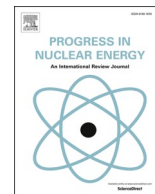




Contents lists available at ScienceDirect

Progress in Nuclear Energy

journal homepage: www.elsevier.com/locate/pnucene

Passive safety systems analysis: A novel approach for inverse uncertainty quantification based on Stacked Sparse Autoencoders and Kriging metamodeling

Giovanni Roma^a, Federico Antonello^a, Francesco Di Maio^{a,*}, Nicola Pedroni^b, Enrico Zio^{a,d}, Andrea Bersano^c, Cristina Bertani^b, Fulvio Mascari^c

^a Energy Department, Politecnico di Milano, Via La Masa 34, Milano, 20156, Italy

^b Energy Department Politecnico di Torino, Corso Duca degli Abruzzi, 24, Torino, 10129, Italy

^c ENEA- BOLOGNA, Via Martiri Monte Sole 4, Bologna, 40129, Italy

^d MINES ParisTech, PSL Research University, CRC, Sophia Antipolis, France

ARTICLE INFO

Keywords:

Nuclear safety
Inverse uncertainty quantification
Bayesian inference
Kriging metamodeling
Autoencoders

ABSTRACT

In passive safety system analysis, it is important to provide the uncertainty quantification of the Thermal-Hydraulic (T-H) code output (e.g., the amount of energy exchanged by the passive safety system during an accidental transient). This requires setting proper Probability Density Functions (PDFs) to represent the uncertainty of selected code inputs and the propagation of this uncertainty through the code. One way to obtain the PDF is by Inverse Uncertainty Quantification (IUQ) methods, which rely directly on experimental data and code simulation results. In this work, we present an innovative IUQ method based on: (i) Stacked Sparse Autoencoders (SSAEs) to reduce the problem dimensionality; and (ii) Kriging metamodels to lower the computational burden associated with the sampling of the uncertain input parameters posterior PDF by Markov Chain Monte Carlo (MCMC) (for which many model simulations are typically required). The novelty stands in the use of SSAEs for dimensionality reduction: this allows using directly the raw data available from experimental facilities or computer codes (typically characterized by small signal-to-noise ratios) without having to resort to filtering techniques, whose choice and setting are nontrivial and bias the results. The proposed approach is applied to the power exchanged by the Heat Exchanger (HX) predicted by the RELAP5-3D model of the PERSEO facility, characterized by a small signal-to-noise ratio (SNR) value. Principal Component Analysis (PCA) and SSAE are compared to explain the application of these methodologies in the context of IUQ and highlight the main advantages and drawbacks while also showing the suitability to deal with non-filtered (raw) data.

* Corresponding author.

E-mail address: francesco.dimaio@polimi.it (F. Di Maio).

¹ When the subscript is within brackets, $z_{(l)}$ indicates the l^{th} hidden layer output of the SSAE (i.e., a vectorial quantity); otherwise, z_k indicates the k^{th} entry of a vector transformed into the p^* -dimensional feature subspace.

<https://doi.org/10.1016/j.pnucene.2022.104209>

Received 17 September 2021; Received in revised form 2 March 2022; Accepted 20 March 2022

Available online 6 April 2022

0149-1970/© 2022 Elsevier Ltd. All rights reserved.

Nomenclature			
Acronyms		LHS	Latin Hypercube Sampling
ANN	Artificial Neural Network	LOOCV	Leave-One-Out Cross Validation
BCs	Boundary conditions	MAP	Maximum A Posteriori
BE	Best Estimate	MCMC	Markov Chain Monte Carlo
BEPU	Best Estimate Plus Uncertainty	MLE	Maximum Likelihood Estimation
CV	Cross Validation	MSE	Mean Squared Error
DAE	Denosing Autoencoder	NPPs	Nuclear Power Plants
DNN	Deep Neural Network	OP	Overall Pool
DOE	Design Of Experiment	PCA	Principal Component Analysis
EM	Expectation-Maximization	PCs	Principal Component Vectors
ENEA	Agenzia nazionale per le nuove tecnologie, l'energia e lo sviluppo economico sostenibile	PDF	Probability Density Function
HX	Heat Exchanger	PERSEO	in-Pool Energy Removal System for Emergency Operation
HXP	Heat Exchanger Pool	PV	Pressure Vessel
ICs	Initial Conditions	SAE	Sparse Autoencoder
IUQ	Inverse Uncertainty Quantification	SSAE	Stacked Sparse Autoencoder
KDE	Kernel Density Estimation	SNR	Signal to Noise Ratio
		T-H	Thermal-Hydraulic
		TV	Triggering Valve
		VAE	Variational Autoencoder

List of symbols

Symbols	Dimension	Description
d	1×1	Number of calibration parameters
q	1×1	Number of design variables
m	1×1	Number of design points
p	1×1	Model output dimension
N_{exp}	1×1	Number of independent experimental measurements
θ	$d \times 1$	Calibration parameters vector
x	$q \times 1$	Design variable vector
$y^M(\theta)$	$p \times 1$	RELAP5-3D computer model output for Test 7 pt.2
y^E	$p \times 1$	Experimental data
\tilde{y}^E	$p \times 1$	Reconstructed experimental data
y		
ϵ	$p \times 1$	Measurement error
$\mathbf{I}_{\epsilon_{exp}}^2$	$p \times p$	Covariance matrix for the measurement error
Σ_{exp}	$p^* \times p^*$	Measurement error covariance matrix in the p^* -dimensional reduced space
Σ_{MM}	$p \times p$	Covariance matrix for the code uncertainty
$\Sigma_{Kriging}$	$p^* \times p^*$	Covariance matrix for the code uncertainty in the reduced space
Θ	$m \times d$	Ensemble of m design points
Y	$p \times m$	Ensemble of m BE computer model outputs
$\bar{\mu}_Y$	$p \times 1$	Column vector of the row means of Y
p^*	1×1	Dimension of the features subspace
Φ	$p^* \times p$	Transformation matrix (PCA)
\tilde{z}^{MM}	$p^* \times 1$	Kriging metamodel prediction in the features subspace
z^{MM}	$p^* \times 1$	Distribution of the Kriging metamodel prediction in the features subspace
z^E	$p^* \times 1$	Experimental data projected in the feature subspace
Z	$p^* \times m$	Ensemble of m RELAP5-3D simulation outputs transformed into the feature subspace
z_k	1×1	k^{th} entry of a vector transformed into the feature subspace
L	1×1	Number of the SSAE hidden layers
K_l	1×1	Dimension of the l^{th} SSAE hidden layer
$z_{(l)}$	$K_l \times 1$	Neurons output of the l^{th} hidden layer ¹
$\sigma(\cdot)$	–	Sigmoid transfer function
$W_{(l)}$	$K_l \times K_{l-1}$	Weight matrix of the l^{th} SSAE layer
$b_{(l)}$	$K_l \times 1$	Bias vector of the l^{th} SSAE layer
E	1×1	SSAE cost function
R_{error}	1×1	Reconstruction error (adopted for both the SSAE and the PCA)
R_{sparse}	1×1	Sparsity regularization term
β	1×1	Coefficient for R_{sparse}
R_{L2}	1×1	L2 regularization term
λ	1×1	Coefficient for R_{L2}
ρ	1×1	Desired averaged neuron activation
N	1×1	Number of RELAP5-3D simulations adopted for the Forward Uncertainty Propagation

1. Introduction

Over the last decades, nuclear safety analysis frameworks based on Best Estimate Plus Uncertainty (BEPU) approaches for thermal-hydraulics transient calculations have gathered great interest (D'Auria et al., 2006; Iaea, 2008). These frameworks stand on Best Estimate (BE) Thermal-Hydraulic (T-H) codes to compute the safety margins of relevant parameters (e.g., fuel pellet maximum centerline temperature) during Nuclear Power Plants (NPPs) accidental scenarios (Iaea, 2008). The computation of the safety margins, taking into account the uncertainty of the calculation, requires identifying the main sources of uncertainty, selecting the most relevant uncertain input parameters and propagating the input uncertainties through the BE T-H code. One of the daunting issues related to this is quantifying the epistemic uncertainty that affects some of the input parameters, i.e., the uncertainty derived from lack of knowledge of the phenomena which these parameters describe (Ferson and Ginzburg, 1996; Durga Rao et al., 2007; Winkler, 1996; Apostolakis, 1994; Ferson et al., 2004). Traditionally, probability constitutes the mathematical structure used to represent epistemic uncertainty and expert judgment is typically used to specify Probability Density Functions (PDFs), nominal values and upper and lower bounds of the parameters (Iaea, 2008; Pourgol-Mohammad, 2009; Helton and Johnson, 2011). Such ad-hoc expert judgment can be aided by Inverse Uncertainty Quantification (IUQ), a powerful tool that determines the uncertainty of the parameters relying on available experimental data and code simulation results (Shrestha and Kozlowski, 2016).

Different approaches have been developed to carry out IUQ. Among them, Bayesian inference is considered the standard approach (Katafygiotis and Beck, 1998a, 1998b), wherein the computation of the posterior PDF is typically tackled through Markov Chain Monte Carlo (MCMC) (Gelman et al., 2015). MCMC is a class of algorithms that allows posterior sampling by running numerous code simulations (up to hundreds of thousands). However, in the context of NPP safety analysis, BE computer models (e.g., RELAP, etc.) are computationally expensive for implementing MCMC sampling that becomes practically unfeasible with the current computational resources available. As an example, MCMC sampling with 10^5 iterations with a BE model that takes 1 h per run would last more than 11 years. To tackle this computational issue, having as a target the prediction of specific phenomena and the related parameters, the BE computer model can be replaced by a computationally cheaper metamodel (Kennedy and O'Hagan, 2001), which is an approximation of the input/output relationship modeled within the behavior of the original BE model (Haftka et al., 2020; Wang and Shan,

2007). For a comprehensive survey of inverse uncertainty quantification methods applied to nuclear system thermal-hydraulics problems, refer to (Wu et al., 2021). Among the various metamodeling approaches (Wang and Shan, 2007), Kriging (Rasmussen and Williams, 2006; Lataniotis et al., 2019) has been successfully applied in various IUQ problems (Kennedy and O'Hagan, 2001; Wilkinson, 2010; Wang et al., 2019; Arendt et al., 2012; Wu et al., 2018a; Wu et al., 2018b). The benefit of Kriging is the uncertainty estimation each time a prediction is performed (Lataniotis et al., 2019).

When the BE model output is a time-dependent scalar quantity (i.e., a time series), at least four different approaches can be adopted to build a metamodel. The first one considers the time-dependent scalar output as a vectorial quantity and builds independent metamodels for each time instance; nevertheless, this approach may result in a significant loss of information, since the multiple outputs can be highly correlated (Fricker et al., 2013). Another method is to treat time as an additional input (Kennedy and O'Hagan, 2001); however, in this approach the number of training points becomes dramatically high in very long time series. For example, if the BE model returns the output at 1000 time instances, with a Design Of Experiment (DOE) of 100 sample points for other inputs such as calibration parameters and design variables, the number of training patterns for the metamodel would be $1000 \cdot 100 = 100000$. The third alternative approach is to use a multi-output emulator (e.g., a Kriging metamodel for predicting a p -dimensional vectorial quantity, where p is the number of time instances) (Fricker et al., 2013; Kleijnen and Mehdad, 2014); this approach outperforms emulators with time as an additional input (Conti and O'Hagan, 2010) but, for very large p , even multi-output emulators may experience a reduction in the metamodeling efficiency (Mohammadi et al., 2019). The fourth alternative approach has been developed to address the abovementioned limitations for high-dimensional time series and consists in performing a dimensionality reduction to retain a relatively small number of significant features to represent the entire output space. In fact, when the output data is redundant (e.g., the BE computer model responses for nearby time instances are strongly correlated) and its dimensionality p is too large to be processed (e.g., through a single multi-output emulator), a reduced set of p^* features containing the most relevant information of the original data can be used instead of the whole set of data and a separate metamodel can be built for each of the p^* extracted features.

Principal Component Analysis (PCA) is one of the most common approaches to carry out dimensionality reduction (Van Der Maaten et al., 2009). PCA performs a linear mapping from a high-dimensional space onto a lower-dimensional space, such that the mapped variables are uncorrelated and keep as much as possible of the original data set variance (Jolliffe, 2002). Higdon et al. (2008) proposed PCA to carry out dimensionality reduction on time series outputs and Kriging metamodels to emulate such reduced outputs. PCA has been used in different Bayesian IUQ/calibration problems with high-dimensional outputs for constructing fast-running metamodels (Wilkinson, 2010; Higdon et al., 2008, 2013; Wu et al., 2018c; Nagel et al., 2020). However, linear dimensionality reduction techniques like PCA cannot deal with complex (real-world) non-linear data (Van Der Maaten et al., 2009). Furthermore, for the specific case of interest here, BE code results may be affected by higher noise (maybe due to numerics or correlation errors (Roma et al., 2021)) in comparison with the experimental data in relation to the power exchanged by the HX (Yang et al., 2018) and typically require pre-processing by the analyst (e.g., filtering of the available raw data) (Roma et al., 2021), which may add a bias. To overcome the limitations of PCA when dealing with time series with small signal-to-noise ratio values, in this work, we explore the use of Autoencoders (AEs) for dimensionality reduction.

An AE is an Artificial Neural Network (ANN) designed to learn new features of the data by reconstructing the input itself (Zhao et al., 2019). It is composed of an encoder and a decoder network. The former maps the high-dimensional input into a small number of features (i.e., into a

lower-dimensional representation), whereas the latter recovers the high-dimensional inputs from the features. AEs are widely used to perform dimensionality reduction (Holden et al., 2006; Wang et al., 2016; Monisha et al., 2019), machine health monitoring (Zhao et al., 2019), and image processing (Mao et al., 2016). Over the last years, several variants of AEs have been developed, such as *sparse* (SAEs), *denoising* (DAEs) and *variational autoencoders* (VAEs) (Olshausen and Fieldt, 1997; Vincent and Larochelle, 2008; Kingma and Welling, 2014). Among them, SAEs, which strive to extract discriminative features avoiding overfitting, are widely used to identify the most relevant and comprehensive set of features for the specific application (Ng, 2011). Moreover, multiple pretrained AEs can be stacked to form a multiple-hidden-layer ANN, called Stacked Sparse Autoencoder (SSAE), improving the representational and modeling power (Utgoff and Stracuzzi, 2002).

In this work, we embed SSAEs for output dimensionality reduction into an IUQ Kriging-based approach, where the Kriging metamodel emulates each of the SSAE extracted features. Up to the authors knowledge, SSAEs (and, in general, non-linear dimensionality reduction techniques) have not yet been applied with Kriging metamodeling in a Bayesian IUQ problem. The rationale for using AEs is related to their capability to: (i) work with raw data without pre-processing; and (ii) deal with nonlinearities (Holden et al., 2006). The proposed approach is applied within a Bayesian IUQ framework aimed at determining the input parameters' PDFs of a RELAP5-3D model of the PERSEO facility (Bandini et al., 2011), for which time series measurements are available. A SSAE is applied to reduce the problem dimensionality, and fast-running Kriging metamodels are implemented to emulate the RELAP5-3D behavior at a lower computational cost. PCA and SSAE are compared to provide some understanding on the application and the use of these methods for IUQ.

