# THE DEEP ARBITRARY POLYNOMIAL CHAOS NEURAL NETWORK OR HOW DEEP ARTIFICIAL NEURAL NETWORKS COULD BENEFIT FROM DATA-DRIVEN HOMOGENEOUS CHAOS THEORY

A PREPRINT

**⊙ Sergey Oladyshkin**
Department of Stochastic Simulation and Safety Research for Hydrosystems,
Institute for Modelling Hydraulic and Environmental Systems, Stuttgart Center for Simulation Science,
University of Stuttgart, Pfaffenwaldring 5a, 70569 Stuttgart, Germany
Sergey.Oladyshkin@iws.uni-stuttgart.de

**Timothy Praditia**
Department of Stochastic Simulation and Safety Research for Hydrosystems,
Institute for Modelling Hydraulic and Environmental Systems, Stuttgart Center for Simulation Science,
University of Stuttgart, Pfaffenwaldring 5a, 70569 Stuttgart, Germany

**⊙ Ilja Kröker**
Department of Stochastic Simulation and Safety Research for Hydrosystems,
Institute for Modelling Hydraulic and Environmental Systems, Stuttgart Center for Simulation Science,
University of Stuttgart, Pfaffenwaldring 5a, 70569 Stuttgart, Germany

**Farid Mohammadi**
Department of Hydromechanics and Modelling of Hydrosystems,
Institute for Modelling Hydraulic and Environmental Systems,
University of Stuttgart, Pfaffenwaldring 61, 70569 Stuttgart, Germany

**⊙ Wolfgang Nowak**
Department of Stochastic Simulation and Safety Research for Hydrosystems,
Institute for Modelling Hydraulic and Environmental Systems, Stuttgart Center for Simulation Science,
University of Stuttgart, Pfaffenwaldring 5a, 70569 Stuttgart, Germany

**Sebastian Otte**
Neuro-Cognitive Modeling, Computer Science Department,
University of Tübingen, Sand 14, 72076 Tübingen, Germany

June 27, 2023

## ABSTRACT

Artificial Intelligence and Machine learning have been widely used in various fields of mathematical computing, physical modeling, computational science, communication science, and stochastic analysis. Approaches based on Deep Artificial Neural Networks (DANN) are very popular in our days. Depending on the learning task, the exact form of DANNs is determined via their multi-layer architecture, activation functions and the so-called loss function. However, for a majority of deep learning approaches based on DANNs, the kernel structure of neural signal processing remains the

same, where the node response is encoded as a linear superposition of neural activity, while the non-linearity is triggered by the activation functions. In the current paper, we suggest to analyze the neural signal processing in DANNs from the point of view of homogeneous chaos theory as known from polynomial chaos expansion (PCE). From the PCE perspective, the (linear) response on each node of a DANN could be seen as a $1^{st}$ degree multi-variate polynomial of single neurons from the previous layer, i.e. linear weighted sum of monomials. From this point of view, the conventional DANN structure relies implicitly (but erroneously) on a Gaussian distribution of neural signals. Additionally, this view revels that by design DANNs do not necessarily fulfill any orthogonality or orthonormality condition for a majority of data-driven applications. Therefore, the prevailing handling of neural signals in DANNs could lead to redundant representation as any neural signal could contain some partial information from other neural signals. To tackle that challenge, we suggest to employ the data-driven generalization of PCE theory known as arbitrary polynomial chaos (aPC) to construct a corresponding multi-variate orthonormal representations on each node of a DANN. Doing so, we generalize the conventional structure of DANNs to Deep arbitrary polynomial chaos neural networks (DaPC NN). They decompose the neural signals that travel through the multi-layer structure by an adaptive construction of data-driven multi-variate orthonormal bases for each layer. Moreover, the introduced DaPC NN provides an opportunity to go beyond the linear weighted superposition of single neurons on each node. Inheriting fundamentals of PCE theory, the DaPC NN offers an additional possibility to account for high-order neural effects reflecting simultaneous interaction in multi-layer networks. Introducing the high-order weighted superposition on each node of the network mitigates the necessity to introduce non-linearity via activation functions and, hence, reduces the room for potential subjectivity in the modeling procedure. Although the current DaPC NN framework has no theoretical restrictions on the use of activation functions. The current paper also summarizes relevant properties of DaPC NNs inherited from aPC as analytical expressions for statistical quantities and sensitivity indexes on each node. We also offer an analytical form of partial derivatives that could be used in various training algorithms. Technically, DaPC NNs require similar training procedures as conventional DANNs, and all trained weights determine automatically the corresponding multi-variate data-driven orthonormal bases for all layers of DaPC NN. The paper makes use of three test cases to illustrate the performance of DaPC NN, comparing it with the performance of the conventional DANN and also with plain aPC expansion. Evidence of convergence over the training data size against validation data sets demonstrates that the DaPC NN outperforms the conventional DANN systematically. Overall, the suggested re-formulation of the kernel network structure in terms of homogeneous chaos theory is not limited to any particular architecture or any particular definition of the loss function. The DaPC NN Matlab Toolbox is available online and users are invited to adopt it for own needs.

**Keywords** Artificial Intelligence · Machine Learning · Deep Artificial Neural Network · Polynomial Chaos Expansion · Arbitrary Polynomial Chaos · Orthogonal decomposition · High-order neural interactions · Deep Arbitrary Polynomial Chaos

# 1 Introduction

During the last decades, Artificial Intelligence (AI) and Machine learning (ML) have been widely used in various fields of mathematical computing, physical modeling, computer science, geosciences, communication science, and stochastic analysis. The terminology AI has been suggested by John McCarthy in 1956 as a neutral title of a Dartmouth workshop [63] to distinguish the research field from cybernetics and also to escape the influence of its originator Norbert Wiener [114]. The closely related term ML has been introduced later in 1959 by Arthur Samuel, where the author explored the logical rules of the game of checkers [90]. Originally, AI and ML have been focused on learning strategies employing logical rules, which were often formalized using an apparatus of discrete mathematics and graph theory. However, with increasing computational power [67] and data availability [69], the fields of AI and ML today employ a much broader spectrum of approaches that originate from stochastic analysis, cybernetics, geosciences, information theory and other disciplines.

