Lambda}$  must be set to the value  $\boldsymbol{\lambda}$  before computing the HSIC. After introducing the function:

$$\begin{aligned} h(\boldsymbol{\theta}^1, \boldsymbol{\theta}^2, \boldsymbol{\theta}^3, \boldsymbol{\theta}^4, \boldsymbol{\lambda}) &= K_{\theta_j}(\theta_j^1, \theta_j^2) K_g(\tilde{f}(\boldsymbol{\theta}^1, \boldsymbol{\lambda}), \tilde{f}(\boldsymbol{\theta}^2, \boldsymbol{\lambda})) \\ &\quad + K_{\theta_j}(\theta_j^1, \theta_j^2) K_g(\tilde{f}(\boldsymbol{\theta}^3, \boldsymbol{\lambda}), \tilde{f}(\boldsymbol{\theta}^4, \boldsymbol{\lambda})) \\ &\quad - 2 K_{\theta_j}(\theta_j^1, \theta_j^2) K_g(\tilde{f}(\boldsymbol{\theta}^1, \boldsymbol{\lambda}), \tilde{f}(\boldsymbol{\theta}^3, \boldsymbol{\lambda})), \end{aligned} \quad (4.3)$$

$H_j(\boldsymbol{\lambda})$  can be rewritten as a single expectation:

$$H_j(\boldsymbol{\lambda}) = \mathbb{E}[h(\boldsymbol{\Theta}^1, \boldsymbol{\Theta}^2, \boldsymbol{\Theta}^3, \boldsymbol{\Theta}^4, \boldsymbol{\lambda})] \quad \text{with} \quad \boldsymbol{\Theta}^1 \perp\!\!\!\perp \boldsymbol{\Theta}^2 \perp\!\!\!\perp \boldsymbol{\Theta}^3 \perp\!\!\!\perp \boldsymbol{\Theta}^4 \sim \mathbb{P}_{\boldsymbol{\Theta}}. \quad (4.4)$$

Assuming that  $K_{\theta_j}$ ,  $K_g$  and  $\tilde{f}$  are measurable,  $h$  is measurable as well, and Tonelli's theorem ensures that the function  $H_j : \mathcal{D}_{\boldsymbol{\Lambda}} \rightarrow \mathbb{R}$  is measurable. This provides a sound theoretical framework for propagating  $\mathbb{P}_{\boldsymbol{\Lambda}}$  to the HSIC index. The random variable  $H_j(\boldsymbol{\Lambda})$  is therefore well-defined. As already pointed out in Remark 3.4, the vectors  $\boldsymbol{\Lambda}$  and  $\boldsymbol{\Theta}$  are independent, which allows us to rewrite  $H_j(\boldsymbol{\Lambda})$  in terms of conditional expectations:

$$H_j(\boldsymbol{\Lambda}) = \text{HSIC}(\boldsymbol{\Theta}_j, G \mid \boldsymbol{\Lambda}) \quad (4.5)$$

$$\begin{aligned} &= \mathbb{E}_{\boldsymbol{\Theta}\boldsymbol{\Theta}'}\left[K_{\theta_j}(\boldsymbol{\Theta}_j, \boldsymbol{\Theta}'_j) K_g(\tilde{f}(\boldsymbol{\Theta}, \boldsymbol{\Lambda}), \tilde{f}(\boldsymbol{\Theta}', \boldsymbol{\Lambda})) \mid \boldsymbol{\Lambda}\right] \\ &\quad + \mathbb{E}_{\boldsymbol{\Theta}\boldsymbol{\Theta}'}\left[K_{\theta_j}(\boldsymbol{\Theta}_j, \boldsymbol{\Theta}'_j)\right] \mathbb{E}_{\boldsymbol{\Theta}\boldsymbol{\Theta}'}\left[K_g(\tilde{f}(\boldsymbol{\Theta}, \boldsymbol{\Lambda}), \tilde{f}(\boldsymbol{\Theta}', \boldsymbol{\Lambda})) \mid \boldsymbol{\Lambda}\right] \\ &\quad - 2 \mathbb{E}_{\boldsymbol{\Theta}\boldsymbol{\Theta}''}\left[K_{\theta_j}(\boldsymbol{\Theta}_j, \boldsymbol{\Theta}'_j) K_g(\tilde{f}(\boldsymbol{\Theta}, \boldsymbol{\Lambda}), \tilde{f}(\boldsymbol{\Theta}'', \boldsymbol{\Lambda})) \mid \boldsymbol{\Lambda}\right]. \end{aligned} \quad (4.6)$$

Equation (4.5) emphasizes that some residual uncertainty, expressed as a measurable function of  $\boldsymbol{\Lambda}$ , remains after the HSIC is computed.

#### 4.1.1. Definition of a $\mathbb{P}_{\boldsymbol{\Lambda}}$ -informed sensitivity measure

To account for the uncertainty coming from the upstream model parameters  $\boldsymbol{\Lambda}$ , a natural idea is to compute the expectation of  $H_j(\boldsymbol{\Lambda})$ . As a consequence, we define:

$$\forall 1 \leq j \leq p, \quad \mathcal{H}_j := \mathbb{E}_{\boldsymbol{\Lambda}}[H_j(\boldsymbol{\Lambda})]. \quad (4.7)$$

Proposition 4.1 shows that  $\mathcal{H}_j$  is a relevant sensitivity measure to detect almost surely (a.s.) any type of statistical dependence between  $\boldsymbol{\Theta}_j$  and  $\tilde{f}(\boldsymbol{\Theta}, \boldsymbol{\Lambda})$ .

**Proposition 4.1.** *Let  $j \in \{1, \dots, p\}$  be fixed. Assume that  $K_{\theta_j}$  and  $K_g$  are two characteristic kernels. It holds that:*

$$\mathcal{H}_j = 0 \iff H_j(\boldsymbol{\Lambda}) = 0 \quad \mathbb{P}_{\boldsymbol{\Lambda}}\text{-a.s.} \iff \boldsymbol{\Theta}_j \perp\!\!\!\perp \tilde{f}(\boldsymbol{\Theta}, \boldsymbol{\Lambda}). \quad (4.8)$$

The detailed proof is provided in Appendix E.1. Since the HSIC is non-negative, the first equivalence is immediate. The second equivalence is slightly more technical. The key argument is that the HSIC, when built using two characteristic kernels, vanishes if and only if the two random objects are independent. Hence, for any  $\boldsymbol{\lambda} \in \mathcal{D}_{\boldsymbol{\Lambda}}$ , if  $H_j(\boldsymbol{\lambda}) = 0$ , then  $\boldsymbol{\Theta}_j \perp\!\!\!\perp G_{\boldsymbol{\lambda}}$ . To prove the direct implication, we verify the factorization property of expectations, while the converse is almost trivial.

**Remark 4.2.** Proposition 4.1 highlights that testing:

$$(H_0) : \Theta_j \perp\!\!\!\perp \tilde{f}(\Theta, \Lambda) \quad (4.9)$$

is not equivalent to testing:

$$(\tilde{H}_0) : \forall \lambda \in \mathcal{D}_\Lambda, \Theta_j \perp\!\!\!\perp \tilde{f}(\Theta, \lambda) \quad (4.10)$$

which is indeed much stronger.  $(\tilde{H}_0)$  implies  $(H_0)$  but the converse is generally false. Having  $\mathcal{H}_j = 0$  only ensures independence between  $\Theta_j$  and  $G_\lambda$  for  $\mathbb{P}_\Lambda$ -almost all  $\lambda \in \mathcal{D}_\Lambda$ . In practice, this property is sufficient and constitutes a solid theoretical guarantee. However, one should be aware that there may exist a  $\mathbb{P}_\Lambda$ -null set  $\mathcal{D}_\Lambda^{\text{null}}$  containing values  $\lambda$  such that  $\Theta_j \not\perp\!\!\!\perp G_\lambda$ .

**Remark 4.3.** Another probabilistic characteristic of  $H_j(\Lambda)$  could be used to define a sensitivity index. For example, one could consider its essential supremum, defined by:

$$\mathcal{H}_j^{\text{sup}} := \text{ess sup } H_j(\Lambda) = \inf \left( \{C > 0 : H_j(\Lambda) \leq C \text{ } \mathbb{P}_\Lambda\text{-a.s.}\} \right). \quad (4.11)$$

Just as  $\mathcal{H}_j$ , the index  $\mathcal{H}_j^{\text{sup}}$  is able to characterize independence between  $\Theta_j$  and  $G$ . Following the same arguments as in the proof of Proposition 4.1, one can show that:

$$\mathcal{H}_j^{\text{sup}} = 0 \iff H_j(\Lambda) = 0 \text{ } \mathbb{P}_\Lambda\text{-a.s.} \iff \Theta_j \perp\!\!\!\perp \tilde{f}(\Theta, \Lambda). \quad (4.12)$$

However, intuitively, it seems harder to estimate accurately  $\mathcal{H}_j^{\text{sup}}$  from a reasonable amount of simulated data.

#### 4.1.2. Estimation from nested simulation data

We now turn to the estimation of  $\mathcal{H}_j$ . Remember that one major constraint of our problem is that inference must be performed from the same data as conditional Bayesian calibration. Consequently, the sensitivity indices  $\mathcal{H}_j$  must be estimable from these specific data. As explained in Section 3.4, the data are organized as a nested Monte Carlo design.

- The outer level  $D_{\Lambda, m} = \{\Lambda^{(k)}\}_{k=1}^m$  consists of an  $m$ -sample of  $\Lambda$ .
- The inner level  $D_{m, n}$  depends on the adopted sampling strategy. For (S1), it contains an  $n$ -sample of all pairs  $(\Theta^{(k)}, \mathbf{Y}^{(k)})$  associated with the values  $\Lambda^{(k)}$ :

$$D_{\Theta^{(k)} \mathbf{Y}^{(k)}, n} = \{D_{\Theta^{(k)}, n}; D_{\mathbf{Y}^{(k)}, n}\} = \left\{ \left( \Theta^{(kl)}, \mathbf{Y}^{(kl)} \right) \right\}_{1 \leq l \leq n}. \quad (4.13)$$

For (S2), it contains an  $n$ -sample of all pairs  $(\Theta, \mathbf{Y}^{(k)})$ :

$$D_{\Theta \mathbf{Y}^{(k)}, n} = \{D_{\Theta, n}; D_{\mathbf{Y}^{(k)}, n}\} = \left\{ \left( \Theta^{(l)}, \mathbf{Y}^{(kl)} \right) \right\}_{1 \leq l \leq n}. \quad (4.14)$$

As already mentioned, a simple post-processing step allows transforming observations of the output vector  $\mathbf{Y}$  into observations of the output of interest  $G$  (depending on the chosen screening approach). In the following,

all introduced notations apply to the post-processed data:

$$\begin{aligned}
 \bullet \quad D_{m,n}^{S_1} &= \{D_{\Theta^{(k)}G^{(k)},n}\}_{1 \leq k \leq m} \quad \text{with} \quad D_{\Theta^{(k)}G^{(k)},n} = \left\{ \left( \Theta^{(kl)}, G^{(kl)} \right) \right\}_{1 \leq l \leq n} ; \\
 \bullet \quad D_{m,n}^{S_2} &= \{D_{\Theta G^{(k)},n}\}_{1 \leq k \leq m} \quad \text{with} \quad D_{\Theta G^{(k)},n} = \left\{ \left( \Theta^{(l)}, G^{(kl)} \right) \right\} .
 \end{aligned} \tag{4.15}$$

Based on these data, the estimation of  $\mathcal{H}_j = \mathbb{E}_{\Lambda} [H_j(\Lambda)]$  is fairly natural. As a first step, the expectation over  $\Lambda$  can be approximated by the empirical average over the values  $\Lambda^{(k)}$ :

$$\widehat{\mathcal{H}}_j := \widehat{\mathcal{H}}_j(D_{\Lambda,m}) = \frac{1}{m} \sum_{k=1}^m H_j(\Lambda^{(k)}) \quad \text{with} \quad H_j(\Lambda^{(k)}) = \text{HSIC}(\Theta_j, G^{(k)} \mid \Lambda^{(k)}) . \tag{4.16}$$

As it embeds the HSIC, the function  $H_j$  is not analytically tractable. The second step therefore consists in estimating all values  $H_j(\Lambda^{(k)})$  from the data provided in  $D_{m,n}$ . Assuming that this estimation step is feasible, the final estimator is given by:

$$\widehat{\widehat{\mathcal{H}}}_j(D_{m,n}) := \frac{1}{m} \sum_{k=1}^m \widehat{H}_j(\Lambda^{(k)}) . \tag{4.17}$$

where  $\widehat{H}_j(\Lambda^{(k)})$  denotes an estimator of  $H_j(\Lambda^{(k)})$ . Two standard estimators of the HSIC are commonly used: the U- and V-statistic estimators (see Appendix C.2). In both cases, the estimate can be computed from an  $n$ -sample of joint observations. In our context, for each quantity  $H_j(\Lambda^{(k)})$  to be estimated, this sample is provided by  $D_{\Theta^{(k)}G^{(k)},n}$  under (S1) and by  $D_{\Theta G^{(k)},n}$  under (S2).

Finally, since each quantity  $H_j(\Lambda^{(k)})$  can be estimated using two different HSIC estimators on two different datasets, four distinct estimators of  $\mathcal{H}_j$  can be constructed:

$$\begin{aligned}
 \text{Estimator 1: } \widehat{\widehat{\mathcal{H}}}_j^{\text{U}}(D_{m,n}^{S_1}) &:= \frac{1}{m} \sum_{k=1}^m \widehat{H}_j^{\text{U}}(D_{\Theta^{(k)}G^{(k)},n}) =: \widehat{\widehat{\mathcal{H}}}_j^1 ; \\
 \text{Estimator 2: } \widehat{\widehat{\mathcal{H}}}_j^{\text{V}}(D_{m,n}^{S_1}) &:= \frac{1}{m} \sum_{k=1}^m \widehat{H}_j^{\text{V}}(D_{\Theta^{(k)}G^{(k)},n}) =: \widehat{\widehat{\mathcal{H}}}_j^2 ; \\
 \text{Estimator 3: } \widehat{\widehat{\mathcal{H}}}_j^{\text{U}}(D_{m,n}^{S_2}) &:= \frac{1}{m} \sum_{k=1}^m \widehat{H}_j^{\text{U}}(D_{\Theta G^{(k)},n}) =: \widehat{\widehat{\mathcal{H}}}_j^3 ; \\
 \text{Estimator 4: } \widehat{\widehat{\mathcal{H}}}_j^{\text{V}}(D_{m,n}^{S_2}) &:= \frac{1}{m} \sum_{k=1}^m \widehat{H}_j^{\text{V}}(D_{\Theta G^{(k)},n}) =: \widehat{\widehat{\mathcal{H}}}_j^4 .
 \end{aligned} \tag{4.18}$$

Above, the superscripts U and V indicate which HSIC estimator is used. Hereafter,  $\widehat{\widehat{\mathcal{H}}}_j$  will refer to any of the four proposed estimators.

**Remark 4.4.** The estimator  $\widehat{\widehat{\mathcal{H}}}_j$  in equation (4.17) is written with a double hat to emphasize the presence of two approximation levels. The estimator  $\widehat{\mathcal{H}}_j$  in equation (4.16) is not tractable, but it plays a crucial role in the analysis of convergence speeds (see Prop. 4.5). In particular, it allows the two sources of error to be disentangled in the proofs of Appendix E.2, thereby reducing the analysis to classical arguments.

In order to find out which estimator(s) should be preferred, the statistical properties (bias, consistency, convergence rate) of all four estimators need to be investigated. For obvious reasons, a minimal requirement for

TABLE 2. Properties of the four proposed estimators.

| Estimators         | Properties     |              |             |  | MSE decay rate  |
|--------------------|----------------|--------------|-------------|--|---|
|                    | Non-negativity | Unbiasedness | Consistency |  |   |
| <b>Estimator 1</b> | ✗              | ✓            | ✓           |  | $\mathcal{O}\left(\frac{1}{m}\right)$                 |
| <b>Estimator 2</b> | ✓              | ✗            | ✓           |  | $\mathcal{O}\left(\frac{1}{m} + \frac{1}{n^2}\right)$ |
| <b>Estimator 3</b> | ✗              | ✓            | ✓           |  | $\mathcal{O}\left(\frac{1}{m} + \frac{1}{n}\right)$   |
| <b>Estimator 4</b> | ✓              | ✗            | ✓           |  | $\mathcal{O}\left(\frac{1}{m} + \frac{1}{n}\right)$   |

the estimator is consistency. It must be proven that  $\widehat{\mathcal{H}}_j$  converges in probability to the target quantity  $\mathcal{H}_j$  as the sample sizes grow ( $m, n \rightarrow \infty$ ). To establish convergence in probability, it is often more convenient to prove convergence in mean square. This amounts to showing that the mean square error (MSE) tends to zero as the sample sizes become arbitrarily large. In our setting, this reads:

$$\text{MSE}\left(\widehat{\mathcal{H}}_j\right) := \mathbb{E}\left[\left|\widehat{\mathcal{H}}_j - \mathcal{H}_j\right|^2\right] \xrightarrow{m, n \rightarrow \infty} 0. \quad (4.19)$$

Convergence in mean square is stronger than convergence in probability, but it is often easier to demonstrate, and it reveals convergence rates more readily. Proposition 4.5 states that all four estimators are consistent, and that two of them are unbiased. In addition, upper bounds on the decay rates of the MSE are provided for the four estimators. These results are particularly valuable, as they can guide the allocation of the simulation budget ( $N_2 = n_{\text{exp}} \times m \times n$ ) between the two levels of sampling.

**Proposition 4.5.** *Let  $j \in \{1, \dots, p\}$  be fixed. Assume that  $K_{\theta_j}$  and  $K_g$  are two bounded kernels. The statistical properties of the four estimators of  $\mathcal{H}_j$  are summarized in Table 2.*

The detailed proofs can be found in Appendix E.2. They are much inspired from ideas presented in a recent work by Fellmann et al. [71].

**Remark 4.6.** It is worth noting that the assumption of bounded kernels is pretty standard in the literature related to HSIC indices, especially to obtain concentration inequalities (see Thm. 4 in [72]). This assumption is not very stringent as it is verified by most translation-invariant kernels used in statistical applications (especially Gaussian and Matérn kernels).

**Remark 4.7.** It should be clearly understood that stating the MSE is  $\mathcal{O}(1/m)$  does not imply that the MSE is asymptotically equivalent to  $1/m$  as  $m, n \rightarrow \infty$ . Rather, this indicates that the MSE can be bounded above by  $C/m$  for some positive constant  $C$ , potentially very large. In a hypothetical scenario where the theoretical values  $\mathcal{H}_j$  are known and the chained model  $y$  can be evaluated without limitation, one could have compared the convergence rates of Proposition 4.5 with the true rates, which might be much faster. Therefore, the rates from Proposition 4.5 should be interpreted as worst-case rates.

Proposition 4.5 provides valuable insights into the statistical properties of the four estimators. All of them are consistent and exhibit convergence rates comparable to those of classical Monte Carlo estimators. This guarantees that the indices  $\mathcal{H}_j$  can be estimated accurately from the available simulation data. Several additional points are worth highlighting.

- **Estimator 1.** It is unbiased, which is natural since it is expressed as the average of  $m$  U-statistics. Its MSE decays at a rate of  $\mathcal{O}(1/m)$ . Accordingly, it appears more efficient to prioritize large values of  $m$

(the number of points used to estimate the expectation over  $\mathbf{\Lambda}$ ) rather than  $n$  (the number of points used to estimate HSIC indices). To reduce the MSE as much as possible, one could be tempted to take  $m \gg n$ . However, as will be seen later, this is not what is done in practice, as doing so would compromise the performance of conditional calibration, which remains the ultimate objective. One typically works with  $n \gg m$  to ensure that a large  $n$ -sample of  $(\Theta, G^{(k)})$  is available for each value  $\mathbf{\Lambda}^{(k)}$ .

- **Estimator 2.** It is biased, as expected due to the use of V-statistics. Unlike Estimator 1, the approximation error cannot be fully eliminated simply by increasing  $m$ . The bias introduced by the V-statistics gives rise to the  $1/n^2$  term in the MSE, which can only be reduced by increasing  $n$ . In light of the previous remark, choosing  $n \gg m$  (for conditional calibration purposes) makes the V-statistic bias practically negligible, and Estimators 1 and 2 converge at essentially the same rate.
- **Estimators 3 and 4.** The remarks made for Estimators 1 and 2 regarding the presence or absence of bias also apply to Estimators 3 and 4. Importantly, their convergence rate of  $\mathcal{O}(1/m + 1/n)$  is slower than that of Estimators 1 and 2. This rules out (S2) in favor of (S1), which therefore emerges as the most effective sampling strategy for achieving maximum accuracy under a limited simulation budget.
- **Estimator 4.** As shown in the consistency proof of  $\widehat{\mathcal{H}}_j^4$  (see Appendix E.2), the MSE satisfies  $\mathcal{O}(1/m + 1/n + 1/n^2)$  where:
  - $1/m$  is due to the empirical mean over the  $m$  observations of  $\mathbf{\Lambda}$ .
  - $1/n$  is due to Strategy (S2);
  - $1/n^2$  is due to V-statistics.

Because of their slower convergence rates, Estimators 3 and 4 are excluded from further analysis. They are neither considered in Section 4.2 nor implemented in Section 5. By contrast, Estimators 1 and 2 are both retained, as neither is inherently superior to the other. Indeed, they inherit the respective advantages and drawbacks of the U- and V-statistic estimators of the HSIC (unbiased but sometimes negative for Estimator 1, biased but always non-negative for Estimator 2).

**Remark 4.8.** The alternative sensitivity index  $\mathcal{H}_j^{\text{null}}$  introduced in equation (4.11) can also be estimated from the nested Monte Carlo design. A counterpart to equation (4.17) is given by:

$$\widehat{\mathcal{H}}_j^{\text{sup}} := \max_{1 \leq k \leq m} \widehat{H}_j(\mathbf{\Lambda}^{(k)}) . \quad (4.20)$$

Once again, four estimators can be constructed by applying the two possible HSIC estimators to the two nested designs. Moreover, it can be easily shown that these estimators are consistent. The proofs are very similar to those presented in Appendix E.2. The keys arguments to demonstrate consistency are the extremal law of large numbers and appropriate concentration inequalities [72]. However, deriving explicit convergence rates may be considerably more challenging.

At this stage, the work is not yet complete, because the estimates of the indices  $\widehat{\mathcal{H}}_j$  are not directly usable for parameter screening. The problem is the same as with HSIC indices, if not worse. Recall that the scale of an HSIC index, measured via the MMD, depends on the marginal distributions of the two random objects and on the kernels chosen to handle these objects. Consequently, the magnitude of the MMD scale is difficult to interpret: given an HSIC estimate, it is impossible to determine whether the underlying HSIC value is zero or not. This is precisely why a proper independence test must be employed to make the final decision.

In the case of  $\mathcal{H}_j$ , the same problem is encountered, as the index is expressed as the average of several HSIC estimates. The difficulty is further compounded by the fact that the related HSIC terms are defined conditional on an event  $\mathbf{\Lambda} = \boldsymbol{\lambda}$ . For each term, dependence is measured between  $\Theta_j$  and  $G_{\boldsymbol{\lambda}} = \tilde{f}(\Theta, \boldsymbol{\lambda})$  where the marginal distribution of  $G_{\boldsymbol{\lambda}}$  naturally varies with  $\boldsymbol{\lambda}$ . For two different values  $\boldsymbol{\lambda}_1 \neq \boldsymbol{\lambda}_2$ , the MMD scales associated with  $H_j(\boldsymbol{\lambda}_1)$  and  $H_j(\boldsymbol{\lambda}_2)$  are therefore not comparable. In short, as with HSIC indices, estimation alone is insufficient: the screening problem must be reformulated as a hypothesis test, and a proper test procedure must be developed.

## 4.2. Testing independence in the bi-level uncertainty framework

### 4.2.1. Problem formulation

The objective is to construct a statistical procedure based on the data  $D_{m,n}^{S_1}$  that allows, for each parameter  $\Theta_j$ , discrimination between the following two hypotheses:

$$(H_0) : \Theta_j \perp\!\!\!\perp G \quad \text{vs.} \quad (H_1) : \Theta_j \not\perp\!\!\!\perp G \quad \text{with} \quad G = \tilde{f}(\Theta, \Lambda). \quad (4.21)$$

By virtue of Proposition 4.1, these statements may be rewritten as:

$$(H_0) : \mathcal{H}_j = 0 \quad \text{vs.} \quad (H_1) : \mathcal{H}_j > 0. \quad (4.22)$$

Naturally, the test statistic will be either Estimator 1 or 2. Given that  $(H_0)$  coincides with  $\mathcal{H}_j = 0$ , the test of independence based on  $\widehat{\mathcal{H}}_j$  is one-sided (upper-tailed), meaning that  $(H_0)$  is rejected when the test statistic takes a sufficiently large value. More precisely, the rejection threshold is the quantile of order  $1 - \alpha$  of the null distribution, with the significance level  $\alpha$  ensuring strict control of Type-I error. Let  $D_{m,n}^{\text{obs}}$  denote the available data. Note that:

- $D_{m,n}^{\text{obs}}$  is the observed realization of the random object  $D_{m,n}$  that generates the data.
- $\widehat{\mathcal{H}}_j(D_{m,n}^{\text{obs}})$  is the observed realization of the test statistic  $\widehat{\mathcal{H}}_j(D_{m,n})$ .