The remainder of the paper is organized as follows. In Section 2 the formulation of the IUQ problem is presented. Section 3 illustrates the proposed approach. The case study regarding the PERSEO experimental facility and the RELAP5-3D model are introduced in Section 4. Section 5 displays the IUQ results for the proposed approach applied to the case study of Section 4 and provides a comparison between PCA and SSAE. Section 6 concludes the work.

2. The formulation of the IUQ problem

Let $y^E(x)$ be a measured experimental quantity and $y^M(x, \theta)$ the corresponding quantity simulated through a computer model. Let $\theta = [\theta_1, \theta_2, \dots, \theta_d]^T$ and $x = [x_1, x_2, \dots, x_q]^T$ be the *calibration parameters* and the *design variables*, respectively (Kennedy and O'Hagan, 2001). Design variables are all the measurable inputs that define the conditions or scenarios under which the experiment is carried out (e.g., Boundary Conditions (BCs) and Initial Conditions (ICs)) (Wu et al., 2018b, 2021), whereas calibration parameters are input to the computer model, but they are unknown or not measurable in the physical experiment (Wu et al., 2018b, 2021). x are uniquely defined by the experiment and, therefore, known a priori (Wu et al., 2018b). The objective of the IUQ is to determine the PDF associated with θ .

The relationship between $y^E(x)$ and $y^M(x, \theta)$ can be described by the *model updating equation* (Wu et al., 2018b):

$$y^E(x) = y^M(x, \theta) + \delta(x) + \varepsilon \quad (1)$$

where $\delta(x)$ and ε are the *model discrepancy* and the *measurement error*, respectively. $\delta(x)$ is due to approximations in $y^M(x, \theta)$, whereas $\varepsilon \sim N(\mu, I\sigma_{exp}^2)$ is an additive measurement error usually assumed to be Gaussian-distributed. The specification of $\delta(x)$ requires complex considerations on the model and is ignored in the current work; thus, the model updating equation reduces to:

$$y^E(x) = y^M(x, \theta) + \varepsilon \quad (2)$$

For a comprehensive discussion about model discrepancy, refer to (Kennedy and O'Hagan, 2001; Arendt et al., 2012; Wu et al., 2018b).

2.1. The Bayesian formulation of the inverse UQ problem

IUQ aims at quantifying the posterior of the calibration parameters θ , $p(\theta|y^E)$ i.e., the distribution after the experimental data is observed. According to Bayes rule, this can be calculated as:

$$p(\theta|y^E) = \frac{p(y^E|\theta)p(\theta)}{\int p(y^E|\theta)p(\theta)d\theta} \quad (3)$$

where $p(\theta)$ is the prior PDF and $p(y^E|\theta)$ is the likelihood function that describes the joint probability of the observed data y^E as a function of the parameters θ .

Let $y^E(x) = [y_1^E(x), \dots, y_p^E(x)]$ and $y^M(x, \theta) = [y_1^M(x, \theta), \dots, y_p^M(x, \theta)]$ be p -dimensional vectors; assuming ε to be zero-mean Gaussian-distributed (i.e., $\varepsilon \sim N(0, I\sigma_{exp}^2)$), the likelihood function can be derived from equation (2):

$$p(y^E|\theta) = \prod_{i=1}^{N_{exp}} \frac{1}{(\sqrt{2\pi})^p \sqrt{|I\sigma_{exp}^2|}} \exp \left[-\frac{1}{2} [y^E(x_i) - y^M(x_i, \theta)]^T (I\sigma_{exp}^2)^{-1} [y^E(x_i) - y^M(x_i, \theta)] \right] \quad (4)$$

where N_{exp} is the number of experiments carried out and $I\sigma_{exp}^2$ is the $p \times p$ covariance matrix of the measurement error. The integral in Eq. (3) is typically analytically intractable; a Markov Chain Monte Carlo (MCMC) algorithm can be implemented to tackle this problem and the MCMC samples are then used to infer $p(\theta|y^E)$. MCMC algorithms may require a significant number (e.g., many thousands) of iterations, each of which needs a code run to evaluate $y^M(x, \theta)$, which can be computationally demanding. Kriging metamodeling is a common choice to overcome this since, with the output estimates, also an estimation of the metamodel uncertainty is provided. Assuming that a Kriging metamodel is used to emulate $y^M(x, \theta)$, the posterior PDF becomes:

$$p(\theta|y^E) \propto p(\theta) \cdot \prod_{i=1}^{N_{exp}} \frac{1}{(\sqrt{2\pi})^p \sqrt{|\Sigma|}} \exp \left[-\frac{1}{2} [y^E(x_i) - \hat{y}(x_i, \theta)]^T \Sigma^{-1} [y^E(x_i) - \hat{y}(x_i, \theta)] \right] \quad (5)$$

where $\hat{y}(x, \theta)$ is the Kriging prediction, whereas the covariance matrix of the likelihood $\Sigma = I\sigma_{exp}^2 + \Sigma_{MM}$ is the sum of the measurement error covariance matrix $I\sigma_{exp}^2$ and the metamodel uncertainty covariance matrix Σ_{MM} .

3. Proposed IUQ approach

To perform IUQ within a Bayesian framework, we propose an approach that comprises two main steps:

1. Dimensionality reduction and Kriging metamodeling (Section 3.1);
2. Bayesian inference (by MCMC sampling) (Section 3.2).

3.1. Dimensionality reduction and Kriging metamodeling

In the context of Bayesian IUQ, the objective of the present Section is to illustrate how to build a metamodel for emulating the time-dependent scalar quantity $y^M(x, \theta; t)$ computed by a BE code. Without loss of generality, let us assume that only a single experimental time series measurement is available for the Bayesian IUQ (i.e., $N_{exp} = 1$, as it is in the case study illustrated in Section 4); then, the dependence of the forward model $y^M(x, \theta; t)$ on x is absorbed into the definition of $y^M = y^M(\theta; t)$. If we assume that $y^M(\theta; t)$ is ticked at p different pre-defined time instances (i.e., $[y^M(\theta; t_1), \dots, y^M(\theta; t_p)]$), the time-dependent scalar output can be treated as multivariate (vectorial), i.e., $y^M(\theta) = [y^M(\theta; t_1), \dots, y^M(\theta; t_p)]$. Let $\theta = [\theta^{(1)}, \dots, \theta^{(m)}]$ be the DOEs and $Y = [y^{(1)}, \dots, y^{(m)}]$ the $p \times m$ matrix containing the relative m p -dimensional BE model responses. A way to handle this data is to transform the p -dimensional output $y \in \mathbb{R}^p$ into a reduced p^* -dimensional features space $Z \subset \mathbb{R}^{p^*}$ (with $p^* \ll p$) through a dimensionality reduction technique (e.g., PCA or SSAEs) and, then, build p^* separate independent metamodels that emulate the p^* extracted features. Actually, it is also possible to build a p^* -dimensional multi-output surrogate model, either by Kriging or ANN, especially convenient

with PCA since the transformed variables are uncorrelated.

In this work, the m BE model responses contained in Y are mapped onto $Z \subset \mathbb{R}^{p^*}$ and arranged in the $p^* \times m$ features matrix $Z = [z^{(1)}, \dots, z^{(m)}]$. Then, the m input-output training patterns (i.e., $\theta = [\theta^{(1)}, \dots, \theta^{(m)}]$ and the corresponding m transformed model responses $[z_j^{(1)}, \dots, z_j^{(m)}]$ are used to train an independent metamodel for each feature j , with $j = 1, 2, \dots, p^*$. Fig. 1 shows a schematic diagram of the metamodel approach adopted.

Among the different methods proposed to perform dimensionality reduction, we present a technique based on SSAEs that is eventually

compared to PCA.

3.1.1. Sparse autoencoders

An *autoencoder* (AE) is a type of Artificial Neural Network (ANN) that is designed to learn new low-dimensional latent features of the data by trying to reconstruct the input data (Zhao et al., 2019). It is composed of an *encoder* network, that maps high-dimensional vector data into a lower-dimensional space of *features* and a *decoder* network that retrieves the original vector from the features (Holden et al., 2006) (Fig. 2).

The encoder transforms a p -dimensional vector y into its K_1 -dimensional hidden representation $z_{(1)} = [z_{(1),1}, z_{(1),2}, \dots, z_{(1),K_1}]$:

$$z_{(1)} = f(\mathbf{W}_{(1)}y + \mathbf{b}_{(1)}) \quad (6)$$

where $z_{(l),j}$ is the j^{th} neuron output of the l^{th} hidden layer (i.e., 1 in the case of basic SAE), f , $\mathbf{W}_{(1)}$, $\mathbf{b}_{(1)}$ are the encoder transfer function, the weight matrix and the bias vector, respectively. The decoder transforms

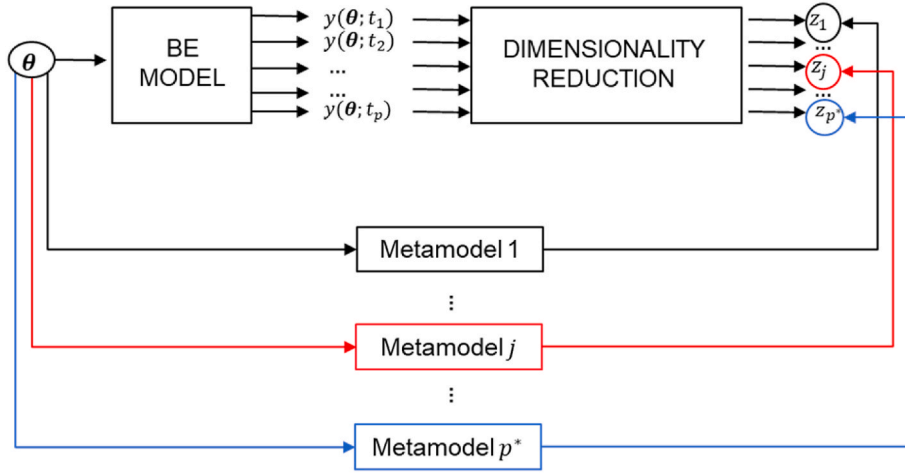


Fig. 1. Diagram of the metamodel methodology adopted for dimensionality reduction.

the hidden representation back into \tilde{y} (i.e., the reconstruction of y):

$$\tilde{y} = g(W_{(2)}z_{(1)} + b_{(2)}) \quad (7)$$

where g , $W_{(2)}$, $b_{(2)}$ are the decoder transfer function, the weight matrix and the bias vector, respectively.

The SAE, which is a variation of the AE, imposes sparsity constraints on the hidden neurons to encourage the identification of discriminative features (Zhao et al., 2019; Ng, 2011). The determination of $W_{(1)}$, $b_{(1)}$, $W_{(2)}$, $b_{(2)}$ is carried out through the minimization of the following cost function:

$$E = R_{error} + \beta R_{sparse} + \lambda R_{L2} \quad (8)$$

where R_{error} , R_{sparse} and R_{L2} are the mean squared error function, the sparsity regularizer and the L_2 regularizer, respectively, whereas β and λ allow tuning the importance of R_{sparse} and R_{L2} in the cost function. The reconstruction error R_{error} allows quantifying the accuracy of the SAE in the reconstruction of the input vectors:

$$R_{error} = \frac{1}{m_{train}} \sum_{i=1}^{m_{train}} \|y^{(i)} - \tilde{y}^{(i)}\|^2 \quad (9)$$

where m_{train} is the number of training patterns.

A hidden layer's neuron is considered "active" when its value is large

(i.e., close to 1.0, in case of sigmoid transfer function) and "inactive" when its value is small (i.e., close to 0.0, in case of sigmoid transfer function). Let $\hat{\rho}_j$ be the average output activation measure of the j^{th} hidden neuron on the training dataset (i.e., for $y^{(i)} \in Y$, $i = 1, 2, \dots, m_{train}$):

$$\hat{\rho}_j = \frac{1}{m_{train}} \sum_{i=1}^{m_{train}} z_{(1)j}^{(i)} \quad (10)$$

where $z_{(1)j}^{(i)}$ is the j^{th} neuron output of the i^{th} hidden representation $z_{(1)}^{(i)} = [z_{(1)1}^{(i)}, \dots, z_{(1)j}^{(i)}, \dots, z_{(1)K_1}^{(i)}]$, $j = 1, \dots, K_1$, $i = 1, 2, \dots, m_{train}$. It has been shown that the extraction of discriminative features $z_{(1)}$ is favored by requiring the sparsity of the AE (Ng, 2011), which imposes the neurons to be inactive most of the time (i.e., all the SAEs hidden neurons are characterized by a small value of $\hat{\rho}_j$, e.g., $\hat{\rho}_j = \rho = 0.05$). To this aim, the sparsity regularization term, R_{sparse} , is included into the expression of E to impose such a sparsity constraint. R_{sparse} penalizes $\hat{\rho}_j$ deviating from ρ through the Kullback-Leibler (KL) divergence measure:

$$R_{sparse} = \sum_{j=1}^{K_1} KL(\rho \| \hat{\rho}_j) = \sum_{j=1}^{K_1} \left[\rho \log \frac{\rho}{\hat{\rho}_j} + (1 - \rho) \log \frac{1 - \rho}{1 - \hat{\rho}_j} \right] \quad (11)$$

It is worth noting that $KL(\rho \| \hat{\rho}_j) = 0$ if $\hat{\rho}_j = \rho$, and it increases

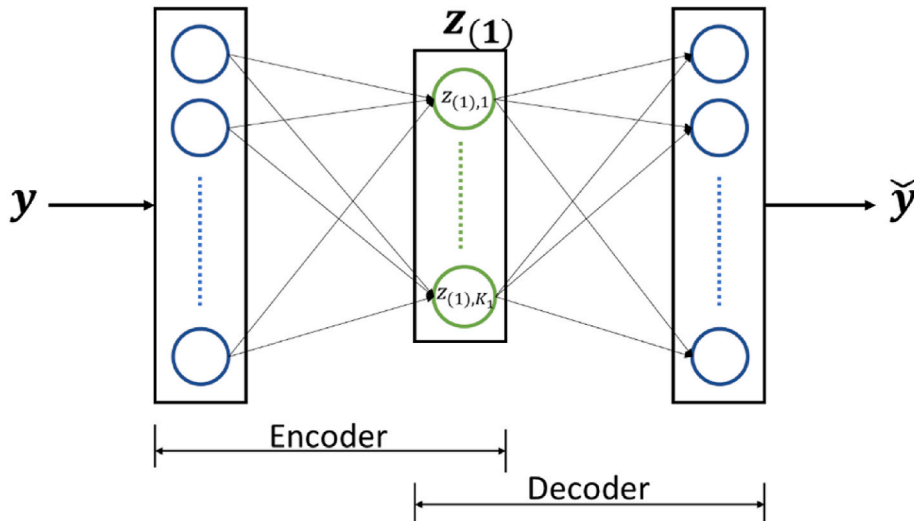


Fig. 2. Structure of a basic SAE.

monotonically as $\hat{\rho}_j$ diverges from ρ . During the training phase, the output value of the hidden neurons may be lowered by increasing the weight values $W_{(1)}$ making R_{sparse} to be small (Olshausen and Fieldt, 1997). To prevent it from happening, the R_{L2} term is added to the cost function:

$$R_{L2} = \frac{1}{2} \|\mathbf{W}\| \quad (12)$$

where \mathbf{W} is the SAE weight matrix.