In particular, approaches based on Deep Artificial Neural Networks (DANN) introduced in cybernetics by Alexey Ivakhnenko and Valentin Lapa in 1967 [43] are currently very popular in AI and ML (see e.g. [91] for a detailed historical recapitulation). DANNs generalize the concept of Artificial Neural Networks (ANN) suggested by Warren McCulloch and Walter Pitts in 1943 [64] to multi-layer structures , i.e. deep ANN. This form of deep learning also gained a strong visibility in society, providing solutions for a broad variety of tasks including recognition of images [105], videos [24], voice [8] and text [95].

Depending on the modelling task, the exact form of the so-called loss function [35] is usually specified to determine the final DANN representation. Various loss functions can be found in literature, suitable for classification tasks [13, 56], Bayesian interpolation [60] or physical regularization [47, 84]. In addition, multiple DANN layers and corresponding neural connections can be customized in various ways, turning DANN into convolutional [19, 85], recurrent [38, 117] or other desired architectures [48]. Due to the magnificent number of works based on the DANN representation, the corresponding literature can hardly be covered in a research paper, and the authors refer to literature [37, 36] and reviews [18, 92] for further information.

The keystone of DANNs can be seen as a certain type of non-linear function that maps from input in $\mathbb{R}^n$ to output in $\mathbb{R}$. To construct such a function, DANNs use a multi-layered approach, where each node of a layer is a linear combination of non-linear univariate functions, known as activation functions. Overall, this yields a chain-rule interaction of neurons in multi-layered architecture [7]. However, there is an alternative branch of approaches that also maps $\mathbb{R}^n$ to $\mathbb{R}$ using linear combinations of non-linear kernel functions or vectors and, hence, could be seen as one-layer approaches. This alternative ML branch consists of polynomial chaos expansions (PCE) introduced in 1938 [113], Kriging introduced in 1951 [55] that is also known as Gaussian process emulator/regression (GPE: [115]) or Wiener–Kolmogorov prediction [27], Support vector regression (SVR) introduced in 1974 [108] and Relevance vector machines (RVM) introduced in 2000 [107]. The common fundamental between PCE, GPE, SVR and RVM can be found in [22].

Regardless of the different mathematical definitions and structures, both one- and multi-layered ML approaches construct the final non-linear representation by determining all their unknown constants. These constants are known as coefficients (terminology for PCE, GPE, SVR and RVN) or weights (terminology for ANN and DANN). For one-layer approaches, the response is established by solving a linear system of equations that encodes the linear combination of non-linear basis functions (polynomials, kernels, vectors, etc.). For multi-layer approaches, the response is constructed by solving a non-linear system of equations that reflects the so-called neural signal processing.

The current paper does not aim at discussing the pros and contras behind various approaches to map $\mathbb{R}^n$ to $\mathbb{R}$. We will rather pay attention to how one-layer findings could be helpful for the multi-layered structure of DANNs. Indeed, let us have a close look into the kernel structure of conventional DANNs. It considers the processing of the neural signal in one node as displayed in Figure 1a. The response of each node (i.e. arrows leaving the central circle) contains linear weighted superposition of single neuron responses (i.e. the arrows entering the central circle) coming from the previous layer. In the figure, the weighting is represented by $w$ and the superposition by $\Sigma$. To obtain the final neuronal response, the superposition is passed through an activation function $\mathcal{A}$ that is usually non-linear. The corresponding weights $w$ of the linear superposition for the input to the featured neurons are to be found by training. However, such a linear weighted superposition of incoming neuron outputs on each node could lead to a redundant representation, as it is not necessarily satisfying orthogonality in signal processing. This aspect has been addressed in the literature on Support Vector Networks in 1995 [26] employing the SVR concept [108]. Also imposing orthogonality within the DANNs has been explored in the context of Recurrent Neural Networks [65] assuring the efficiency of the training procedure. The paper [111] highlights the benefit of using orthogonal weight matrices in Recurrent Neural Networks, as they preserve gradient norm during backpropagation making them highly desirable for optimization purposes. However, the imposition of hard constraints on orthogonality within Recurrent Neural Networks may negatively affect convergence speed [111] as applying Gaussian prior regularization may not be appropriate for many applications. Additionally, in order to overcome the difficulty of training Recurrent Neural Networks caused by vanishing and exploding gradients, a new approach have been addressed that learns a unitary weight matrix with eigenvalues [116, 9], enabling optimization in the complex domain. Nevertheless, accounting for orthogonality seems to be very promising, as it could mitigate redundancy in DANN representation [112] and potentially could provide a better ability for generalization [46].

Additionally, the actual non-linearity of DANNs is triggered by the non-linearity of the choosen activation functions. However, the choice of the activation functions is an extremely non-trivial task itself [66] and can be very subjective [94], posing additional challenges for DANN users. Very recently, the work [23] suggested to replace non-linear activation functions via non-linear polynomial representations, introducing co-called $\Pi$-Nets. The introduced $\Pi$-Nets consider high-order polynomial terms, which consistently improves the performance in discriminative and generative tasks for images and audio processing [23]. Nevertheless, similar to conventional DANNs, $\Pi$-Nets do not yet consider orthogonality in processing the neural signal and, hence, could also lead to redundancy in representation. We argue that employing orthogonal (or even better orthonormal) decompositions for processing neural signals could be extremely relevant, especially for data-poor applications.

In the current work, we will take the reader on a journey to the early stages of ML. We will pay special attention to the PCE, introduced via the homogeneous chaos theory by Norbert Wiener in 1938 [113] as already mentioned above. The so-called non-intrusive version of the PCE [33, 59] and its advanced extensions towards sparse quadrature [50], sparse regression [4, 16] or multi-element approaches [6, 57] gained popularity for surrogate building in computationally demanding modelling tasks [31, 73]. They all employ the idea of multi-variate polynomial representation. The

Figure 1: Sketches of the neural processing in (a) DANN and (b) DaPC NN

fundamental concept of PCE theory lies in projection of functions onto a space spanned by orthonormal polynomials that capture the non-linear dependencies.

Using the DANN vocabulary, the PCE could be seen as a one-layer ML approach, where the response is approximated as a linear combination of non-linear multivariate orthonormal polynomial basis functions. At the same time, employing the PCE vocabulary, the response of one node in a DANN could be seen as a $1^{st}$ degree multi-variate polynomial of single neurons from the previous layer, i.e. linear weighted superposition, then passed through an activation function. To be more specific, the $1^{st}$ degree multi-variate polynomial representation in one node of a DANN consists of $1^{st}$ degree monomials that reflect the incoming neural signals and the $0^{th}$ degree term that is known as bias. Obviously, this representation via linear weighted superposition of monomials does not necessarily fulfill the orthogonality or orthonormality condition that is targeted by PCE approach for the majority of applications.