The  $p$ -value of the test is defined as the probability, under  $(H_0)$ , of observing a value of the test statistic at least as extreme as the one actually observed. In our context, it is given by:

$$p_j := \mathbb{P}_{H_0} \left( \widehat{\mathcal{H}}_j(D_{m,n}) > \widehat{\mathcal{H}}_j(D_{m,n}^{\text{obs}}) \right). \quad (4.23)$$

For the HSIC-based test of independence, three distinct test procedures can be used to estimate the  $p$ -value, each suited to a particular range of sample sizes (see Appendix C.3).

- **Large sample sizes.** For  $n \geq 500$ , the asymptotic test aims to approximate the spectral distribution ruling the asymptotic behavior under  $(H_0)$  of the rescaled test statistic. The  $p$ -value is estimated at the cost of  $\mathcal{O}(n^2)$  operations.
- **Small sample sizes.** For  $n \leq 100$ , only the permutation-based test is applicable. This procedure involves randomly permuting the output data and recomputing the test statistic for each permuted sample, thereby generating an empirical approximation of the null distribution. The  $p$ -value is then estimated as the proportion of simulated values that exceed the original observed value. The computational complexity of this approach is  $\mathcal{O}(B n^2)$ , where  $B$  denotes the number of permutations.
- **Intermediate sample sizes.** For  $100 \leq n \leq 500$ , the non-asymptotic Gamma test enables retrieving the  $p$ -value of the permutation-based test, but without performing any permutation. The computational complexity is therefore reduced to  $\mathcal{O}(n^2)$ .

We now aim to adapt these test procedures to our specific screening problem. From a purely technical standpoint, studying the asymptotic behavior of the (rescaled) test statistic appears to be challenging. Indeed, because of the bi-level sampling scheme, classical asymptotic results for U- and V-statistics cannot be directly applied. In any case, given that conditional calibration requires  $N_2 = n_{\text{exp}} \times m \times n$  evaluations of the downstream model, the asymptotic framework is clearly unrealistic. By contrast, the two alternative test procedures are relatively easy to implement. A permutation-based test is developed in the next section, while a non-asymptotic Gamma test is presented in Appendix D.2.

#### 4.2.2. Permutation-based test procedure

At this stage, the objective is to test  $(H_0) : \Theta_j \perp\!\!\!\perp G$  using the nested Monte Carlo design  $D_{m,n}^{S_1}$ . Let  $D_{m,n}^{\text{obs}}$  denote the available data. From a practical viewpoint,  $D_{m,n}^{\text{obs}}$  is the generated nested design, which will later be used for conditional Bayesian calibration. From a probabilistic viewpoint,  $D_{m,n}^{\text{obs}}$  is the only observed realization of  $D_{m,n}^{S_1}$ , which represents the underlying data generation process. As explained in Section 3.4,  $D_{m,n}^{\text{obs}}$  is composed of an  $m$ -sample of  $\Lambda$ :

$$D_{\Lambda,m}^{\text{obs}} := \left\{ \boldsymbol{\lambda}^{(k)} \right\}_{1 \leq k \leq m} , \quad (4.24)$$

together with an  $n$ -sample of each pair  $(\Theta^{(k)}, G^{(k)})$ :

$$D_k^{\text{obs}} := D_{\Theta^{(k)}G^{(k)},n}^{\text{obs}} = \left\{ (\boldsymbol{\theta}^{(kl)}, g^{(kl)}) \right\}_{1 \leq l \leq n} \quad \text{with} \quad g^{(kl)} = \tilde{f}(\boldsymbol{\theta}^{(kl)}, \boldsymbol{\lambda}^{(k)}) . \quad (4.25)$$

In these equations, lowercase letters indicate that the data are fixed during the test procedure.

To simulate the distribution of the test statistic under  $(H_0)$ , we draw inspiration from the permutation-based test (see Appendix C.3), which is arguably the simplest strategy to do so in a finite-sample setting. Recall the equivalence stated in Proposition 4.1:

$$\Theta_j \perp\!\!\!\perp G \iff \exists \tilde{\mathcal{D}}_{\Lambda} \subseteq \mathcal{D}_{\Lambda} \text{ such that } \mathbb{P}_{\Lambda}(\tilde{\mathcal{D}}_{\Lambda}) = 1 \text{ and } \forall \boldsymbol{\lambda} \in \tilde{\mathcal{D}}_{\Lambda}, \Theta_j \perp\!\!\!\perp G_{\boldsymbol{\lambda}} . \quad (4.26)$$

While not fully rigorous, it is reasonable to assume that the values  $\boldsymbol{\lambda}^{(k)}$  belong to  $\tilde{\mathcal{D}}_{\Lambda}$ . Under this assumption, one has:

$$(H_0) \implies \forall 1 \leq k \leq m, \Theta_j^{(k)} \perp\!\!\!\perp G^{(k)} \quad \text{where} \quad G^{(k)} = G_{\boldsymbol{\lambda}^{(k)}} = \tilde{f}(\Theta^{(k)}, \boldsymbol{\lambda}^{(k)}) . \quad (4.27)$$

A necessary condition to simulate  $(H_0)$  is therefore to enforce independence within each pair  $(\Theta^{(k)}, G^{(k)})$ . To achieve this, it suffices to apply random permutations to the output data within each sample  $D_k^{\text{obs}}$ . Let  $\tau_k$  denote the permutation applied to the  $k$ -th sample. In practice,  $\tau_k$  is randomly selected from the set  $\mathbb{S}_n$  of all permutations of  $\{1, \dots, n\}$ . The corresponding permuted sample is obtained as follows:

$$D_k^{\tau_k} := D_{\Theta^{(k)}G^{(k)},n}^{\tau_k} = \left\{ \left( \boldsymbol{\theta}^{(kl)}, g^{(k\tau_k(l))} \right) \right\}_{1 \leq l \leq n} . \quad (4.28)$$

Let  $\boldsymbol{\tau} = (\tau_1, \dots, \tau_m)$  denote the collection of all permutations. Accordingly,  $D_{m,n}^{\boldsymbol{\tau}}$  denotes the nested design derived from  $D_{m,n}^{\text{obs}}$  after applying the permutations  $\tau_k$  to the samples  $D_k^{\text{obs}}$ . Then, the reasoning is straightforward.

- Computing  $\widehat{H}_j(D_k^{\tau_k})$  produces a realization of  $\widehat{H}_j(\boldsymbol{\lambda}^{(k)})$  under  $\Theta_j \perp\!\!\!\perp \tilde{f}(\Theta^{(k)}, \boldsymbol{\lambda}^{(k)})$ .
- Computing  $\widehat{\mathcal{H}}_j(D_{m,n}^{\boldsymbol{\tau}})$  produces a realization of  $\widehat{\mathcal{H}}_j(D_{m,n})$  under  $(H_0)$ .

The complete procedure used to produce a pseudo-realization of the null distribution is illustrated in Figure 4. Repeating this procedure  $B$  times, each time with a new collection of permutations  $\boldsymbol{\tau}_b = (\tau_{b1}, \dots, \tau_{bm}) \in (\mathbb{S}_n)^m$ , yields an empirical approximation of the null distribution in the form of a  $B$ -sample. The  $p$ -value is then estimated by:

$$\widehat{p}_j = \frac{1}{B} \sum_{b=1}^B \mathbf{1}_{\left\{ \widehat{\mathcal{H}}_j(D_{m,n}^{\boldsymbol{\tau}_b}) > \widehat{\mathcal{H}}_j(D_{m,n}^{\text{obs}}) \right\}} . \quad (4.29)$$

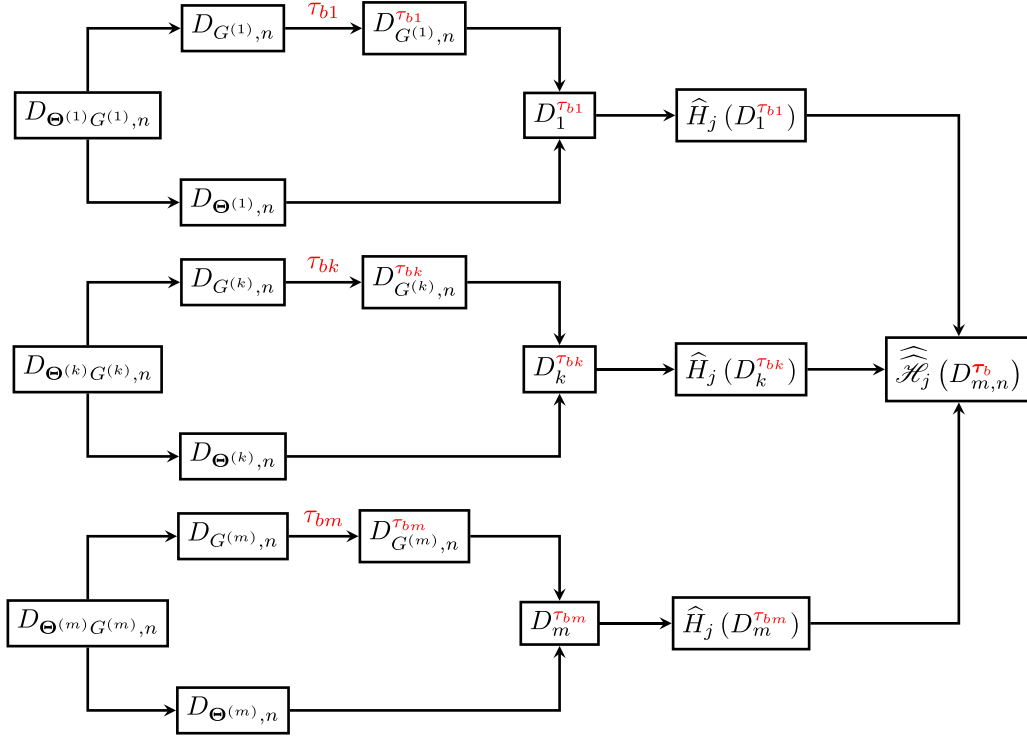


FIGURE 4. Schematic of the procedure used to simulate the  $b$ -th realization of the test statistic  $\widehat{\mathcal{H}}_j$  under the null hypothesis  $(H_0) : \Theta_j \perp\!\!\!\perp \tilde{f}(\Theta, \Lambda)$ . To this end, a realization of each estimator  $\widehat{H}_j(\lambda^{(k)})$  must be simulated under independence between  $\Theta_j$  and  $\tilde{f}(\Theta, \lambda^{(k)})$ . This is achieved by applying a random permutation  $\tau_{bk}$  to each sample  $D_{G^{(k),n}}$ . The permuted output sample  $D_{G^{(k),n}}^{\tau_{bk}}$  and the input sample  $D_{\Theta^{(k),n}}$  are then reassembled to form a sample  $D_k^{\tau_{bk}}$  which is used to estimate  $H_j(\lambda^{(k)})$ . The empirical average of these values yields one realization of  $\widehat{\mathcal{H}}_j$  under  $(H_0)$ .

The overall computational complexity of this test procedure is  $\mathcal{O}(mBn^2)$ . Since the number of permutations is often set to  $B = 500$ , this can become costly, even for intermediate sample sizes ( $100 \leq n \leq 500$ ). Using the non-asymptotic Gamma test described in Appendix D.2 reduces the complexity to  $\mathcal{O}(mn^2)$ .

### 4.3. Summary of contributions

First, we have introduced a sensitivity measure that takes into account the bi-level parametric uncertainty arising from the chaining between Models 1 and 2 (see Fig. 3). This framework leads to consider a set of  $\mathbb{P}_\Lambda$ -informed HSIC indices computed between the parameters  $\Theta_j$  of Model 2 and the output of interest  $G = \tilde{f}(\Theta, \Lambda)$ . A key result (see Prop. 4.1) is that  $\mathcal{H}_j$  provides an exact characterization of the probabilistic independence between  $\Theta_j$  and  $G$ . This makes  $\mathcal{H}_j$  a particularly relevant criterion for conditional calibration since the influence exerted by  $\Theta_j$  is quantified by integrating over the full uncertainty on  $\Lambda$ .

Second, we have shown that four distinct estimators of  $\mathcal{H}_j$  can be constructed from the nested Monte Carlo design intended for conditional calibration. Their statistical properties have been thoroughly investigated (see Prop. 4.5) in order to identify the most reliable estimators. In particular, we have demonstrated that two estimators are unbiased, and that all four estimators are consistent, with convergence rates comparable to those

of classical Monte Carlo estimators. These findings have confirmed that  $\mathcal{H}_j$  can be accurately estimated using a reasonable simulation budget. They have also enabled us to discard two estimators based on their MSE decay rates.

Finally, we have highlighted that point estimation alone is insufficient for parameter screening due to the kernel-based metric underlying  $\mathcal{H}_j$ . Thanks to Proposition 4.1, the screening problem has been reformulated as a hypothesis testing problem, where the null hypothesis corresponds to  $\mathcal{H}_j = 0$ . Two test procedures, largely inspired by what exists for HSIC indices, have been developed to assess the statistical significance of the estimated indices.

## 5. REFINED ANALYSIS OF THE FISSION GAS BEHAVIOR MODEL

In this section, we revisit the case study introduced in Section 3.1. The goal remains to screen the parameters  $\boldsymbol{\theta}$  of the fission gas behavior model, but this time taking into account the residual uncertainty in the thermal model conductivity  $\lambda$ . To do so, we apply the methodology developed in Section 4. The following procedure needs to be executed.

- Generate the design  $D_{\Lambda, m}$  with  $m = 20$  and  $\mathbb{P}_{\Lambda}$  estimated in [32].
- Generate all the designs  $D_{\boldsymbol{\Theta}^{(k)}, n}$  with  $n = 200$  and  $\mathbb{P}_{\boldsymbol{\Theta}}$  derived from expert judgment.
- Evaluate the chained model  $y$  to compute the output data  $D_{\mathbf{Y}^{(k)}, n}$ .
- Estimate the sensitivity indices with respect to  $G = \tilde{f}(\boldsymbol{\Theta}, \Lambda)$ .
- Run the test procedures to estimate the  $p$ -values.
- Select the most influential parameters using the  $p$ -values ( $\alpha = 5\%$ ) or the detection rates ( $\beta = 80\%$ ).

The definition of the analyzed output  $G$  depends on the chosen screening approach (see Rem. 3.3). The final nested design  $D_{m, n}^{S_1}$  is obtained by gathering the input and output data according to equation (3.18). This sampling strategy entails the subsequent use of Estimators 1 and 2. The two test procedures developed in Section 4.2 apply to Estimator 2, whereas only the permutation-based test applies to Estimator 1.

**Remark 5.1.** Looking at Table 2, one might wonder why  $m$  is chosen much smaller than  $n$ , whereas the convergence rates of the MSE of Estimators 1 and 2 suggest the opposite. This situation is dictated by the conditional calibration algorithm, which requires large designs in the parameter space  $\mathcal{D}_{\boldsymbol{\Theta}}$  to work efficiently. This illustrates a situation where estimator performance is suboptimal because the data were not originally intended for GSA.

The results obtained for Approach A are reported in Figure 5, which is arranged in the same way as Figure 2. It displays the test results for  $(H_0) : \Theta_j \perp\!\!\!\perp Y_i$  across all  $(i, j)$  pairs. This is done for the two possible estimators and the permutation-based test. For Estimator 1, the detection rates  $\rho_j$  summarizing the  $n_{\text{exp}}$  tests performed for each parameter  $\Theta_j$  are reported in Column A of Table 3. The detection rates for Estimator 2 are not shown as they exactly coincide with those of Estimator 1.

We observe that parameters  $\Theta_5, \Theta_6, \Theta_7$ , and  $\Theta_{11}$  are influential on the RGF of almost all fuel rods. For these four parameters, the detection rates exceed 90%. Then comes  $\Theta_{10}$  with  $\rho_{10} = 65\%$ . If the decision threshold is set to  $\beta = 60\%$ , only these five parameters are selected. Nevertheless, it is worth noting that the influence of  $\Theta_2$  and  $\Theta_3$  is not negligible, with detection rates  $\rho_2 = 22.5\%$  and  $\rho_3 = 47.5\%$ , respectively.

Comparing these results with those from Section 3.3 reveals two important insights. First, the four most influential parameters are identified in the same way. Second, our screening methodology appears to be more conservative: the selected set of parameters is larger, and the detection rates are higher for all parameters. This behavior was expected, as the null hypothesis is much stronger. Indeed, for  $(H_0) : \Theta_j \perp\!\!\!\perp \tilde{f}(\boldsymbol{\Theta}, \Lambda)$  to hold,

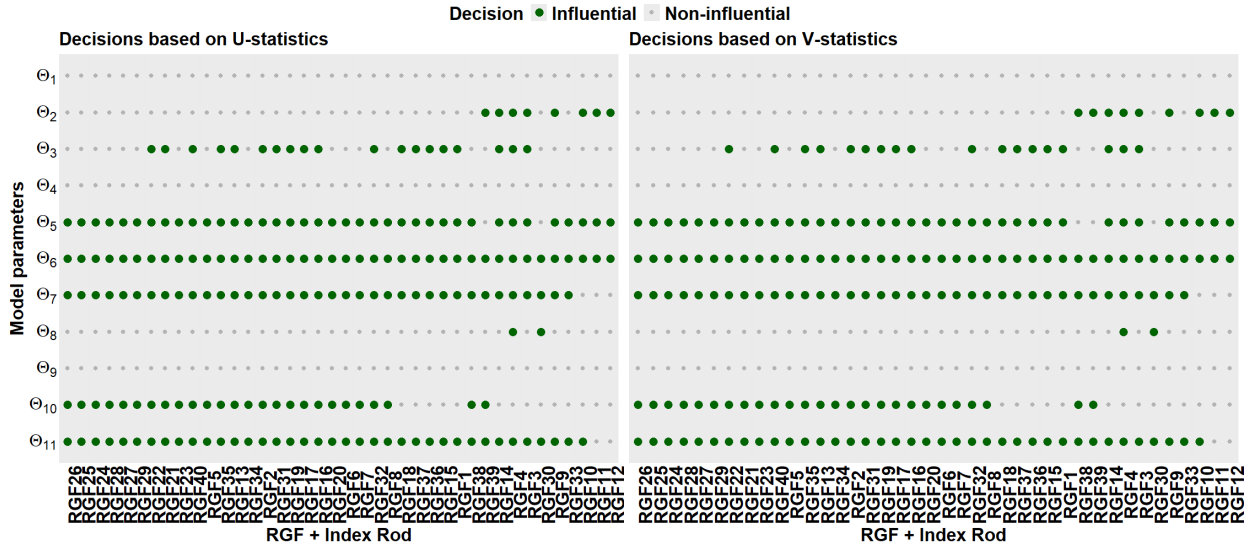


FIGURE 5. Results obtained for Approach A after incorporating the uncertainty on  $\Lambda$ . For each fuel rod, the influence of each parameter  $\Theta_j$  is assessed using the test procedure developed in Section 4.2.2. The test statistic is either Estimator 1 (left) or Estimator 2 (right). The significance level is set to  $\alpha = 5\%$ . The associated detection rates are reported in Table 3.

TABLE 3. Comparison of the screening results obtained for the three considered approaches after incorporating the uncertainty on  $\Lambda$ . Only the results obtained with Estimator 1 are shown.

| Results with U-statistics |  | Detection rates |       | $p$ -values |  |
|---------------------------|--|-----------------|-------|-------------|--|
| Approach                  |  | A               | B     | C           |  |
| Parameters                |  |                 |       |             |  |
| $\Theta_1$                |  | 0%              | 1     | 1           |  |
| $\Theta_2$                |  | 20%             | 0     | 0.997       |  |
| $\Theta_3$                |  | 47.5%           | 0.011 | 0           |  |
| $\Theta_4$                |  | 0%              | 1     | 1           |  |
| $\Theta_5$                |  | 95%             | 0     | 0           |  |
| $\Theta_6$                |  | 100%            | 0     | 0           |  |
| $\Theta_7$                |  | 92.5%           | 0     | 0           |  |
| $\Theta_8$                |  | 5%              | 0.963 | 0.903       |  |
| $\Theta_9$                |  | 0%              | 1     | 1           |  |
| $\Theta_{10}$             |  | 65%             | 0.658 | 0.002       |  |
| $\Theta_{11}$             |  | 95%             | 0     | 0           |  |

the parameter  $\Theta_j$  must be independent of  $\tilde{f}(\Theta, \lambda)$  for almost all  $\lambda \in \mathcal{D}_\Lambda$ , not merely for the nominal value  $\lambda = \lambda_{\text{nom}}$ . Consequently,  $(H_0)$  is rejected more frequently, which increases the power of the test procedure.

The  $p$ -values obtained for Approaches B and C are also reported in Table 3. The results from the three approaches are highly consistent, each leading to the selection of the same four most influential parameters, namely  $\Theta_5$ ,  $\Theta_6$ ,  $\Theta_7$ , and  $\Theta_{11}$ . Despite this overall agreement, a few differences between the three approaches can still be observed. Approach C, which targets the least-squares criterion  $L = \tilde{f}^C(\Theta, \Lambda)$ , readily detects  $\Theta_{10}$  but discards  $\Theta_2$ . This suggests that  $\Theta_{10}$  plays an essential role in the conditional calibration procedure, unlike  $\Theta_2$ .

Interestingly, for Approach B, which operates on the vector  $\mathbf{Y} = \tilde{f}^B(\boldsymbol{\Theta}, \Lambda)$  of all RGFs, the parameter  $\Theta_{10}$  is not detected. Recall that in this approach the multivariate output is handled via a multivariate kernel obtained by tensorization of univariate Gaussian kernels. This construction implicitly standardizes the output components, thereby balancing their relative contributions. Such an effect is often desirable for training a metamodel, but in our setting it tends to penalize parameters whose influence is confined to a small subset of outputs. This may partly explain why  $\Theta_{10}$  is not selected under Approach B. To facilitate the detection of  $\Theta_{10}$ , it may be worth considering the weighted PCA-kernel (see Prop. 2 in [58]). It is constructed as a weighted linear combination (using the PCA eigenvalues as weights) of univariate kernels applied to the principal components. It has the drawback of not being characteristic (see Appendix C.5.3), but it offers improved performance in terms of detection capabilities. This alternative kernel has been experimented and indeed allows the successful detection of  $\Theta_{10}$ .

Overall, the three approaches converge towards the selection of six parameters:  $\Theta_3$ ,  $\Theta_5$ ,  $\Theta_6$ ,  $\Theta_7$ ,  $\Theta_{10}$ , and  $\Theta_{11}$ . Thus, the proposed screening methodology corroborates the findings of Section 3.2 and allows for a more broader selection of parameters. Additionally, all conclusions are robust with respect to both the choice of estimator and test procedure. In particular, the  $p$ -values returned by the non-asymptotic Gamma test closely match those obtained from the permutation-based test. In Figure E.1, the histogram of pseudo-realizations of Estimator 2 under  $(H_0)$  is compared with the parametric Gamma approximation, revealing a good fit.

## 6. CONCLUSIONS AND PERSPECTIVES

This paper proposes a screening methodology aimed at reducing the number of parameters involved in the conditional Bayesian calibration of the downstream component of a computational chain. The proposed approach properly accounts for the residual uncertainty arising from the calibration of the upstream component. Besides, it does not induce any additional computational cost, as it can be implemented directly from the nested Monte Carlo design required to carry out conditional calibration.

The proposed sensitivity measure is defined as an integrated version of the HSIC. It explicitly incorporates the upstream uncertainty level and is able to correctly identify the most influential parameters of the downstream model, regardless of the values taken by the upstream model parameters. The proposed estimators are consistent, and their convergence rates are comparable to those of standard Monte Carlo estimators, allowing for accurate inference even with a limited simulation budget. Decision making relies on a test of independence that can be performed using the provided nested design. Two test procedures have been developed, both inspired by existing approaches for HSIC indices. While the permutation-based test entails a high computational complexity, this burden can be alleviated by resorting to the non-asymptotic Gamma test.

The methodology has been successfully applied to an industrial case study from the nuclear industry, involving the chaining between a thermal model and a fission gas behavior model within the ALCYONE code. The results confirm that the proposed approach fulfills its intended purpose: it leads to the selection of a broader set of influential parameters, as a direct consequence of the rigorous treatment of upstream uncertainty.

In its current form, the methodology presents several notable strengths. Since it is built upon the HSIC framework, it inherits most of its advantages, in particular the ability to handle any type of model output. In addition, the approach naturally extends to more complex computational chains involving more than two components. To further align the methodology with operational constraints, future work could focus on extending the theoretical developments to alternative experimental designs, such as Latin hypercube sampling (LHS).