Once the SSAE training is completed, its encoder is used to perform dimensionality reduction. The encoder transforms a generic time series $y^{(i)}$ into $z_{(L)}^{(i)}$, that is a K_L -dimensional vector, where K_L is the number of neurons of the innermost layer.

3.1.2. Stacked Sparse Autoencoders

Training non-linear autoencoders with multiple hidden layers is a difficult task (Holden et al., 2006). To tackle this problem, Hinton and Salakhutdinov (Holden et al., 2006) proposed the breakthrough approach adopted in this work, which consists of a *pre-training phase* and a *fine-tuning phase*. Let us consider an L -hidden-layer SAE (Fig. 3):

1. The pre-training regards a consecutive training of L (basic) SAEs. Initially, the first basic SAE is trained using the input vectors $y^{(i)} \in \mathbf{Y}$; then, the corresponding extracted features $z_{(1)}^{(i)}$ are used as training input vectors for the next basic SAE, which transforms $z_{(1)}^{(i)}$ into $z_{(2)}^{(i)}$ (Fig. 3a). This step is carried out up to the last SAE training.
2. Then, the SSAE is built by stacking all the basic SAE (Yang et al., 2018) (Fig. 3b).
3. In the fine-tuning phase, the SSAE obtained from the pre-training phase is fine-tuned using the backpropagation of error derivatives (Rumelhart et al., 1986).

3.1.3. SSAE performance assessment

Given a set of DOE points $\boldsymbol{\theta} = [\boldsymbol{\theta}^{(1)}, \dots, \boldsymbol{\theta}^{(m)}]$, typically built by a Latin Hypercube Sampling (LHS) (McKay et al., 2000) according to the ranges of the prior distributions, $\boldsymbol{\theta}$ are simulated through the BE model and the results collected into the $p \times m$ data matrix $\mathbf{Y} = [y^{(1)}, \dots, y^{(m)}]$. The m simulated time series collected in \mathbf{Y} are split into two groups: \mathbf{Y}^{TRAIN} and \mathbf{Y}^{TEST} . The training set contains m_{train} time series stored in the $p \times m_{train}$ data matrix \mathbf{Y}^{TRAIN} , whereas the test set contains $(m - m_{train})$ time series stored in the $p \times (m - m_{train})$ data matrix \mathbf{Y}^{TEST} . The training of the SSAE starts by defining its architecture (i.e., the number of hidden layers L and the number of neurons per each layer, that is, K_1, \dots, K_L) and setting the hyperparameters values (i.e., ρ , β , λ); then, L basic SAEs are trained, stacked and fine-tuned according to the procedure described in Section 3.1.2, and the fine-tuned SSAE is obtained. The SSAE capability in reconstructing time series that differ from the training one is tested on \mathbf{Y}^{TEST} through the reconstruction error:

$$R_{error} = \frac{1}{m - m_{train}} \sum_{i=1}^{m - m_{train}} \left\| y_{test}^{(i)} - \hat{y}_{test}^{(i)} \right\|^2 \quad (13)$$

Once the SSAE training is completed, K_L independent Kriging metamodels are built (one for each component of $z_{(L)}$) using the m inputs $\boldsymbol{\theta} = [\boldsymbol{\theta}^{(1)}, \dots, \boldsymbol{\theta}^{(m)}]$ and the respective transformed time series $\mathbf{Z}_{(L)} = [z_{(L)}^{(1)}, \dots, z_{(L)}^{(m)}]$.

On the other hand, in order to assess the predictive accuracy of each of the $p^* = K_L$ Kriging metamodels, the normalized Leave One Out Cross Validation (LOOCV) error is computed as follows (refer to the Appendix for further details):

$$\varepsilon_{LOOCVj} = \frac{\frac{1}{m} \sum_{i=1}^m \left(z_j^{(i)} - \hat{z}_{(-i)}^{MM}(z_j^{(i)}) \right)^2}{\frac{1}{m} \sum_{i=1}^m \left(z_j^{(i)} - \bar{z}_j \right)^2} \quad (14)$$

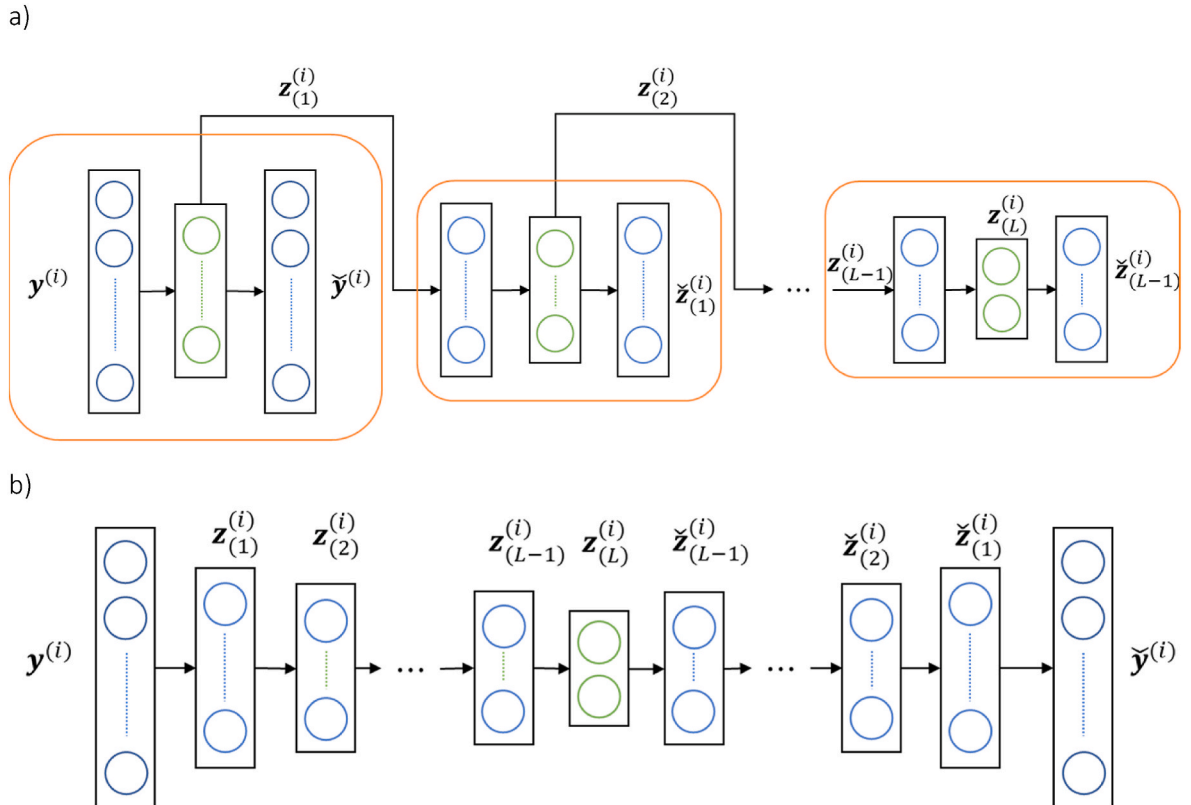


Fig. 3. Steps of the pre-training phase for an L -hidden-layer SAE: (a) training of L basic SAEs; (b) stacking of the basic SAEs.

where m is the Kriging training dataset size, $j = 1, \dots, p^*$, $\hat{z}_{(-i)}^{MM}(\theta^{(i)})$ is the Kriging metamodel obtained using all the points of θ , except $\theta^{(i)}$, $\bar{z}_j = \frac{1}{m} \sum_{i=1}^m z_j^{(i)}$. It should be noticed that in (14), the LOOCV error (i.e., $\frac{1}{m} \sum_{i=1}^m (z_j^{(i)} - \hat{z}_{(-i)}^{MM}(\theta^{(i)}))^2$) is normalized with respect to the training output sample variance (i.e., $\frac{1}{m} \sum_{i=1}^m (z_j^{(i)} - \bar{z}_j)^2$).

The setting of the SSAE hyperparameters (i.e., ρ, β, λ) is performed by trial-and-error, considering the values of ε_{LOOCV} for each metamodel and R_{error} . When a SSAE is trained, the corresponding ε_{LOOCV} errors are evaluated and, if acceptable, the SSAE is retained; otherwise, different hyperparameters are set and the entire procedure is repeated. There is no predefined threshold below which R_{error} and ε_{LOOCV} are considered acceptable; further details about this last point are given in Section 5.

3.2. Bayesian inference (by MCMC sampling)

In line with (Wu et al., 2018c), also in this work, the IUQ is carried out in a reduced space: y^E is mapped to the feature space (to obtain z^E). The distribution of $p(\theta)$, typically set on the basis of expert judgement when scarce information is available for θ , is taken with as large as possible prior ranges, letting the data speak for themselves. The posterior $p(\theta|z^E)$ is proportional to $p(\theta)$ multiplied by the likelihood $p(z^E|\theta)$:

$$p(\theta|z^E) \propto p(\theta)p(z^E|\theta) \quad (15)$$

The challenging objective addressed in this Section is to formulate the likelihood $p(z^E|\theta)$ in the case of non-linear dimensionality reduction of the output (e.g., by SSAE). In this work, we derive an expression for the likelihood $p(z^E|\theta)$ by propagating $y^E \sim N(y(\theta), I\sigma_{exp}^2)$ through the SSAE encoder by means of an Extended Kalman Filter (EKF) (Welch and Bishop, 2006). The expression of $p(z^E|\theta)$, in the case of a single independent experimental measurement (i.e., $N_{exp} = 1$), results:

$$p\left(z^E|\theta\right) = \frac{1}{(\sqrt{2\pi})^{p^*} \sqrt{|\Sigma_{exp}(\theta) + \Sigma_{Kriging}(\theta)|}} \exp\left[-\frac{1}{2}\left[z^E - \hat{z}_{(L)}^{MM}(\theta)\right]^T \left(\Sigma_{exp}(\theta) + \Sigma_{Kriging}(\theta)\right)^{-1} \left[z^E - \hat{z}_{(L)}^{MM}(\theta)\right]\right] \quad (16)$$

where $\hat{z}_{(L)}^{MM}(\theta)$ is the Kriging prediction in the reduced space; $\Sigma_{exp}(\theta)$ is the experimental uncertainty covariance matrix in the p^* -dimensional reduced space; and $\Sigma_{Kriging}(\theta)$ is the covariance matrix of the Kriging prediction uncertainty, i.e., a $p^* \times p^*$ matrix having the mean squared errors of each feature prediction on the diagonal entries:

$$\Sigma_{Kriging} = \begin{bmatrix} \sigma_{z_1(\theta)}^2 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \sigma_{z_{p^*}(\theta)}^2 \end{bmatrix} \quad (17)$$

Details on the hypothesis adopted for the derivation of $p(z^E|\theta)$ are given in the Appendix. Once the likelihood is formulated, a MCMC algorithm can be adopted to sample from $p(\theta|z^E)$ and the samples are used to find the PDF.

4. Case study

We show the application of the IUQ methodology to a TH RELAP5-3D model (Idaho National Laboratory, 2015) developed by Politecnico di Torino (Bersano et al., 2020) for the PERSEO test facility (sketched in Fig. 4) (Mascari et al., 2019). For more details about the facility, please refer to (Mascari et al., 2019; Ferri et al., 2005).

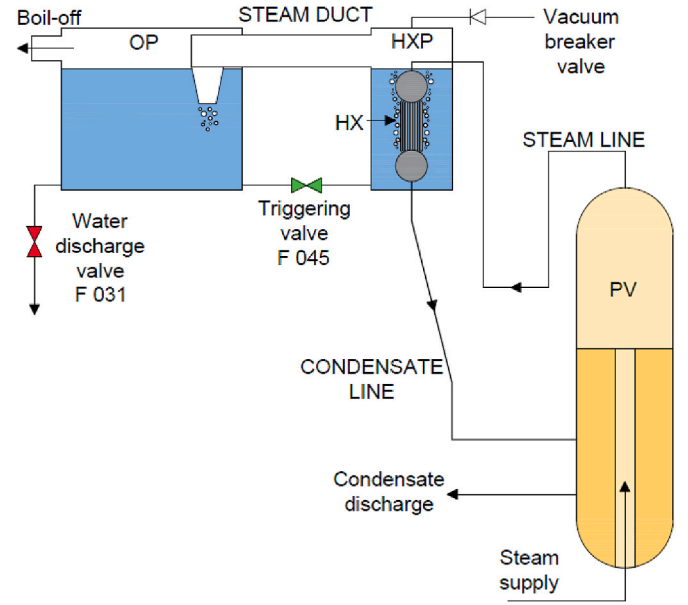


Fig. 4. Scheme of the PERSEO facility (Bersano et al., 2020).

Among the nine tests of the PERSEO experimental campaign (Mascari et al., 2019; Ferri et al., 2005), experimental data from the second part of Test 7 are here used for IUQ. According to the classification of input parameters made in Section 2, the RELAP5-3D code contains different model input parameters that could be included in the IUQ process as calibration variables (i.e., θ). However, in the present analysis, only some of them are considered: in fact, the objective of the present study is to show how the IUQ could be performed by taking advantage of the adopted methodology instead of carrying out a com-

plete uncertainty analysis. A uniform prior distribution is set for each parameter: $p(\theta) = \frac{1}{U-L}$ for $\theta \in [L, U]$ and $p(\theta) = 0$ otherwise, where U and L are the upper and lower bounds, respectively (as reported in Table 1). It is worth mentioning that the values reported in Table 1 are rescaled factors (i.e., all the parameters have been normalized with respect to their prior nominal values). For more details about the description of the

Table 1

The RELAP5-3D model input parameters selected for the IUQ.