Such a non-orthogonal linear weighted superposition of incoming neuron outputs is a possibly redundant representation of the overall signal flow: any neural signal could contain some partial information from other neural signals. Similarly, the non-linear polynomial representation in Π-Nets [23] could also profit from orthonormal decomposition of neural signals. Therefore, employing multi-variate orthonormal polynomial bases to process the neural signal on each node, as in the PCE approach, could be beneficial for both conventional DANN structures and for advanced Π-Nets. Another approach that combines variational autoencoder and PCE methods is expected to appear soon in the literature forming the PCE-Net model [97].

Unfortunately, due to the data-driven nature of all these networks, a direct transfer of classical PCE theory or its generalized extension [120] for constructing orthonormal representations on each node is not possible. The reason is that PCE requires knowledge of probability density functions for all inputs, and this knowledge is not available in the context of DANNs, where the probability density function for all inputs to all nodes in all layers would have to be known.

Therefore, the current paper suggests to account for the data-driven nature of neural signals in DANNs and thus uses the data-driven generalization of PCE known as arbitrary polynomial chaos (aPC: [75]) to construct corresponding multi-variate orthonormal representations on each node of DANNs. Doing so, we generalize the conventional structure of DANNs to Deep arbitrary polynomial chaos neural networks (DaPC NN)[1]. They decompose the neural signals traveling through the multi-layer structure through an adaptive construction of data-driven multi-variate orthonormal bases for each layer. Doing so, we provide an opportunity to go beyond the linear weighted superposition of single (univariate) neurons on each node that is traditionally employed in various DANN architectures.

Inheriting fundamentals of the PCE, the DaPC NN offers an additional possibility to account for high-order neural effects that reflect simultaneous interaction of neurons in the multi-layer network. Thus, the modeler is prompted to specify the DaPCE NN architecture through not only the number of layers, the number of nodes per layer and the activation function as in conventional DANNs, but also through the desired polynomial degree of non-linearity per layer. Figure 1b schematically illustrates the neural processing in the DaPC NN through a multi-variate orthonormal basis $\Psi$ and a corresponding high-order weighted superposition $\sum_{\forall d}$. Similar to the conventional DANN, unknown weights of the DaPC NN should be determined by training.

---

[1] The authors of the current paper came across the interesting work under revision [123] denoted as Deep aPCE that seems to be developed in parallel with the current DaPC NN work. In order to avoid any confusion for the reader, we would like to clarify that both works employ the data-driven fundamentals of aPC [75]. However, regardless the similarity in the abbreviations, the Deep aPCE makes use of DANNs to construct the aPC expansion coefficients, but the DaPC NN employs aPC representations to construct a re-formulated DANN. We refer the reader to the original paper [123] for more details.

In this sense, the paper does not aim to contribute to the DANN architecture, to any particular definition of loss functions, or to optimal training schemes. Instead, we offer a mathematical reformulation of the kernel DANN structure in terms of homogeneous chaos theory. Moreover, introducing the high-order weighted superposition on each node may remove the necessity to introduce non-linearity via activation functions if desired. Hence, one can reduce the room for potential subjectivity in the modeling procedure. We expect that joining the multi-layered structure of neural networks together with the theory of polynomial chaos expansion could be beneficial for ML tasks and also creates a foundation for further investigations.

The rest of the paper is structured as follows: Section 2 summarizes the necessary background on the data-driven aPC in Section 2.1 and the conventional DANN structure in Section 2.2. Section 2.2 also points out the implicit Gaussian assumptions of neural signal processing in conventional DANN structures from the view point of PCE theory. Section 2.3 offers our novel DaPC NN formulation. It introduces the adaptive, deep, multi-variate orthonormal polynomial representation via a multi-layered structure. Section 2.3 also provides relevant properties that the DaPC NN inherits from aPC. Section 2.4 briefly outlines an example of the conventional training procedure to compute the weights by providing the required partial derivatives. Section 3 illustrates the performance of DaPC NN together with a conventional DANN and the aPC. In this comparison, we use exactly the same training data and loss function for DANN, DaPC NN and aPC in a total of three different test cases. Additionally, Section 3 shows evidence of convergence against reference validation data sets, comparing how the size of the training data sets affects the performance of the considered ML approaches.

## 2 Machine learning with non-redundant decomposition

### 2.1 Arbitrary Polynomial Chaos Expansion

Theory of polynomial chaos expansion (PCE) was originally introduced by Norbert Wiener [113] in 1938. In PCE, the dependence of the model response (output) on all inputs is expressed using projection onto an orthogonal or orthonormal multi-variate polynomial basis [33, 104]. For the current paper, we will employ a purely data-driven generalization of the PCE introduced in 2012 by Oladyshkin and Nowak [75] that will open the pathway for a data-driven deep ML structure.

#### 2.1.1 Homogeneus chaos theory

Let us consider a model response $\mathcal{R}(\boldsymbol{\omega})$ that depends on some multi-dimensional input $\boldsymbol{\omega} = \{\omega_1, \ldots, \omega_n\}$ from the input space $\Omega$, where $n$ is the number of inputs. Let $L^2(\Omega)$ denote the $L^2$-space on $\Omega$, weighted by the assumed probability distribution. According to PCE theory [21, 113], the model response $\mathcal{R}(\boldsymbol{\omega}) \in L^2(\Omega)$ can be expanded in $\boldsymbol{\omega}$ in the following manner to map $\mathbb{R}^n$ to $\mathbb{R}$:

$$\mathcal{R}(\boldsymbol{\omega}) = \sum_{i=0}^{\infty} w_i \Psi_i(\boldsymbol{\omega}) \approx \sum_{i=0}^{M} w_i \Psi_i(\boldsymbol{\omega}), \tag{1}$$

where $\Psi_i(\boldsymbol{\omega})$ are basis functions from a multi-variate orthonormal (w.r.t. the assumed probability distribution) polynomial basis $\{\Psi_0(\boldsymbol{\omega}), \ldots, \Psi_M(\boldsymbol{\omega})\}$ defined on the input space $\Omega$, and $w_i$ are corresponding coefficients that determine the form of the expansion in equation (1). The total number of expansion terms $M$ depends on the number of inputs $N$ and on the desired degree $d$ of the polynomial representation as $M = (n + d)!/(n!d!)$. This formulation of $M$ is based on a total-degree truncation, and other alternatives exist [104]. For a comprehensive discussion of the requirements for existence and completeness of an orthonormal polynomial basis $\{\Psi_0(\boldsymbol{\omega}), \ldots, \Psi_M(\boldsymbol{\omega})\}$ of $L^2(\Omega)$, we refer to [30, 32, 104].