## ACKNOWLEDGMENTS

The major part of this work was done when Oumar Baldé was a PhD student in CEA, DES, Service de Génie Logiciel pour la Simulation.

## REFERENCES

- [1] T.J. Santner, B.J. Williams, W.I. Notz and B.J. Williams, *The Design and Analysis of Computer Experiments*, vol. 1. Springer (2003).
- [2] K. Campbell, Statistical calibration of computer simulations. *Reliabil. Eng. Syst. Saf.* **91** (2006) 1358–1363.
- [3] T.G. Trucano, L.P. Swiler, T. Igusa, W.L. Oberkampf and M. Pilch, Calibration, validation, and sensitivity analysis: What's what. *Reliabil. Eng. Syst. Saf.* **91** (2006) 1331–1357.
- [4] E. de Rocquigny, N. Devictor and S. Tarantola, *Uncertainty in Industrial Practice: A Guide to Quantitative Uncertainty Management*. John Wiley & Sons (2008).
- [5] A. Bouloré, Importance of uncertainty quantification in nuclear fuel behaviour modelling and simulation. *Nucl. Eng. Des.* **355** (2019) 110311.
- [6] M. Gu and L. Wang, Scaled Gaussian stochastic process for computer model calibration and prediction. *SIAM/ASA J. Uncertainty Quantif.* **6** (2018) 1555–1583.
- [7] R.K. Wong, C.B. Storlie and T.C. Lee, A frequentist approach to computer model calibration. *J. Roy. Statist. Soc. B Statist. Methodol.* **79** (2017) 635–648.
- [8] M.C. Kennedy and A. O'Hagan, Bayesian calibration of computer models. *J. Roy. Statist. Soc. B (Statist. Methodol.)* **63** (2001) 425–464.
- [9] X. Wu, T. Kozłowski, H. Meidani and K. Shirvan, Inverse uncertainty quantification using the modular Bayesian approach based on Gaussian process. Part 1. Theory. *Nucl. Eng. Des.* **335** (2018) 339–355.
- [10] R.E. Kass and L. Wasserman, The selection of prior distributions by formal rules. *J. Am. Statist. Assoc.* **91** (1996) 1343–1370.
- [11] S. Chib and E. Greenberg, Understanding the Metropolis–Hastings algorithm. *Am. Statist.* **49** (1995) 327–335.
- [12] B. Currin, T. Mitchell, M. Morris and D. Ylvisaker, Bayesian prediction of deterministic functions, with applications to the design and analysis of computer experiments. *J. Am. Statist. Assoc.* **86** (1991) 953–963.
- [13] J. Sacks, W.J. Welch, T.J. Mitchell and H.P. Wynn, Design and analysis of computer experiments. *Statist. Sci.* **4** (1989) 409–423.
- [14] C. Andrieu and J. Thoms, A tutorial on adaptive MCMC. *Statist. Comput.* **18** (2008) 343–373.
- [15] V.V. Dighe, M. Becker, T. Göçmen, B. Sanderse and J.-W. van Wingerden, Sensitivity analysis and Bayesian calibration of a dynamic wind farm control model: FLORIDyn, in *Journal of Physics: Conference Series*, vol. 2265. IOP Publishing (2022) 022062.
- [16] J. Gou, C. Miao, Q. Duan, Q. Tang, Z. Di, W. Liao, J. Wu and R. Zhou, Sensitivity analysis-based automatic parameter calibration of the VIC model for streamflow simulations over China. *Water Resources Res.* **56** (2020) e2019WR025968.
- [17] J.B. Nagel, J. Rieckermann and B. Sudret, Principal component analysis and sparse polynomial chaos expansions for global sensitivity analysis and model calibration: application to urban drainage simulation. *Reliabil. Eng. Syst. Saf.* **195** (2020) 106737.
- [18] G. Perret, D. Wicaksono, I.D. Clifford and H. Ferroukhi, Global sensitivity analysis and Bayesian calibration on a series of refflood experiments with varying boundary conditions. *Nucl. Technol.* **208** (2022) 711–722.
- [19] D.M. Hamby, A review of techniques for parameter sensitivity analysis of environmental models. *Environ. Monitor. Assessment* **32** (1994) 135–154.
- [20] M.D. Morris, Factorial sampling plans for preliminary computational experiments. *Technometrics* **33** (1991) 161–174.
- [21] J. Morio, Global and local sensitivity analysis methods for a physical system. *Eur. J. Phys.* **32** (2011) 1577.
- [22] A. Saltelli, Making best use of model evaluations to compute sensitivity indices. *Comput. Phys. Commun.* **145** (2002) 280–297.
- [23] I.M. Sobol, Global sensitivity indices for nonlinear mathematical models and their Monte Carlo estimates. *Math. Comput. Simul.* **55** (2001) 271–280.
- [24] S. Kucherenko, M. Rodriguez-Fernandez, C. Pantelides and N. Shah, Monte Carlo evaluation of derivative-based global sensitivity measures. *Reliabil. Eng. Syst. Saf.* **94** (2009) 1135–1148.
- [25] L. Clouvel, B. Iooss, V. Chabridon, M.I. Idrissi and F. Robin, An overview of variance-based importance measures in the linear regression context: comparative analyses and numerical tests. *Socioenviron. Syst. Model.* **7** (2025) 18681–18681.

- [26] B. Iooss and C. Prieur, Shapley effects for sensitivity analysis with correlated inputs: Comparisons with Sobol' indices, numerical estimation and applications. *Int. J. Uncertain. Quantif.* **9** (2019) 493–514.
- [27] A.B. Owen and C. Prieur, On Shapley value for measuring importance of dependent inputs. *SIAM/ASA J. Uncertain. Quantif.* **5** (2017) 986–1002.
- [28] S. Da Veiga, Global sensitivity analysis with dependence measures. *J. Statist. Computat. Simul.* **85** (2015) 1283–1305.
- [29] M. De Lozzo and A. Marrel, New improvements in the use of dependence measures for sensitivity analysis and screening. *J. Statist. Computat. Simul.* **86** (2016) 3038–3058.
- [30] A. Gretton, K. Fukumizu, C. Teo, L. Song, B. Schölkopf and A. Smola, A kernel statistical test of independence. *Adv. Neural Inform. Process. Syst.* **20** (2007) 585–592.
- [31] C. Introïni, I. Ramière, J. Sercombe, B. Michel, T. Helfer and J. Fauque, ALCYONE: the fuel performance code of the PLEIADES platform dedicated to PWR fuel rods behavior. *Ann. Nucl. Energy* **207** (2024) 110711.
- [32] A. Bouloré, C. Struzik, V. Bouineau, F. Gaudier, G. Damblin and S. Bernaud, Modelling of UO<sub>2</sub> thermal conductivity: improvement of the irradiation defects contribution and uncertainty quantification. *Nucl. Eng. Des.* **407** (2023) 112304.
- [33] G. Damblin, P. Barbillon, M. Keller, A. Pasanisi and É. Parent, Adaptive numerical designs for the calibration of computer codes. *SIAM/ASA J. Uncertain. Quantif.* **6** (2018) 151–179.
- [34] C.P. Robert, G. Casella and G. Casella, *Monte Carlo Statistical Methods*, vol. 2. Springer (1999).
- [35] G. Damblin and P. Gaillard, Bayesian inference and non-linear extensions of the CIRCE method for quantifying the uncertainty of closure relationships integrated into thermal-hydraulic system codes. *Nucl. Eng. Des.* **359** (2020) 110391.
- [36] S. Barde, Bayesian estimation of large-scale simulation models with Gaussian process regression surrogates. *Computat. Statist. Data Anal.* **196** (2024) 107972.
- [37] M. Van Oijen, D. Cameron, K. Butterbach-Bahl, N. Farahbakhshazad, P.-E. Jansson, R. Kiese, K.H. Rahn, C. Werner and J. Yeluripati, A Bayesian framework for model calibration, comparison and analysis: application to four models for the biogeochemistry of a Norway spruce forest. *Agric. For. Meteorol.* **151** (2011) 1609–1621.
- [38] X. Wu, T. Mui, G. Hu, H. Meidani and T. Kozłowski, Inverse uncertainty quantification of TRACE physical model parameters using sparse grid stochastic collocation surrogate model. *Nucl. Eng. Des.* **319** (2017) 185–200.
- [39] A. Marrel, B. Iooss and V. Chabridon, The ICSCREAM methodology: identification of penalizing configurations in computer experiments using screening and metamodel: applications in thermal-hydraulics. *Nucl. Sci. Eng.* **196** (2022) 301–321.
- [40] A. Marrel and B. Iooss, Probabilistic surrogate modeling by Gaussian process: a new estimation algorithm for more robust prediction. *Reliabil. Eng. Syst. Saf.* **247** (2024) 110120.
- [41] B. Zhou, Z. Shi, S. Kucherenko and H. Zhao, A unified approach for global sensitivity analysis based on active subspace and Kriging. *Reliabil. Eng. Syst. Saf.* **217** (2022) 108080.
- [42] M. Binois and N. Wycoff, A survey on high-dimensional Gaussian process modeling with application to Bayesian optimization. *ACM Trans. Evol. Learn. Optim.* **2** (2022) 1–26.
- [43] H. Haario, E. Saksman and J. Tamminen, An adaptive Metropolis algorithm. *Bernoulli* **7** (2001) 223–242.
- [44] F. Liu, M. Bayarri and J. Berger, Modularization in Bayesian analysis, with emphasis on analysis of computer models. *Bayesian Anal.* **4** (2009) 119–150.
- [45] X. Wu, K. Shirvan and T. Kozłowski, Demonstration of the relationship between sensitivity and identifiability for inverse uncertainty quantification. *J. Computat. Phys.* **396** (2019) 2–30.
- [46] X. Xu, C. Sun, G. Huang and B.P. Mohanty, Global sensitivity analysis and calibration of parameters for a physically-based agro-hydrological model. *Environ. Model. Softw.* **83** (2016) 88–102.
- [47] M. Zambrano-Bigiarini, Z. Zając and S. Tarantola, Global sensitivity analysis for the calibration of a fully-distributed hydrological model, in *7th International Conference on Sensitivity Analysis of Model Output, MASCOT-SAMO* (2013).
- [48] A.B. Owen, Sobol' indices and Shapley value. *SIAM/ASA J. Uncertain. Quantif.* **2** (2014) 245–251.
- [49] M. Hérin, M. Il Idrissi, V. Chabridon and B. Iooss, Proportional marginal effects for global sensitivity analysis. *SIAM/ASA J. Uncertain. Quantif.* **12** (2024) 667–692.

- [50] F. Gamboa, A. Janon, T. Klein and A. Lagnoux, Sensitivity analysis for multidimensional and functional outputs. *Electron. J. Statist.* **8** (2014) 575–603.
- [51] F. Gamboa, T. Klein, A. Lagnoux and L. Moreno, Sensitivity analysis in general metric spaces. *Reliabil. Eng. Syst. Saf.* **212** (2021) 107611.
- [52] A. Janon, T. Klein, A. Lagnoux, M. Nodet and C. Prieur, Asymptotic normality and efficiency of two Sobol index estimators. *ESAIM Probab. Statist.* **18** (2014) 342–364.
- [53] A. Marrel, B. Iooss, B. Laurent and O. Roustant, Calculations of Sobol’ indices for the Gaussian process metamodel. *Reliabil. Eng. Syst. Saf.* **94** (2009) 742–751.
- [54] S. Da Veiga, F. Gamboa, A. Lagnoux, T. Klein and C. Prieur, New estimation of Sobol’ indices using kernels (2023). <https://arxiv.org/pdf/2303.17832>.
- [55] E. Borgonovo, A new uncertainty importance measure. *Reliabil. Eng. Syst. Saf.* **92** (2007) 771–784.
- [56] S. Rahman, The f-sensitivity index. *SIAM/ASA J. Uncertain. Quantif.* **4** (2016) 130–162.
- [57] M.R. El Amri and A. Marrel, Optimized HSIC-based tests for sensitivity analysis: application to thermohydraulic simulation of accidental scenario on nuclear reactor. *Qual. Reliabil. Eng. Int.* **38** (2022) 1386–1403.
- [58] M.R. El Amri and A. Marrel, More powerful HSIC-based independence tests, extension to space-filling designs and functional data. *Int. J. Uncertain. Quantif.* **14** (2024) 69–98.
- [59] G.-K. Delipei, J. Garnier, J. Le Pallec and B. Normand, Uncertainty analysis methodology for multi-physics coupled rod ejection accident, in *International Conference on Mathematics and Computational Methods Applied to Nuclear Science and Engineering (M&C 2019)* (2019).
- [60] N. Marie, S. Li, A. Marrel, M. Marquès, S. Bajard, A. Tosello, J. Perez, B. Grosjean, A. Gerschenfeld, M. Anderhuber et al., VVUQ of a thermal-hydraulic multi-scale tool on unprotected loss of flow accident in SFR reactor. *EPJ N Nucl. Sci. Technol.* **7** (2021) 3.
- [61] Á. Rollón de Pinedo, M. Couplet, B. Iooss, N. Marie, A. Marrel, E. Merle and R. Sueur, Functional outlier detection by means of h-mode depth and dynamic time warping. *Appl. Sci.* **11** (2021) 11475.
- [62] T. Wang, X. Dai and Y. Liu., Learning with Hilbert–Schmidt independence criterion: A review and new perspectives. *Knowl. Based Syst.* **234** (2021) 107567.
- [63] T. Wang, Z. Hu and H. Liu, A unified view of feature selection based on Hilbert–Schmidt independence criterion. *Chemometrics Intell. Lab. Syst.* **236** (2023) 104807.
- [64] T. Wang, J. Lu and G. Zhang, Two-stage fuzzy multiple kernel learning based on Hilbert–Schmidt independence criterion. *IEEE Trans. Fuzzy Syst.* **26** (2018) 3703–3714.
- [65] A. Gretton, O. Bousquet, A. Smola and B. Schölkopf, Measuring statistical dependence with Hilbert–Schmidt norms, in *International Conference on Algorithmic Learning Theory*. Springer (2005) 63–77.
- [66] S. Bernaud, I. Ramière, G. Latu and B. Michel, PLEIADES: A numerical framework dedicated to the multiphysics and multiscale nuclear fuel behavior simulation. *Ann. Nucl. Energy* **205** (2024) 110577.
- [67] B. Iooss and A. Marrel, Advanced methodology for uncertainty propagation in computer experiments with large number of inputs. *Nucl. Technol.* **205** (2019) 1588–1606.
- [68] C. Sriperumbudur, K. Fukumizu and G. Lanckriet, On the relation between universality, characteristic kernels and RKHS embedding of measures, in *Proceedings of the 13th International Conference on Artificial Intelligence and Statistics*. JMLR Workshop and Conference Proceedings (2010) 773–780.
- [69] J. Ziegel, D. Ginsbourger and L. Dümbgen, Characteristic kernels on Hilbert spaces, Banach spaces, and on sets of measures. *Bernoulli* **30** (2024) 1441–1457.
- [70] O. Baldé, G. Damblin, A. Marrel, A. Bouloré and L. Giraldi, Nonparametric Bayesian approach for quantifying the conditional uncertainty of input parameters in chained numerical models (2023). <https://arxiv.org/pdf/2307.01111>.
- [71] N. Fellmann, C. Blanchet-Scalliet, C. Helbert, A. Spagnol and D. Sinoquet, Kernel-based sensitivity analysis for (excursion) sets. *Technometrics* **66** (2024) 575–587.
- [72] L. Song, A. Smola, A. Gretton, J. Bedo and K. Borgwardt, Feature selection via dependence maximization. *J. Mach. Learn. Res.* **13** (2012) 1393–1434.
- [73] B.K. Sriperumbudur, A. Gretton, K. Fukumizu, B. Schölkopf and G.R. Lanckriet, Hilbert space embeddings and metrics on probability measures. *J. Mach. Learn. Res.* **11** (2010) 1517–1561.
- [74] Z. Szabó and B.K. Sriperumbudur, Characteristic and universal tensor product kernels. *J. Mach. Learn. Res.* **18** (2017) 1–29.

- [75] R. Serfling, Approximation Theorems of Mathematical Statistics. John Wiley & Sons (1980).
- [76] K. Zhang, J. Peters, D. Janzing and B. Schölkopf, Kernel-based conditional independence test and application in causal Discovery, in *Proceedings of the Twenty-Seventh Conference on Uncertainty in Artificial Intelligence* (2011) 804–813.
- [77] F. Kazi-Aoual, S. Hitier, R. Sabatier and J.-D. Lebreton, Refined approximations to permutation tests for multivariate inference. *Computat. Statist. Data Anal.* **20** (1995) 643–656.
- [78] M. Yamada, W. Jitkrittum, L. Sigal, E.P. Xing and M. Sugiyama, High-dimensional feature selection by feature-wise kernelized lasso. *Neural Computat.* **26** (2014) 185–207.
- [79] M. Kanagawa, P. Hennig, D. Sejdinovic and B.K. Sriperumbudur, Gaussian processes and kernel methods: a review on connections and Equivalences (2018). <https://arxiv.org/pdf/1807.02582>.
- [80] A. Gretton and L. Györfi, Consistent nonparametric tests of independence. *J. Mach. Learn. Res.* **11** (2010) 1391–1423.



**Please help to maintain this journal in open access!**

This journal is currently published in open access under the Subscribe to Open model (S2O). We are thankful to our subscribers and supporters for making it possible to publish this journal in open access in the current year, free of charge for authors and readers.

Check with your library that it subscribes to the journal, or consider making a personal donation to the S2O programme by contacting [subscribers@edpsciences.org](mailto:subscribers@edpsciences.org).

More information, including a list of supporters and financial transparency reports, is available at <https://edpsciences.org/en/subscribe-to-open-s2o>.

APPENDIX A. PRELIMINARIES ON CHARACTERISTIC KERNELS

**A.1 Characteristic kernels**

Throughout this section,  $\mathcal{Z}$  is a separable metric space. The set of all Borel probability measures on  $\mathcal{Z}$  is denoted by  $\mathcal{M}_1^+(\mathcal{Z})$ .

**Definition A.1** (Kernel mean embedding). Let  $K : \mathcal{Z} \times \mathcal{Z} \rightarrow \mathbb{R}$  be a kernel. For any probability distribution  $\nu \in \mathcal{M}_1^+(\mathcal{Z})$  such that  $\mathbb{E}_\nu [\sqrt{K(Z, Z)}] < \infty$ , the *kernel mean embedding* of  $\nu$  is the function defined as:

$$\begin{aligned} \mu_\nu : \mathcal{Z} &\longrightarrow \mathbb{R} \\ z &\longmapsto \mathbb{E}_\nu[K(z, Z)] = \int_{\mathcal{Z}} K(z, z') \, d\nu(z'). \end{aligned} \tag{A.1}$$

In particular,  $\mu_\nu$  belongs to the reproducing kernel Hilbert space (RKHS)  $\mathcal{H}$  induced by  $K$ .

**Definition A.2** (Characteristic kernel). Let  $K : \mathcal{Z} \times \mathcal{Z} \rightarrow \mathbb{R}$  be a bounded kernel. This assumption ensures that  $\mu_\nu$  is well-defined for any probability measure  $\nu \in \mathcal{M}_1^+(\mathcal{Z})$ . The kernel  $K$  is said to be characteristic on  $\mathcal{Z}$  if the mapping  $\nu \in \mathcal{M}_1^+(\mathcal{Z}) \mapsto \mu_\nu \in \mathcal{H}$  is injective. This means that:

$$\forall \nu_1, \nu_2 \in \mathcal{M}_1^+(\mathcal{Z}), \quad \mu_{\nu_1} = \mu_{\nu_2} \implies \nu_1 = \nu_2. \tag{A.2}$$

**Definition A.3** (Maximum mean discrepancy). Let  $K : \mathcal{Z} \times \mathcal{Z} \rightarrow \mathbb{R}$  be a bounded kernel, and let  $\mathcal{H}$  denote the associated RKHS. For any probability measures  $\nu_1, \nu_2 \in \mathcal{M}_1^+(\mathcal{Z})$ , the maximum mean discrepancy (MMD) is defined by:

$$\text{MMD}(\nu_1, \nu_2) := \|\mu_{\nu_1} - \mu_{\nu_2}\|_{\mathcal{H}}. \tag{A.3}$$

If  $K$  is characteristic on  $\mathcal{Z}$ , the MMD induces a proper metric on  $\mathcal{M}_1^+(\mathcal{Z})$ , that is:

$$\text{MMD}(\nu_1, \nu_2) = 0 \iff \mu_{\nu_1} = \mu_{\nu_2} \iff \nu_1 = \nu_2 . \quad (\text{A.4})$$

**Definition A.4** (Integrally strictly positive definite kernel). Let  $K : \mathcal{Z} \times \mathcal{Z} \rightarrow \mathbb{R}$  be a bounded kernel. It is said to be integrally strictly positive definite on  $\mathcal{Z}$  if:

$$\forall \nu \in \mathcal{M}_1^+(\mathcal{Z}), \quad \int_{\mathcal{Z}} \int_{\mathcal{Z}} K(z, z') d\nu(z) d\nu(z') > 0 . \quad (\text{A.5})$$

**Theorem A.5** (Theorem 7 in [73]). *A bounded kernel  $K : \mathcal{Z} \times \mathcal{Z} \rightarrow \mathbb{R}$  is characteristic on  $\mathcal{Z}$  if and only if it is integrally strictly positive definite.*

This result establishes that being integrally strictly positive definite is a sufficient condition for being characteristic. As shown in [69], a specific class of kernels is integrally strictly positive definite, and therefore characteristic.

**Proposition A.6** (Proposition 5.2 in [69]). *Let  $\mathcal{F}$  be a separable Hilbert space. Let  $T : \mathcal{Z} \rightarrow \mathcal{F}$  and  $\varphi : [0, \infty) \rightarrow \mathbb{R}$  be two measurable maps. Assume that  $T$  is injective and that the function:*

$$\forall \mathbf{x}, \mathbf{x}' \in \mathbb{R}^q, \quad k_\varphi(\mathbf{x}, \mathbf{x}') = \varphi(\|\mathbf{x} - \mathbf{x}'\|_{\mathbb{R}^q}^2) \quad (\text{A.6})$$

defines a characteristic kernel on  $\mathbb{R}^q$ . Then, the function:

$$\forall z, z' \in \mathcal{Z}, \quad K(z, z') = \varphi(\|T(z) - T(z')\|_{\mathcal{F}}^2) \quad (\text{A.7})$$

defines a characteristic kernel on  $\mathcal{Z}$ .

## A.2 Translation-invariant kernels

**Definition A.7** (Translation-invariant function). A function  $K : \mathbb{R}^q \times \mathbb{R}^q \rightarrow \mathbb{R}$  is said to be translation-invariant if there exists a shift function  $\psi : \mathbb{R}^q \rightarrow \mathbb{R}$  such that:

$$\forall \mathbf{x}, \mathbf{x}' \in \mathbb{R}^q, \quad K(\mathbf{x}, \mathbf{x}') = \psi(\mathbf{x} - \mathbf{x}') . \quad (\text{A.8})$$

**Assumption A.8.** Let  $K : \mathbb{R}^q \times \mathbb{R}^q \rightarrow \mathbb{R}$  be a translation-invariant function. It is assumed that its shift function  $\psi : \mathbb{R}^q \rightarrow \mathbb{R}$  is continuous and integrable on  $\mathbb{R}^q$ :

$$\psi \in C^0(\mathbb{R}^q) \cap L^1(\mathbb{R}^q) . \quad (\text{A.9})$$

Note that Assumption A.8 is not stringent, as it is verified by most kernels used in practice (especially Gaussian and Matérn kernels).

**Theorem A.9** (Bochner). *Let  $K : \mathbb{R}^q \times \mathbb{R}^q \rightarrow \mathbb{R}$  be a translation-invariant function satisfying Assumption A.8. Then,  $K$  is a kernel if and only if the Fourier transform (FT) of  $\psi$ , defined by:*

$$\forall \mathbf{w} \in \mathbb{R}^q, \quad \widehat{\psi}(\mathbf{w}) = \frac{1}{(\sqrt{2\pi})^q} \int_{\mathbb{R}^q} e^{-i\mathbf{w}^t \mathbf{x}} \psi(\mathbf{x}) d\mathbf{x} , \quad (\text{A.10})$$

is a real-valued non-negative function. Moreover, when  $K$  is a kernel, the associated shift function  $\psi$  is more commonly called the signature of  $K$ .