θ_i	Parameter (multiplication factor)	Parameter Name	Lower bound	Upper bound
θ_1	Inner fouling factor	Inner FF	0.5	1.5
θ_2	Outer fouling factor	Outer FF	1.0	1.5
θ_3	Injector K factor	K injector	0.5	1.5
θ_4	Sum of the steam line's K factors	K sum steam	0.5	1.5
θ_5	Sum of condensate line's K factors	K sum condensate	0.5	1.5
θ_6	Diaphragm K factor	K diaphragm	0.5	1.5
θ_7	Rockwool thermal conductivity	K rockwool	1.0	1.5
θ_8	HXP first pipe flow area	A effective	0.5	1.5

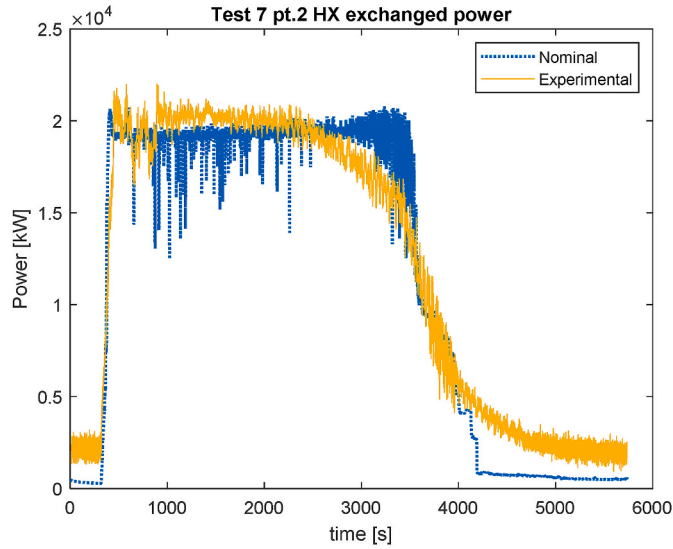


Fig. 5. HX exchanged power (Test 7-Part 2). The experimental measurement, in yellow solid line; the RELAP5-3D simulation for the nominal values, in blue dotted line.

parameters and the selection of the prior ranges, please refer to (Roma et al., 2021).

The power exchanged by the HX is used as the experimental data to perform the IUQ, because it carries significant information regarding the transient; thus, it is one of the most representative output. Fig. 5 compares the RELAP5-3D HX exchanged power computed for the nominal values of the prior distribution and the corresponding experimental data.

The dependence on x of the forward model $y^M(x, \theta)$ is absorbed into the definition of $y^M = y^M(\theta)$ that represents the RELAP5-3D model for the only experiment (i.e., Test 7 part 2) considered. Such a model predicts the vector $y^M = [y^M(\theta; t_1), \dots, y^M(\theta; t_p)]$, whose entries $y^M(\theta; t)$ are the power exchanged during Test 7-Part 2 at times t . Each RELAP5-3D simulation employs an Intel Core i7-7500U processor and takes around 2 h. Simulations are carried out at the Energy Department of Politecnico di Torino by faculty members.

4.1. Data description

Let $y^E = [y_1^E, \dots, y_p^E]^T$ be the available experimental HX exchanged power time series measured during Test 7-Part 2. We adopt Latin Hypercube Sampling (LHS) to build the DOE $\theta = [\theta^{(1)}, \dots, \theta^{(m)}]$, that contains $m = 180$ training inputs, in line with the distributions described in Table 1. The corresponding RELAP5-3D output realizations $y^{(i)}$ are, then, stored into the $p \times m$ training output data matrix $Y = [y^{(1)}, \dots, y^{(m)}]$ with $p = 5723$. Fig. 6 shows that the ensemble of the $m = 180$ RELAP5-3D output realizations (i.e., time series) contained in Y are affected by oscillations with a small signal-to-noise ratio. We apply the novel IUQ approach, based on SSAEs, on the dataset described here and compare the results with the standard approach based on PCA.

In this work, the design variables vector x (that defines the conditions under which the PERSEO experiment is carried out) is fixed since $N_{exp} = 1$; for this reason, the trends of the output training data reported in Fig. 6 look very similar. If $N_{exp} > 1$, the trends of the training output could be very different; consequently, the number of features required to map the outputs in the reduced space may increase, whereas the Kriging performances associated with each feature should be assessed case by case, but this is out of the scope of the present paper.

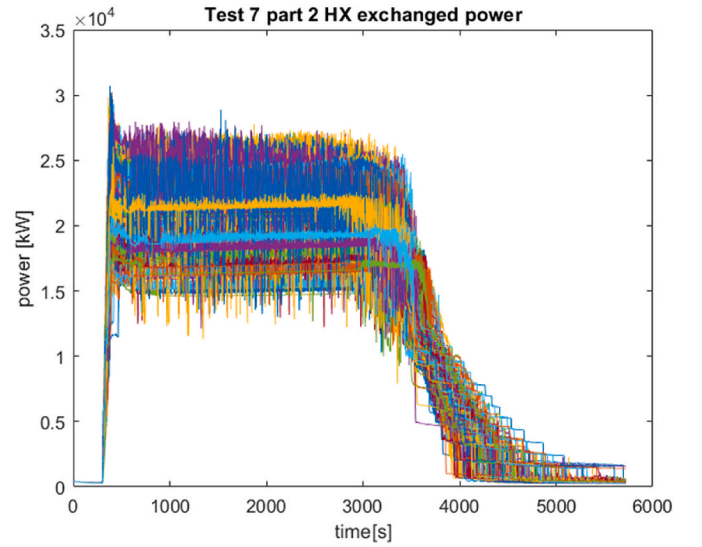


Fig. 6. The ensemble of the $m = 180$ RELAP5-3D output realizations contained in Y , characterized by a small signal-to-noise ratio.

5. Results

To analyze the effectiveness of the proposed SSAE approach, in this Section we compare the SSAE results with those obtained by a standard PCA-based dimensionality reduction approach of literature (Wilkinson, 2010; Higdon et al., 2008, 2013; Wu et al., 2018c; Nagel et al., 2020; Roma et al., 2021). In particular, Section 5.1 shows 1) how PCA is not capable of dealing with non-filtered (raw) data and 2) proposes some reflections on the main limitations involved by data filtering as a possible solution to deal with such raw data. Moreover, to assess the capability of the SSAE to deal with raw data, in Sections 5.2.1 and 5.2.2 the SSAE is fed with non-filtered (raw) and filtered time series, respectively. Section 5.3 shows the results of the forward uncertainty propagation performed through both the metamodel and the RELAP5-3D model.

5.1. PCA-based dimensionality reduction

Let $\theta = [\theta^{(1)}, \dots, \theta^{(m)}]$ and $Y = [y^{(1)}, \dots, y^{(m)}]$ be the $m = 180$ Design of experiment (DOE) points and the corresponding RELAP5-3D model outputs reported in Section 4.1, respectively. To perform the PCA, the $p \times m$ data matrix Y is centered, obtaining $Y_{centered}$; then, the Singular Value Decomposition (SVD) is carried out for $Y_{centered}$. More detail about the PCA are provided in (Roma et al., 2021). A number $p^* = 31$ of PCs gives a cumulative percentage of variation explained at least to 95%. Then, p^* independent Kriging metamodels are built, one for each feature (i.e., for each PC), employing the m input-output training patterns given

$$\text{by } \theta = [\theta^{(1)}, \dots, \theta^{(m)}] \text{ and } Z = \Phi(Y - \bar{\mu}_Y) = \begin{bmatrix} z_1^{(1)} & \dots & z_1^{(m)} \\ \vdots & \ddots & \vdots \\ z_{p^*}^{(1)} & \dots & z_{p^*}^{(m)} \end{bmatrix}, \text{ where}$$

$\bar{\mu}_Y = \frac{1}{m} \sum_{i=1}^m y^{(i)}$ and Φ is the PCA transformation matrix. Each Kriging metamodel is trained to map from $\theta \subset \mathbb{R}^d$ to $z_j \in \mathbb{R}$ (with $j = 1, \dots, p^*$). To evaluate the capability of PCA of reconstructing the noisy transients, we transform back Z into the original space

$$\hat{Y} = \bar{\mu}_Y + \Phi^T Z \quad (18)$$

and compute the reconstruction error R_{error} through Eq. (9). We find that $R_{error} = 2.2647 \times 10^5 \text{ kW}$.

In the current research, Matérn 5/2 correlation kernel, constant

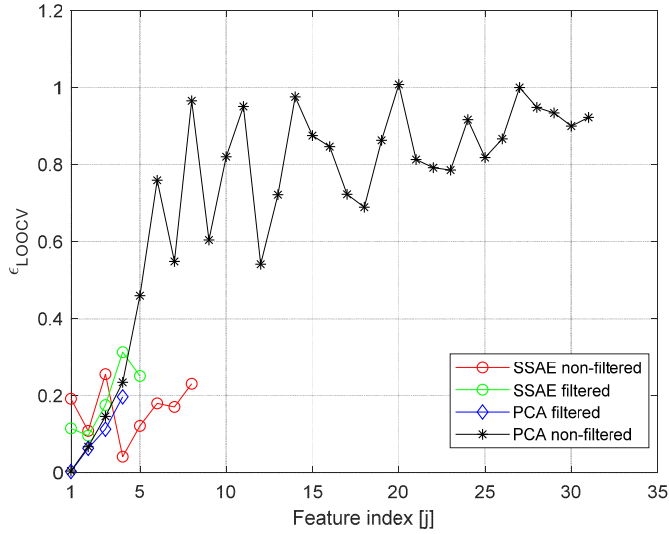


Fig. 7. ε_{LOOCV} in the case of PCA with non-filtered data, (2) PCA with filtered data (3) SSAAE with non-filtered data and (4) SSAAE with filtered data.

trend function and CV hyperparameter estimation are adopted for each Kriging metamodel, as done in (Roma et al., 2021) where data filtering is applied on the same training dataset, before PCA. Refer to (Roma et al., 2021) for further details on Kriging metamodeling. In order to assess the predictive accuracy of each of the $p^* = 31$ Kriging metamodels, the normalized Leave One Out Cross Validation (LOOCV) error ε_{LOOCV} is computed and reported in Fig. 7: ε_{LOOCV} increases for PCs of higher-order, which means that the metamodel prediction capability decreases for higher-order PCs. Likely, the first PCs are those associated with the genuine signal variation (due to physical phenomena); on the contrary, PCs of higher-order relate to the higher noise in comparison with the experimental data of the HX exchanged power predicted by RELAP5-3D model that is significant for the case study proposed. Consequently, training a metamodel for the first PCs is easier since there is an underlying function (given by the physics of the phenomena) that links the inputs of θ and the transformed output of Z .

According to its definition, an ε_{LOOCV} close to 1.0 relates to scarce metamodel predictive capability; indeed, it implies that the LOOCV error of the j^{th} Kriging metamodel

$$\frac{1}{m} \sum_{i=1}^m \left(z_j^{(i)} - \hat{z}_{(-i)}^{MM}(\theta^{(i)}) \right)^2 \quad (19)$$

is of the same order of magnitude of the transformed output sample variance $\text{var}[z_j]$:

$$\text{var}[z_j] = \frac{1}{m} \sum_{i=1}^m \left(z_j^{(i)} - \bar{z}_j \right)^2 \quad (20)$$

Because of the scarce performances that characterize most of the Kriging metamodels (for many of them $\varepsilon_{LOOCV} \cong 1.0$ (Fig. 7)), they cannot be used to replace the RELAP5-3D model in the IUQ process. This shows that in this context, PCA cannot cope with raw (non-filtered) data; thus, this analysis is not carried out.

Given the limitations of the PCA to cope with raw data (i.e., ε_{LOOCV} close to 1.0 for most of the features), a possible solution is to filter the HX power exchanged predicted by RELAP5-3D (noisy) raw time series collected in Y before applying PCA and Kriging metamodels. Although this allows removing the part of output variability due to oscillations, data filtering gives rise to a nontrivial issue, i.e., choosing a proper filtering technique (that should consider the smallest timescale on which physical phenomena take place during the transient). Moreover, even

though filtering (raw) data affected by numerical oscillations is common practice, (1) defending the choice of a particular filtering technique rather than another is nontrivial, and (2) the IUQ results may be affected by the selection of the specific filtering method adopted. In Ref. (Roma et al., 2021), we apply a moving median filter (i.e., the MATLAB function) on the columns of the same dataset Y before performing PCA and Kriging metamodeling, and finding that a smaller number of PCs (i.e., $p^* = 4$, rather than 31) is able to explain the same percentage of the variance of the dataset (i.e., 95%). These results are found in (Roma et al., 2021) and compared, in Section 5.2, to those obtained for the SSAE.

5.2. SSAE-based dimensionality reduction

5.2.1. Non-filtered (raw) data

In this Section, we take advantage of the SSAE properties to perform dimensionality reduction without applying data filtering. The SSAE architecture and the pre-training hyperparameters are set according to Table 2. The hyper-parameters K_1, K_2, K_3 and L are tuned following a trial-and-error approach based on the SSAE performances (i.e., R_{error} and the ε_{LOOCV} for each feature) (given that $L = 3; p^* = K_3$). More powerful tuning approaches (e.g., extensive grid search and evolutionary optimization) can be used, but at a much higher computational cost.

Unlike PCA, there is no rule of thumb to define p^* , i.e., K_3 , which is here set to $K_3 = 8$. The $m = 180$ simulated time series collected in Y are split into two groups: Y^{TRAIN} , that contains $m_{train} = 162$ time series, and Y^{TEST} that contains $m_{test} = m - m_{train} = 18$ time series (i.e., 10% of the available data). $L = 3$ basic SAEs are pre-trained using Y^{TRAIN} , with the hyperparameter reported in Table 2; then, they are stacked and fine-tuned according to the procedure described in Section 3.1.2. The fine-tuned SSAE thereby obtained is tested on Y^{TEST} obtaining a reconstruction error equal to $R_{error} = 8.4935 \times 10^5 kW$. A separate independent Kriging is built for each feature and ε_{LOOCV} is computed for each of them through (14). Therefore, $K_3 = 8$ independent Kriging metamodels are built adopting the Matérn 5/2 correlation kernel, constant trend function and CV hyperparameter estimation. In this case, R_{error} is computed on Y^{TEST} , whereas in the case of PCA R_{error} is computed on Y . We can notice that, even though R_{error} is higher, it has the same order of magnitude of that obtained by PCA with filtered and raw data (Table 6). Fig. 7 shows that the SSAE, without any user-experience-based data filtering approach, obtains ε_{LOOCV} values that are comparable (i.e., between 0.040 and 0.260) to those of PCA applied to filtered data.

Once the SSAE and the Kriging metamodels are trained, the algorithm recalled in Section 3.2 (and presented in the Appendix) is implemented to carry out Bayesian inference. The measurement error standard deviation σ_{exp} is set to 500 kW (Ferri et al., 2005) and an adaptive Metropolis algorithm is employed to produce 8 parallel chains with $1 \cdot 10^5$ iterations. It took nearly 16 hours to calculate the posterior using an Intel Core i7-7500U processor. According to (Gelman et al., 2015), we post-process the samples by discarding, for each chain, the first half for burn-in to diminish the influence of the starting samples, as a conservative choice. The MCMC convergence is examined through the approach proposed in (Gelman et al., 2015). The KDE of the posterior marginals PDFs obtained through the SSAE with non-filtered data are displayed in Fig. 8. Some summary statistics of the posterior distribution are reported in Table 3, whereas Table 4 shows the correlation among

Table 2
SSAAE architecture and hyperparameters (non-filtered data).