The multi-variate polynomial basis $\{\Psi_0(\boldsymbol{\omega}), \ldots, \Psi_M(\boldsymbol{\omega})\}$ is comprised of the tensor product of univariate orthonormal polynomials $\{\phi_j^{(0)}, \ldots, \phi_j^{(d)}\}$ of degree $d$ for the inputs $\omega_j$, assuming that the inputs are statistically independent:

$$\Psi_\alpha(\boldsymbol{\omega}) = \prod_{j=1}^{N} \phi_j^{(\alpha_j)}(\omega_j), \quad \sum_{j=1}^{N} \alpha_j \leq d. \tag{2}$$

Here $\alpha = (\alpha_1, \ldots, \alpha_N) \in \mathbb{N}_0^N$ is a multivariate index describing polynomial degree that contains the combinatoric information how to enumerate all possible products of individual univariate basis functions and contains the corresponding polynomial degree for input $\omega_j$ within the univariate polynomials $\phi_j^{(\alpha_j)}(\omega_j)$. The set of polynomials

$\left\{\phi_j^{(\alpha_j)} \mid \alpha_j = 0, \ldots, d\right\}$ forms an orthonormal basis of polynomial degree at most $d$ for each input $\omega_j$:

$$\left\langle \phi_j^{(\alpha_j)}(\omega_j), \phi_j^{(\alpha'_j)}(\omega_j) \right\rangle_{L^2(\Omega)} = \delta_{\alpha_j, \alpha'_j}, \tag{3}$$

where $\delta_{\alpha_j, \alpha'_j}$ represents the Kronecker delta $\forall \alpha_j, \alpha'_j = 0, \ldots, d$.

### 2.1.2 Data-driven orthonormal representation

The original theory of homogeneous chaos [113] is based on orthonormal Hermite polynomials $\phi_j^{(\alpha_j)}(\omega_j)$ satisfying eq. (3), which are optimal [75, 104] for Gaussian distributed inputs $\boldsymbol{\omega}$. Further extensions to a number of parametric statistical distributions (Gamma, Beta, Uniform, etc.) have been suggested in [119] and [118], based on the Askey scheme [10] of orthonormal polynomials.

Such approaches assume an exact knowledge of the involved probability density functions [104], which is not available in various applications or often requires additional assumptions [86]. To assure a purely data-driven representation for ML, in the current paper we will consider the data-driven generalization of the polynomial chaos expansion known as the arbitrary polynomial chaos (aPC) introduced in [75]. The necessity to adapt to arbitrary distributions in practical tasks is discussed in more detail in [72].

The aPC technique adapts to arbitrary probability distribution shapes of the inputs and can be inferred from limited input data through a few statistical moments [75]. Formally, univariate orthonormal polynomials bases $\phi_j^{(\alpha_j)}(\omega_j)$ of polynomial degree $\alpha_j$ can be written as a sum of the following monomials:

$$\phi_j^{(\alpha_j)}(\omega_j) = \frac{1}{\sqrt{\kappa_{\alpha_j}}} \sum_{i=0}^{\alpha_j} m_i^{(\alpha_j)} \omega_j^i, \quad \alpha_j = \overline{0, d}, \tag{4}$$

where $\kappa_{\alpha_j} = m_{\alpha_j}^{(\alpha_j)}$ is a constant representing the norm of the univariate polynomial. The corresponding monomial coefficients $m_i^{(\alpha_j)}$ can be defined according to the (empirical or theoretical) raw moments of the inputs $\omega_j$ as follows [75, 79]:

$$\begin{bmatrix} \mu_0(\omega_j) & \mu_1(\omega_j) & \ldots & \mu_{\alpha_j}(\omega_j) \\ \vdots & \vdots & \ddots & \vdots \\ \mu_{\alpha_j-1}(\omega_j) & \mu_{\alpha_j}(\omega_j) & \ldots & \mu_{2\alpha_j-1}(\omega_j) \\ \mu_{\alpha_j}(\omega_j) & \mu_{\alpha_j+1}(\omega_j) & \ldots & \mu_{2\alpha_j}(\omega_j) \end{bmatrix} \begin{bmatrix} m_0^{(\alpha_j)} \\ \vdots \\ m_{\alpha_j-1}^{(\alpha_j)} \\ m_{\alpha_j}^{(\alpha_j)} \end{bmatrix} = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix}, \tag{5}$$

where $\mu_k(\omega_j)$ denotes the $k$-th raw stochastic moment of the input $\omega_j$.

Therefore, constructing the multi-variate orthonormal polynomial basis $\left\{\Psi_0(\boldsymbol{\omega}), \ldots, \Psi_M(\boldsymbol{\omega})\right\}$ of degree $d$ for an aPC representation $\mathcal{R}(\boldsymbol{\omega})$ in equation (1) is based on the data-driven formulation in equations (2), (4) and (5). In particular, the aPC approach is based on $2d$ raw stochastic moments only. These moments can be evaluated directly from an available data set of limited size. Matrix on the left-hand side of equation (5) is known as the Hankel matrix of moments [49] and its properties have been analysed in [98]. The resulting polynomials are real if, and only if, the Hankel matrix of moments is positive definite, see also the related Hamburger moment problem, e.g. [5, 30, 96, 104]. From the practical point of view, the solution of the linear system of equations in (5) can be obtained directly numerically, via lower-order moments representation [75], via recursive relations [1], via Gram-Schmidt orthogonalization [45] or via the Stieltjes procedure [102].

Overall, the polynomial representation in equation (1) is one of the oldest ML approaches to map $\mathbb{R}^n$ to $\mathbb{R}$. It quantifies the response $\mathcal{R}$ to the inputs $\boldsymbol{\omega}$ through an orthonormal basis $\left\{\Psi_0(\boldsymbol{\omega}), \ldots, \Psi_M(\boldsymbol{\omega})\right\}$ and computes the expansion coefficients $w_i$ [2]. Each coefficient $w_i$ for $i > 0$ indicates how much variance one or another term brings into the overall composition (see also relation to global sensitivity analysis in [74, 122]). The unknown expansion coefficients $w_i$ can be determined using Galerkin projection [54, 75], numerical integration, regression or collocation approaches [58, 73, 109].