**Proposition A.10** (Proposition 8 in [68]). *Let  $K : \mathbb{R}^q \times \mathbb{R}^q \rightarrow \mathbb{R}$  be a translation-invariant kernel satisfying Assumption A.8. Then,  $K$  is characteristic on  $\mathbb{R}^q$  if and only if the FT of  $\psi$  is positive on  $\mathbb{R}^q$ , which can be equivalently*

expressed as:

$$\text{supp}(\widehat{\psi}) := \overline{\{\mathbf{w} \in \mathbb{R}^q : \widehat{\psi}(\mathbf{w}) > 0\}} = \mathbb{R}^q. \quad (\text{A.11})$$

### A.3 Tensor product of univariate kernels

In this section, the objective is to construct a characteristic kernel on  $\mathbb{R}^q$  that could be used to compute the HSIC on a vector output (see Approach B in Sect. 3.2). Let  $\{K_i\}_{i=1}^q$  be a set of  $q$  univariate kernels. The *tensor-product kernel* is defined by:

$$\forall \mathbf{x}, \mathbf{x}' \in \mathbb{R}^q, \quad K_{\text{tens}}(\mathbf{x}, \mathbf{x}') := \left[ \bigotimes_{i=1}^q K_i \right](\mathbf{x}, \mathbf{x}') = \prod_{i=1}^q K_i(x_i, x'_i). \quad (\text{A.12})$$

The question is whether taking characteristic kernels  $K_i$  on  $\mathbb{R}$  is sufficient to obtain a characteristic kernel  $K_{\text{tens}}$  on  $\mathbb{R}^q$ . In general, the answer is no. Indeed, a counterexample for  $q = 3$  may be found in [74]. However, under additional assumptions on the kernels  $K_i$ , it is possible to show that  $K_{\text{tens}}$  remains characteristic.

We first establish that the result holds when tensorizing Gaussian kernels (see Prop. A.11). The proof follows almost immediately from Proposition A.6. We then extend the result to the case where all kernels  $K_i$  are translation-invariant (see Prop. A.12). In this case, the proof consists in computing the FT of the signature  $\psi_{\text{tens}}$  and showing that its support is  $\mathbb{R}^q$ .

**Proposition A.11.** *Assume that the kernels  $\{K_i\}_{i=1}^q$  are all Gaussian kernels:*

$$\forall x_i, x'_i \in \mathbb{R}, \quad K_i(x_i, x'_i) = \exp\left[-\frac{1}{2} \left(\frac{x_i - x'_i}{\gamma_i}\right)^2\right]. \quad (\text{A.13})$$

Then, the kernel  $K_{\text{tens}}$  is characteristic on  $\mathbb{R}^q$ .

*Proof.* If taking Gaussian kernels, equation (A.12) becomes:

$$K_{\text{tens}}(\mathbf{x}, \mathbf{x}') = \prod_{i=1}^q \exp\left[-\frac{1}{2} \left(\frac{x_i - x'_i}{\gamma_i}\right)^2\right] = \exp\left[-\frac{1}{2} \sum_{i=1}^q \left(\frac{x_i - x'_i}{\gamma_i}\right)^2\right] \quad (\text{A.14})$$

$$= \varphi(\|T(\mathbf{x}) - T(\mathbf{x}')\|^2), \quad (\text{A.15})$$

where the last equation is obtained by identification with:

- $\varphi : z \mapsto \exp(-|z|)$  ;
- $T : \mathbf{x} = (x_1, \dots, x_d) \mapsto (x_1/\gamma_1, \dots, x_d/\gamma_d)/\sqrt{2}$ .

On the one hand,  $T$  is a homothety and is therefore injective. On the other hand, with this choice of  $\varphi$ , the kernel  $k_\varphi$  from equation (A.6) is a multivariate Laplacian kernel, which is known to be characteristic on  $\mathbb{R}^q$ . Since all assumptions of Proposition A.6 are satisfied, one can conclude that  $K_{\text{tens}}$  is integrally strictly positive definite, hence characteristic on  $\mathbb{R}^q$ .  $\square$

**Proposition A.12.** *Assume that each kernel  $K_i$  is translation-invariant, satisfies Assumption A.8, and is characteristic on  $\mathbb{R}$ . Then, the kernel  $K_{\text{tens}}$  is characteristic on  $\mathbb{R}^q$ .*

*Proof.* First, observe that  $K_{\text{tens}}$  is a translation-invariant kernel on  $\mathbb{R}^q$ :

$$K_{\text{tens}}(\mathbf{x}, \mathbf{x}') = \prod_{i=1}^q K_i(x_i, x'_i) = \prod_{i=1}^q \psi_i(x_i - x'_i) = \psi_{\text{tens}}(\mathbf{x} - \mathbf{x}'). \quad (\text{A.16})$$

Its signature  $\psi_{\text{tens}}$  is given by the tensor product of the signatures  $\psi_i$  of the kernels  $K_i$ :

$$\forall \mathbf{r} \in \mathbb{R}^q, \quad \psi_{\text{tens}}(\mathbf{r}) := \left[ \bigotimes_{i=1}^q \psi_i \right] (\mathbf{r}) = \prod_{i=1}^q \psi_i(r_i). \quad (\text{A.17})$$

According to Proposition A.10, in order to prove that  $K_{\text{tens}}$  is characteristic on  $\mathbb{R}^q$ , it is sufficient to prove that the FT of  $\psi_{\text{tens}}$  is positive on  $\mathbb{R}^q$ . A direct consequence of the Fubini-Tonelli theorem is the separability of the FT:

$$\begin{aligned} \widehat{\psi}_{\text{tens}}(\mathbf{w}) &= \frac{1}{(\sqrt{2\pi})^q} \int_{\mathbb{R}^q} e^{-i\mathbf{w}^t \mathbf{x}} \psi_{\text{tens}}(\mathbf{x}) \, d\mathbf{x} \\ &= \prod_{i=1}^q \left( \frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} e^{-iw_i x_i} \psi_i(w_i) \, dw_i \right) = \prod_{i=1}^q \widehat{\psi}_i(w_i). \end{aligned} \quad (\text{A.18})$$

As all kernels  $K_i$  are characteristic on  $\mathbb{R}$ , the FTs  $\widehat{\psi}_i$  are positive on  $\mathbb{R}$ . It follows that  $\widehat{\psi}$  is positive on  $\mathbb{R}^q$ , which completes the proof by Proposition A.10.  $\square$

#### APPENDIX B. TECHNICAL RESULTS ON U- AND V-STATISTICS

For  $s \geq 2$ , let  $Z_1, \dots, Z_s$  be independent and identically distributed (i.i.d.) random variables with common distribution  $\mathbb{P}_Z \in \mathcal{M}_1^+(\mathcal{Z})$ . Let  $\eta : \mathcal{Z}^s \rightarrow \mathbb{R}$  be a measurable function. Throughout Appendix B, the target quantity is the following expectation:

$$\zeta := \mathbb{E}_{\mathbb{P}_Z} [\eta(Z_1, \dots, Z_s)] = \int_{\mathcal{Z}} \dots \int_{\mathcal{Z}} \eta(z_1, \dots, z_s) \, d\mathbb{P}_Z(z_1) \dots d\mathbb{P}_Z(z_s). \quad (\text{B.1})$$

In this multivariate setting, the function  $\eta$  is said to be *symmetric* if it is invariant under any permutation of its arguments:

$$\forall \sigma \in \mathbb{S}_s, \quad \forall \mathbf{z} \in \mathcal{Z}^s, \quad \eta(z_{\sigma(1)}, \dots, z_{\sigma(s)}) = \eta(z_1, \dots, z_s), \quad (\text{B.2})$$

where  $\mathbb{S}_s$  denotes the set of all permutations of  $\{1, \dots, s\}$ . If the function  $\eta$  is not symmetric, it must be replaced by its symmetrized version  $\widetilde{\eta}$ , defined by:

$$\forall \mathbf{z} \in \mathcal{Z}^s, \quad \widetilde{\eta}(z_1, \dots, z_s) := \frac{1}{s!} \sum_{\sigma \in \mathbb{S}_s} \eta(z_{\sigma(1)}, \dots, z_{\sigma(s)}). \quad (\text{B.3})$$

By construction,  $\widetilde{\eta}$  is symmetric. Note that  $\eta = \widetilde{\eta}$  if and only if  $\eta$  is already symmetric. Moreover, replacing  $\eta$  with  $\widetilde{\eta}$  leaves the target quantity unchanged:

$$\mathbb{E}_{\mathbb{P}_Z} [\widetilde{\eta}(Z_1, \dots, Z_s)] = \mathbb{E}_{\mathbb{P}_Z} [\eta(Z_1, \dots, Z_s)] = \zeta, \quad (\text{B.4})$$

because the vectors  $(Z_1, \dots, Z_s)$  and  $(Z_{\sigma(1)}, \dots, Z_{\sigma(s)})$  have the same joint distribution for any permutation  $\sigma \in \mathbb{S}_s$ .

**Definition B.1** (U-statistic). Let  $(Z^{(1)}, \dots, Z^{(n)})$  be an i.i.d. sample from the distribution  $\mathbb{P}_Z$ , with  $n \geq s$ . The U-statistic estimator of  $\zeta$  is constructed as follows:

$$\widehat{\zeta}^U = \binom{n}{s}^{-1} \sum_{\mathbf{i} \in \mathcal{A}_n^s} \widetilde{\eta}(Z^{(i_1)}, \dots, Z^{(i_s)}), \quad (\text{B.5})$$

where the summation set is:

$$\mathcal{A}_n^s = \{\mathbf{i} := (i_1, \dots, i_s) : 1 \leq i_1 < \dots < i_s \leq n\} \quad \text{with} \quad \text{Card}(\mathcal{A}_n^s) = \binom{n}{s} = \frac{n!}{s!(n-s)!}. \quad (\text{B.6})$$

**Definition B.2** (V-statistic). Let  $(Z^{(1)}, \dots, Z^{(n)})$  be an i.i.d. sample from the distribution  $\mathbb{P}_Z$ , with  $n \geq s$ . The V-statistic estimator of  $\zeta$  is constructed as follows:

$$\widehat{\zeta}^V = \frac{1}{n^s} \sum_{\mathbf{i} \in \mathcal{B}_n^s} \widetilde{\eta}(Z^{(i_1)}, \dots, Z^{(i_s)}), \quad (\text{B.7})$$

where the summation set is:

$$\mathcal{B}_n^s = \{\mathbf{i} := (i_1, \dots, i_s) : 1 \leq i_1, \dots, i_s \leq n\} \text{ with } \text{Card}(\mathcal{B}_n^s) = n^s. \quad (\text{B.8})$$

It is straightforward to see that  $\widehat{\zeta}^U$  is unbiased. The constraint  $\mathbf{i} \in \mathcal{A}_n^s$  on the multi-index  $\mathbf{i} = (i_1, i_2, \dots, i_s)$  ensures that  $i_1 < i_2 < \dots < i_s$ , so that each term  $\widetilde{\eta}(Z^{(i_1)}, \dots, Z^{(i_s)})$  has expectation exactly equal to  $\zeta$ . On the contrary,  $\widehat{\zeta}^V$  is biased because its summation set  $\mathcal{B}_n^s$  includes some multi-indices  $\mathbf{i}$  for which  $i_a = i_b$  with  $(a, b) \in \{1, \dots, s\}^2$  and  $a \neq b$ .

$\widehat{\zeta}^U$  is the unbiased estimator of  $\zeta$  with the smallest variance. The following lemma provides a simple upper bound for its variance.

**Lemma B.3** (Lemma A of Sect. 5.2.1 in [75]). *Assume that  $\mathbb{E}_{\mathbb{P}_Z} [|\widetilde{\eta}(Z_1, \dots, Z_s)|^2] < \infty$ . Then, the variance of  $\widehat{\zeta}^U$  satisfies:*

$$0 \leq \mathbb{V}(\widehat{\zeta}^U) \leq \frac{s}{n} \mathbb{V}_{\mathbb{P}_Z}(\widetilde{\eta}(Z_1, \dots, Z_s)). \quad (\text{B.9})$$

An immediate consequence of this lemma is that the variance of  $\widehat{\zeta}^U$  vanishes as  $n \rightarrow \infty$ .

Asymptotically,  $\widehat{\zeta}^U$  and  $\widehat{\zeta}^V$  behave similarly. This is formalized by the following lemma, which provides convergence rates for the moments of their difference.

**Lemma B.4** (Lemma of Sect. 5.7.3 in [75]). *Let  $r \geq 1$  be fixed. Assume that:*

$$\forall 1 \leq i_1, \dots, i_s \leq s, \quad \mathbb{E}_{\mathbb{P}_Z} [|\widetilde{\eta}(Z_{i_1}, \dots, Z_{i_s})|^r] < \infty. \quad (\text{B.10})$$

Then, one has:

$$\mathbb{E}_{\mathbb{P}_Z} \left[ \left| \widehat{\zeta}^U - \widehat{\zeta}^V \right|^r \right] = \mathcal{O} \left( \frac{1}{n^r} \right). \quad (\text{B.11})$$

In this work, this Big-O convergence rate is insufficient. Instead, we require a moment bound with an explicit constant. From this perspective, Lemma B.5 fills the gap left by Lemma B.4.

**Lemma B.5.** *Assuming that the conditions of Lemma B.4 hold, one has:*

$$\mathbb{E} \left[ \left| \widehat{\zeta}^U - \widehat{\zeta}^V \right|^r \right] \leq C_r \left( \frac{s(s-1)}{n} \right)^r \quad (\text{B.12})$$

with:

$$C_r = C_r(\widetilde{\eta}, \mathbb{P}_Z) := \max_{1 \leq i_1, \dots, i_s \leq s} \mathbb{E} [|\widetilde{\eta}(Z_{i_1}, \dots, Z_{i_s})|^r]. \quad (\text{B.13})$$

The book by Serfling [75] only provides a sketch of the proof of Lemma B.4. Following the line of argument suggested there, one can recover Lemma B.5. For completeness, we include below a detailed and self-contained proof of Lemma B.5.

*Proof.* For technical reasons, we introduce the following two index sets:

$$\mathcal{E}_n^s := \{\mathbf{i} = (i_1, \dots, i_s) : 1 \leq i_1 \neq \dots \neq i_s \leq n\} \quad (\text{B.14})$$

$$\mathcal{F}_n^s := \{\mathbf{i} = (i_1, \dots, i_s) : \exists (a, b) \in \{1, \dots, s\}^2 \text{ such that } i_a = i_b \text{ and } a \neq b\}. \quad (\text{B.15})$$

The elements of  $\mathcal{E}_n^s$  are obtained by applying all permutations in  $\mathbb{S}_s$  to each element of  $\mathcal{A}_n^s$ . The set  $\mathcal{F}_n^s$  is the complement of  $\mathcal{E}_n^s$  in  $\mathcal{B}_n^s$ . Consequently, their cardinalities are given by:

$$\text{Card}(\mathcal{E}_n^s) = s! \text{Card}(\mathcal{A}_n^s) = s! \binom{n}{s} = \frac{n!}{(n-s)!} = n(n-1) \dots (n-s+1) =: (n)_s \quad (\text{B.16})$$

$$\text{Card}(\mathcal{F}_n^s) = \text{Card}(\mathcal{B}_n^s) - \text{Card}(\mathcal{E}_n^s) = n^s - (n)_s. \quad (\text{B.17})$$

The key element of the proof is the statistic:

$$\widehat{\zeta}^{\text{W}} := \frac{1}{n^s - (n)_s} \sum_{\mathbf{i} \in \mathcal{F}_n^s} \tilde{\eta} \left( Z^{(i_1)}, \dots, Z^{(i_s)} \right), \quad (\text{B.18})$$

which is defined analogously to  $\widehat{\zeta}^{\text{U}}$  and  $\widehat{\zeta}^{\text{V}}$ , except that the averaging is performed over the index set  $\mathcal{F}_n^s$  (instead of  $\mathcal{A}_n^s$  and  $\mathcal{B}_n^s$ ). From now, in order to simplify the notation throughout the proof, we write  $\tilde{\eta}_i$  for the terms appearing in  $\widehat{\zeta}^{\text{U}}$ ,  $\widehat{\zeta}^{\text{V}}$ , and  $\widehat{\zeta}^{\text{W}}$ .

To prove the inequality in equation (B.12), we proceed in several steps, as recommended in [75]. First, we show that:

$$n^s \left( \widehat{\xi}^{\text{U}} - \widehat{\xi}^{\text{V}} \right) = (n^s - (n)_s) \left( \widehat{\xi}^{\text{U}} - \widehat{\xi}^{\text{W}} \right). \quad (\text{B.19})$$

Next, we establish that:

$$n^s - (n)_s = \mathcal{O}(n^{s-1}). \quad (\text{B.20})$$

Combining equations (B.19) and (B.20), the desired inequality follows easily.

Let us see how to establish equation (B.19). From the definitions of the index sets, one has:

$$n^s \widehat{\zeta}^{\text{V}} = \sum_{\mathbf{i} \in \mathcal{B}_n^s} \tilde{\eta}_i = \sum_{\mathbf{i} \in \mathcal{E}_n^s} \tilde{\eta}_i + \sum_{\mathbf{i} \in \mathcal{F}_n^s} \tilde{\eta}_i \quad (\text{B.21})$$

$$= s! \sum_{\mathbf{i} \in \mathcal{A}_n^s} \tilde{\eta}_i + \sum_{\mathbf{i} \in \mathcal{F}_n^s} \tilde{\eta}_i \quad (\text{B.22})$$

$$= (n)_s \widehat{\zeta}^{\text{U}} + (n^s - (n)_s) \widehat{\zeta}^{\text{W}}. \quad (\text{B.23})$$

Equation (B.21) follows from the disjoint union  $\mathcal{B}_n^s = \mathcal{E}_n^s \sqcup \mathcal{F}_n^s$ . Equation (B.22) uses the symmetry of the function  $\tilde{\eta}$ : under this assumption, summing over  $\mathcal{E}_n^s$  is equivalent to summing over  $\mathcal{A}_n^s$  and multiplying by  $s! = \text{Card}(\mathbb{S}_s)$ . Equation (B.23) directly stems from Definition B.1 and equation (B.18). Finally, equation (B.23) yields equation (B.19) after straightforward manipulations.

To demonstrate equation (B.20), we use induction on  $s$ . Initialization for  $s = 1$  is trivial:

$$n^s - (n)_s = n - n = 0 = \mathcal{O}(1).$$

Then, assume that equation (B.20) holds for some  $s \geq 1$ . It must be shown that it also holds at order  $s + 1$ . One can write:

$$n^{s+1} - (n)_{s+1} = n^{(s+1)} - n(n-1)(n-2) \dots (n-s) \quad (\text{B.24})$$

$$= n [n^s - (n-1)(n-2) \dots (n-s)] \quad (\text{B.25})$$

$$= n[n^s - (n)_s] + n[(n)_s - (n-1)_s]. \quad (\text{B.26})$$

On the one hand, by the induction hypothesis, there exists  $\gamma_s > 0$  such that:

$$n^s - (n)_s \leq \gamma_s n^{s-1}. \quad (\text{B.27})$$

On the other hand, straightforward algebra yields:

$$(n)_s - (n-1)_s = n(n-1)(n-2)\dots(n-s+1) - (n-1)(n-2)\dots(n-s) \quad (\text{B.28})$$

$$= [(n-1)(n-2)\dots(n-s+1)][n - (n-s)] \quad (\text{B.29})$$

$$= (n-1)(n-2)\dots(n-s+1)s \quad (\text{B.30})$$

$$\leq n^{s-1}s. \quad (\text{B.31})$$

To obtain the last inequality, note that each of the  $s-1$  factors in  $(n-1)(n-2)\dots(n-s+1)$  is bounded above by  $n$ . Then, by substituting equations (B.27) and (B.31) in equation (B.26), we obtain:

$$n^{s+1} - (n)_{s+1} \leq (s + \gamma_s)n^s = \gamma_{s+1}n^s \quad \text{with} \quad \gamma_{s+1} = s + \gamma_s. \quad (\text{B.32})$$

This shows that equation (B.20) holds at order  $s+1$ , and therefore completes the induction. Moreover, equation (B.32) allows the constants  $\gamma_s$  to be determined explicitly. From the initialization step, one has  $\gamma_1 = 0$ . The recursion  $\gamma_{s+1} = s + \gamma_s$  then gives:

$$\gamma_s = (s-1) + \gamma_{s-1} = (s-1) + (s-2) + \dots + 2 + 1 + \gamma_1 = \sum_{k=1}^{s-1} k = \frac{s(s-1)}{2}. \quad (\text{B.33})$$

Consequently, equation (B.20) can be written explicitly as:

$$n^s - (n)_s \leq \left(\frac{s(s-1)}{2}\right)n^{s-1}. \quad (\text{B.34})$$

Now that equations (B.19) and (B.20) have been established, we can derive equation (B.12).

In Lemma B.4, for the chosen exponent  $r \geq 1$ , it was assumed that:

$$\forall 1 \leq i_1, \dots, i_s \leq s, \quad \mathbb{E}_{\mathbb{P}_Z} [|\tilde{\eta}(Z_{i_1}, \dots, Z_{i_s})|^r] < \infty. \quad (\text{B.35})$$

This assumption ensures finiteness of the constant:

$$C_r = \max_{1 \leq i_1, \dots, i_s \leq s} \mathbb{E}_{\mathbb{P}_Z} [|\tilde{\eta}(Z_{i_1}, \dots, Z_{i_s})|^r]. \quad (\text{B.36})$$

In the context of Lemma B.5,  $C_r$  is a constant with respect to the given i.i.d. sample, but it depends on  $\tilde{\eta}$  and  $\mathbb{P}_Z$ .

Combining equations (B.19) and (B.34) yields the following inequalities:

$$(\hat{\zeta}^U - \hat{\zeta}^V) \leq \left(\frac{\gamma_s}{n}\right) (\hat{\zeta}^U - \hat{\zeta}^W) \quad (\text{B.37})$$

$$|\hat{\zeta}^U - \hat{\zeta}^V|^r \leq \left(\frac{\gamma_s}{n}\right)^r |\hat{\zeta}^U - \hat{\zeta}^W|^r \quad (\text{B.38})$$

$$\mathbb{E} [|\hat{\zeta}^U - \hat{\zeta}^V|^r] \leq \left(\frac{\gamma_s}{n}\right)^r \mathbb{E} [|\hat{\zeta}^U - \hat{\zeta}^W|^r]. \quad (\text{B.39})$$

Using the triangle inequality in the space  $\mathbb{L}^r(\Omega)$  of all random variables with finite  $r$ -th moment, we obtain:

$$\left(\mathbb{E}\left[|\widehat{\zeta}^{\mathbb{U}} - \widehat{\zeta}^{\mathbb{W}}|^r\right]\right)^{1/r} = \|\widehat{\zeta}^{\mathbb{U}} - \widehat{\zeta}^{\mathbb{W}}\|_{\mathbb{L}^r} \leq \|\widehat{\zeta}^{\mathbb{U}}\|_{\mathbb{L}^r} + \|\widehat{\zeta}^{\mathbb{W}}\|_{\mathbb{L}^r} . \quad (\text{B.40})$$

Raising both sides of the inequality to the power  $r$  results in:

$$\|\widehat{\zeta}^{\mathbb{U}} - \widehat{\zeta}^{\mathbb{W}}\|_{\mathbb{L}^r}^r \leq \left(\|\widehat{\zeta}^{\mathbb{U}}\|_{\mathbb{L}^r} + \|\widehat{\zeta}^{\mathbb{W}}\|_{\mathbb{L}^r}\right)^r \leq 2^{r-1} \left(\|\widehat{\zeta}^{\mathbb{U}}\|_{\mathbb{L}^r}^r + \|\widehat{\zeta}^{\mathbb{W}}\|_{\mathbb{L}^r}^r\right) . \quad (\text{B.41})$$

The two inequalities above follow from the properties of the power function  $\varphi_r(t) = t^r$  on  $[0, +\infty)$ , which is notably increasing and convex for  $r \geq 1$ . Specifically, the second inequality is an instance of Jensen's inequality applied to a sum of two non-negative terms with equal weights  $1/2$ . After replacing  $\mathbb{L}^r$ -norms with expectations, equation (B.41) writes:

$$\mathbb{E}\left[|\widehat{\zeta}^{\mathbb{U}} - \widehat{\zeta}^{\mathbb{W}}|^r\right] \leq 2^{r-1} \left(\mathbb{E}\left[|\widehat{\zeta}^{\mathbb{U}}|^r\right] + \mathbb{E}\left[|\widehat{\zeta}^{\mathbb{W}}|^r\right]\right) . \quad (\text{B.42})$$

It now remains to bound the two expectations appearing on the right-hand side. With more convenient notations, the U-statistic defined in equation (B.1) can be written as:

$$\widehat{\zeta}^{\mathbb{U}} = \frac{1}{\text{Card}(\mathcal{A}_n^s)} \sum_{i \in \mathcal{A}_n^s} \widetilde{\eta}_i . \quad (\text{B.43})$$

In the same spirit as equations (B.40) and (B.41), by successively applying the triangle inequality, the monotonicity of  $\varphi_r$ , and Jensen's inequality, we obtain:

$$|\widehat{\zeta}^{\mathbb{U}}| \leq \frac{1}{\text{Card}(\mathcal{A}_n^s)} \sum_{i \in \mathcal{A}_n^s} |\widetilde{\eta}_i| \quad (\text{B.44})$$

$$|\widehat{\zeta}^{\mathbb{U}}|^r \leq \left(\frac{1}{\text{Card}(\mathcal{A}_n^s)} \sum_{i \in \mathcal{A}_n^s} |\widetilde{\eta}_i|\right)^r \leq \frac{1}{\text{Card}(\mathcal{A}_n^s)} \sum_{i \in \mathcal{A}_n^s} |\widetilde{\eta}_i|^r . \quad (\text{B.45})$$

Taking the expectation on both sides leads to:

$$\mathbb{E}\left[|\widehat{\zeta}^{\mathbb{U}}|^r\right] \leq \frac{1}{\text{Card}(\mathcal{A}_n^s)} \sum_{i \in \mathcal{A}_n^s} \mathbb{E}\left[|\widetilde{\eta}_i|^r\right] \leq C_r \times \frac{1}{\text{Card}(\mathcal{A}_n^s)} \sum_{i \in \mathcal{A}_n^s} 1 = C_r . \quad (\text{B.46})$$

The same reasoning applies to the other expectation. Injecting the bound  $C_r$  in equation (B.42) yields:

$$\mathbb{E}\left[|\widehat{\zeta}^{\mathbb{U}} - \widehat{\zeta}^{\mathbb{W}}|^r\right] \leq 2^r C_r . \quad (\text{B.47})$$

Finally, bringing together equations (B.39) and (B.47) gives:

$$\mathbb{E}\left[|\widehat{\zeta}^{\mathbb{U}} - \widehat{\zeta}^{\mathbb{V}}|^r\right] \leq \left(\frac{\gamma_s}{n}\right)^r 2^r C_r \leq \left(\frac{s(s-1)}{n}\right)^r C_r . \quad (\text{B.48})$$

□

## APPENDIX C. A CLOSER LOOK AT THE HSIC

This section focuses on the Hilbert-Schmidt independence criterion (HSIC). First, Appendix C.1 introduces the HSIC between two random objects and presents several equivalent formulations. Next, Appendix C.2 describes how to estimate the HSIC from a sample of joint observations of the two objects. Then, Appendix C.3 explains how the HSIC can be used to test independence between two objects. Following this, Appendix C.4 provides several statistical examples illustrating the use of HSIC tests. Finally, Appendix C.5 addresses the particular case of parameter screening for model calibration.