Architecture	Hyperparameters (SAE pre-training)		
L	3	ρ	0.05
K_1	200	β	1
K_2	20	λ	0.001
K_3	8		

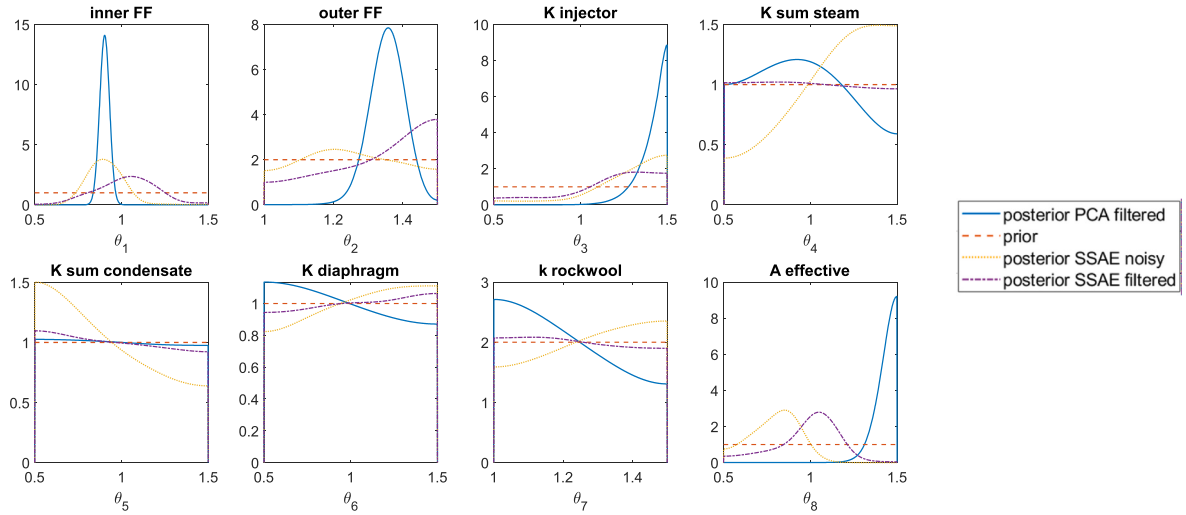


Fig. 8. Prior distributions and posteriors' kernel density estimations obtained adopting three different output dimensionality reduction techniques (i.e., PCA with filtered data and SSAE with both filtered and non-filtered data) in the IUQ process.

Table 3

Posterior summaries (SSAE with non-filtered data).

θ_i	Parameter	Mean value	Mode	5th percentile	95th percentile
θ_1	Inner FF	0.91	0.84	0.76	1.06
θ_2	Outer FF	1.25	1.22	1.04	1.46
θ_3	K injector	1.25	0.97	0.74	1.48
θ_4	K sum steam	1.13	1.49	0.64	1.47
θ_5	K sum condensate	0.91	0.61	0.53	1.42
θ_6	K diaphragm	1.03	1.35	0.57	1.46
θ_7	K rockwool	1.27	1.46	1.03	1.48
θ_8	A effective	0.80	0.84	0.57	0.98

Table 4

Correlation Matrix computed using the MCMC samples (SSAE with non-filtered data).

	θ_1	θ_2	θ_3	θ_4	θ_5	θ_6	θ_7	θ_8
θ_1	1,0000	-	-	-	-	-	-	-
θ_2	0,0758	1,0000	-	-	-	-	-	-
θ_3	-0,0815	-0,4110	1,0000	-	-	-	-	-
θ_4	-0,1261	-0,0720	-0,1443	1,0000	-	-	-	-
θ_5	0,0006	-0,0943	0,0772	-0,1062	1,0000	-	-	-
θ_6	-0,0723	-0,0788	-0,0279	0,0156	0,0249	1,0000	-	-
θ_7	-0,0731	0,0623	-0,1287	-0,0523	-0,0490	-0,0092	1,0000	-
θ_8	0,2947	-0,4725	-0,0484	0,1701	-0,0192	0,0345	-0,0055	1,0000

the calibration parameters.

5.2.2. Filtered data

For comparison purposes, we feed the SSAE with filtered data. The SSAE architecture is set according to the parameters reported in Table 5. Also in this case, a trial-and-error approach based on the SSAE performances (i.e., R_{error} and the ϵ_{LOOCV} for each feature) is adopted.

The same train-test split procedure of Section 5.2.1 is implemented on Y , and a reconstruction error equal to $R_{error} = 7.9071 \times 10^5 kW$ is obtained. $K_2 = 5$ independent Kriging metamodells are built adopting the Matérn 5/2 correlation kernel, constant trend function and CV hyperparameter estimation. Table 6 and Fig. 7 show, respectively, R_{error} and ϵ_{LOOCV} of the Kriging metamodells compared in the cases of: (1) PCA with non-filtered raw data, (2) PCA with filtered data, (3) SSAE with non-filtered raw data and (4) SSAE with filtered data.

Notice that performing data filtering on Y before applying PCA allows reducing: (1) the number of PCs required to have the same cu-

mulative percentage of variation explained (i.e., $p^* = 4$, instead of 31); and (2) the highest value of ϵ_{LOOCV} (that is reduced below 0.20). In contrast, data filtering, for the SSAE architecture and hyperparameter setting proposed in Table 5, does not bring a significant improvement, with respect to the case of SSAE with non-filtered data in terms of R_{error} on test data and ϵ_{LOOCV} . In our case study, output training data have similar trends (as illustrated in Fig. 6). If, on the contrary, training time series would have been very different, the number of extracted features (i.e., the number of neurons in the hidden layer) required to get the same SSAE reconstruction error R_{error} could be higher (i.e., the ANN complexity increases). In this case, since we build an independent Kriging metamodel for each extracted feature (as shown Fig. 1), the number of Kriging increases and the efficiency of each metamodel should be re-assessed (since it can not be a priori predicted). More sophisticated approaches than the simple trial-and-error (e.g., evolutionary optimization techniques) can be implemented to find the set of SSAE's hyperparameters that minimize both R_{error} and ϵ_{LOOCV} . However,

Table 5
SSAE architecture and hyperparameters (filtered data).

Architecture		Hyperparameters (SAE pre-training)	
L	2	ρ	0.05
K_1	50	β	1
K_2	5	λ	0.0001

Table 6
 R_{error} in the case of: (1) PCA with non-filtered data, (2) PCA with filtered data (3) SSAE with non-filtered data and (4) SSAE with filtered data.

Reconstruction error R_{error}	
PCA non-filtered data	$2.2647 \times 10^5 kW$
PCA filtered data	$2.1308 \times 10^5 kW$
SSAE non-filtered data (computed for Y^{TEST})	$8.4935 \times 10^5 kW$
SSAE filtered data (computed for Y^{TEST})	$7.9071 \times 10^5 kW$

Table 7
Posterior summaries (SSAE with filtered data).

θ_i	Parameter	Mean value	Mode	5th percentile	95th percentile
θ_1	Inner FF	1.03	1.18	0.75	1.26
θ_2	Outer FF	1.32	1.46	1.05	1.49
θ_3	K injector	1.17	1.24	0.63	1.47
θ_4	K sum steam	1.00	1.16	0.55	1.45
θ_5	K sum condensate	0.98	0.61	0.54	1.44
θ_6	K diaphragm	1.01	1.08	0.55	1.45
θ_7	k rockwool	1.25	1.00	1.00	1.50
θ_8	A effective	0.98	1.14	0.63	1.20

this would be cumbersome considering that: (1) PCA is easier to implement and (2) at lower computational cost, it performs better in terms of R_{error} and ϵ_{LOOCV} with filtered data. Moreover, for PCA the Σ_{exp} matrix does not need to be computed at each MCMC iteration (unlike the SSAE case, as shown in the Appendix), and this significantly reduces the computational cost required by the algorithm. In fact, in the case of SSAE with filtered data, the adaptive Metropolis algorithm takes almost 10.5 hours for $2 \cdot 10^5$ iterations on an Intel Core i7-7500U and, considering that the MCMC algorithm developed for the PCA takes 2.5 hours for 10^5 iterations (on the same processor), the computational cost required by the MCMC algorithm in this case of SSAE is twice that required for PCA. Thus, we can conclude that using the SSAE is conveniently applied to raw data (as expected in real applications), avoiding filtering.

From each chain, the first half of the samples are discarded for burn-in. A common practice adopted to reduce autocorrelation is thinning (Gelman et al., 2015). It consists in keeping every k^{th} sample from each sequence and discarding the rest: we kept every 2000^{th} (out of $5 \cdot 10^5$), 4000^{th} (out of $1 \cdot 10^5$) and 8000^{th} (out of $2 \cdot 10^5$) sample from each chain for PCA with filtered data, SSAE with raw data and SSAE with filtered data, respectively. The KDE of the posterior marginal PDFs obtained through the SSAE with non-filtered data are displayed in Fig. 8. The

Table 8
Correlation Matrix computed using the MCMC samples (SSAE with filtered data).

	θ_1	θ_2	θ_3	θ_4	θ_5	θ_6	θ_7	θ_8
θ_1	1,0000	-	-	-	-	-	-	-
θ_2	0,3468	1,0000	-	-	-	-	-	-
θ_3	-0,1763	-0,1013	1,0000	-	-	-	-	-
θ_4	-0,0319	0,0062	0,0441	1,0000	-	-	-	-
θ_5	-0,0018	-0,0038	-0,0048	-0,0074	1,0000	-	-	-
θ_6	0,0384	0,1051	-0,0768	-0,0130	-0,0158	1,0000	-	-
θ_7	0,0319	-0,0184	0,0331	-0,0270	0,0101	-0,0010	1,0000	-
θ_8	0,4391	0,2914	-0,3216	0,0207	0,0446	0,0317	-0,0491	1,0000

statistics of the posterior PDF are summarized in Table 7. The marginal posterior distributions in Fig. 8 cannot be used to draw samples, because the calibration parameters are not independent (see Table 8).

In all the cases examined, the marginal posteriors of $\theta_4, \theta_5, \theta_6, \theta_7$ do not differ very much from their priors and are defined on the same supports of their priors, whereas for θ_1 and θ_8 a significant update can be seen with respect to the priors. This is in line with the results found in (Roma et al., 2021), where the sensitivity analysis, carried out through first-order Sobol' indices (Saltelli et al., 2010), revealed that $\theta_4, \theta_5, \theta_6, \theta_7$ are less influential. In particular, for the PCA-based approach, the marginal posteriors of $\theta_1, \theta_2, \theta_3, \theta_8$ show a considerable modification regarding their priors and the posterior of θ_8 is peaked at the upper bound of the prior, whereas in both SSAE cases the marginal posterior of θ_2 is quite similar to the prior, and the posterior of θ_8 is peaked near the prior mean value.

As a general result, the posterior PDFs obtained through the SSAE, with filtered and raw data, are wider than those obtained through PCA. In this view, the PCA-based approach, providing sharper posterior PDFs (i.e., characterized by smaller variances), seems to allow reducing epistemic uncertainty about calibration parameters more than the SSAE-based approach. However, it is difficult to assess and comment on the consistency of such posterior PDFs; in this regard, Section 5.3 proposes the propagation of such uncertainty to check the consistency of the results to experimental data.

5.3. Forward uncertainty quantification: comparison of SSAE and PCA

To compare and evaluate the relevance of the proposed approaches, we perform forward uncertainty propagation. In particular, (1) we feed the T-H model with the posterior samples obtained from the MCMC after the thinning and (2) we compare the ensemble of time series obtained to y^E ; this allows comparing the posterior PDFs obtained applying both PCA and SSAE. Both the RELAP5-3D model and the Kriging metamodels are used to propagate the posteriors' uncertainty. The Kriging-based forward uncertainty propagation is proposed in Section 5.3.1, whereas Section 5.3.3 shows the results obtained by propagating through the RELAP5-3D model.

5.3.1. Kriging metamodels

For each IUQ approach proposed, 100 posterior samples are simulated through the Kriging metamodels to obtain the predictions in the p^* -dimensional reduced space. Such predictions are, then, transformed into the p -dimensional space, adopting transformation matrix Φ in the case of PCA and the SSAE decoders in the other case. In particular, Fig. 9a, Fig. 9b and c compare the reconstructed experimental data y^E to the reconstructed Kriging outcomes of 100 posterior samples, the prior nominal value, the posterior mode and the posterior mean value.

In each case, the simulated posterior samples envelop y^E . Moreover, in the case of PCA, the posterior's prediction ensemble (Fig. 9a) shows smaller variance than that of the SSAE (Fig. 9b and c); this is in line with the narrower PDFs that characterize the posterior marginal KDE obtained applying PCA (Fig. 8). In general, the reconstructed Kriging predictions of the posterior samples obtained adopting PCA reproduce

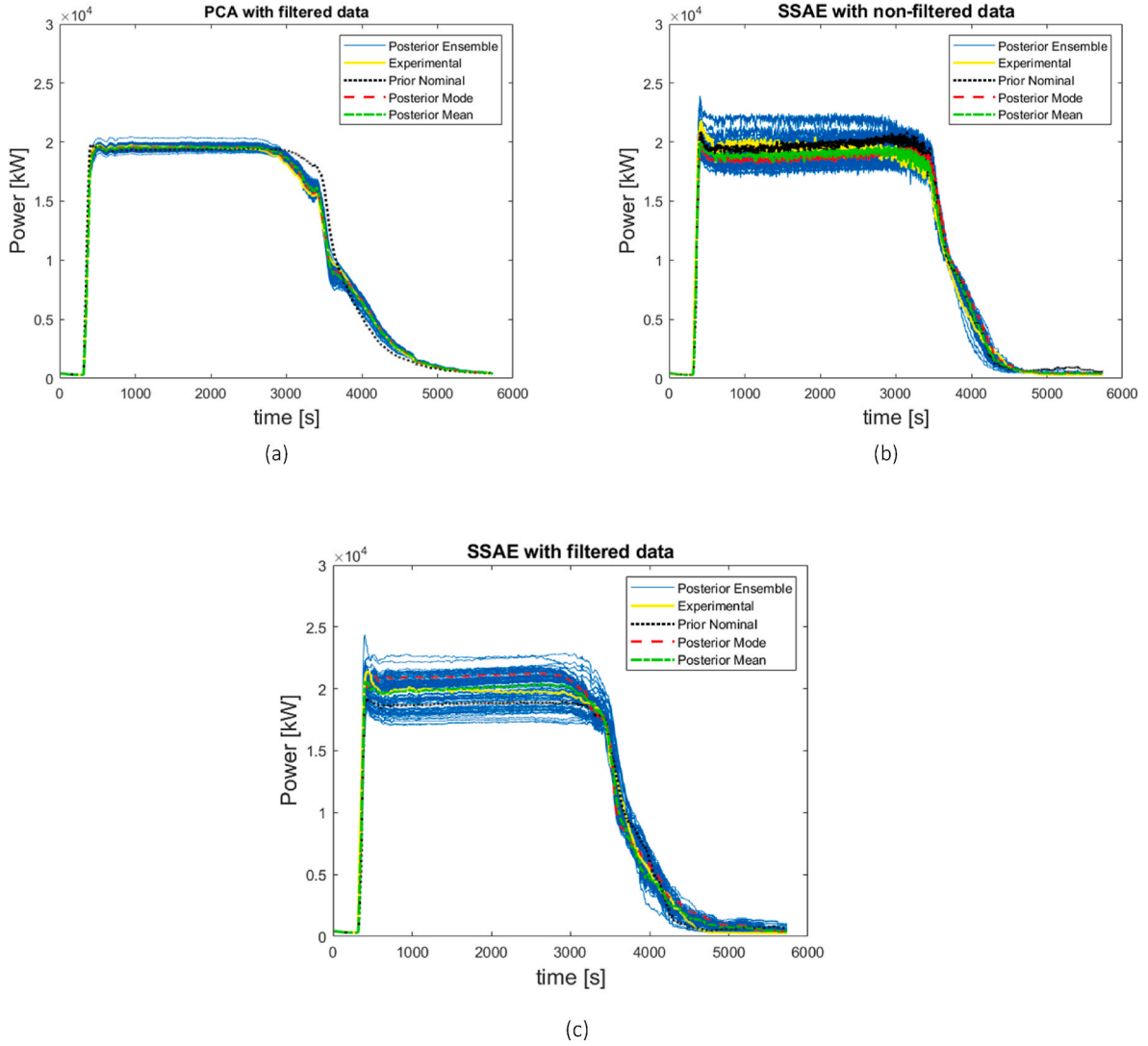


Fig. 9. Reconstructed Kriging outcomes of 100 posterior samples, the prior nominal value, the posterior mode and the posterior mean value, compared with the reconstructed experimental data. The results obtained applying PCA with filtered data are reported in Fig. 9a, whereas Fig. 9b and c report the results obtained applying the SSAE with non-filtered and filtered data, respectively.

y^E better than the reconstructed Kriging predictions of the posterior samples obtained adopting the SSAE with both filtered and raw data. Finally, we can notice that, when the SSAE is applied, data filtering does not bring a significant improvement in terms of agreement of the simulated posterior samples with y^E . It is worth mentioning that the low number of simulated posterior samples (i.e., 100) is due to the fact that, as explained in Section 5.2.2, thinning has been performed on the MCMC samples to reduce autocorrelation, reducing the effective number of MCMC posterior samples to 1000 and 100, in PCA and SSAE cases, respectively.