---

[2]Traditionally, the response $\mathcal{R}$ has often been used as an approximation of some full-complexity physical model $\mathcal{M}$ in order to learn about the non-linear dependence of model output on modelling inputs $\boldsymbol{\omega}$, i.e. $\mathcal{R}(\boldsymbol{\omega}) \approx \mathcal{M}(\boldsymbol{\omega})$. In that sense the PCE projection in equation (1) is often used a surrogate (response surface or reduced model), that is considered to be a special case of supervised machine learning.

## 2.2   Deep Artificial Neural Networks

Let us shortly summarize the key aspects of conventional DANN structures. They all rely on the multi-layer concept introduced by Alexey Ivakhnenko and Valentin Lapa in 1967 [43]. Similar to the method introduced in Section 2.1, DANNs maps $\mathbb{R}^n$ to $\mathbb{R}$ by providing an approximation of the response (output). In DANNs, information flows from input nodes, through co-called hidden layers, and then to output [35]. The term hidden refers to the fact that the computations that occur within the hidden layer are not visible from the outside of the network. Each hidden layer in a neural network contains several hidden nodes, also known as hidden neurons. These nodes are responsible for performing computations on the input data and transmitting the results to other nodes in the network, ultimately leading to the output. It has gained significant popularity recently because it is a universal approximator [39] and due to technological advances in computing power.

### 2.2.1   Conventional DANN structure

We will consider a conventional Deep Artificial Neural Network (DANN) with $L$ hidden layers and corresponding numbers $N^{(\mathcal{L})}$ of hidden neurons for each layer $\mathcal{L}$ ($\mathcal{L} = 1, \ldots, L$). Similar to the ML approach in Section 2.1, DANNs provide the response $\mathcal{R}(\boldsymbol{\omega})$ as a non-linear dependence [7] on the same multidimensional input $\boldsymbol{\omega} = \{\omega_1, \ldots, \omega_n\}$ from the input space $\Omega$, where $n$ is the number of inputs. We will denote the response on each hidden layer $\mathcal{L}$ as a vector $\boldsymbol{\mathcal{R}}^{(\mathcal{L})} = \left\{ \mathcal{R}^{(\mathcal{L},1)}, \ldots, \mathcal{R}^{(\mathcal{L},N^{\mathcal{L}})} \right\}$. This vector represents the responses (outputs) from the corresponding hidden node $1, \ldots, N^{\mathcal{L}}$ on the current layer $\mathcal{L}$. Then, the response of the DANN representation $\mathcal{R}(\boldsymbol{\omega})$ could be seen [35] as recursive encapsulation of responses from hidden layers that contains $N^{(\mathcal{L})}$ hidden neurons for each layer $\mathcal{L}$. Formally, representation $\mathcal{R}(\boldsymbol{\omega})$ could be written in the following form:

$$\mathcal{R}(\boldsymbol{\omega}) = \boldsymbol{\mathcal{R}}^{(L)}\left( \boldsymbol{\mathcal{R}}^{(L-1)}\left( \ldots \left( \boldsymbol{\mathcal{R}}^{(1)}\left( \boldsymbol{\omega} \right) \right) \right) \right), \tag{6}$$

where the input for each hidden layer $\boldsymbol{\mathcal{R}}^{(\mathcal{L})}$ is the response from the previous layer $\boldsymbol{\mathcal{R}}^{(\mathcal{L}-1)}$ ($\mathcal{L} = 2, \ldots, L$). The input for the first layer $\boldsymbol{\mathcal{R}}^{(1)}$ is the overall input $\boldsymbol{\omega} = \{\omega_1, \ldots, \omega_n\}$

Each response $\mathcal{R}^{(\mathcal{L},\mathcal{N})}$ of the node $\mathcal{N}$ in the hidden layer $\mathcal{L}$ is defined according to the ANN representation [35] as follows:

$$\begin{aligned} \mathcal{R}^{(\mathcal{L},\mathcal{N})}\left( \boldsymbol{\mathcal{R}}^{(\mathcal{L}-1)} \right) = & w_0^{(\mathcal{L},\mathcal{N})} \\ & + \sum_{i=1}^{M^{(\mathcal{L})}} w_i^{(\mathcal{L},\mathcal{N})} \mathcal{A}^{(\mathcal{L})}(\mathcal{R}^{(\mathcal{L}-1,i)}). \end{aligned} \tag{7}$$

Here, $\boldsymbol{\mathcal{R}}^{(\mathcal{L}-1)} = \left\{ \mathcal{R}^{(\mathcal{L}-1,1)}, \ldots, \mathcal{R}^{(\mathcal{L}-1,N^{\mathcal{L}-1})} \right\}$ is the response from the previous layer $\mathcal{L}-1$ (where $\boldsymbol{\mathcal{R}}^{(1)} = \boldsymbol{\omega}$), $M^{(\mathcal{L})}$ is the number of weights for node $\mathcal{N}^{(\mathcal{L}-1)}$, $\mathcal{A}^{(\mathcal{L})}$ is the activation function[3] for the layer $\mathcal{L}$, $w_0^{(\mathcal{L},\mathcal{N})}$ is the bias in the node $\mathcal{N}$ of the layer $\mathcal{L}$ and $w_i^{(\mathcal{L},\mathcal{N})}$ ($i = 1, \ldots, M^{(\mathcal{L})}$) are the weights in the node $\mathcal{N}$ of layer $\mathcal{L}$. The notation $\mathcal{R}^{(\mathcal{L},\mathcal{N})}$ in equation (7) represents the non-activated response in the current paper and, hence, $\mathcal{A}^{(\mathcal{L})}(\mathcal{R}^{(\mathcal{L},\mathcal{N})})$ corresponds to the activated response. We would like to clarify to the reader, that there are two common ways [7] to write the equation for the response on a hidden node of a DANN using post and pre-activation function. The post-activation function applies directly to the inputs from the previous layer before computing the weighted sum (as in equation (7)), while the pre-activation function applies to the weighted sum of inputs from the previous layer. However, once the pre-activation formulation is utilized, the output of DANN is generated without use of an activation function. Despite these two different ways of writing the equation, they ultimately lead to the same formal representation of the response as a function of the inputs. This is because the two formulations are mathematically equivalent and can be transformed into each other using simple algebraic manipulations.