Throughout this section,  $X$  and  $Y$  denote the two random objects under consideration. It is only assumed that they take values in separable spaces  $\mathcal{X}$  and  $\mathcal{Y}$ . Their marginal distributions are denoted by  $\mathbb{P}_X$  and  $\mathbb{P}_Y$ , respectively. The pair  $Z := (X, Y)$  takes values in the space  $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$  and follows the joint distribution  $\mathbb{P}_Z = \mathbb{P}_{XY}$ .

**C.1 Definition of the HSIC**

Any dependence measure aims to quantify the discrepancy between the joint distribution  $\mathbb{P}_{XY}$ , representing the actual dependence within the pair  $Z = (X, Y)$  and the product of the marginal distributions  $\mathbb{P}_X \otimes \mathbb{P}_Y$ , representing a hypothetical situation of independence between the two objects. The HSIC is one such dependence measure, which captures statistical dependence using covariance kernels.

The idea behind the HSIC is to transport the distributions  $\mathbb{P}_{XY}$  and  $\mathbb{P}_X \otimes \mathbb{P}_Y$  into an RKHS of functions  $h : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ , where their discrepancy can be more easily quantified. To this end, a kernel is assigned to both objects. Let  $K_x : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  and  $K_y : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  denote the kernels associated with  $X$  and  $Y$ , respectively. The corresponding RKHSs are denoted by  $\mathcal{H}_x$  and  $\mathcal{H}_y$ . The HSIC is then defined as the square of the MMD between  $\mathbb{P}_{XY}$  and  $\mathbb{P}_X \otimes \mathbb{P}_Y$  in the RKHS  $\mathcal{H}_x \otimes \mathcal{H}_y$  induced by the tensor-product kernel  $K_x \otimes K_y$  [65]. This can be written as:

$$\text{HSIC}(X, Y) := \text{MMD}^2(\mathbb{P}_{XY}, \mathbb{P}_X \otimes \mathbb{P}_Y) = \|\mu_{\mathbb{P}_{XY}} - \mu_{\mathbb{P}_X \otimes \mathbb{P}_Y}\|_{\mathcal{H}_x \otimes \mathcal{H}_y}^2. \quad (\text{C.1})$$

Let us assume that both kernels are characteristic. Under this assumption, the HSIC equals zero if and only if the two random objects are independent:

$$\text{HSIC}(X, Y) = 0 \iff X \perp\!\!\!\perp Y, \quad (\text{C.2})$$

where the symbol  $\perp\!\!\!\perp$  indicates an independence relationship. As will be seen later in Appendix C.3, this property is essential for independence testing [30].

By using the integral expression of the kernel mean embedding given in equation (A.1), the HSIC can be rewritten entirely in terms of kernel-based moments:

$$\begin{aligned} \text{HSIC}(X, Y) &= \mathbb{E}[K_x(X_1, X_2) K_y(Y_1, Y_2)] \\ &\quad + \mathbb{E}[K_x(X_1, X_2)] \mathbb{E}[K_y(Y_1, Y_2)] \\ &\quad - 2 \mathbb{E}[K_x(X_1, X_2) K_y(Y_1, Y_3)]. \end{aligned} \quad (\text{C.3})$$

Here,  $Z_1 = (X_1, Y_1)$ ,  $Z_2 = (X_2, Y_2)$  and  $Z_3 = (X_3, Y_3)$  are three independent pairs, with common distribution  $\mathbb{P}_Z = \mathbb{P}_{XY}$ . This reformulation makes the HSIC amenable to inference. The next section explains how to construct an estimator of the HSIC from a given sample of observations.

**C.2 Estimation of the HSIC**

Throughout Appendix C, the data consist of an  $n$ -sample of  $Z$ . It is composed of  $n$  independent and identically distributed (i.i.d.) observations of  $Z$ :

$$D_n := \left\{ \left( Z^{(i)} \right) \right\}_{1 \leq i \leq n} = \left\{ \left( X^{(i)}, Y^{(i)} \right) \right\}_{1 \leq i \leq n} \quad \text{with} \quad Z^{(i)} \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}_Z. \quad (\text{C.4})$$

The idea is to express the HSIC in the canonical form given in equation (B.1), enabling the use of a U- or V-statistic estimator. By merging all expectations, the HSIC expression established in equation (C.3) can be compacted into a single

expectation:

$$\text{HSIC}(X, Y) = \mathbb{E}_{\mathbb{P}_Z} [\eta(Z_1, Z_2, Z_3, Z_4)] , \quad (\text{C.5})$$

where the variables  $Z_i$  are independent copies of  $Z \sim \mathbb{P}_Z$ , and  $\eta$  is the four-argument function defined by:

$$\eta(z_1, z_2, z_3, z_4) = K_x(x_1, x_2) K_y(y_1, y_2) + K_x(x_1, x_2) K_y(y_3, y_4) - 2 K_x(x_1, x_2) K_y(y_1, y_3) . \quad (\text{C.6})$$

While the kernels  $K_x$  and  $K_y$  are symmetric, the function  $\eta$  itself is not. It is therefore necessary to introduce its symmetrized version. Specializing equation (B.3) yields:

$$\tilde{\eta}(z_1, z_2, z_3, z_4) = \frac{1}{4!} \sum_{\sigma \in \mathbb{S}_4} \eta(z_{\sigma(1)}, z_{\sigma(2)}, z_{\sigma(3)}, z_{\sigma(4)}) , \quad (\text{C.7})$$

where  $\mathbb{S}_4$  denotes the set of all permutations of  $\{1, \dots, 4\}$ . Since any permutation of the vector  $(Z_1, Z_2, Z_3, Z_4)$  has the same distribution, one still has:

$$\text{HSIC}(X, Y) = \mathbb{E}_{\mathbb{P}_Z} [\tilde{\eta}(Z_1, Z_2, Z_3, Z_4)] , \quad (\text{C.8})$$

but the four-argument function is now symmetric. Then, the HSIC can be estimated by using either a U- or V-statistic:

$$\bullet \hat{H}_n^{\text{U}} := \binom{n}{4}^{-1} \sum_{1 \leq i_1 < \dots < i_4 \leq n} \tilde{\eta}(Z^{(i_1)}, \dots, Z^{(i_4)}) , \quad (\text{C.9})$$

$$\bullet \hat{H}_n^{\text{V}} := \frac{1}{n^4} \sum_{1 \leq i_1, \dots, i_4 \leq n} \tilde{\eta}(Z^{(i_1)}, \dots, Z^{(i_4)}) . \quad (\text{C.10})$$

These two estimators can then be instantiated by substituting  $\tilde{\eta}$  by its expression from equation (C.7). In both cases, the resulting quadruple sum can be decomposed into three terms, each of which can be further simplified by accounting for repeated terms through combinatorial counting. The final expression of the U-statistic is given by:

$$\begin{aligned} \hat{H}_n^{\text{U}} &= \frac{1}{(n)_2} \sum_{(i_1, i_2) \in \mathcal{I}_n^2} K_x(X^{(i_1)}, X^{(i_2)}) K_y(Y^{(i_1)}, Y^{(i_2)}) \\ &\quad + \frac{1}{(n)_4} \sum_{(i_1, i_2, i_3, i_4) \in \mathcal{I}_n^4} K_x(X^{(i_1)}, X^{(i_2)}) K_y(Y^{(i_3)}, Y^{(i_4)}) \\ &\quad - \frac{2}{(n)_3} \sum_{(i_1, i_2, i_3) \in \mathcal{I}_n^3} K_x(X^{(i_1)}, X^{(i_2)}) K_y(Y^{(i_1)}, Y^{(i_3)}) . \end{aligned} \quad (\text{C.11})$$

Regarding the notations above,  $(n)_k := n!/(n-k)!$  and the summation set  $\mathcal{I}_n^k$  is defined by:

$$\mathcal{I}_n^k = \{\mathbf{i} = (i_1, i_2, \dots, i_k) : 1 \leq i_1 \neq i_2 \neq \dots \neq i_k \leq n\} \quad \text{with} \quad \text{Card}(\mathcal{I}_n^k) = (n)_k . \quad (\text{C.12})$$

Likewise, the final expression of the V-statistic is given by:

$$\begin{aligned} \hat{H}_n^{\text{V}} &= \frac{1}{n^2} \sum_{(i_1, i_2) \in \mathcal{J}_n^2} K_x(X^{(i_1)}, X^{(i_2)}) K_y(Y^{(i_1)}, Y^{(i_2)}) \\ &\quad + \frac{1}{n^4} \sum_{(i_1, i_2, i_3, i_4) \in \mathcal{J}_n^4} K_x(X^{(i_1)}, X^{(i_2)}) K_y(Y^{(i_3)}, Y^{(i_4)}) \end{aligned} \quad (\text{C.13})$$

$$-\frac{2}{n^3} \sum_{(i_1, i_2, i_3) \in \mathcal{J}_n^3} K_x \left( X^{(i_1)}, X^{(i_2)} \right) K_y \left( Y^{(i_1)}, Y^{(i_3)} \right),$$

where the summation set  $\mathcal{J}_n^k$  is defined by:

$$\mathcal{J}_n^k = \{\mathbf{i} = (i_1, i_2, \dots, i_k) : 1 \leq i_1, i_2, \dots, i_k \leq n\} \quad \text{with} \quad \text{Card}(\mathcal{J}_n^k) = n^k. \quad (\text{C.14})$$

The V-statistic estimator admits a simple matrix representation:

$$\widehat{H}_n^V = \frac{1}{n^2} \text{Tr}(\mathbf{H}\mathbf{L}_x\mathbf{H}\mathbf{L}_y) = \frac{1}{n^2} \text{Tr}(\widetilde{\mathbf{L}}_x\widetilde{\mathbf{L}}_y) = \frac{1}{n^2} \mathbf{1}_n^t (\widetilde{\mathbf{L}}_x \odot \widetilde{\mathbf{L}}_y) \mathbf{1}_n, \quad (\text{C.15})$$

where:

- $\text{Tr}(\cdot)$  denotes the matrix trace operator.
- $\odot$  denotes the elementwise matrix product.
- $\mathbf{L}_x = \left[ K_x \left( X^{(i)}, X^{(j)} \right) \right]_{1 \leq i, j \leq n}$  and  $\mathbf{L}_y = \left[ K_y \left( Y^{(i)}, Y^{(j)} \right) \right]_{1 \leq i, j \leq n}$  are the Gram matrices.
- $\mathbf{H} = \mathbf{I}_n - (1/n)\mathbf{1}_n\mathbf{1}_n^t$  is the centering matrix, with  $\mathbf{I}_n \in \mathbb{R}^{n \times n}$  the identity matrix, and  $\mathbf{1}_n \in \mathbb{R}^n$  a vector of ones.
- $\widetilde{\mathbf{L}}_x := \mathbf{H}\mathbf{L}_x\mathbf{H}$  and  $\widetilde{\mathbf{L}}_y := \mathbf{H}\mathbf{L}_y\mathbf{H}$  are the double-centered Gram matrices.

The second equality in equation (C.15) allows rewriting the HSIC in terms of two double-centered symmetric matrices, as required by the non-asymptotic Gamma test (see Appendix D). The third equality highlights that  $\mathcal{O}(n^2)$  operations are sufficient to compute the V-statistic.

As a U-statistic,  $\widehat{H}_n^U$  is the unbiased estimator of  $\text{HSIC}(X, Y)$  with the smallest variance [75]. However, this estimator may take negative values, which is somehow unsettling (because the HSIC is a non-negative quantity). In contrast, the V-statistic  $\widehat{H}_n^V$  has complementary properties: it is biased, asymptotically unbiased and always non-negative.

### C.3 Independence testing with the HSIC

The HSIC can be used to test independence between the random objects  $X$  and  $Y$ . The idea is to make a decision between the two following hypotheses:

$$(\text{H}_0) : X \perp\!\!\!\perp Y \quad \text{vs.} \quad (\text{H}_1) : X \not\perp\!\!\!\perp Y, \quad (\text{C.16})$$

on the only basis of the observed data:

$$D_n^{\text{obs}} := \left\{ z^{(i)} \right\}_{1 \leq i \leq n} = \left\{ \left( x^{(i)}, y^{(i)} \right) \right\}_{1 \leq i \leq n}. \quad (\text{C.17})$$

The data in  $D_n^{\text{obs}}$  are written in lowercase to emphasize that  $D_n^{\text{obs}}$  is a fixed dataset, as opposed to  $D_n$  which models the underlying random mechanism generating the data.  $D_n^{\text{obs}}$  is used to perform the test, whereas  $D_n$  is only needed for theoretical considerations.

Assuming that  $K_x$  and  $K_y$  are characteristic kernels, recall the fundamental property stated in equation (C.2). The HSIC vanishes if and only if  $X$  and  $Y$  are independent objects. Consequently, the two hypotheses in equation (C.16) can be reformulated as:

$$\text{HSIC}(X, Y) = 0 \quad \text{vs.} \quad \text{HSIC}(X, Y) > 0. \quad (\text{C.18})$$

Two important points follow from this reformulation of the problem. First, a natural choice for the test statistic is to take either the U- or V-statistic estimator of the HSIC. In the following, it is denoted by  $\widehat{H}^{\text{stat}}(D_n)$ , allowing both cases to be handled simultaneously. Second, the test is one-sided (upper-tailed), meaning that  $(\text{H}_0)$  is rejected when the observed value  $\widehat{H}^{\text{stat}}(D_n^{\text{obs}})$  of the test statistic  $\widehat{H}^{\text{stat}}(D_n)$  exceeds the quantile of order  $(1 - \alpha)$  of the null distribution.

The significance level, often set to  $\alpha = 5\%$ , ensures strict control of Type-I error. Alternatively, the  $p$ -value can be used to quantify how extreme the observed value is relative to the null distribution:

$$p\text{-val} := \mathbb{P}_{H_0} \left( \widehat{H}^{\text{stat}}(D_n) \geq \widehat{H}^{\text{stat}}(D_n^{\text{obs}}) \right). \quad (\text{C.19})$$

The null hypothesis ( $H_0$ ) is then rejected when the  $p$ -value is smaller than  $\alpha$ . Generally speaking, estimating a test  $p$ -value requires knowledge of the null distribution. It may be known exactly, approximated by a parametric model, or simulated using a routine applied to the data. Depending on the sample size, three procedures are available for the HSIC-based independence test.

- **The asymptotic test.** When  $n$  is large ( $n \geq 500$ ), asymptotic theory can be used. The rescaled HSIC estimator  $n \widehat{H}^{\text{stat}}(D_n)$  converges to a degenerate limiting distribution. More precisely, this limit is a spectral distribution characterized by the eigenvalues of the centered kernels  $\widetilde{K}_x$  and  $\widetilde{K}_y$  (obtained by centering  $K_x$  and  $K_y$  with respect to  $\mathbb{P}_X$  and  $\mathbb{P}_Y$ ). Several techniques have been proposed to approximate the spectral distribution, including asymptotic moment formulas to fit a Gamma distribution [30] or kernel PCA to estimate the eigenvalues [76]. Regardless of the approach, the  $p$ -value is obtained at the cost of  $\mathcal{O}(n^2)$  operations.
- **The permutation-based test.** When  $n$  is small ( $n \leq 100$ ), the asymptotic test cannot be applied. The null distribution of the HSIC estimator must be simulated with a permutation mechanism. Random permutations are sequentially applied to the output data and the test statistic is recomputed on each permuted sample [29, 57]. For any given permutation  $\tau \in \mathbb{S}_n$ , the permuted sample  $D_n^\tau$  is derived from  $D_n^{\text{obs}}$  by permuting the outputs  $y^{(i)}$  according to  $\tau$ :

$$D_n^\tau = \left\{ \left( x^{(i)}, y^{(\tau(i))} \right) \right\}_{1 \leq i \leq n}. \quad (\text{C.20})$$

The resulting HSIC estimate  $\widehat{H}^{\text{stat}}(D_n^\tau)$  can be regarded as a realization of the test statistic under ( $H_0$ ). Repeating this procedure for  $B$  permutations yields  $B$  realizations from the null distribution. Using an empirical estimator of the probability in equation (C.19), the  $p$ -value is approximated as:

$$\widehat{p}\text{-val} = \frac{1}{B} \sum_{b=1}^B \mathbf{1}_{\{\widehat{H}^{\text{stat}}(D_n^{\tau_b}) > \widehat{H}^{\text{stat}}(D_n^{\text{obs}})\}}. \quad (\text{C.21})$$

Overall, the computational complexity of this procedure amounts to  $\mathcal{O}(B n^2)$ .

- **The non-asymptotic Gamma test.** When  $n$  is intermediate ( $100 \leq n \leq 500$ ), the permutation-based test becomes computationally expensive. To alleviate this burden, a parametric estimate of the null distribution can be constructed. A computational trick allows the first two moments of the test statistic under the permutation distribution to be derived analytically, without explicitly generating permutations [58, 77]. These moments are used to fit a Gamma distribution, from which the  $p$ -value is obtained. With an algorithmic complexity of only  $\mathcal{O}(n^2)$ , this procedure can be viewed as a low-complexity version of the permutation-based test. Appendix D.1 provides all technical details needed to implement this procedure.

#### C.4 Statistical applications of the HSIC

As highlighted in Appendix C.3, the HSIC provides a nonparametric independence test between two random objects. Importantly, it only requires a sample of  $n$  joint observations.

In the seminal work by Gretton et al. [30], the HSIC test was first used as a validation tool for *independent component analysis* (ICA). In that setting, the data consist of signals obtained from mixtures of mutually independent sources. The goal of ICA is to estimate the mixing matrix and to reconstruct the source signals. At the end of the procedure, independence between the reconstructed sources can be assessed through HSIC tests.

Subsequently, the HSIC was promoted as a powerful tool for *feature selection* in supervised learning [72]. A natural idea in this context is to prioritize explanatory variables that exhibit strong statistical dependence with the target variable. By relying on HSIC tests, one can easily detect variables that have a strongly nonlinear effect on the target. However, this strategy may lead to select variables with similar explanatory power, thereby unnecessarily complicating

the predictive model. To mitigate this redundancy issue, the HSIC-Lasso leverages an  $\ell^1$ -penalty to reduce intra-feature dependencies [78].

In most safety assessment studies, the output variable  $Y = f(\mathbf{X})$  is computed from a set of input variables  $\mathbf{X} = (X_1, \dots, X_d)$  by evaluating a deterministic simulation code  $f : \mathcal{X} \rightarrow \mathcal{Y}$ . In high-dimensional settings, where  $f$  takes as input several dozens of parameters, many numerical methods in learning and optimization suffer from the so-called *curse of dimensionality*. To remedy this problem, a preliminary step of *screening* is often performed in order to identify the inputs  $X_i$  with negligible influence on  $Y$ , so that they can be fixed at their nominal values. As explained by Da Veiga [28], this can be achieved using HSIC tests, since they allow identifying the inputs for which it is reasonable to accept the null hypothesis  $(H_0) : X_i \perp\!\!\!\perp Y$ .

**Remark C.1.** Let us assume that  $X_i$  is not a dummy input, meaning that  $x_i \mapsto f(\mathbf{x})$  is not a constant function. In most cases, this functional relationship induces a probabilistic dependence between  $X_i$  and  $Y$ . Hence, from a theoretical standpoint, the alternative hypothesis  $(H_1) : X_i \not\perp\!\!\!\perp Y$  holds, even if the influence of  $X_i$  on  $Y$  is negligible. Moreover, when characteristic kernels are used, the HSIC test is asymptotically consistent [30]. This property ensures that the probability of detecting an input variable  $X_i$  for which  $(H_1)$  is true converges to one as the sample size  $n$  increases. Therefore, the HSIC test is expected to detect all input variables. In practice, this is not what happens because  $n$  is limited (typically  $10^2 \leq n \leq 10^3$  for most industrial numerical simulators). In this regime, it is difficult to distinguish weak dependence from true independence using the available data. As a result, the HSIC test leads to reject the null hypothesis  $(H_0)$  for weakly influential input variables, making it a powerful screening tool.

### C.5 Using the HSIC for parameter screening

Parameter screening prior to Bayesian calibration is a specific instance of the general screening problem, with a few particularities. Using the notations introduced in Section 3.2, the HSIC test can be employed to select the most influential calibration parameters  $\Theta_j$  with respect to the output of interest  $G_{\text{nom}} = f(\Theta)$ . The function  $f$  represents the deterministic relationship between the calibration parameters  $\Theta = (\Theta_1, \dots, \Theta_p)$  and the output under study  $G_{\text{nom}}$ . It encapsulates the full simulation chain  $y$  and depends on the chosen screening strategy. Recall the three possible definitions of  $G_{\text{nom}}$ :

- **Approach A:**  $G_{\text{nom}} = f_i^A(\Theta) = Y_i = y(\mathbf{x}_i, \lambda_{\text{nom}}, \Theta)$  ;
- **Approach B:**  $G_{\text{nom}} = f^B(\Theta) = \mathbf{Y} = [Y_i]_{1 \leq i \leq n_{\text{exp}}}$  ;
- **Approach C:**  $G_{\text{nom}} = f^C(\Theta) = L = (\mathbf{z} - \mathbf{Y})^t \Sigma_{\epsilon}^{-1} (\mathbf{z} - \mathbf{Y})$  .

The estimation of HSIC indices (see Appendix C.2) and the execution of independence tests (see Appendix C.3) is straightforward, provided that all variables are equipped with appropriate kernels. All calibration parameters  $\Theta_j$  and the outputs  $Y_i$  and  $L$  are scalar and continuous. Consequently, these variables can be equipped with standard univariate kernels (such as Gaussian or Matérn kernels). In our numerical experiments (see Sects. 3.3 and 5), Gaussian kernels are employed, and their bandwidth parameters are chosen using the median heuristic. The only critical point is the choice of a kernel for the multivariate output  $\mathbf{Y} \in \mathbb{R}^q$  (with  $q = n_{\text{exp}}$ ). Several options can be considered.

#### C.5.1 RBF kernels

The simplest choice is to consider multivariate kernels with radial isotropic structure, commonly referred to as *radial basis function* (RBF) kernels. For example, the multivariate Gaussian kernel is defined as:

$$\forall \mathbf{y}, \mathbf{y}' \in \mathbb{R}^q, \quad K_{\mathcal{N}}(\mathbf{y}, \mathbf{y}') = \exp\left(-\frac{\|\mathbf{y} - \mathbf{y}'\|^2}{2\gamma}\right). \quad (\text{C.22})$$

According to Proposition A.10,  $K_{\mathcal{N}}$  is characteristic on  $\mathbb{R}^q$ . The FT of its signature  $\psi$ , required to establish equation (A.11), may be found in [79] (see Ex. 2.7). Note that  $K_{\mathcal{N}}$  is *radial* (as it only depends on the Euclidean distance between  $\mathbf{y}$  and  $\mathbf{y}'$ ) and *isotropic* (as it involves a unique bandwidth parameter  $\gamma$ ). It is therefore not recommended when the components of  $\mathbf{Y}$  have very different marginal distributions.