5.3.2. Safety margin calculation

In BEPU methodologies, the results are expressed in terms of uncertainty ranges for the calculated Figure of Merit (FOM); this allows to compute safety margins with respect to safety threshold values. The current BEPU methods can be subdivided in (Iaea, 2008): (1) probabilistic approaches (e.g., CSAU, GRS and ASTRUM), (2) deterministic methods (e.g., AEAW and EDF-Framatome), and (3) methods based on the extrapolation of output uncertainty (e.g., UMAE). Following the GRS method (Glaeser, 2008), Wilk's formula (Zio et al., 2010) can be used to

compute the number of code runs N that ensures the confidence level β and the probability content γ (Di Maio et al., 2016), and this allows computing margins to a safety threshold values. According to the Wilk's formula, the one-sided confidence level is given by the following expression:

$$1 - \gamma^N \geq \beta \quad (21)$$

This expression is valid for first-order Wilks' formula, that is the case in which the highest (lowest) outcome is inside the upper (lower) 5% range with at least 95% confidence.

Following (Roma et al., 2021), the HX exchanged energy over the mission time $T_{mission} = 5736$ s is selected as FOM for the current analysis:

Table 9

The lowest simulated values of HX exchanged energy and the respective margins with respect to the threshold value $0.9E_{nominal}$.

	E_{lowest}	M
PCA with filtered data	6.716×10^7 kJ	7.120×10^6 kJ
SSAE with non-filtered data	6.366×10^7 kJ	3.618×10^6 kJ
SSAE with filtered data	6.388×10^7 kJ	3.844×10^6 kJ

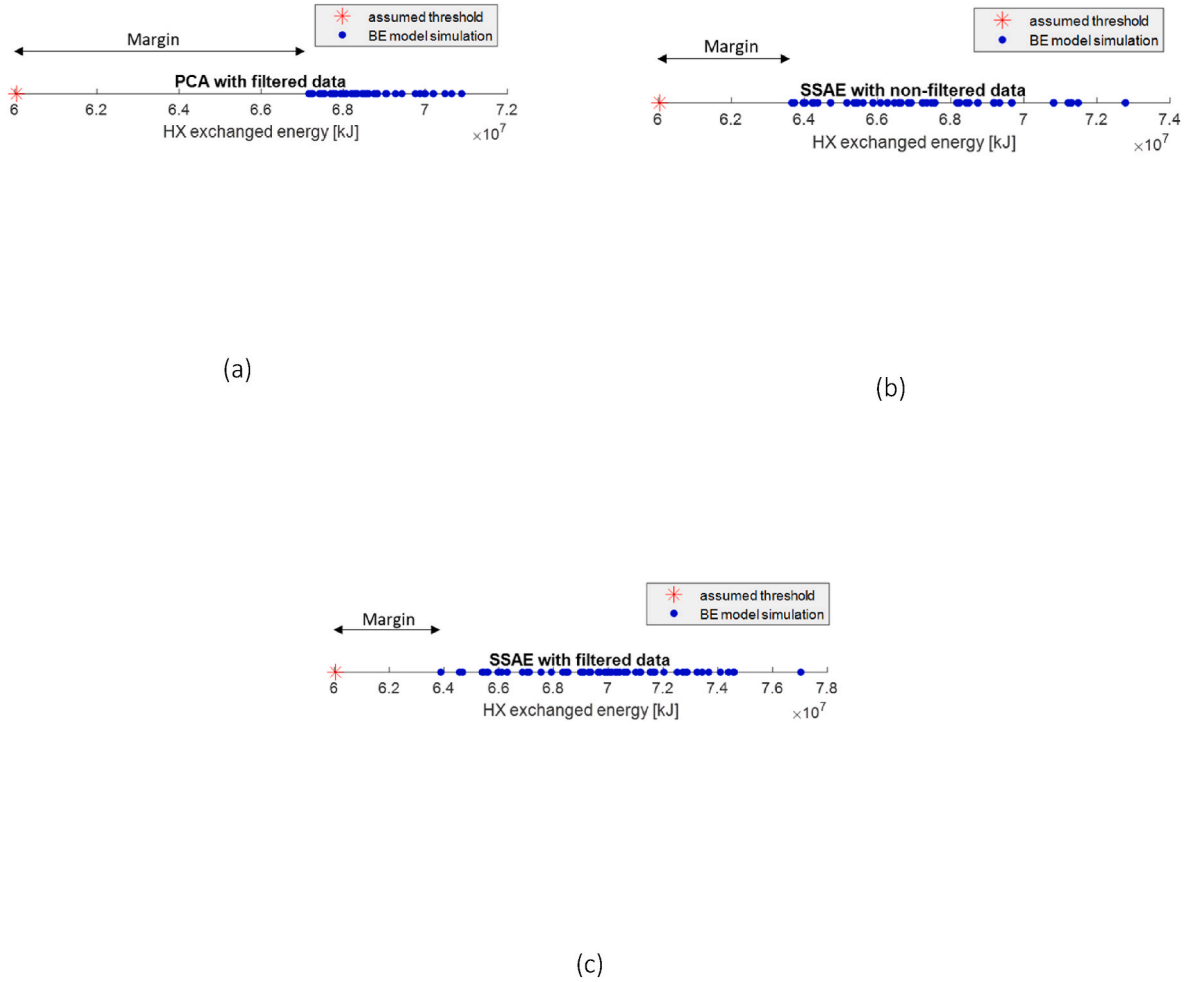


Fig. 10. Safety margins M in the case of PCA with filtered data (Fig. 10a), SSAE with non-filtered data (Fig. 10b) and SSAE with filtered data (Fig. 10c).

$$E = \int_0^{T_{mission}} P(t) dt \quad (22)$$

For demonstration purposes, in the three cases analyzed (i.e., PCA with filtered data and SSAE with both filtered and non-filtered data), the failure criterion for E is set to $0.9E_{nominal}$, i.e., the system fails if $E < E_{threshold} = 0.9 \int_0^{T_{mission}} P_{nominal}(t) dt = 6.004 \times 10^7 \text{ kJ}$, where $P_{nominal}$ is the HX exchanged power simulated through RELAP5-3D utilizing as input the prior nominal values. According to Wilks' formula (Zio et al., 2010), for each case, $N = 59$ RELAP5-3D simulations are carried out to calculate the one-sided statistical tolerance limit with $\beta = 95\%$ confidence level and $\gamma = 95\%$ probability content (see Eq. (21)). Among the N simulated values of E , according to Wilks' formula, the smallest one is contained by the lower 5% range with at least 95% confidence; thus, its margin is given by $M = E_{lowest} - 0.9E_{nominal}$. Table 9 reports the values of E_{lowest} and the respective margins M (with respect to the threshold value $0.9E_{nominal}$); the latter values are also graphically reported, along with the $N = 59$ computed values of E , in Fig. 10a, Fig. 10b and c. It is important to remark that the safety margin analysis performed in this Section is presented only for demonstration reasons and regards only epistemic uncertainty.

It can be noticed that the safety margins computed for the SSAE are lower than those obtained for the PCA with filtered data; this is a direct consequence of the larger variance that characterizes the posterior PDFs found through the SSAE-based approach of IUQ. In this regard, if aleatory uncertainty is also considered, the margin M can even be lower; thus, for a realistic estimate of, it is recommended to properly characterize such aleatory uncertainty and repeat the safety margin calculation.

5.3.3. RELAP5-3D

The good agreement observed for the PCA-based approach (applied to filtered data), between \hat{y}^{-E} and the simulated posterior samples in the Kriging-based forward uncertainty propagation should be examined for the RELAP5-3D BE model. The $N = 59$ available RELAP5-3D simulations, carried out to calculate the safety margins as in Section 5.3.2, are used to assess and compare the PCA and the SSAE from 1) a qualitative point of view and 2) a quantitative point of view, by computing the Signal to Noise Ratio. More precisely, Fig. 11a, Figs. 11b and c compare the Test 7-Part 2 experimental data with respect to the RELAP5-3D predictions of the following input parameters: $N = 59$ posterior samples, the prior nominal value, the posterior mode $\theta_{posterior}^{mode}$ and posterior

mean value $\theta_{posterior}^{mean}$ (these two latter obtained from the thousands of MCMC samples after burn-in).

One can notice that, when PCA is applied to filtered data, the RELAP5-3D simulated posterior mode $y^M(\theta_{posterior.PCA}^{mode})$ and mean value $y^M(\theta_{posterior.PCA}^{mean})$ (and, in general, all the $N = 59$ simulated posterior samples) display a much wider noise with respect to the case in which the SSAE is applied to non-filtered data. Since the oscillations affecting the RELAP5-3D simulated HX power exchanged are higher in comparison with the experimental data, their signal-to-noise ratio can be considered a characteristic to evaluate their physical consistency: the lower the signal-to-noise ratio, the lower the physical consistency. In this regard, for the PCA and the SSAE with/without filter results shown in Fig. 11, we compute the SNR (see Eq. (23)) of the $N = 59$ RELAP5-3D output time series: we model the (unknown) noise-free output $p \times N$ matrix S with a moving median filter (with a 50 s sliding window to

accommodate sudden power oscillations, i.e. ~ 1 MW) and we compute the $p \times N$ noise matrix ξ by subtracting S from the non-filtered output (i.e., $\xi_{ij} = y_j^{RELAP}(\theta_i) - S_{ij}$ with $i = 1, 2, \dots, N$ and $j = 1, 2, \dots, p$):

$$SNR = 20 \log_{10} \left(\frac{\sqrt{\sum_{i=1}^N \sum_{j=1}^p |S_{ij}|^2}}{\sqrt{\sum_{i=1}^N \sum_{j=1}^p |\xi_{ij}|^2}} \right) \quad (23)$$

For PCA and SSAE with/without filtering, the SNR values are equal to: $SNR_{PCA} = 23.04$ dB, $SNR_{SSAE, noisy} = 31.12$ dB and $SNR_{SSAE, filtered} = 26.15$ dB. Concerning this last point, the SSAE-based approach applied to raw data outperforms the PCA-based approach. Moreover, comparing Fig. 11b and c, we can notice that data-filtering does not provide any benefits both in terms of posterior mode accuracy (in predicting y^E) and

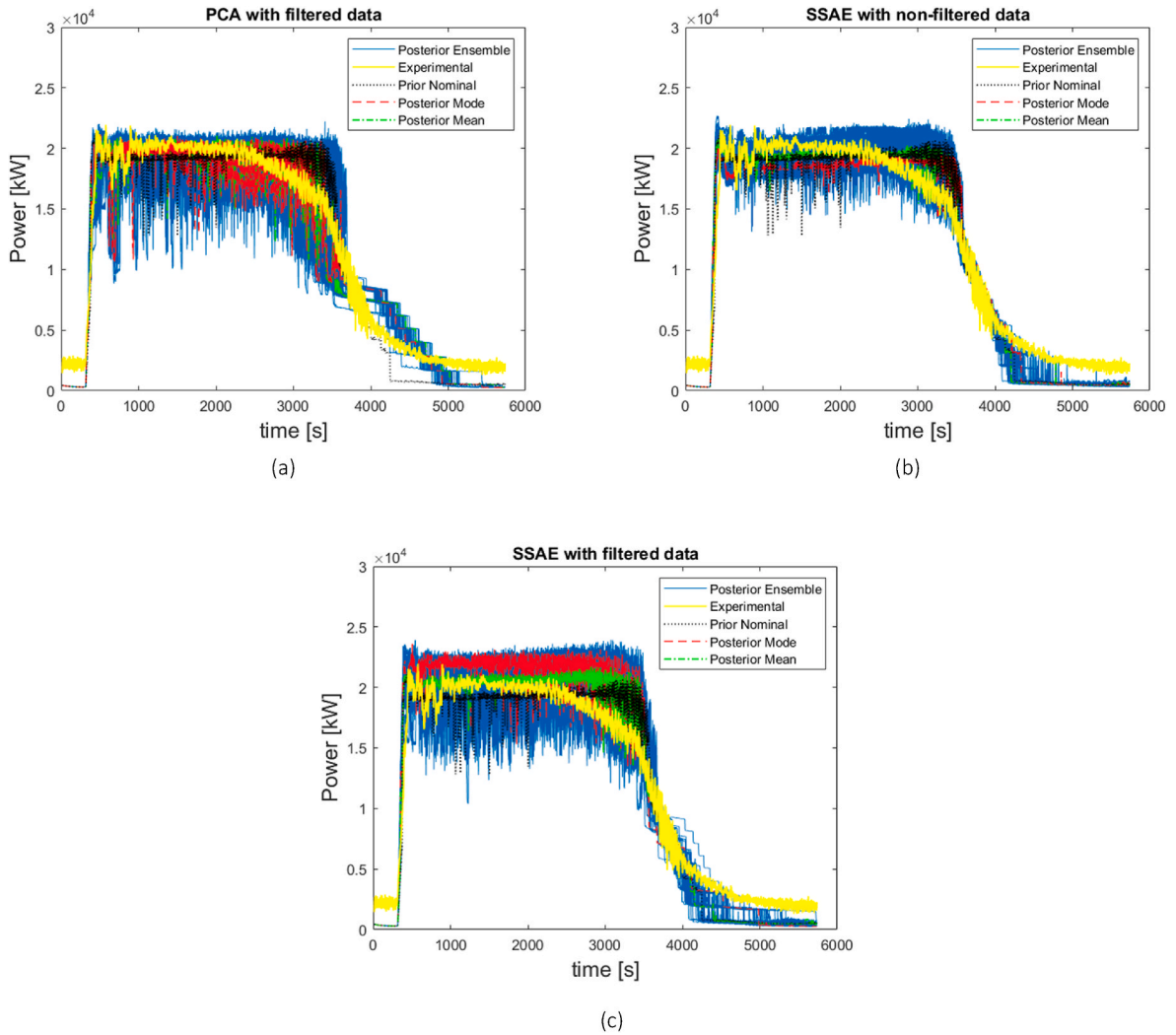


Fig. 11. RELAP5-3D outcomes of $N = 59$ posterior samples with respect to y^E and the RELAP5-3D outcomes of the prior nominal value, the posterior mode and posterior mean value. The results obtained applying PCA with filtered data are reported in Fig. 11a, whereas Figs. 11b and c report the results obtained applying the SSAE with non-filtered and filtered data, respectively.

signal-to-noise ratio (characterizing the $N = 59$ simulated posterior samples), when the SSAE-based approach is applied. A possible reason is that filtered data may bias autoencoders and Kriging training. This may impair the metamodel generalization capability when fed with noisy data. This finally suggests that filtering should be avoided in favour of the use of raw data, whose information carried is significative and can be capitalized by the autoencoder in SSAE-based IUQ approaches.