There is a variety of non-linear (and also linear) functions that are commonly applied as activation functions $\mathcal{A}^{(\mathcal{L})}$ for an arbitrary input $\mathcal{I}$. The most popular activation functions choice are the sigmoid in equation (8), the hyperbolic tangent in equation (9) and the rectified linear unit in equation (10):

$$\mathcal{A}_{sig}^{(\mathcal{L})}(\mathcal{I}) = \frac{1}{1 + e^{-\mathcal{I}}}. \tag{8}$$

---

[3]The activation function could be also specified for each individual node $\mathcal{N}$ in the layer $\mathcal{L}$ as $\mathcal{A}^{(\mathcal{L},\mathcal{N})}$, but in order to keep transparency for the reader we keep the formulation where the activation function $\mathcal{A}^{(\mathcal{L})}$ is the same for all nodes of the layer

$$\mathcal{A}^{(\mathcal{L})}_{tanh}(\mathcal{I}) = \frac{e^{\mathcal{I}} - e^{-\mathcal{I}}}{e^{\mathcal{I}} + e^{-\mathcal{I}}}. \tag{9}$$

$$\mathcal{A}^{(\mathcal{L})}_{ReLU}(\mathcal{I}) = \max(\mathcal{I}, 0). \tag{10}$$

The decision to choose a specific function depends heavily on the prediction task and the data type [94]. Furthermore, care has to be taken when choosing the activation functions, as it can lead to the so-called vanishing or exploding gradient problems during DANN training [3].

The weights can be seen as a vector $\boldsymbol{w} = \left\{ w_i^{(\mathcal{L},\mathcal{N})}, i = 1, \ldots, N_w \right\}$ with the total number of weights $N_w$ that depends on the number of layers $L$, the number of hidden neurons $N^{(\mathcal{L})}$ per layer $\mathcal{L}$ ($\mathcal{L} = 1, \ldots, L$) and the number of weights $M^{(\mathcal{L})}$ per hidden neurons $N^{(\mathcal{L})}$ of the layer $\mathcal{L}$. The weights $w_i^{(\mathcal{L},\mathcal{N})}$ of the DANN define the final form of the representation $\mathcal{R}(\boldsymbol{\omega})$ in equation (6), and they are determined via a training procedure. More details of the training procedure are discussed in Section 2.4.

### 2.2.2 Connection between DANN and aPC

According to equation (7), the conventional structure of DANNs propagates the signal from inputs to the response through deep layers, where each node of the layer employs a linear combination of zeroth and first-order monomials. The monomials represent the activated output from the nodes of the previous layer. Considering the definitions in Section 2.1, the linear representation via monomials is a particular case of the PCE theory, where the polynomial basis of $0^{th}$ degree $\phi_j^{(0)}$ and $1^{st}$ degree $\phi_j^{(1)}$ is defined explicitly as:

$$\phi_j^{(0)}(\omega_*) = m_0^{(0)}, \quad \phi_j^{(1)}(\omega_*) = m_0^{(1)} + m_1^{(1)}\omega_i, \tag{11}$$

where $\omega_i$ is some input of a hidden node $\mathcal{N}$ in a layer $\mathcal{L}$.

As it have been stated in Section 2.1.1, the PCE theory requires that the set of polynomials (as well in equation 11)) forms an orthonormal basis for each input $\omega_j$ in the input space $\Omega$, i.e. satisfying equation (3). For example, the original theory of homogeneous chaos [113] is based on Hermite polynomials that are orthonormal for Gaussian distributed inputs. The use of any other basses for Gaussian distributed inputs will result in an erroneous non-orthogonal decomposition and considered to be not optimal [75, 104]. The aPC theory [75] generalizes the original PCE theory by allowing for arbitrary input distributions and constructs the orthonormal polynomial basis from the available data-driven input distribution that is encoded in raw moments. However, the conventional DANN imposes the basis via a particular form of $0^{th}$ and $1^{st}$ order polynomial in equation (11), but we do not know for which underlying distribution that polynomial basis would satisfy the orthonormality conditions, i.e. whether it is optimal. To find out the underlying orthonormal distribution for the conventional DANN that satisfies equation (3), we will explore equations (18) and (22) from the paper [79], that are written as following:

$$\begin{bmatrix} m_0^{(0)} & 0 & \cdots & 0 \\ m_0^{(1)} & m_1^{(1)} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ m_0^{(\alpha_j)} & m_1^{(\alpha_j)} & \cdots & m_{\alpha_j}^{(\alpha_j)} \end{bmatrix} \begin{bmatrix} \mu_0 \\ \mu_1 \\ \vdots \\ \mu_{\alpha_j} \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

and

$$\begin{bmatrix} m_0^{(0)} & 0 & \cdots & 0 \\ m_0^{(1)} & m_1^{(1)} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ m_0^{(\alpha_j)} & m_1^{(\alpha_j)} & \cdots & m_{\alpha_j}^{(\alpha_j)} \end{bmatrix} \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_{\alpha_j+1} \end{bmatrix} = \begin{bmatrix} m_0^{(0)}\mu_1 \\ 1 \\ \vdots \\ 0 \end{bmatrix}.$$

Solving the system of linear equations presented above, we can reconstruct the first two raw moments behind the polynomial basis in equation (11) that are used in the conventional DANN. This process entails:

$$\mu_1 = 0, \mu_2 = 1. \tag{12}$$

Therefore, we can conclude that conventional DANN structures implicitly assume a Gaussian distribution with zero mean and unit variance (i.e. standard Gaussian distribution) for propagating the signal through hidden layers. In other

words, the conventional DANN representation optimally preserves orthonormality, if and only if, all neural inputs of all hidden layers were standard Gaussian. However, if the inputs of all hidden layers do not follow a standard Gaussian distribution, then the conventional DANN structure commonly used in machine learning may not be the most optimal way to propagate the signals through the hidden layers. Applied ML tasks could often be used in situations where the propagation of neural signals from layer to layer is not necessarily Gaussian, let alone standardized to unit variance. Therefore, based on results presented in [75] and [79], we conclude that the linear representation via non-orthonormal monomials is not optimal. Not even batch normalization [41] of each node's input could mitigate the mentioned effect due to its linear nature. Theoretically, some general non-linear transformation could map the distributions of all neural inputs onto the Gaussian, but this is not feasible due to the data-driven nature of such neural inputs.