### C.5.2 Tensor-product kernels

An alternative kernel is the tensor-product kernel  $K_{\text{tens}}$  defined in equation (A.12). The idea is to associate a univariate kernel  $K_i$  with each component  $Y_i$ , and take their tensor product. In the case of Gaussian kernels, this yields:

$$\forall \mathbf{y}, \mathbf{y}' \in \mathbb{R}^q, \quad K_{\text{tens}}(\mathbf{y}, \mathbf{y}') = \prod_{i=1}^q K_i(y_i, y'_i) = \prod_{i=1}^q \exp \left[ \frac{1}{2} \left( \frac{y_i - y'_i}{\gamma_i} \right)^2 \right]. \quad (\text{C.23})$$

If all univariate kernels  $K_i$  are translation-invariant and characteristic on  $\mathbb{R}$ , then  $K_{\text{tens}}$  is also translation-invariant and characteristic on  $\mathbb{R}^q$  (see Prop. A.12).  $K_{\text{tens}}$  is considerably more flexible than  $K_{\mathcal{N}}$ , as each bandwidth parameter  $\gamma_i$  can be adapted to the scale of  $Y_i$ . Empirically, however,  $K_{\text{tens}}$  tends to disadvantage parameters whose influence is confined to a small number of output components.

### C.5.3 Weighted PCA kernel

To address this limitation, we consider the weighted PCA kernel proposed by El Amri and Marrel [58]. For simplicity, the random vector  $\mathbf{Y}$  is assumed to be centered and scaled. Let  $\Sigma$  denote its covariance matrix:

$$\Sigma = [\sigma_{kl}]_{1 \leq k, l \leq q} \quad \text{with} \quad \sigma_{kl} = \text{Cov}(Y_k, Y_l). \quad (\text{C.24})$$

In practice, this matrix is estimated empirically from centered and scaled observations of  $\mathbf{Y}$ . Let  $\{\lambda_i\}_{i=1}^q$  and  $\{\mathbf{v}_i\}_{i=1}^q$  denote the eigenvalues and eigenvectors of  $\Sigma$ . Using first the canonical basis  $(\mathbf{e}_1, \dots, \mathbf{e}_q)$  and then the eigenbasis  $(\mathbf{v}_1, \dots, \mathbf{v}_q)$ , the vector  $\mathbf{Y}$  can be decomposed as follows:

$$\mathbf{Y} = \sum_{i=1}^q Y_i \mathbf{e}_i = \sum_{i=1}^q C_i \mathbf{v}_i \quad \text{with} \quad C_i = \langle \mathbf{Y}, \mathbf{v}_i \rangle_{\mathbb{R}^q} = \mathbf{v}_i^t \mathbf{Y}. \quad (\text{C.25})$$

The total variance of  $\mathbf{Y}$  is defined as the sum of the variances of its individual components:

$$\text{Tr}(\Sigma) = \sum_{i=1}^q \mathbb{V}(Y_i) = \sum_{i=1}^q \lambda_i. \quad (\text{C.26})$$

We retain the smallest number  $r$  of principal components that ensures a loss of information smaller than  $\epsilon = 5\%$  of the total variability:

$$r := \min \left\{ 1 \leq k \leq q : \sum_{i=1}^k \lambda_i / \sum_{i=1}^q \lambda_i \geq 1 - \epsilon \right\}. \quad (\text{C.27})$$

The projection  $\mathcal{P} : \mathbb{R}^q \rightarrow \mathbb{R}^r$  from the physical space to the latent space of principal components is given by:

$$\mathcal{P}(\mathbf{y}) := [c_i]_{1 \leq i \leq r} = \mathbf{V}_r^t \mathbf{y} \quad \text{with} \quad \mathbf{V}_r = [ \mathbf{v}_1 \mid \mathbf{v}_2 \mid \dots \mid \mathbf{v}_r ] \in \mathbb{R}^{q \times r} \quad (\text{C.28})$$

The weighted PCA kernel is then defined as:

$$\forall \mathbf{y}, \mathbf{y}' \in \mathbb{R}^q, \quad K_{\text{PCA}}^r(\mathbf{y}, \mathbf{y}') = \sum_{i=1}^r \lambda_i K_i(c_i, c'_i), \quad (\text{C.29})$$

where the functions  $K_i$  are univariate kernels (for instance Gaussian kernels) assigned to the first  $r$  principal components. Whenever  $r < q$ ,  $K_{\text{PCA}}^r$  ceases to be characteristic on  $\mathbb{R}^q$ . Indeed, one can find two distributions ( $\mathbb{P}_1 \neq \mathbb{P}_2$ ) that share identical marginal distributions for the first  $r$  principal components but differ on the remaining components. As the resulting kernel mean embeddings are identical ( $\mu_{\mathbb{P}_1} = \mu_{\mathbb{P}_2}$ ),  $K_{\text{PCA}}^r$  cannot be characteristic. Nevertheless, empirical

evidence shows that this kernel provides increased detection power for inputs whose influence is localized on a small number of output components.

#### APPENDIX D. NON-ASYMPTOTIC GAMMA TEST PROCEDURE

##### D.1 Standard procedure

Let us consider the classical setting, in which the HSIC is used to make a decision between the two following hypotheses:

$$(H_0) : X \perp\!\!\!\perp Y \quad \text{vs.} \quad (H_1) : X \not\perp\!\!\!\perp Y . \quad (\text{D.1})$$

An overview of existing test procedures is provided in Appendix C.3. In particular, it is reported that the non-asymptotic Gamma test is well suited to intermediate sample sizes ( $100 \leq n \leq 500$ ). In this regime, both the asymptotic test and the permutation-based test suffer from limitations. On the one hand, the sample size  $n$  is not sufficiently large for the asymptotic test to be used reliably. On the other hand, the computational complexity of the permutation-based test, namely  $\mathcal{O}(Bn^2)$ , can already be substantial, especially if the number of permutations  $B$  is set very high. In this context, the non-asymptotic Gamma test lies midway, thanks to a Gamma approximation of the null distribution obtained without relying on any asymptotic result. It combines two advantages: it can be used in the non-asymptotic regime (contrary to the asymptotic test) and it has a low computational complexity (compared to that of the permutation-based test).

This test procedure applies only to the V-statistic estimator of the HSIC because its key element is the associated matrix representation:

$$\widehat{H}^V(D_n^{\text{obs}}) = \frac{1}{n^2} \text{Tr}(\mathbf{A}\mathbf{W}) \quad \text{with} \quad \begin{cases} \mathbf{A} = \widetilde{\mathbf{L}}_x = \mathbf{H}\mathbf{L}_x\mathbf{H} \\ \mathbf{W} = \widetilde{\mathbf{L}}_y = \mathbf{H}\mathbf{L}_y\mathbf{H} \end{cases} . \quad (\text{D.2})$$

All details regarding the construction of this estimator are provided in Appendix C.2. The available data  $D_n^{\text{obs}}$  is a Monte Carlo sample of  $n$  joint observations, as in equation (C.17). To simulate  $(H_0)$ , the permutation-based test randomly permutes the output data and recomputes the test statistic for each permuted sample. Using a sequence of  $B$  permutations, denoted by  $\{\tau_b\}_{b=1}^B$ , the simulated values of the V-statistic are:

$$\left\{ \widehat{H}^V(D_n^{\tau_b}) \right\}_{1 \leq b \leq B} \quad \text{with} \quad \forall 1 \leq b \leq B, \quad D_n^{\tau_b} = \left\{ \left( x^{(i)}, y^{(\tau_b(i))} \right) \right\}_{1 \leq i \leq n} . \quad (\text{D.3})$$

The starting point is to notice that the histogram of these values appears to be consistent with a Gamma distribution. This is unsurprising, as  $n$  is already moderately large, and it is well established that the spectral distribution governing the asymptotic behavior of the V-statistic can be closely approximated by a Gamma distribution [30, 80].

A Gamma distribution could be fitted to the sample of HSIC values provided in equation (D.3), but this would not help reduce the computational complexity of the permutation-based test. Conveniently, as explained in [58], a Gamma approximation can be constructed directly, without relying on the values simulated by the permutation-based test.

Before continuing, we clarify some elements of the probabilistic modeling. Let  $\tau$  be a random variable uniformly distributed over the set  $\mathbb{S}_n$  of all possible permutations. In the following, the probability measure  $\mathbb{P}_\tau$  will be referred to as the *permutation distribution*. Since the sample  $D_n^{\text{obs}}$  is fixed, the random variable  $\widehat{H}^V(D_n^\tau)$  is the image of  $\tau$  under the operation of permuting the output data and computing the V-statistic.

When a permutation  $\tau \in \mathbb{S}_n$  is applied to the output data, only the matrix  $\mathbf{W}$  is modified: its rows and columns are permuted according to  $\tau$ . The new value of the V-statistic admits the following matrix representation:

$$\mathcal{S}(\tau) := n^2 \widehat{H}^V(D_n^\tau) = \text{Tr}(\mathbf{A}\mathbf{W}^{\tau\tau}) \quad \text{with} \quad \mathbf{W}^{\tau\tau} := [W_{\tau(i)\tau(j)}]_{1 \leq i, j \leq n} . \quad (\text{D.4})$$

The first two moments of  $\mathcal{S}(\tau)$  under  $\mathbb{P}_\tau$  can be computed exactly. This result can be stated in a general setting. For any square matrix  $\mathbf{M} \in \mathbb{R}^{n \times n}$ , we first introduce the following three functionals:

$$D_1^{\mathbf{M}} := \text{Tr}(\mathbf{M}) = \sum_{i=1}^n M_{ii} \quad (\text{D.5})$$

$$D_2^M := \text{Tr}(\mathbf{M}^{\odot 2}) = \sum_{i=1}^n M_{ii}^2 \quad (\text{D.6})$$

$$S_2^M := \text{Tr}(\mathbf{M}^2) = \text{Sum}(\mathbf{M}^{\odot 2}) = \sum_{i=1}^n \sum_{j=1}^n M_{ij}^2. \quad (\text{D.7})$$

For a given matrix  $\mathbf{M}$ , computing these three quantities requires  $n$ ,  $2n$  and  $2n^2$  operations, respectively. To simplify the formulas to come, we also introduce some auxiliary notations:

$$f_1(\mathbf{M}) := (n-1)S_2^M - (D_1^M)^2 \quad (\text{D.8})$$

$$f_2(\mathbf{M}) := n(n+1)D_2^M - (n-1)\left[(D_1^M)^2 + 2S_2^M\right] \quad (\text{D.9})$$

$$p_n := (n-1)^2(n+1)(n-2) \quad (\text{D.10})$$

$$q_n := (n+1)n(n-1)(n-2)(n-3). \quad (\text{D.11})$$

The main result, stated in the following lemma, can now be written in a compact form.

**Lemma D.1.** *Assume that  $\mathbf{A}$  and  $\mathbf{W}$  are two double-centered symmetric matrices of size  $n$ . Let  $\mathcal{S}(\tau)$  be the random variable defined in equation (D.4). Its first two moments under the permutation distribution  $\mathbb{P}_\tau$  are given by:*

$$\mu_{\text{perm}} := \mathbb{E}_{\mathbb{P}_\tau} [\text{Tr}(\mathbf{A}\mathbf{W}^{\tau\tau})] = \frac{1}{n!} \sum_{\tau \in \mathbb{S}_n} \text{Tr}(\mathbf{A}\mathbf{W}^{\tau\tau}) \quad (\text{D.12})$$

$$= \frac{1}{n-1} D_1^{\mathbf{A}} D_1^{\mathbf{W}} =: E_0(\mathbf{A}, \mathbf{W}) \quad (\text{D.13})$$

$$\sigma_{\text{perm}}^2 := \mathbb{V}_{\mathbb{P}_\tau} (\text{Tr}(\mathbf{A}\mathbf{W}^{\tau\tau})) = \mathbb{E}_{\mathbb{P}_\tau} [\text{Tr}(\mathbf{A}\mathbf{W}^{\tau\tau})^2] - (\mu_{\text{perm}})^2 \quad (\text{D.14})$$

$$= \frac{2}{p_n} f_1(\mathbf{A}) f_1(\mathbf{W}) + \frac{1}{q_n} f_2(\mathbf{A}) f_2(\mathbf{W}) =: V_0(\mathbf{A}, \mathbf{W}). \quad (\text{D.15})$$

The calculations are detailed in [77] (see Sect. 3). In particular, the formulas in Lemma D.1 show that the exact values of  $\mu_{\text{perm}}$  and  $\sigma_{\text{perm}}^2$  can be obtained at the cost of  $\mathcal{O}(n^2)$  operations.

Back to the HSIC framework, applying Lemma D.1 to equation (D.2) yields:

$$\mu_H := \mathbb{E}_\tau [\widehat{H}^V(D_n^\tau)] = \frac{1}{n^2} E_0(\mathbf{A}, \mathbf{W}) \quad (\text{D.16})$$

$$\sigma_H^2 := \mathbb{V}_\tau (\widehat{H}^V(D_n^\tau)) = \frac{1}{n^4} V_0(\mathbf{A}, \mathbf{W}). \quad (\text{D.17})$$

Then, we resort to the method of moments to construct a Gamma approximation of the V-statistic  $\widehat{H}^V(D_n^\tau)$  under the permutation distribution  $\mathbb{P}_\tau$ . Recall that, for a random variable  $S \sim \text{Ga}(a, b)$  following a Gamma distribution with shape parameter  $a$  and scale parameter  $b$ , the first two moments satisfy:

$$\begin{cases} \mu_S := \mathbb{E}[S] = ab \\ \sigma_S^2 := \mathbb{V}(S) = ab^2 \end{cases} \quad \text{and thus} \quad \begin{cases} a = \mu_S^2 / \sigma_S^2 \\ b = \sigma_S^2 / \mu_S \end{cases}. \quad (\text{D.18})$$

Consequently, the parameters of the Gamma distribution used to approximate the null distribution can be estimated by:

$$\widehat{a}_H = \mu_H^2 / \sigma_H^2 \quad \text{and} \quad \widehat{b}_H = \sigma_H^2 / \mu_H. \quad (\text{D.19})$$

Finally, the test  $p$ -value can be computed directly from the survival function of the estimated Gamma distribution.

Throughout the test procedure, it is never necessary to permute the output data and recompute the V-statistic. In terms of algorithm complexity, the construction of the Gamma approximation entails  $\mathcal{O}(n^2)$  operations, which is much cheaper than the permutation-based test. However, as with the asymptotic test, using a Gamma approximation remains heuristic [80]. Even though this parametric choice was validated empirically in the numerical study proposed in [58], there is no theoretical guarantee that it strictly controls Type-I and Type-II errors.

## D.2 New procedure in the bi-level uncertainty framework

The permutation-based test procedure developed in Section 4.2 allows one to decide between the two following hypotheses:

$$(H_0) : \Theta_j \perp\!\!\!\perp G \quad \text{vs.} \quad (H_1) : \Theta_j \not\perp\!\!\!\perp G \quad \text{with} \quad G = \tilde{f}(\Theta, \Lambda). \quad (\text{D.20})$$

In this context, any of the four estimators of  $\mathcal{H}_j$  proposed in equation (4.18) can be used as test statistic. Interestingly, the permutation-based test applies to all estimators. However, only Estimators 1 and 2 are considered in Section 4.2, since they lead to faster MSE convergence rates, as reported in Table 2.

The computational complexity of the permutation-based test is  $\mathcal{O}(B m n^2)$ . Once again, this raises the classical issue of the high computational complexity associated with permutation-based test procedures. To reduce this complexity, our idea is to leverage the computational trick of the non-asymptotic Gamma test. Since this procedure is specific to the V-statistic estimator of the HSIC, we will focus exclusively on Estimator 2:

$$\widehat{\mathcal{H}}_j^2(D_{m,n}^{\text{obs}}) = \frac{1}{m} \sum_{k=1}^m \widehat{H}_j^{\text{V}}(D_k^{\text{obs}}). \quad (\text{D.21})$$

As detailed in Section 3.4,  $D_{m,n}^{\text{obs}}$  is a nested Monte Carlo design generated according to (S1), and each design  $D_k^{\text{obs}}$  is an  $n$ -sample of the pair  $(\Theta^{(k)}, G^{(k)})$ . Replacing the V-statistics by their matrix representations leads to:

$$\widehat{\mathcal{H}}_j^2(D_{m,n}^{\text{obs}}) = \frac{1}{m n^2} \sum_{k=1}^m \text{Tr}(\mathbf{A}_{jk} \mathbf{W}_k) \quad \text{with} \quad \begin{cases} \mathbf{A}_{jk} := \tilde{\mathbf{L}}_{\theta_j}^{(k)} = \mathbf{H} \mathbf{L}_{\theta_j}^{(k)} \mathbf{H} \\ \mathbf{W}_k := \tilde{\mathbf{L}}_g^{(k)} = \mathbf{H} \mathbf{L}_g^{(k)} \mathbf{H} \end{cases}, \quad (\text{D.22})$$

where  $\mathbf{L}_{\theta_j}^{(k)}$  and  $\mathbf{L}_g^{(k)}$  denote the Gram matrices built from the observations of  $\Theta_j^{(k)}$  and  $G^{(k)}$  found in the sample  $D_k^{\text{obs}}$ . To simulate the null distribution, a permutation  $\tau_k$  is applied to the output data of each design  $D_k^{\text{obs}}$ . A realization of the null distribution is thus obtained by randomly selecting a set of  $m$  permutations  $\tau = (\tau_1, \dots, \tau_m)$ , applying them to the samples  $D_k^{\text{obs}}$ , and recomputing the test statistic:

$$\mathcal{S}_j(\tau) := \widehat{\mathcal{H}}_j^2(D_{m,n}^{\tau}) = \frac{1}{m} \sum_{k=1}^m \widehat{H}_j^{\text{V}}(D_k^{\tau_k}) = \frac{1}{m n^2} \sum_{k=1}^m \text{Tr}(\mathbf{A}_{jk} \mathbf{W}_k^{\tau_k}). \quad (\text{D.23})$$

If the process is repeated  $B$  times, a  $B$ -sample of the null distribution is produced. When plotting the associated histogram (see Fig. E.1), it appears fully consistent with a Gamma distribution. This suggests that a Gamma approximation could be constructed.

As for the non-asymptotic Gamma test, the idea is to calculate the exact values of the first two moments, from which a Gamma distribution is fitted. The only difference is that the permutation scheme here involves a set of permutations  $\tau \in (\mathbb{S}_n)^m$ , instead of a single permutation  $\tau \in \mathbb{S}_n$ . Fortunately, Lemma D.1 allows for a straightforward derivation of the first two moments of  $\mathcal{S}_j(\tau)$  under  $\mathbb{P}_{\tau}$ :

$$\mu_j = \mathbb{E}_{\tau}[\mathcal{S}_j(\tau)] = \frac{1}{m n^2} \sum_{k=1}^m \mathbb{E}_{\tau_k}[\text{Tr}(\mathbf{A}_{jk} \mathbf{W}_k^{\tau_k})] = \frac{1}{m n^2} \sum_{k=1}^m E_0(\mathbf{A}_{jk}, \mathbf{W}_k) \quad (\text{D.24})$$

$$\sigma_j^2 = \mathbb{V}_\tau(\mathcal{S}_j(\boldsymbol{\tau})) = \frac{1}{m^2 n^4} \sum_{k=1}^m \mathbb{V}_{\tau_k}(\text{Tr}(\mathbf{A}_{jk} \mathbf{W}_k^{\tau_k})) = \frac{1}{m^2 n^4} \sum_{k=1}^m V_0(\mathbf{A}_{jk}, \mathbf{W}_k). \quad (\text{D.25})$$

The variance formula holds because of the mutual independence of the random permutations in the set  $\boldsymbol{\tau}$ . Once  $\mu_j$  and  $\sigma_j^2$  are computed, the parameters  $(\hat{a}_j, \hat{b}_j)$  of the Gamma approximation and the final  $p$ -value follow directly.

Figure E.1 reveals that the Gamma approximation closely matches the histogram of the null distribution obtained from the permutation-based test, which provides empirical support for this test procedure.

With an overall computational complexity of  $\mathcal{O}(m n^2)$ , instead of  $\mathcal{O}(B m n^2)$ , this non-asymptotic Gamma test stands as a low-complexity alternative to the permutation-based test.

#### APPENDIX E. TECHNICAL PROOFS

We begin with a simple lemma that will be used in a proof of this section.

**Lemma E.1.** *Let  $(Z_1, \dots, Z_m)$  be an i.i.d. sample from the distribution  $\mathbb{P}_Z$ . The second moment of the empirical mean  $\bar{Z}_m := \frac{1}{m} \sum_{k=1}^m Z_k$  is given by:*

$$\mathbb{E}[\bar{Z}_m^2] = \frac{1}{m} \mathbb{E}[Z_1^2] + \frac{m-1}{m} (\mathbb{E}[Z_1])^2 \leq \mathbb{E}[Z_1^2] \quad (\text{E.1})$$

*Proof.* Let  $\mu_1$ ,  $\mu_2$  and  $\sigma^2$  be the first moment, the second moment, and the variance of  $\mathbb{P}_Z$ , respectively. Then, one has  $\sigma^2 = \mu_2 - \mu_1^2$  (by definition) and  $\mu_1^2 \leq \mu_2$  (by convexity). For the empirical mean  $\bar{Z}_m$ , it is well-known that:

$$\mathbb{E}[\bar{Z}_m] = \mu_1 \quad \text{and} \quad \mathbb{V}(\bar{Z}_m) = \frac{\sigma^2}{m}. \quad (\text{E.2})$$

Therefore, the second moment satisfies:

$$\mathbb{E}[\bar{Z}_m^2] = \mathbb{V}(\bar{Z}_m) + \mathbb{E}[\bar{Z}_m]^2 = \frac{\sigma^2}{m} + \mu_1^2 = \frac{\mu_2 - \mu_1^2}{m} + \mu_1^2 = \frac{1}{m} \mu_2 + \frac{m-1}{m} \mu_1^2 \leq \mu_2. \quad (\text{E.3})$$

□

#### E.1 Proof of Proposition 4.1

*Proof.* Since the kernels  $K_\theta_j$  and  $K_g$  are characteristic, the fundamental property of the HSIC stated in equation (C.2) holds. In the present setting, it reads:

$$\forall \boldsymbol{\lambda} \in \mathcal{D}_\Lambda, \quad H_j(\boldsymbol{\lambda}) = 0 \iff \text{HSIC}(\Theta_j, G_\lambda) = 0 \iff \Theta_j \perp\!\!\!\perp G_\lambda. \quad (\text{E.4})$$

As a consequence, the following equivalences hold:

$$\mathcal{H}_j = 0 \iff \mathbb{E}_\Lambda[H_j(\boldsymbol{\Lambda})] = 0 \quad (\text{E.5})$$

$$\iff H_j(\boldsymbol{\Lambda}) = 0 \quad \mathbb{P}_\Lambda\text{-a.s.} \quad (\text{E.6})$$

$$\iff \exists \tilde{\mathcal{D}}_\Lambda \subseteq \mathcal{D}_\Lambda \text{ such that } \mathbb{P}_\Lambda(\tilde{\mathcal{D}}_\Lambda) = 1 \text{ and } \forall \boldsymbol{\lambda} \in \tilde{\mathcal{D}}_\Lambda, H_j(\boldsymbol{\lambda}) = 0 \quad (\text{E.7})$$

$$\iff \exists \tilde{\mathcal{D}}_\Lambda \subseteq \mathcal{D}_\Lambda \text{ such that } \mathbb{P}_\Lambda(\tilde{\mathcal{D}}_\Lambda) = 1 \text{ and } \forall \boldsymbol{\lambda} \in \tilde{\mathcal{D}}_\Lambda, \Theta_j \perp\!\!\!\perp G_\lambda. \quad (\text{E.8})$$

Equation (E.6) is immediate since  $H_j(\boldsymbol{\lambda}) = \text{HSIC}(\Theta_j, G_\lambda) \geq 0$  for all  $\boldsymbol{\lambda} \in \mathcal{D}_\Lambda$ . Equation (E.7) is obtained by expliciting the almost sure equality with respect to  $\mathbb{P}_\Lambda$ . Finally, using equation (E.4) for all  $\boldsymbol{\lambda} \in \tilde{\mathcal{D}}_\Lambda \subseteq \mathcal{D}_\Lambda$  yields equation (E.8). This last equation indicates that  $\Theta_j$  and  $G_\lambda = \tilde{f}(\boldsymbol{\Theta}, \boldsymbol{\lambda})$  are independent for almost all values  $\boldsymbol{\lambda} \in \mathcal{D}_\Lambda$ . This is in fact sufficient