6. Conclusions

In this work, we propose a novel IUQ approach for T-H model code input parameters in case of noisy (i.e., small signal-to-noise ratio) time-dependent outputs, that paves the way for a novel dimensionality reduction method capable of dealing with raw data in the context of Bayesian IUQ. The approach is developed within a Bayesian framework and adopts a SSAE-based dimensionality reduction technique to extract significant features from the simulated time series and involves implementing Kriging metamodels for the quick emulation of such features.

The effectiveness of the proposed approach has been firstly demonstrated considering (noisy) raw data simulated by a time-dependent RELAP5-3D model the PERSEO facility in relation to the power exchanged by the HX. The results show *i*) the capability of the adopted SSAE to reduce the input data dimensionality while preserving its most significant characteristics and *ii*) the ability of the proposed IUQ approach in dealing with (noisy) raw data. Moreover, the comparison of the proposed approach with a standard dimensionality reduction method (i.e., PCA), exhibit the inability of the latter in dealing with (noisy) raw data. This highlights the novel characteristic of the SSAE-based approach, which, in contrast to the standard dimensionality reduction methods, allows going through the IUQ without resorting to filtering techniques, which are based on expert judgment and can affect the IUQ results. Moreover, the PDFs obtained applying the SSAE-based approach to raw data, when propagated through the RELAP5-3D code, give power exchanged by the HX time series characterized by a higher

signal-to-noise ratio than in the PCA-based approach. This can be considered an element to assess the physical consistency of the uncertainty propagated to the code output, which theoretically should not be affected by large noise/oscillations.

Also, the proposed approach has been applied to a filtered data set, still related with the RELAP5-3D model of the PERSEO facility. It can be concluded that performing data filtering before applying the SSAE does not bring any benefits in terms of metamodel accuracy and consistency of the propagated results with respect to experimental data. This is in line with the fact that SSAEs, being able to perform denoising, should not be affected by the noise. The comparison with the PCA-based approach (applied to filtered data) shows that *i*) PCA allows reducing epistemic uncertainty more than the SSAE-based approach since the former provides sharper posterior PDFs (i.e., characterized by minor variance) and *ii*) the MCMC sampling is computationally more expensive for SSAE than for PCA.

Future research lies on the possibility of exploring: 1) more powerful tuning approaches to optimize the SSAE architecture (e.g., extensive grid search and evolutionary optimization); 2) new approaches of uncertainty propagation through DNNs (e.g., the Monte Carlo sampling, the entire-DNN unscented transform and the piecewise exponential approximation of the transfer function) taking, also, into account their computational cost; 3) new approaches, such as multivariate Kriging metamodels (Kleijnen and Mehdad, 2014), to take into account dependencies among the output components. In fact, another limitation that should be further investigated is the effect of building a distinct independent metamodel for each feature extracted, i.e., assuming the features to be independent.

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Appendix

The challenging objective of this Section is to formulate the likelihood $p(z^E|\theta)$ in the case of non-linear dimensionality reduction of the output (e.g., SSAE). This problem can be tackled by propagating the uncertainty of $y^E \sim N(y(\theta), I\sigma_{exp}^2)$ through the SSAE's decoder. Monte Carlo (MC) sampling and the Unscented Transform (Abdelaziz et al., 2015) have been employed in past works to propagate the uncertainty through Deep Neural Networks; that is, given a multivariate distribution of the input layer (e.g., $y^E \sim N(y(\theta), I\sigma_{exp}^2)$), the mean vector $z_{(L)}$ and the covariance matrix $\Sigma_{(L)}$ of the output layer L are estimated (Abdelaziz et al., 2015; Hadjajmadi and Homayounpour, 2015); however, these methods (particularly the MC sampling) can be computationally expensive (Abdelaziz et al., 2015; Titensky et al., 2018). Following the approach proposed in (Titensky et al., 2018), in this work, we use Extended Kalman Filtering (EKF) (Welch and Bishop, 2006) to propagate $y^E \sim N(y(\theta), I\sigma_{exp}^2)$ through the SSAE encoder in order to derive an expression for the likelihood $p(z^E|\theta)$. EKF is an algorithm that is used to estimate the state of non-linear discrete-time dynamic systems when the system state cannot be directly measured. In the EKF, the system's state at time-step l is treated as an uncertain quantity characterized by a mean vector $z_{(l)}$ and a covariance matrix $\Sigma_{(l)}$. The EKF algorithm consists of two steps: the prediction step and the update step. In the prediction step the system's state $z_{(l)}$ is predicted along with its error covariance $\Sigma_{(l)}$ starting from (1) the process noise, (2) the control input and (3) the previous step's state. In the updating step, the prior estimates computed in the prediction step are updated (through indirect measurement of the system state) to find the posterior estimate of the state and its error covariance. For further details about Kalman filtering refer to (Welch and Bishop, 2006).

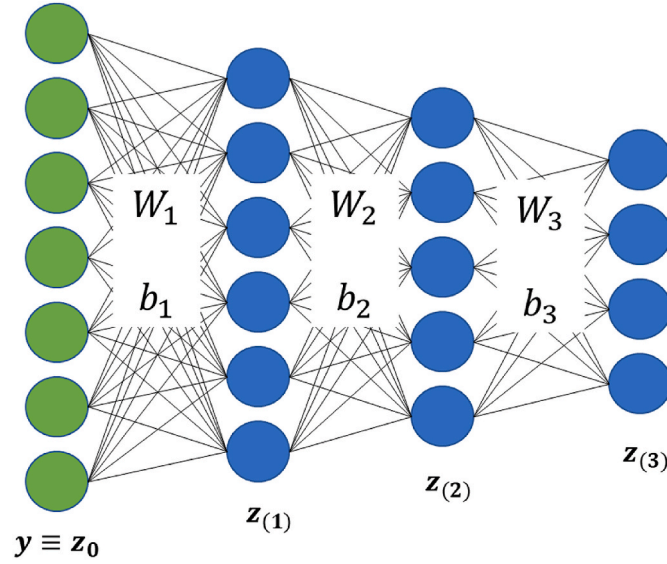


Fig. 12. Example of DNN where $L = 3$.

In the context of EKF applied for the uncertainty propagation in Deep Neural Networks, the concept of “system” does not refer to any real physical system; in fact, we treat the layers of the SSAE encoders as the states of a fictitious unforced (i.e., without a control input) system at different time-steps (i.e., the input layer of the SSAE $z_{(0)}$ represents the system state at time $l = 0$, the second layer $z_{(1)}$ represents the systems state at time $l = 1$ et cetera). Moreover, only the prediction step of the EKF algorithm is applied since there is no real physical system that allows measurements (used in the update step). Let us assume that the SSAE encoder is composed of $L + 1$ layers and that $z_{(l)}$ and $\Sigma_{(l)}$ are, respectively, the vector representing the state estimate (i.e., the mean value) and the covariance matrix of the layer l , such that:

$$\Sigma_{(l),j,k} = \text{cov}(z_{(l),j}, z_{(l),k}) \quad (\text{C.1})$$

The system evolves from the state (read: layer) $l - 1$ to the state l through a non-linear transformation:

$$z_{(l)} = \sigma(\mathbf{W}_{(l)}z_{(l-1)} + \mathbf{b}_{(l)}) \quad (\text{C.2})$$

where $\sigma(\bullet) = \frac{e^{\bullet}}{e^{\bullet} + 1}$ is the sigmoid transfer function of the SSAE, $\mathbf{W}_{(l)}$ is the SSAE weight matrix between layer $l - 1$ and layer l , $\mathbf{b}_{(l)}$ is the bias vector for layer l . Using the equation of the EKF's prediction step, we can find the state estimate $z_{(l)}$ and the covariance matrix $\Sigma_{(l)}$ for each layer:

$$z_{(l)} = \sigma(\mathbf{W}_{(l)}z_{(l-1)} + \mathbf{b}_{(l)}) \quad (\text{C.3})$$

$$\Sigma_{(l)} = \mathbf{F}_{(l)}\Sigma_{(l-1)}\mathbf{F}_{(l)} + \mathbf{Q}_{(l)} \quad (\text{C.4})$$

where $\mathbf{F}_l = \nabla_{z_{(l-1)}}z_l$ is the Jacobian matrix and \mathbf{Q}_l is the process noise covariance matrix that takes into account the inherent error introduced by the dimensionality reduction itself (i.e., the error related to the fact that the SSAE is not a perfect model). In this work, under the hypothesis that the SSAE is a perfect model, \mathbf{Q}_l is neglected. If a sigmoid transfer function is adopted, It can be shown that:

$$F_{(l),j,k} = \sigma\left(\sum_r W_{(l),j,r} \cdot z_{(l-1),r} + b_{(l),j}\right) \left(1 - \sigma\left(\sum_r W_{(l),j,r} \cdot z_{(l-1),r} + b_{(l),j}\right)\right) \cdot W_{(l),j,k} \quad (\text{C.5})$$

where $F_{(l),j,k}$ and $W_{(l),j,k}$ represent, respectively, the element at the j^{th} row and k^{th} column of the $\mathbf{F}_{(l)}$ and $\mathbf{W}_{(l)}$ matrices; $z_{(l-1),r}$ is the r^{th} entry of $z_{(l-1)}$. Iteratively applying equations (C.3) and (C.4), one can propagate $y^E \sim N(y(\theta), I\sigma_{\text{exp}}^2)$ through the SSAE encoder in order to derive the distribution of z^E that has mean value $z_{(L)}$ and the covariance matrix $\Sigma_{(L)}$. It should be noted that, since the sigmoid function is non-linear, the Jacobian matrices $\mathbf{F}_{(l)}$, and, in turn $\Sigma_{(L)}$, both depend on the input vector $z_{(0)} \equiv y(\theta)$ that is a priori unknown since the Kriging metamodels predict directly $\hat{z}_{(L)}^{MM}(\theta)$. To address this problem, we reconstruct $\hat{z}_{(L)}^{MM}(\theta)$ (through the SSAE decoder) obtaining $y^K(\theta)$ that is eventually used to derive $\Sigma_{(L)}$ through the procedure shown above. This procedure, unlike in the case of PCA, is repeated for each iteration of the MCMC algorithm because $\Sigma_{(L)}$ depends on θ . It is unlikely to assume that $p(z^E|\theta)$ is exactly multivariate Gaussian because of the non-linear transformations introduced by the sigmoid transfer function; however, it still can be approximated by Gaussian distribution:

$$p(z^E|\theta) = N(z^E|z_{(L)}, \Sigma_{(L)}) \quad (C.6)$$

Considering the observation reported above; for each iteration of the MCMC sampling, we propose the following algorithm to determine and compute the analytical expression of the likelihood:

1. Compute the Kriging prediction $\hat{z}_{(L)}^{MM}(\theta)$ in the p^* -dimensional features space;
2. Reconstruct the Kriging prediction through the SSAE decoder in order to obtain $y^K(\theta)$;
3. Impose $z_{(0)} = y^K(\theta)$ and $\Sigma_{(0)} = I\sigma_{exp}^2$, then compute $\Sigma_{(L)}$ applying the EKF;
4. Assume the likelihood $p(z^E|\theta)$ to be Gaussian distributed:

$$z^E \sim N(z^E|z_{(L)}^{MM}(\theta), \Sigma_{(L)}(\theta)) = z_{(L)}^{MM}(\theta) + N(0, \Sigma_{(L)}(\theta)) \quad (C.7)$$

according to the Kriging theory:

$$z_{(L)}^{MM}(\theta) \sim N(\hat{z}_{(L)}^{MM}(\theta), \Sigma_{Kriging}(\theta)) = \hat{z}_{(L)}^{MM}(\theta) + N(0, \Sigma_{Kriging}(\theta)) \quad (C.8)$$

where $\Sigma_{Kriging}$ is the covariance matrix associated with the Kriging prediction uncertainty, that is a $p^* \times p^*$ matrix having the mean square errors of each feature prediction as diagonal entries:

$$\Sigma_{Kriging} = \begin{bmatrix} \sigma_{z_1(\theta)}^2 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \sigma_{z_{p^*}(\theta)}^2 \end{bmatrix} \quad (C.9)$$

5. Finally, substitute equation (C.8) into (C.7) and assume that $N(0, \Sigma_{(L)})$ and $N(0, \Sigma_{Kriging})$ are statistically independent, then $p(z^E|\theta)$ can be written as:

$$p(z^E|\theta) = N(\hat{z}_{(L)}^{MM}(\theta), \Sigma_{(L)}(\theta) + \Sigma_{Kriging}(\theta)) \quad (C.10)$$

Since in this work $N_{exp} = 1$, considering equation (C.10), the posterior PDF reduces to:

$$p(\theta|z^E) \propto p(\theta) \frac{1}{(\sqrt{2\pi})^p \sqrt{|\Sigma|}} \exp\left[-\frac{1}{2} [z^E - \hat{z}_{(L)}^{MM}(\theta)]^T \Sigma^{-1} [z^E - \hat{z}_{(L)}^{MM}(\theta)]\right] \quad (C.11)$$

where $\Sigma = \Sigma_{(L)}(\theta) + \Sigma_{Kriging}(\theta)$. Note that $\Sigma_{(L)}(\theta)$, in the expression (16) of Section 3.2, coincides with $\Sigma_{exp}(\theta)$.