## 2.3 Deep Arbitrary Polynomial Chaos Neural Network

Let us generalize the structure of conventional DANNs in Section 2.2 to overcome their redundancy caused by the non-orthonormal representation of neural signals. To do so, we will employ the orthogonal representation via the data-driven theory of polynomial chaos expansion introduced in Section 2.1. This also introduces the possibility to consider additional high-order interactions between neurons through the non-linear multivariate terms from the PCE representation in the conventional DANN structure in equation (7).

### 2.3.1 Deep orthonormal polynomial representation

We will consider the number of deep layers $L$ and the corresponding number of neurons $N^{(\mathcal{L})}$ for each layer $\mathcal{L}$ as in Section 2.2. Similar to Section 2.1 and Section 2.2, we will map the multi-dimensional input $\boldsymbol{\omega} = \{\omega_1, \ldots, \omega_n\}$ to a response $\mathcal{R}(\boldsymbol{\omega})$. Combining the deep ML representation in equation (6) with the orthonormal expansion in equation (1), we will construct a generalized representation denoted as Deep Arbitrary Polynomial Chaos Neural Network (DaPC NN). Here, the response $\mathcal{R}^{(\mathcal{L},\mathcal{N})}$ of the node $\mathcal{N}$ in the hidden layer $\mathcal{L}$ forms a vector of layer responses $\boldsymbol{\mathcal{R}}^{(\mathcal{L})} = \left\{ \mathcal{R}^{(\mathcal{L},1)}, \ldots, \mathcal{R}^{(\mathcal{L},N^{\mathcal{L}})} \right\}$ as follows:

$$
\begin{aligned}
\mathcal{R}^{(\mathcal{L},\mathcal{N})} \left( \boldsymbol{\mathcal{R}}^{(\mathcal{L}-1)} \right) = \sum_{i=0}^{M^{(\mathcal{L})}} w_i^{(\mathcal{L},\mathcal{N})} \Psi_i^{(\mathcal{L})} \Big[ \mathcal{A}^{(\mathcal{L})}(\mathcal{R}^{(\mathcal{L}-1,1)}), \\
..., \mathcal{A}^{(\mathcal{L})}(\mathcal{R}^{(\mathcal{L}-1,N^{\mathcal{L}-1})}) \Big],
\end{aligned}
\tag{13}
$$

where $\boldsymbol{\mathcal{R}}^{(\mathcal{L}-1)} = \left\{ \mathcal{R}^{(\mathcal{L}-1,1)}, \ldots, \mathcal{R}^{(\mathcal{L}-1,N^{\mathcal{L}-1})} \right\}$ is the neural response from the previous layer $\mathcal{L}-1$ (again with $\boldsymbol{\mathcal{R}}^{(1)} = \boldsymbol{\omega}$), $M^{(\mathcal{L})}$ is the total number of terms for each node on the layer $\mathcal{L}$, $\mathcal{A}^{(\mathcal{L})}$ is an activation function for the layer $\mathcal{L}$, $\Psi_i^{(\mathcal{L})}$ is a multivariate orthonormal (w.r.t. the probability distribution given by response of the previous layer) polynomial from the basis $\left\{ \Psi_0^{(\mathcal{L})}, \ldots, \Psi_{M^{(\mathcal{L})}}^{(\mathcal{L})} \right\}$ of degree $d^{(\mathcal{L})}$ for the layer $\mathcal{L}$ and $w_i^{(\mathcal{L},\mathcal{N})}$ are the weights of the node $\mathcal{N}$ in layer $\mathcal{L}$, where the term $w_0^{(\mathcal{L},\mathcal{N})}$ represents bias as in the conventional DANN definition in equation (7).

The DaPC NN representation in equation (13) reflects the non-linear interaction between the neurons using high-order multivariate terms and orthonormal representation in contrast to DANNs in equation (7). According to Section 2.1, the total number of weights $M^{(\mathcal{L})}$ on each node $\mathcal{N}$ of the layer $\mathcal{L}$ depends on the number of layer inputs $N^{(\mathcal{L}-1)}$ from the previous layer and on the desired degree $d^{(\mathcal{L})}$ of the polynomial representation for the layer $\mathcal{L}$ as:

$$
M^{(\mathcal{L})} = \frac{(N^{(\mathcal{L}-1)} + d^{(\mathcal{L})})!}{N^{(\mathcal{L}-1)}! d^{(\mathcal{L})}!}.
\tag{14}
$$

To ensure the optimal transfer of neural signals from layer to layer and to mitigate redundancy within each node by the DaPC NN representation, we will construct the multivariate orthonormal polynomial basis $\left\{ \Psi_0^{(\mathcal{L})}, \ldots, \Psi_{M^{(\mathcal{L})}}^{(\mathcal{L})} \right\}$ of degree $d^{(\mathcal{L})}$ for each layer $\mathcal{L}$ depending on the layer input, i.e. depending on the activated response from the previous layer $\mathcal{A}^{(\mathcal{L})}(\boldsymbol{\mathcal{R}}^{(\mathcal{L}-1)})$.

As the input layer ($\mathcal{L} = 1$) directly corresponds to the inputs $\boldsymbol{\omega}$ (i.e. $\boldsymbol{\mathcal{R}}^{(1)} = \boldsymbol{\omega}$) and the inputs $\boldsymbol{\omega}$ follow some arbitrary (but given by each specific application) data-driven distribution, the response $\boldsymbol{\mathcal{R}}^{(1)}$ follows the exactly same distribution. For example, the input training data set could be employed in a purely data-driven way to serve as empirical distribution. After that, the responses $\boldsymbol{\mathcal{R}}^{(\mathcal{L})}$ ($\mathcal{L} = 2, .., L$) on the all nodes of the multi-layered structure will follow distributions that result from all previous weights, biases, polynomials and activation functions. That

means, the procedure of orthonormalization proceeds sequentially (forward) through the layers. To maintain the data-driven approach of the orthonormalization procedure, we will utilize the aPC representation for computing the data-driven multi-variate orthonormal basis. This basis is introduced in equations (2), (4) and (5), which will be used to construct the corresponding orthonormal bases $\left\{\Psi_0^{(\mathcal{L})}, \ldots, \Psi_{M^{(\mathcal{L})}}^{(\mathcal{L})}\right\}$ in each layer. Automatically, by ensuring the orthogonal decomposition, the weights $w_i^{(\mathcal{L},\mathcal{N})}$ of a particular node $\mathcal{N}$ and a layer $\mathcal{L}$ gain a meaning according to global sensitivity analysis [74]: the weights reflects the partial contribution of each single neuron (linear univariate terms) or simultaneous combination of neurons (non-linear multivariate terms) to the total variance of the response $\mathcal{R}^{(\mathcal{L},\mathcal{N})}$ for the node $\mathcal{N}$ and the layer $\mathcal{L}$.