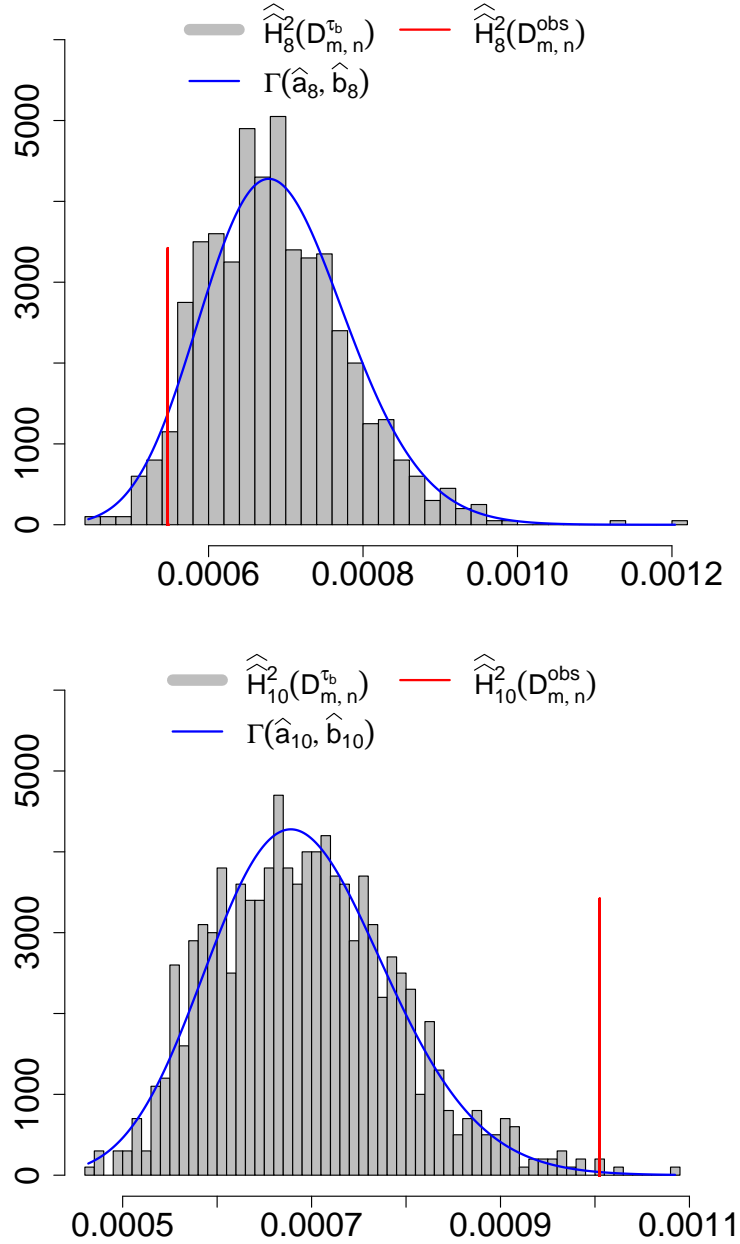


FIGURE E.1. Approximation of the null distribution of Estimator 2 by the permutation-based test and the non-asymptotic Gamma test. The red vertical line indicates the observed value of the test statistic  $\widehat{\mathcal{H}}_j^2(D_{m,n}^{\text{obs}})$ . The histogram shows the values  $\left\{ \widehat{\mathcal{H}}_j^2(D_{m,n}^{\tau_b}) \right\}_{1 \leq b \leq B}$  simulated by the permutation-based test (with  $B = 10^3$ ). The blue curve represents the parametric estimate constructed by the non-asymptotic Gamma test. Results of Approach C are shown for  $\Theta_8$  (top) and  $\Theta_{10}$  (bottom).

to conclude that  $\Theta_j \perp\!\!\!\perp G = \tilde{f}(\Theta, \Lambda)$ . This can be easily verified using the factorization property of expectations. For any pair  $(\varphi, \psi)$  of bounded Borelian functions, one must check that:

$$\mathbb{E}[\varphi(\Theta_j) \psi(G)] = \mathbb{E}[\varphi(\Theta_j)] \mathbb{E}[\psi(G)]. \quad (\text{E.9})$$

Let us successively rewrite the mixed moment:

$$\mathbb{E}[\varphi(\Theta_j) \psi(G)] = \mathbb{E}\left[\varphi(\Theta_j) \psi\left(\tilde{f}(\Theta, \Lambda)\right)\right] \quad (\text{E.10})$$

$$= \mathbb{E}\left[\varphi(\Theta_j) \psi\left(\tilde{f}(\Theta, \Lambda)\right) \mathbf{1}_{\tilde{\mathcal{D}}_\Lambda}(\Lambda)\right] \quad (\text{E.11})$$

$$= \mathbb{E}_\Lambda\left[\mathbb{E}\left[\varphi(\Theta_j) \psi\left(\tilde{f}(\Theta, \Lambda)\right) \mathbf{1}_{\tilde{\mathcal{D}}_\Lambda}(\Lambda)\right] \middle| \Lambda\right] \quad (\text{E.12})$$

$$= \mathbb{E}_\Lambda[g(\Lambda)]. \quad (\text{E.13})$$

Equation (E.11) holds because  $\mathbb{P}_\Lambda(\tilde{\mathcal{D}}_\Lambda) = 1$ . Equation (E.12) follows directly from the tower property of conditional expectation. Furthermore, since  $\Theta \perp\!\!\!\perp \Lambda$ , the function  $g : \mathcal{D}_\Lambda \rightarrow \mathbb{R}$  introduced in equation (E.13) can be expressed as:

$$g(\lambda) = \mathbb{E}\left[\varphi(\Theta_j) \psi\left(\tilde{f}(\Theta, \Lambda)\right) \mathbf{1}_{\tilde{\mathcal{D}}_\Lambda}(\Lambda)\right] \middle| \Lambda = \lambda \quad (\text{E.14})$$

$$= \mathbb{E}\left[\varphi(\Theta_j) \psi\left(\tilde{f}(\Theta, \lambda)\right) \mathbf{1}_{\tilde{\mathcal{D}}_\Lambda}(\lambda)\right] = \mathbb{E}\left[\varphi(\Theta_j) \psi(G_\lambda) \mathbf{1}_{\tilde{\mathcal{D}}_\Lambda}(\lambda)\right] \quad (\text{E.15})$$

To factorize the last expectation, two cases must be distinguished. On the one hand, if  $\lambda \in \tilde{\mathcal{D}}_\Lambda$ , one has:

$$\begin{aligned} \mathbb{E}\left[\varphi(\Theta_j) \psi(G_\lambda) \mathbf{1}_{\tilde{\mathcal{D}}_\Lambda}(\lambda)\right] &= \mathbb{E}[\varphi(\Theta_j) \psi(G_\lambda)] && \text{since } \mathbf{1}_{\tilde{\mathcal{D}}_\Lambda}(\lambda) = 1 ; \\ &= \mathbb{E}[\varphi(\Theta_j)] \mathbb{E}[\psi(G_\lambda)] && \text{since } \Theta_j \perp\!\!\!\perp G_\lambda \\ &= \mathbb{E}[\varphi(\Theta_j)] \mathbb{E}\left[\psi(G_\lambda) \mathbf{1}_{\tilde{\mathcal{D}}_\Lambda}(\lambda)\right] && \text{since } \mathbf{1}_{\tilde{\mathcal{D}}_\Lambda}(\lambda) = 1 . \end{aligned} \quad (\text{E.16})$$

On the other hand, if  $\lambda \notin \tilde{\mathcal{D}}_\Lambda$ , one has:

$$\underbrace{\mathbb{E}\left[\varphi(\Theta_j) \psi(G_\lambda) \mathbf{1}_{\tilde{\mathcal{D}}_\Lambda}(\lambda)\right]}_{=0} = 0 = \mathbb{E}[\varphi(\Theta_j)] \underbrace{\mathbb{E}\left[\psi(G_\lambda) \mathbf{1}_{\tilde{\mathcal{D}}_\Lambda}(\lambda)\right]}_{=0}. \quad (\text{E.17})$$

This shows that the factorization holds for all  $\lambda \in \mathcal{D}_\Lambda$ . Hence, the function  $g$  can now be expressed as:

$$g(\lambda) = \mathbb{E}[\varphi(\Theta_j)] \mathbb{E}\left[\psi(G_\lambda) \mathbf{1}_{\tilde{\mathcal{D}}_\Lambda}(\lambda)\right] \quad (\text{E.18})$$

$$= \mathbb{E}[\varphi(\Theta_j)] \mathbb{E}\left[\psi\left(\tilde{f}(\Theta, \Lambda)\right) \mathbf{1}_{\tilde{\mathcal{D}}_\Lambda}(\Lambda)\right] \middle| \Lambda = \lambda. \quad (\text{E.19})$$

Injecting this last expression in equation (E.13) leads to:

$$\mathbb{E}[\varphi(\Theta_j) \psi(G)] = \mathbb{E}_\Lambda\left[\mathbb{E}[\varphi(\Theta_j)] \mathbb{E}\left[\psi\left(\tilde{f}(\Theta, \Lambda)\right) \mathbf{1}_{\tilde{\mathcal{D}}_\Lambda}(\Lambda)\right] \middle| \Lambda\right] \quad (\text{E.20})$$

$$= \mathbb{E}[\varphi(\Theta_j)] \mathbb{E}_\Lambda\left[\mathbb{E}\left[\psi\left(\tilde{f}(\Theta, \Lambda)\right) \mathbf{1}_{\tilde{\mathcal{D}}_\Lambda}(\Lambda)\right] \middle| \Lambda\right] \quad (\text{E.21})$$

$$= \mathbb{E}[\varphi(\Theta_j)] \mathbb{E}\left[\psi\left(\tilde{f}(\Theta, \Lambda)\right) \mathbf{1}_{\tilde{\mathcal{D}}_\Lambda}(\Lambda)\right] \quad (\text{E.22})$$

$$= \mathbb{E}[\varphi(\Theta_j)] \mathbb{E}[\psi(G)]. \quad (\text{E.23})$$

To move from equation (E.20) to equation (E.21), the constant  $\mathbb{E}[\varphi(\Theta_j)]$  is taken outside the expectation operator. The remaining steps follow the same reasoning as above, namely the tower property and the fact that  $\mathbb{P}_\Lambda(\tilde{\mathcal{D}}_\Lambda) = 1$ . Since equation (E.9) holds for all pairs  $(\varphi, \psi)$  of bounded Borelian functions, this proves that  $\Theta_j \perp\!\!\!\perp G$ .

Now that  $\mathcal{H}_j = 0 \Rightarrow \Theta_j \perp\!\!\!\perp G$  has been proven, we turn to the other implication. It is assumed that  $\Theta_j \perp\!\!\!\perp G$ . Recall how  $\mathcal{H}_j$  is expressed in terms of kernel-based moments:

$$\begin{aligned} \mathcal{H}_j &= \mathbb{E}_\Lambda [H_j(\Lambda)] = \mathbb{E} \left[ K_{\theta_j}(\Theta_j, \Theta'_j) K_g \left( \tilde{f}(\Theta, \Lambda), \tilde{f}(\Theta', \Lambda) \right) \right] \\ &\quad + \mathbb{E} \left[ K_{\theta_j}(\Theta_j, \Theta'_j) \right] \mathbb{E} \left[ K_g \left( \tilde{f}(\Theta, \Lambda), \tilde{f}(\Theta', \Lambda) \right) \right] \\ &\quad - 2 \mathbb{E} \left[ K_{\theta_j}(\Theta_j, \Theta'_j) K_g \left( \tilde{f}(\Theta, \Lambda), \tilde{f}(\Theta'', \Lambda) \right) \right], \end{aligned} \quad (\text{E.24})$$

where  $\Theta \perp\!\!\!\perp \Theta' \perp\!\!\!\perp \Theta'' \sim \mathbb{P}_\Theta$ . The expectations appearing in the first and third lines can be factorized and coincide exactly with the second line:

$$\begin{aligned} &\mathbb{E} \left[ K_{\theta_j}(\Theta_j, \Theta'_j) K_g \left( \tilde{f}(\Theta, \Lambda), \tilde{f}(\Theta', \Lambda) \right) \right] \\ &= \mathbb{E} \left[ K_{\theta_j}(\Theta_j, \Theta'_j) K_g \left( \tilde{f}(\Theta, \Lambda), \tilde{f}(\Theta'', \Lambda) \right) \right] \\ &= \mathbb{E} \left[ K_{\theta_j}(\Theta_j, \Theta'_j) \right] \mathbb{E} \left[ K_g \left( \tilde{f}(\Theta, \Lambda), \tilde{f}(\Theta'', \Lambda) \right) \right]. \end{aligned} \quad (\text{E.25})$$

Let us justify how to factorize the first expectation. On the one hand, by hypothesis, one has:

$$\Theta_j \perp\!\!\!\perp \tilde{f}(\Theta, \Lambda) \quad \text{and} \quad \Theta'_j \perp\!\!\!\perp \tilde{f}(\Theta', \Lambda). \quad (\text{E.26})$$

On the other hand, Remark 3.4 provides:

$$\Theta \perp\!\!\!\perp \Theta' \perp\!\!\!\perp \Theta'' \perp\!\!\!\perp \Lambda, \quad (\text{E.27})$$

which in turn gives:

$$\Theta_j \perp\!\!\!\perp \tilde{f}(\Theta', \Lambda) \quad \text{and} \quad \Theta'_j \perp\!\!\!\perp \tilde{f}(\Theta, \Lambda) \quad (\text{E.28})$$

Equations (E.26) and (E.28) can be summarized as:

$$\begin{bmatrix} \Theta_j \\ \Theta'_j \end{bmatrix} \perp\!\!\!\perp \begin{bmatrix} \tilde{f}(\Theta, \Lambda) \\ \tilde{f}(\Theta', \Lambda) \end{bmatrix}, \quad (\text{E.29})$$

which eventually yields:

$$K_{\theta_j}(\Theta_j, \Theta'_j) \perp\!\!\!\perp K_g \left( \tilde{f}(\Theta, \Lambda), \tilde{f}(\Theta', \Lambda) \right). \quad (\text{E.30})$$

This justifies the factorization for the first expectation in equation (E.25). The same reasoning holds for the other expectation. Injecting the two factorized expressions in equation (E.24) leads to  $\mathcal{H}_j = 0$ . For this implication, note that characteristic kernels are not required.  $\square$

## E.2 Proof of Proposition 4.5

### E.2.0 Preparatory results

For a fixed value of  $\lambda \in \mathcal{D}_\Lambda$ , the HSIC index between  $\Theta_j$  and  $G_\lambda = \tilde{f}(\Theta_j, \lambda)$  is given by:

$$\begin{aligned} H_j(\lambda) &= \mathbb{E} \left[ K_{\theta_j}(\Theta_j^1, \Theta_j^2) K_g \left( \tilde{f}(\Theta^1, \lambda), \tilde{f}(\Theta^2, \lambda) \right) \right] \\ &\quad + \mathbb{E} \left[ K_{\theta_j}(\Theta_j^1, \Theta_j^2) \right] \mathbb{E} \left[ K_g \left( \tilde{f}(\Theta^3, \lambda), \tilde{f}(\Theta^4, \lambda) \right) \right] \\ &\quad - 2 \mathbb{E} \left[ K_{\theta_j}(\Theta_j^1, \Theta_j^2) K_g \left( \tilde{f}(\Theta^1, \lambda), \tilde{f}(\Theta^3, \lambda) \right) \right], \end{aligned} \quad (\text{E.31})$$

where  $\Theta^1, \Theta^2, \Theta^3$  and  $\Theta^4$  are four independent and identically distributed random vectors, sharing the same joint distribution  $\mathbb{P}_\Theta$ . Following the same approach as in Appendix C.2,  $H_j(\lambda)$  can be written as a single expectation:

$$H_j(\lambda) = \mathbb{E}_{\mathbb{P}_\Theta} [\eta_\lambda(\Theta^1, \Theta^2, \Theta^3, \Theta^4)], \quad (\text{E.32})$$

where the function  $\eta_\lambda$  is defined as:

$$\begin{aligned} \eta_\lambda : \mathcal{D}_\Theta \times \mathcal{D}_\Theta \times \mathcal{D}_\Theta \times \mathcal{D}_\Theta &\longrightarrow \mathbb{R} \\ (\theta^1, \theta^2, \theta^3, \theta^4) &\longmapsto K_{\theta_j}(\theta_j^1, \theta_j^2) K_g \left( \tilde{f}(\theta^1, \lambda), \tilde{f}(\theta^2, \lambda) \right) \\ &\quad + K_{\theta_j}(\theta_j^1, \theta_j^2) K_g \left( \tilde{f}(\theta^3, \lambda), \tilde{f}(\theta^4, \lambda) \right) \\ &\quad - 2 K_{\theta_j}(\theta_j^1, \theta_j^2) K_g \left( \tilde{f}(\theta^1, \lambda), \tilde{f}(\theta^3, \lambda) \right). \end{aligned} \quad (\text{E.33})$$

As  $\eta_\lambda$  is not symmetric, it must be symmetrized:

$$\tilde{\eta}_\lambda(\theta^1, \theta^2, \theta^3, \theta^4) = \frac{1}{4!} \sum_{\sigma \in \mathbb{S}_4} \eta_\lambda(\theta^{\sigma(1)}, \theta^{\sigma(2)}, \theta^{\sigma(3)}, \theta^{\sigma(4)}). \quad (\text{E.34})$$

This symmetrization does not modify the target expectation:

$$H_j(\lambda) = \mathbb{E}_{\mathbb{P}_\Theta} [\tilde{\eta}_\lambda(\Theta^1, \Theta^2, \Theta^3, \Theta^4)], \quad (\text{E.35})$$

but is required by the theory of U- and V-statistics.

The boundedness assumption on the kernels implies that:

$$\begin{aligned} \exists M_{\theta_j} > 0 : \quad \forall (\theta_j, \theta'_j) \in (\mathcal{D}_{\Theta_j})^2, \quad |K_{\theta_j}(\theta_j, \theta'_j)| &\leq M_{\theta_j} < \infty \\ \exists M_g > 0 : \quad \forall (g, g') \in (\mathcal{D}_G)^2, \quad |K_g(g, g')| &\leq M_g < \infty. \end{aligned} \quad (\text{E.36})$$

As a consequence, both  $\eta_\lambda$  and its symmetrized version  $\tilde{\eta}_\lambda$  are bounded. The corresponding bounds admit explicit expressions in terms of  $M_{\theta_j}$  and  $M_g$ :

$$\forall \lambda \in \mathcal{D}_\Lambda, \quad \forall (\theta^1, \theta^2, \theta^3, \theta^4) \in (\mathcal{D}_\Theta)^4, \quad |\eta_\lambda(\theta^1, \theta^2, \theta^3, \theta^4)| \leq 4 M_{\theta_j} M_g \quad (\text{E.37})$$

$$|\tilde{\eta}_\lambda(\theta^1, \theta^2, \theta^3, \theta^4)| \leq 4 M_{\theta_j} M_g. \quad (\text{E.38})$$

This ensures that the function  $\tilde{\eta}_\lambda$  satisfies the assumptions required by the lemmas stated in Appendix B. Specifically, using equation (E.38), one can easily show that:

$$H_j(\lambda) = \mathbb{E} [\tilde{\eta}_\lambda(\Theta^1, \Theta^2, \Theta^3, \Theta^4)] \leq 4 M_{\theta_j} M_g \quad (\text{E.39})$$

$$\mathbb{V}_{\mathbb{P}_{\Theta}}(\tilde{\eta}_{\lambda}(\Theta^1, \Theta^2, \Theta^3, \Theta^4)) \leq (4 M_{\theta_j} M_g)^2 < \infty \quad (\text{E.40})$$

$$C_r(\tilde{\eta}_{\lambda}, \mathbb{P}_{\Theta}) := \max_{1 \leq i_1, \dots, i_4 \leq 4} \mathbb{E} \left[ \left| \tilde{\eta}_{\lambda}(\Theta^{i_1}, \dots, \Theta^{i_4}) \right|^r \right] \leq (4 M_{\theta_j} M_g)^r < \infty. \quad (\text{E.41})$$

Let us now consider a Monte Carlo sample composed of  $n$  joint observations of  $(\Theta, G_{\lambda})$  where  $G_{\lambda} = \tilde{f}(\Theta, \lambda)$ . The dataset can be written as:

$$D_{\Theta G_{\lambda}, n} = \left\{ \left( \Theta^{(l)}, G_{\lambda}^{(l)} \right) \right\}_{1 \leq l \leq n} \quad \text{with} \quad G_{\lambda}^{(l)} = \tilde{f}(\Theta^{(l)}, \lambda). \quad (\text{E.42})$$

Let  $\hat{H}_j^U(\lambda)$  and  $\hat{H}_j^V(\lambda)$  denote the U- and V-statistic estimators of  $H_j(\lambda)$  based on the data from  $D_{\Theta G_{\lambda}, n}$ . Using equations (E.40) and (E.41), Lemma B.3 and Lemma B.5 can be directly specialized, which yields the following corollary.

**Corollary E.2.** *Under the boundedness assumption on kernels, one has:*

$$\mathbb{V}(\hat{H}_j^U(\lambda)) \leq \frac{\alpha}{n} = \mathcal{O}\left(\frac{1}{n}\right) \quad \text{with} \quad \alpha = 64 M_{\theta_j}^2 M_g^2 \quad (\text{E.43})$$

$$\mathbb{E} \left[ \left| \hat{H}_j^U(\lambda) - \hat{H}_j^V(\lambda) \right|^r \right] \leq \frac{\beta_r}{n^r} = \mathcal{O}\left(\frac{1}{n^r}\right) \quad \text{with} \quad \beta_r = \left(48 M_{\theta_j}^2 M_g^2\right)^r. \quad (\text{E.44})$$

Note that the constants  $\alpha$  and  $\beta_r$  depend only on the kernel bounds  $M_{\theta_j}$  and  $M_g$ , and are therefore independent of  $\lambda$  and  $\mathbb{P}_{\Theta}$ .

### E.2.1 Consistency of Estimator 1

Estimator 1 is given by:

$$\widehat{\mathcal{H}}_j^{-1} \left( D_{m,n}^{S_1} \right) = \frac{1}{m} \sum_{k=1}^m \hat{H}_j^U \left( D_{\Theta^{(k)} G^{(k)}, n} \right), \quad (\text{E.45})$$

where  $D_{m,n}^{S_1}$  is a nested Monte Carlo design generated according to Strategy (S1). To generate  $D_{m,n}^{S_1}$ , proceed as follows.

1. Generate an  $m$ -sample  $D_{\Lambda, m} = \left\{ \Lambda^{(k)} \right\}_{1 \leq k \leq m}$  from  $\mathbb{P}_{\Lambda}$ .
2. For each  $\Lambda^{(k)}$  in  $D_{\Lambda, m}$ :
  - Generate an  $n$ -sample  $D_{\Theta^{(k)}, n} = \left\{ \Theta^{(kl)} \right\}_{1 \leq l \leq n}$  from  $\mathbb{P}_{\Theta}$ .
  - Compute the output values  $G^{(kl)} = \tilde{f}(\Theta^{(kl)}, \Lambda^{(k)})$  and store them in  $D_{G^{(k)}, n}$ .
  - Merge  $D_{\Theta^{(k)}, n}$  and  $D_{G^{(k)}, n}$  to form  $D_{\Theta^{(k)} G^{(k)}, n}$ .

This sampling process guarantees the following probabilistic relationships between the data:

$$D_{\Theta^{(1)}, n} \perp\!\!\!\perp \dots \perp\!\!\!\perp D_{\Theta^{(m)}, n} \sim (\mathbb{P}_{\Theta})^{\otimes n} \quad (\text{E.46})$$

$$D_{G^{(1)}, n} \perp\!\!\!\perp \dots \perp\!\!\!\perp D_{G^{(m)}, n} \sim \mathbb{Q}_A \quad (\text{E.47})$$

$$D_{\Theta^{(1)} G^{(1)}, n} \perp\!\!\!\perp \dots \perp\!\!\!\perp D_{\Theta^{(m)} G^{(m)}, n} \sim \mathbb{Q}_B. \quad (\text{E.48})$$

Above,  $\mathbb{Q}_A$  denotes the joint distribution of the data contained in each sample  $D_{G^{(k)}, n}$ . The corresponding cumulative distribution function (CDF) is given by:

$$\forall (g_1, \dots, g_n) \in (\mathcal{D}_G)^n, \quad F_A(g_1, \dots, g_n) = \int_{\mathcal{D}_{\Lambda}} F_{G_{\lambda}}(g_1) \dots F_{G_{\lambda}}(g_n) d\mathbb{P}_{\Lambda}(\lambda), \quad (\text{E.49})$$

where  $F_{G_\lambda}$  is the CDF of  $G_\lambda = \tilde{f}(\Theta, \lambda)$ . Similarly,  $\mathbb{Q}_B$  is the joint distribution of the data contained in each sample  $D_{\Theta^{(k)}G^{(k)},n}$ . Its CDF is:

$$\begin{aligned} \forall ((\theta_1, g_1), \dots, (\theta_n, g_n)) \in (\mathcal{D}_\Theta \times \mathcal{D}_G)^n, \\ F_B((\theta_1, g_1), \dots, (\theta_m, g_m)) = \int_{\mathcal{D}_\Lambda} F_{\Theta G_\lambda}(\theta_1, g_1) \dots F_{\Theta G_\lambda}(\theta_n, g_n) d\mathbb{P}_\Lambda(\lambda), \end{aligned} \quad (\text{E.50})$$

where  $F_{\Theta G_\lambda}$  is the CDF of  $(\Theta, G_\lambda)$ .