References

- Abdelaziz, A.H., Watanabe, S., Hershey, J.R., Kolossa, D., Abdelaziz, A.H., Watanabe, S., et al., 2015. Uncertainty Propagation through Deep Neural Networks to Cite This Version : Others.
- Apostolakis, G., 1994. A Commentary on Model Uncertainty.
- Arendt, P.D., Apley, D.W., Chen, W., 2012. Quantification of model uncertainty: calibration, model discrepancy, and identifiability. *J. Mech. Des. Trans. ASME* 134, 1–12. <https://doi.org/10.1115/1.4007390>.
- Bandini, G., Meloni, P., Polidori, M., Lombardo, C., 2011. Validation of CATHARE V2.5 thermal-hydraulic code against full-scale PERSEO tests for decay heat removal in LWRs. *Nucl. Eng. Des.* 241, 4662–4671. <https://doi.org/10.1016/j.nucengdes.2011.02.034>.
- Bersano, A., Bertani, C., Falcone, N., de Salve, M., Mascari, F., Meloni, P., 2020. Qualification of RELAP5-3D code against the in-pool passive energy removal system PERSEO data. In: 30th Eur Saf Reliab Conf ESREL 2020 15th Probabilistic Saf Assess Manag Conf PSAM 2020, 1150–7.
- Conti, S., O'Hagan, A., 2010. Bayesian emulation of complex multi-output and dynamic computer models. *J. Stat. Plann. Inference* 140, 640–651. <https://doi.org/10.1016/j.jspi.2009.08.006>.
- Di Maio, F., Rai, A., Zio, E., 2016. A dynamic probabilistic safety margin characterization approach in support of Integrated Deterministic and Probabilistic Safety Analysis. *Reliab. Eng. Syst. Saf.* 145, 9–18. <https://doi.org/10.1016/j.res.2015.08.016>.
- Durga Rao, K., Kushwaha, H.S., Verma, A.K., Srividya, A., 2007. Quantification of epistemic and aleatory uncertainties in level-1 probabilistic safety assessment studies. *Reliab. Eng. Syst. Saf.* 92, 947–956. <https://doi.org/10.1016/j.res.2006.07.002>.
- D'Auria, F., Bousbia-salah, A., Petrucci, A., del Nevo, A., 2006. State of the art in using best estimate calculation tools in nuclear technology. *Nucl. Eng. Technol.* 38, 11–32.
- Ferri, R., Achilli, A., Cattadori, G., Bianchi, F., Meloni, P., 2005. Design, experiments and Relap5 code calculations for the perseo facility. *Nucl. Eng. Des.* 235, 1201–1214. <https://doi.org/10.1016/j.nucengdes.2005.02.011>.
- Ferson, S., Ginzburg, L.R., 1996. Different methods are needed to propagate ignorance and variability. *Reliab. Eng. Syst. Saf.* 54, 133–144. [https://doi.org/10.1016/S0951-8320\(96\)00071-3](https://doi.org/10.1016/S0951-8320(96)00071-3).
- Ferson, S., Joslyn, C.A., Helton, J.C., Oberkampf, W.L., Sentz, K., 2004. Summary from the epistemic uncertainty workshop: consensus amid diversity. *Reliab. Eng. Syst. Saf.* 85, 355–369. <https://doi.org/10.1016/j.res.2004.03.023>.
- Fricker, T.E., Oakley, J.E., Urban, N.M., 2013. Multivariate Gaussian process emulators with nonseparable covariance structures. *Technometrics* 55, 47–56. <https://doi.org/10.1080/00401706.2012.715835>.
- Gelman, A., Carlin, J.B., Stern, H.S., Dunson, D.B., Vehtari, A., Rubin, D.B., 2015. *Bayesian Data Analysis*, third ed., vol. 1542. <https://doi.org/10.1017/CBO9781107415324.004>.
- Glaeser, H., 2008. GRS method for uncertainty and sensitivity evaluation of code results and applications. *Sci. Technol. Nucl. Install.* 2008 <https://doi.org/10.1155/2008/798901>.
- Hadjajmadi, A.H., Homayounpour, M.M., 2015. Uncertainty propagation through neural network bottleneck features. *ICEE 2015 - Proc. 23rd Iran Conf. Electr. Eng.* 10, 567–575. <https://doi.org/10.1109/IranianCEE.2015.7146280>.
- Haftka, T., Shyy, W., Tucker, P.K., 2020. *Surrogate-based Analysis and Optimization*.
- Helton, J.C., Johnson, J.D., 2011. Quantification of margins and uncertainties: alternative representations of epistemic uncertainty. *Reliab. Eng. Syst. Saf.* 96, 1034–1052. <https://doi.org/10.1016/j.res.2011.02.013>.
- Higdon, D., Gattiker, J., Williams, B., Rightley, M., 2008. Computer model calibration using high-dimensional output. *J. Am. Stat. Assoc.* 103, 570–583. <https://doi.org/10.1198/01621450700000888>.
- Higdon, D., Geelhood, K., Williams, B., Unal, C., 2013. Calibration of tuning parameters in the FRAPCON model. *Ann. Nucl. Energy* 52, 95–102. <https://doi.org/10.1016/j.anucene.2012.06.018>.
- Holden, A.J., Robbins, D.J., Stewart, W.J., Smith, D.R., Schultz, S., Wegener, M., et al., 2006. Reducing the Dimensionality of Data with Neural Networks. *Science* 313, 504–507.
- Iaea, 2008. *Best estimate safety analysis for nuclear power Plants: uncertainty evaluation*. *Saf. Rep. Ser.* 52, 1–211.

- Idaho National Laboratory, 2015. RELAP5-3D Code Manual Volume I: Code Structure, System Models and Solution Methods, 2015.
- Jolliffe, I.T., 2002. *Principal Component Analysis*, second ed. Springer-Verlag, New York.
- Katafygiotis, L.S., Beck, J.L., 1998a. Updating models and their uncertainties. II: model identifiability. *J. Eng. Mech.* 124, 463–467. [https://doi.org/10.1061/\(asce\)0733-9399\(1998\)124:4\(463\)](https://doi.org/10.1061/(asce)0733-9399(1998)124:4(463)).
- Katafygiotis, L.S., Beck, J.L., 1998b. Updating models and their uncertainties. II: model identifiability. *J. Eng. Mech.* 124, 463–467.
- Kennedy, M.C., O'Hagan, A., 2001. Bayesian calibration of computer models. *J. Res. Stat. Soc. Ser. B (Statist. Methodol.)* 63, 425–464. <https://doi.org/10.1111/1467-9868.00294>.
- Kingma, D.P., Welling, M., 2014. Auto-encoding variational bayes. 2nd Int. Conf. Learn Represent ICLR 2014 - Conf. Track Proc. 1–14.
- Kleijnen, J.P.C., Mehdad, E., 2014. Multivariate versus univariate Kriging metamodels for multi-response simulation models. *Eur. J. Oper. Res.* 236, 573–582. <https://doi.org/10.1016/j.ejor.2014.02.001>.
- Lataniotis, C., Wicaksono, D., Marelli, S., Sudret, B., 2019. UQLab user manual: kriging (Gaussian process modeling) report # UQLab-V1.3-105. Chair risk. Saf Uncertain Quantif ETH Zurich, Switz 1–18.
- Mao, X.-J., Shen, C., Yang, Y.-B., 2016. Image Restoration Using Convolutional Auto-Encoders with Symmetric Skip Connections, pp. 1–17.
- Mascari, F., Lombardo, C., De Salve, M., Bertani, C., Bersano, A., Falcone, N., et al., 2019. Description of PERSEO Test N. 7 for International Open Benchmark Exercise. ADPFISS-LP1-126.
- McKay, M.D., Beckman, R.J., Conover, W.J., 2000. A comparison of three methods for selecting values of input variables in the analysis of output from a computer code. *Technometrics* 42, 55–61. <https://doi.org/10.1080/00401706.2000.10485979>.
- Mohammadi, H., Challenor, P., Goodfellow, M., 2019. Emulating dynamic non-linear simulators using Gaussian processes. *Comput. Stat. Data Anal.* 139, 178–196. <https://doi.org/10.1016/j.csda.2019.05.006>.
- Monisha, R., Mrinalini, R., Britto, M.N., Ramakrishnan, R., Rajinikanth, V., 2019. Smart Intell. Comput. Appl. 104 <https://doi.org/10.1007/978-981-13-1921-1>.
- Nagel, J.B., Rieckermann, J., Sudret, B., 2020. Principal component analysis and sparse polynomial chaos expansions for global sensitivity analysis and model calibration: application to urban drainage simulation. *Reliab. Eng. Syst. Saf.* 195, 106737. <https://doi.org/10.1016/j.res.2019.106737>.
- Ng, A., 2011. Sparse Autoencoder, CS294A Lecture Notes, pp. 1–19.
- Olshausen, B.A., Fieldt, D.J., 1997. Sparse coding with an overcomplete basis set: a strategy employed by V1 ? coding V1 gabor-wavelet natural images. *Vis. Res.* 37, 3311–3325.
- Pourgol-Mohammad, M., 2009. Thermal-hydraulics system codes uncertainty assessment: a review of the methodologies. *Ann. Nucl. Energy* 36, 1774–1786. <https://doi.org/10.1016/j.anucene.2009.08.018>.
- Rasmussen, C.E., Williams, C.K.I., 2006. *Gaussian Processes for Machine Learning*, vol. 38, 2006.
- Roma, G., Di Maio, F., Bersano, A., Pedroni, N., Bertani, C., Mascari, F., et al., 2021. A Bayesian framework of inverse uncertainty quantification with principal component analysis and Kriging for the reliability analysis of passive safety systems. *Nucl. Eng. Des.* 379, 111230. <https://doi.org/10.1016/j.nucengdes.2021.111230>.
- Rumelhart, D.E., Hinton, G.E., Williams, R.J., 1986. Learning representations by back-propagating errors. *Nature* 323, 533–536. <https://doi.org/10.1038/323533a0>.
- Saltelli, A., Annoni, P., Azzini, I., Campolongo, F., Ratto, M., Tarantola, S., 2010. Variance based sensitivity analysis of model output. Design and estimator for the total sensitivity index. *Comput. Phys. Commun.* 181, 259–270. <https://doi.org/10.1016/j.cpc.2009.09.018>.
- Shrestha, R., Kozlowski, T., 2016. Inverse uncertainty quantification of input model parameters for thermal-hydraulics simulations using expectation-maximization under Bayesian framework. *J. Appl. Stat.* 43, 1011–1026. <https://doi.org/10.1080/02664763.2015.1089220>.
- Titensky, J.S., Jananthan, H., Kepner, J., 2018. Uncertainty propagation in Deep neural networks using extended kalman filtering. In: 2018 IEEE MIT Undergrad Res Technol Conf URTC 2018. <https://doi.org/10.1109/URTC45901.2018.9244804>.
- Utgoff, P.E., Stracuzzi, D.J., 2002. Many-layered learning. Proc - 2nd Int Conf Dev Learn ICDL 2002. <https://doi.org/10.1109/DEVLRN.2002.1011824>, 141–6.
- Van Der Maaten, L.J.P., Postma, E.O., Van Den Herik, H.J., 2009. Dimensionality reduction: a comparative review. *J. Mach. Learn. Res.* 10, 1–41. <https://doi.org/10.1080/13506280444000102>.
- Vincent, P., Larochelle, H., 2008. Extracting and Composing Robust Features with Denoising.Pdf, pp. 1096–1103.
- Wang, G.G., Shan, S., 2007. Review of metamodeling techniques in support of engineering design optimization. *J. Mech. Des. Trans. ASME* 129, 370–380. <https://doi.org/10.1115/1.2429697>.
- Wang, Y., Yao, H., Zhao, S., 2016. Auto-encoder based dimensionality reduction. *Neurocomputing* 184, 232–242. <https://doi.org/10.1016/j.neucom.2015.08.104>.
- Wang, C., Wu, X., Kozlowski, T., 2019. Gaussian process-based inverse uncertainty quantification for TRACE physical model parameters using steady-state PSBT benchmark. *Nucl. Sci. Eng.* 193, 100–114. <https://doi.org/10.1080/00295639.2018.1499279>.
- Welch, G., Bishop, G., 2006. An introduction to the kalman filter. In *Pract. 7*, 1–16.
- Wilkinson, R.D., 2010. Bayesian calibration of expensive multivariate computer experiments. In: *Large-Scale Inverse Probl Quantif Uncertain*, pp. 195–215. <https://doi.org/10.1002/9780470685853.ch10>.
- Winkler, R.L., 1996. Uncertainty in probabilistic risk assessment. *Reliab. Eng. Syst. Saf.* 54, 127–132. [https://doi.org/10.1016/S0951-8320\(96\)00070-1](https://doi.org/10.1016/S0951-8320(96)00070-1).
- Wu, X., Kozlowski, T., Meidani, H., Shirvan, K., 2018a. Inverse uncertainty quantification using the modular Bayesian approach based on Gaussian Process, Part 2: application to TRACE. *Nucl. Eng. Des.* 335, 417–431. <https://doi.org/10.1016/j.nucengdes.2018.06.003>.
- Wu, X., Kozlowski, T., Meidani, H., Shirvan, K., 2018b. Inverse uncertainty quantification using the modular Bayesian approach based on Gaussian process, Part 1: Theory. *Nucl. Eng. Des.* 335, 339–355. <https://doi.org/10.1016/j.nucengdes.2018.06.004>.
- Wu, X., Kozlowski, T., Meidani, H., 2018c. Kriging-based inverse uncertainty quantification of nuclear fuel performance code BISON fission gas release model using time series measurement data. *Reliab. Eng. Syst. Saf.* 169, 422–436. <https://doi.org/10.1016/j.res.2017.09.029>.
- Wu, X., Xie, Z., Alsafadi, F., Kozlowski, T., 2021. A comprehensive survey of inverse uncertainty quantification of physical model parameters in nuclear system thermal-hydraulics codes. *Nucl. Eng. Des.* 384. <https://doi.org/10.1016/j.nucengdes.2021.111460>.
- Yang, Z., Baraldi, P., Zio, E., 2018. Automatic extraction of a health indicator from vibrational data by sparse autoencoders. In: Proc - 2018 3rd Int Conf Syst Reliab Safety. ICSRS. <https://doi.org/10.1109/ICSRS.2018.8688720>, 2019:328–32.
- Zhao, R., Yan, R., Chen, Z., Mao, K., Wang, P., Gao, R.X., 2019. Deep learning and its applications to machine health monitoring. *Mech. Syst. Signal Process.* 115, 213–237. <https://doi.org/10.1016/j.ymsp.2018.05.050>.
- Zio, E., Di Maio, F., Tong, J., 2010. Safety margins confidence estimation for a passive residual heat removal system. *Reliab. Eng. Syst. Saf.* 95, 828–836. <https://doi.org/10.1016/j.res.2010.03.006>.