*Proof.* We first prove that  $\widehat{\mathcal{H}}_j^1$  is unbiased.

$$\mathbb{E} \left[ \widehat{\mathcal{H}}_j^1 \right] = \frac{1}{m} \sum_{k=1}^m \mathbb{E} \left[ \widehat{H}_j^U (D_{\Theta^{(k)}G^{(k)},n}) \right] \quad (\text{E.51})$$

$$= \mathbb{E} \left[ \widehat{H}_j^U (D_{\Theta^{(1)}G^{(1)},n}) \right] \text{ with equation (E.48)} \quad (\text{E.52})$$

$$= \mathbb{E}_{\Lambda^{(1)}} \left[ \mathbb{E} \left[ \widehat{H}_j^U (D_{\Theta^{(1)}G^{(1)},n}) \middle| \Lambda^{(1)} \right] \right] \text{ with the tower property.} \quad (\text{E.53})$$

In particular, the conditional expectation can be computed as follows:

$$\mathbb{E} \left[ \widehat{H}_j^U (D_{\Theta^{(1)}G^{(1)},n}) \middle| \Lambda^{(1)} = \lambda \right] = \mathbb{E} \left[ \widehat{H}_j^U (D_{\Theta^{(1)}G_\lambda,n}) \right] = H_j(\lambda). \quad (\text{E.54})$$

The first equality holds because  $\Lambda^{(1)} \perp\!\!\!\perp \Theta^{(1)}$ . The second equality is obtained by recalling that a U-statistic is unbiased. Then, injecting equation (E.54) in equation (E.53) gives:

$$\mathbb{E} \left[ \widehat{\mathcal{H}}_j^1 \right] = \mathbb{E}_\Lambda [H_j(\Lambda^{(1)})] = \mathcal{H}_j, \quad (\text{E.55})$$

which confirms that Estimator 1 is unbiased.

Now, let us prove that this estimator is consistent. To this end, we prove that Estimator 1 converges in mean square to  $\mathcal{H}_j$ . In fact, the derivations are simpler and provide a convergence rate in the process. The mean square error (MSE) can be written as:

$$\text{MSE} \left( \widehat{\mathcal{H}}_j^1 \right) = \mathbb{E} \left[ \left( \widehat{\mathcal{H}}_j^1 - \mathcal{H}_j \right)^2 \right] = \mathbb{V} \left( \widehat{\mathcal{H}}_j^1 \right) \text{ as } \widehat{\mathcal{H}}_j^1 \text{ is unbiased} \quad (\text{E.56})$$

$$= \frac{1}{m} \mathbb{V} \left( \widehat{H}_j^U (D_{\Theta^{(1)}G^{(1)},n}) \right) \text{ with equation (E.48)} \quad (\text{E.57})$$

The total variance formula yields:

$$\mathbb{V} \left( \widehat{H}_j^U (D_{\Theta^{(1)}G^{(1)},n}) \right) = \mathbb{V}_{\Lambda^{(1)}} \left( \mathbb{E} \left[ \widehat{H}_j^U (D_{\Theta^{(1)}G^{(1)},n}) \middle| \Lambda^{(1)} \right] \right) + \mathbb{E}_{\Lambda^{(1)}} \left[ \mathbb{V} \left( \widehat{H}_j^U (D_{\Theta^{(1)}G^{(1)},n}) \middle| \Lambda^{(1)} \right) \right]. \quad (\text{E.58})$$

On the one hand, the bound provided in equation (E.39) can be used in equation (E.54) to obtain:

$$\mathbb{V}_{\Lambda^{(1)}} \left( \mathbb{E} \left[ \widehat{H}_j^U (D_{\Theta^{(1)}G^{(1)},n}) \middle| \Lambda^{(1)} \right] \right) = \mathbb{V}_{\Lambda^{(1)}} \left( H_j(\Lambda^{(1)}) \right) \leq (4M_{\theta_j} M_g)^2. \quad (\text{E.59})$$

On the other hand, the conditional variance can be bounded using Corollary E.2:

$$\mathbb{V} \left( \widehat{H}_j^U (D_{\Theta^{(1)}G^{(1)},n}) \middle| \Lambda^{(1)} = \lambda \right) = \mathbb{V} \left( \widehat{H}_j^U (D_{\Theta^{(1)}G_\lambda,n}) \right) \leq \frac{\alpha}{n}, \quad (\text{E.60})$$

which then yields:

$$\mathbb{E}_{\Lambda^{(1)}} \left[ \mathbb{V} \left( \widehat{H}_j^U (D_{\Theta^{(1)}G^{(1)},n}) \middle| \Lambda^{(1)} \right) \right] \leq \frac{\alpha}{n}. \quad (\text{E.61})$$

Combining equations (E.59) and (E.61) leads to:

$$\text{MSE} \left( \widehat{\mathcal{H}}_j^1 \right) \leq \frac{1}{m} \left( (4M_x M_y)^2 + \frac{\alpha}{n} \right) = \mathcal{O} \left( \frac{1}{m} \right). \quad (\text{E.62})$$

Therefore, Estimator 1 converges in mean square to  $\mathcal{H}_j$ , with the MSE decaying at a rate of  $1/m$ .  $\square$

### E.2.2 Consistency of Estimator 2

Estimator 2 is given by:

$$\widehat{\mathcal{H}}_j^2 (D_{m,n}^{\text{S1}}) = \frac{1}{m} \sum_{k=1}^m \widehat{H}_j^V (D_{\Theta^{(k)}G^{(k)},n}), \quad (\text{E.63})$$

where  $D_{m,n}^{\text{S1}}$  is a nested Monte Carlo design generated according to Strategy (S1). Compared to Estimator 1, the only difference lies in the use of V-statistics instead of U-statistics to estimate HSIC indices. Since V-statistics are not unbiased in general, Estimator 2 is typically biased. However, because the experimental design remains  $D_{m,n}^{\text{S1}}$ , the probabilistic relationships among the data highlighted in equations (E.46)–(E.48) still hold.

*Proof.* For Estimator 2, the proof strategy relies on a decomposition of the MSE obtained by introducing Estimator 1 as an intermediate quantity. Using the classical inequality  $(a+b)^2 \leq 2(a^2+b^2)$ , the MSE of Estimator 2 can be bounded as follows:

$$\frac{1}{2} \text{MSE} \left( \widehat{\mathcal{H}}_j^2 \right) = \frac{1}{2} \mathbb{E} \left[ \left( \widehat{\mathcal{H}}_j^2 - \mathcal{H}_j \right)^2 \right] = \frac{1}{2} \mathbb{E} \left[ \left( \widehat{\mathcal{H}}_j^2 - \widehat{\mathcal{H}}_j^1 + \widehat{\mathcal{H}}_j^1 - \mathcal{H}_j \right)^2 \right] \quad (\text{E.64})$$

$$\leq \mathbb{E} \left[ \left( \widehat{\mathcal{H}}_j^2 - \widehat{\mathcal{H}}_j^1 \right)^2 \right] + \mathbb{E} \left[ \left( \widehat{\mathcal{H}}_j^1 - \mathcal{H}_j \right)^2 \right]. \quad (\text{E.65})$$

The second term is the MSE of Estimator 1, and it was established in equation (E.62) that:

$$\mathbb{E} \left[ \left( \widehat{\mathcal{H}}_j^1 - \mathcal{H}_j \right)^2 \right] = \mathcal{O} \left( \frac{1}{m} \right). \quad (\text{E.66})$$

Regarding the first term, it is the mean squared difference between Estimators 1 and 2. Note that the difference may be written as:

$$\widehat{\mathcal{H}}_j^2 - \widehat{\mathcal{H}}_j^1 = \frac{1}{m} \sum_{k=1}^m \left( \widehat{H}_j^V (D_{\Theta^{(k)}G^{(k)},n}) - \widehat{H}_j^U (D_{\Theta^{(k)}G^{(k)},n}) \right) = \frac{1}{m} \sum_{k=1}^m \Delta_{jk}. \quad (\text{E.67})$$

Equation (E.48) indicates that the datasets  $\{D_{\Theta^{(k)}G^{(k)},n}\}_{1 \leq k \leq m}$  form an i.i.d. sample from  $\mathbb{Q}_B$ . As a consequence, the variables  $\{\Delta_{jk}\}_{1 \leq k \leq m}$  also form an i.i.d. sample as well. Applying Lemma E.1 to this sample, we obtain:

$$\mathbb{E} \left[ \left( \widehat{\mathcal{H}}_j^2 - \widehat{\mathcal{H}}_j^1 \right)^2 \right] = \mathbb{E} \left[ \left( \frac{1}{m} \sum_{k=1}^m \Delta_{jk} \right)^2 \right] \leq \mathbb{E} [\Delta_{1j}^2]. \quad (\text{E.68})$$

The expectation on the right-hand side can reformulated as:

$$\mathbb{E} [\Delta_{1j}^2] = \mathbb{E} \left[ \left( \widehat{H}_j^V (D_{\Theta^{(1)}G^{(1)},n}) - \widehat{H}_j^U (D_{\Theta^{(1)}G^{(1)},n}) \right)^2 \right] \quad (\text{E.69})$$

$$= \mathbb{E}_{\Lambda^{(1)}} \left[ \mathbb{E} \left[ \left( \widehat{H}_j^V (D_{\Theta^{(1)}G^{(1)},n}) - \widehat{H}_j^U (D_{\Theta^{(1)}G^{(1)},n}) \right)^2 \middle| \Lambda^{(1)} \right] \right], \quad (\text{E.70})$$

and Corollary E.2 (with  $r = 2$ ) can be used to bound the conditional moment:

$$\begin{aligned} & \mathbb{E} \left[ \left( \widehat{H}_j^V (D_{\Theta^{(1)}G^{(1)},n}) - \widehat{H}_j^U (D_{\Theta^{(1)}G^{(1)},n}) \right)^2 \middle| \Lambda^{(1)} = \lambda \right] \\ &= \mathbb{E} \left[ \left( \widehat{H}_j^V (D_{\Theta^{(1)}G_{\lambda,n}}) - \widehat{H}_j^U (D_{\Theta^{(1)}G_{\lambda,n}}) \right)^2 \right] \leq \frac{\beta_2}{n^2} = \mathcal{O} \left( \frac{1}{n^2} \right). \end{aligned} \quad (\text{E.71})$$

Since the constant  $\beta_2$  does not depend on  $\lambda$ , the converge rate remains unchanged after taking the expectation with respect to  $\Lambda^{(1)}$ , which gives  $\mathbb{E} [\Delta_{j1}^2] = \mathcal{O}(1/n^2)$ . Finally, putting together the elements of equations (E.65), (E.66), (E.68) and (E.71) leads to:

$$\text{MSE} \left( \widehat{\mathcal{H}}_j^2 \right) = \mathcal{O} \left( \frac{1}{m} + \frac{1}{n^2} \right). \quad (\text{E.72})$$

□

### E.2.3 Consistency of Estimator 3

The consistency proofs for Estimators 3 and 4 follow the same line of reasoning as those for Estimators 1 and 2. The MSE is split into two well-chosen expectation terms. The tower property of conditional expectation is then used to handle the two levels of uncertainty separately.

Estimator 3 is given by:

$$\widehat{\mathcal{H}}_j^3 (D_{m,n}^{S_2}) = \frac{1}{m} \sum_{k=1}^m \widehat{H}_j^U (D_{\Theta G^{(k)},n}), \quad (\text{E.73})$$

where  $D_{m,n}^{S_2}$  is a nested Monte Carlo design generated according to Strategy (S2). To generate  $D_{m,n}^{S_2}$ , proceed as follows.

1. Generate an  $m$ -sample  $D_{\Lambda,m} = \left\{ \Lambda^{(k)} \right\}_{1 \leq k \leq m}$  from  $\mathbb{P}_{\Lambda}$ .
2. Generate an  $n$ -sample  $D_{\Theta,n} = \left\{ \Theta^{(l)} \right\}_{1 \leq l \leq n}$  from  $\mathbb{P}_{\Theta}$ .
3. For each  $\Lambda^{(k)}$  in  $D_{\Lambda,m}$ :
  - Compute the output values  $G^{(kl)} = \tilde{f} \left( \Theta^{(l)}, \Lambda^{(k)} \right)$  and store them in  $D_{G^{(k)},n}$ .
  - Merge  $D_{\Theta,n}$  and  $D_{G^{(k)},n}$  to form  $D_{\Theta G^{(k)},n}$ .

The same design  $D_{\Theta,n}$  provides the input observations of  $\Theta$  used to construct all output designs  $D_{G^{(k)},n}$ . Therefore, these designs are no longer independent, which constitutes a major difference from the data generated under (S1). However, they still share the same distribution. More precisely, using the notations introduced in equations (E.47) and (E.48), one can write:

$$\forall 1 \leq k \leq m, \quad D_{G^{(k)},n} \sim \mathbb{Q}_A \quad \text{and} \quad D_{\Theta G^{(k)},n} \sim \mathbb{Q}_B. \quad (\text{E.74})$$

Similarly, all pairs  $\left( \Lambda^{(k)}, D_{\Theta G^{(k)},n} \right)$  share the same distribution  $\mathbb{Q}_C$ , yet are not independent for the same reason as explained above. The distribution  $\mathbb{Q}_C$  can be defined as the push-forward measure of  $\mathbb{P}_{\Lambda} \otimes (\mathbb{P}_{\Theta})^{\otimes n}$  through the measurable

mapping:

$$\begin{aligned} \psi : \mathcal{D}_\Lambda \times (\mathcal{D}_\Theta)^{\otimes n} &\longrightarrow \mathcal{D}_\Lambda \times (\mathcal{D}_\Theta)^{\otimes n} \times (\mathcal{D}_G)^{\otimes n} \\ \left( \lambda, (\theta_1, \dots, \theta_n) \right) &\longmapsto \left( \lambda, (\theta_1, \dots, \theta_n), (\tilde{f}(\lambda, \theta_1), \dots, \tilde{f}(\lambda, \theta_n)) \right). \end{aligned} \quad (\text{E.75})$$

Therefore, in the same spirit as equation (E.74), one can write:

$$\forall 1 \leq k \leq m, \quad \left( \Lambda^{(k)}, D_{\Theta_{G^{(k)}}} \right) \sim \mathbb{Q}_C. \quad (\text{E.76})$$

*Proof.* To establish the unbiasedness of Estimator 3, we proceed exactly as for Estimator 1. The only difference is that equation (E.48) can no longer be invoked because it is not satisfied by the data generated under (S2). However, equation (E.74) suffices. In fact, independence between the output samples  $D_{G^{(k)},n}$  is not required to prove unbiasedness.

It remains to show consistency. For Estimator 3, the quantity used to split the MSE is the pseudo-estimator  $\widehat{\mathcal{H}}_j$  from equation (4.16):

$$\widehat{\mathcal{H}}_j := \widehat{\mathcal{H}}_j(D_{\Lambda,m}) = \frac{1}{m} \sum_{k=1}^m H_j(\Lambda^{(k)}) \quad \text{with} \quad H_j(\Lambda^{(k)}) = \text{HSIC} \left( \Theta_j, G^{(k)} \mid \Lambda^{(k)} \right). \quad (\text{E.77})$$

It is useful from a theoretical perspective, although it can be computed in practice since the theoretical values of HSIC indices are not tractable. The MSE can therefore be split as follows:

$$\frac{1}{2} \text{MSE} \left( \widehat{\mathcal{H}}_j^3 \right) = \frac{1}{2} \mathbb{E} \left[ \left( \widehat{\mathcal{H}}_j^3 - \mathcal{H}_j \right)^2 \right] = \frac{1}{2} \mathbb{E} \left[ \left( \widehat{\mathcal{H}}_j^3 - \widehat{\mathcal{H}}_j + \widehat{\mathcal{H}}_j - \mathcal{H}_j \right)^2 \right] \quad (\text{E.78})$$

$$\leq \mathbb{E} \left[ \left( \widehat{\mathcal{H}}_j^3 - \widehat{\mathcal{H}}_j \right)^2 \right] + \mathbb{E} \left[ \left( \widehat{\mathcal{H}}_j - \mathcal{H}_j \right)^2 \right]. \quad (\text{E.79})$$

By convexity of the square function, the quantity inside the first expectation satisfies:

$$\left( \widehat{\mathcal{H}}_j^3 - \widehat{H}_j \right)^2 = \left( \frac{1}{m} \sum_{k=1}^m \left( \widehat{H}_j^U(D_{\Theta_{G^{(k)},n}}) - H_j(\Lambda^{(k)}) \right) \right)^2 \quad (\text{E.80})$$

$$\leq \frac{1}{m} \sum_{k=1}^m \left( \widehat{H}_j^U(D_{\Theta_{G^{(k)},n}}) - H_j(\Lambda^{(k)}) \right)^2. \quad (\text{E.81})$$

Then, taking expectations gives:

$$\mathbb{E} \left[ \left( \widehat{\mathcal{H}}_j^3 - \widehat{H}_j \right)^2 \right] \leq \frac{1}{m} \sum_{k=1}^m \mathbb{E} \left[ \left( \widehat{H}_j^U(D_{\Theta_{G^{(k)},n}}) - H_j(\Lambda^{(k)}) \right)^2 \right] \quad (\text{E.82})$$

$$= \mathbb{E} \left[ \left( \widehat{H}_j^U(D_{\Theta_{G^{(1)},n}}) - H_j(\Lambda^{(1)}) \right)^2 \right] \quad (\text{E.83})$$

$$= \mathbb{E}_{\Lambda^{(1)}} \left[ \mathbb{E} \left[ \left( \widehat{H}_j^U(D_{\Theta_{G^{(1)},n}}) - H_j(\Lambda^{(1)}) \right)^2 \mid \Lambda^{(1)} \right] \right]. \quad (\text{E.84})$$

Equation (E.83) is a consequence of all terms in the sum sharing the same distribution  $\mathbb{Q}_C$ , as stated in equation (E.76). The conditional expectation in equation (E.84) can be bounded using the simple arguments introduced earlier ( $\Lambda^{(1)} \perp\!\!\!\perp \Theta$ ,

the unbiasedness of U-statistics, and Corollary E.2):

$$\mathbb{E} \left[ \left( \widehat{H}_j^U (D_{\Theta_{G^{(1)},n}}) - H_j(\Lambda^{(1)}) \right)^2 \middle| \Lambda^{(1)} = \lambda \right] = \mathbb{E} \left[ \left( \widehat{H}_j^U (D_{\Theta_{G_{\lambda},n}}) - H_j(\lambda) \right)^2 \right] \quad (\text{E.85})$$

$$= \mathbb{V} \left( \widehat{H}_j^U (D_{\Theta_{G_{\lambda},n}}) \right) \quad (\text{E.86})$$

$$\leq \frac{\alpha}{n}. \quad (\text{E.87})$$

Since the constant  $\alpha$  does not depend on  $\lambda$ , taking expectations with respect to  $\Lambda^{(1)}$  provides:

$$\mathbb{E} \left[ \left( \widehat{\mathcal{H}}_j^3 - \widehat{\mathcal{H}}_j \right)^2 \right] = \mathcal{O} \left( \frac{1}{n} \right). \quad (\text{E.88})$$

The Big-O rate associated with the second expectation in equation (E.78) can easily be derived:

$$\mathbb{E} \left[ \left( \widehat{\mathcal{H}}_j - \mathcal{H}_j \right)^2 \right] = \mathbb{V} \left( \widehat{\mathcal{H}}_j \right) = \frac{1}{m} \mathbb{V}_{\Lambda^{(1)}} \left( H_j(\Lambda^{(1)}) \right) \leq \frac{1}{m} (4M_x M_y)^2 = \mathcal{O} \left( \frac{1}{m} \right). \quad (\text{E.89})$$

The first equality holds because  $\mathbb{E} \left[ \widehat{\mathcal{H}}_j \right] = \mathcal{H}_j$ . Moreover, the bound on the variance of  $H_j(\Lambda^{(1)})$  follows from equation (E.39). Finally, combining equations (E.88) and (E.89) yields:

$$\text{MSE} \left( \widehat{\mathcal{H}}_j^3 \right) = \mathcal{O} \left( \frac{1}{m} + \frac{1}{n} \right). \quad (\text{E.90})$$

□

#### E.2.4 Consistency of Estimator 4

Estimator 4 is given by:

$$\widehat{\mathcal{H}}_j^4 \left( D_{m,n}^{\text{S}_2} \right) = \frac{1}{m} \sum_{k=1}^m \widehat{H}_j^V \left( D_{\Theta^{(k)}_{G^{(k)},n}} \right), \quad (\text{E.91})$$

where  $D_{m,n}^{\text{S}_2}$  is a nested Monte Carlo design generated according to Strategy (S2). Compared to Estimator 3, the only difference lies in the use of V-statistics, which typically introduces bias. The probabilistic relationships among the data highlighted in equations (E.74) and (E.75) remain valid.

*Proof.* Up to equation (E.85), we proceed exactly as for Estimator 3. The analogous expression is:

$$\mathbb{E} \left[ \left( \widehat{H}_j^V (D_{\Theta_{G^{(1)},n}}) - H_j(\Lambda^{(1)}) \right)^2 \middle| \Lambda^{(1)} = \lambda \right] = \mathbb{E} \left[ \left( \widehat{H}_j^V (D_{\Theta_{G_{\lambda},n}}) - H_j(\lambda) \right)^2 \right]. \quad (\text{E.92})$$

The expectation on the right-hand side can be split as follows:

$$\begin{aligned} & \frac{1}{2} \mathbb{E} \left[ \left( \widehat{H}_j^V (D_{\Theta_{G_{\lambda},n}}) - H_j(\lambda) \right)^2 \right] \\ & \leq \mathbb{E} \left[ \left( \widehat{H}_j^U (D_{\Theta_{G_{\lambda},n}}) - H_j(\lambda) \right)^2 \right] + \mathbb{E} \left[ \left( \widehat{H}_j^V (D_{\Theta_{G_{\lambda},n}}) - \widehat{H}_j^U (D_{\Theta_{G_{\lambda},n}}) \right)^2 \right]. \quad (\text{E.93}) \end{aligned}$$

Bounds for these two terms are provided directly by Corollary E.2:

$$\mathbb{E} \left[ \left( \widehat{H}_j^U (D_{\Theta_{G,\lambda},n}) - H_j(\boldsymbol{\lambda}) \right)^2 \right] = \mathbb{V} \left( \widehat{H}_j^U (D_{\Theta_{G,\lambda},n}) \right) \leq \frac{\alpha}{n} \quad (\text{E.94})$$

$$\mathbb{E} \left[ \left( \widehat{H}_j^V (D_{\Theta_{G,\lambda},n}) - \widehat{H}_j^U (D_{\Theta_{G,\lambda},n}) \right)^2 \right] \leq \frac{\beta_2}{n^2}. \quad (\text{E.95})$$

Reassembling the full line of reasoning, we arrive at:

$$\begin{aligned} \frac{1}{4} \text{MSE} \left( \widehat{\mathcal{H}}_j^4 \right) &\leq \mathbb{E}_{\boldsymbol{\Lambda}^{(1)}} \left[ \mathbb{E} \left[ \left( \widehat{H}_j^U (D_{\Theta_{G^{(1)},n})} - H_j(\boldsymbol{\Lambda}^{(1)}) \right)^2 \middle| \boldsymbol{\Lambda}^{(1)} \right] \right] \\ &\quad + \mathbb{E}_{\boldsymbol{\Lambda}^{(1)}} \left[ \mathbb{E} \left[ \left( \widehat{H}_j^V (D_{\Theta_{G^{(1)},n})} - \widehat{H}_j^U (D_{\Theta_{G^{(1)},n})} \right)^2 \middle| \boldsymbol{\Lambda}^{(1)} \right] \right] \\ &\quad + \mathbb{E} \left[ \left( \widehat{\mathcal{H}}_j - \mathcal{H}_j \right)^2 \right]. \end{aligned} \quad (\text{E.96})$$

Finally, combining the Big-O rates obtained for all three expectations yields:

$$\text{MSE} \left( \widehat{\mathcal{H}}_j^4 \right) = \mathcal{O} \left( \frac{1}{n} \right) + \mathcal{O} \left( \frac{1}{n^2} \right) + \mathcal{O} \left( \frac{1}{m} \right) = \mathcal{O} \left( \frac{1}{m} + \frac{1}{n} \right). \quad (\text{E.97})$$

□