The training procedure itself will be discussed in Section 2.4, where the weights $\boldsymbol{w} = \left\{w_i^{(\mathcal{L},\mathcal{N})}, i = 1, \ldots, N_w\right\}$ in equation (13) will be obtained via a similar training procedure as for DANNs. However, independent from any particular training procedure, the weights $\boldsymbol{w}$ determine uniquely the corresponding data-driven orthonormal bases $\left\{\Psi_0^{(\mathcal{L})}, \ldots, \Psi_{M^{(\mathcal{L})}}^{(\mathcal{L})}\right\}$ for each layer $\mathcal{L}$ ($\mathcal{L} = 1, \ldots, L$) of DaPC NNs. Indeed, the orthonormal basis on a particular layers $\mathcal{L}$ depends on the neural response from the previous layer $\boldsymbol{\mathcal{R}}^{(\mathcal{L}-1)}$ only. Remarking that $\boldsymbol{\mathcal{R}}^{(1)} = \boldsymbol{\omega}$, the corresponding data-driven orthonormal bases $\left\{\Psi_0^{(2)}, \ldots, \Psi_{M^{(2)}}^{(2)}\right\}$ for the second layer (i.e. $\mathcal{L} = 2$) can be constructed via equations (2), (4) and (5). The neural response $\boldsymbol{\mathcal{R}}^{(2)}$ on the second layer itself is again fully determined by the corresponding weights through equation (13). Taking into consideration the recursive encapsulation of neural responses in equation (6), it is easy to see that orthonormal bases on all layers are dictated by the weights only. Therefore, all trained weights $\boldsymbol{w}$ determine uniquely the corresponding data-driven orthonormal bases $\left\{\Psi_0^{(\mathcal{L})}, \ldots, \Psi_{M^{(\mathcal{L})}}^{(\mathcal{L})}\right\}$ for each layer $\mathcal{L}$ ($\mathcal{L} = 1, \ldots, L$) of DaPC NN without any additionally actions.

Figure 2 schematically illustrates the structure of the introduced DaPC NN. Similar to DANNs, the structure of DaPC NNs is specified via hidden layers $\mathcal{L}$ ($\mathcal{L} = 1, \ldots, L$), hidden nodes $\mathcal{N}$ ($\mathcal{N} = 1, \ldots, N^{(\mathcal{L})}$) and activation functions $\mathcal{A}^{\mathcal{L}}$. Each layer $\mathcal{L}$ is equipped with an orthonormal basis $\left\{\Psi_0^{(\mathcal{L})}, \ldots, \Psi_{M^{(\mathcal{L})}}^{(\mathcal{L})}\right\}$.Basically, neural signals traveling from layer to layer should pass through a sort of data-driven filter that constructs an optimal orthonormal representation for each layer as illustrated in Figure 2. Such an orthonormal basis is the same for all nodes $\mathcal{N}$ ($1, \ldots, N^{(\mathcal{L})}$) of a given layer $\mathcal{L}$ and employs the neural signal coming as response $\mathcal{R}^{(\mathcal{L}-1,\mathcal{N})}$ from the previous layer $\mathcal{L} - 1$. Moreover, the suggested DaPC NN structure offers flexibility to specify the degree of non-linearity $d^{(\mathcal{L})}$ for each particular layer $\mathcal{L}$ ($\mathcal{L} = 1, \ldots, L$) in order to go beyond a linear representation of univariate neurons, if desired.

### 2.3.2 Properties of the DaPC NN

Let us introduce relevant properties of the DaPC NN that could be useful for practical applications.

**Property 1:** Due to the orthonormal DaPC NN representation, the expected value $\mu\left(\mathcal{R}^{(\mathcal{L},\mathcal{N})}\right)$ and total variance $\sigma^2\left(\mathcal{R}^{(\mathcal{L},\mathcal{N})}\right)$ of the response $\mathcal{R}^{(\mathcal{L},\mathcal{N})}$ on each particular node $\mathcal{N}$ of any layer $\mathcal{L}$ (including the last layer forming the total response $\mathcal{R}$) can be quantified analytically [75] using the following explicit form:

$$\mu\left(\mathcal{R}^{(\mathcal{L},\mathcal{N})}\right) = w_0^{(\mathcal{L},\mathcal{N})}, \tag{15}$$

$$\sigma^2\left(\mathcal{R}^{(\mathcal{L},\mathcal{N})}\right) = \sum_{i=1}^{M^{(\mathcal{L})}} \left(w_i^{(\mathcal{L},\mathcal{N})}\right)^2, \tag{16}$$

where the explicit analytical relations in equation (15) are written with respect to the response $\boldsymbol{\mathcal{R}}^{(\mathcal{L}-1)}$ from the previous layer $\mathcal{L}$ and not with respect to the inputs $\boldsymbol{\omega}$.

**Property 2:** Due to the orthonormal DaPC NN representation, Sobol [101] sensitivity indices $S_I^{(\mathcal{L},\mathcal{N})}$ of a particular node $\mathcal{N}$ and layer $\mathcal{L}$ can be explicitly computed as follows according to global sensitivity analysis [74, 103]:

$$S_I^{(\mathcal{L},\mathcal{N})} = \frac{\left(w_I^{(\mathcal{L},\mathcal{N})}\right)^2}{\sum_{i=1}^{M^{(\mathcal{L})}} \left(w_i^{(\mathcal{L},\mathcal{N})}\right)^2}, \tag{17}$$

where the Sobol index $S_I^{(\mathcal{L},\mathcal{N})}$ reflects the relative partial contribution of each single neuron (linear univariate terms) or simultaneous combination of neurons (non-linear multivariate terms) to the total variance of the response $\mathcal{R}^{(\mathcal{L},\mathcal{N})}$ for that particular node $\mathcal{N}$ in a layer $\mathcal{L}$.

Figure 2: An illustration of the DaPC NN structure. The input of layer $\mathcal{L}$ is formed by the output of layer $\mathcal{L}-1$ and the orthonormal basis $\left\{\Psi_0^{(\mathcal{L})}, \ldots, \Psi_{M^{(\mathcal{L})}}^{(\mathcal{L})}\right\}$ of the polynomial degree $d^{(\mathcal{L})}$.

**Property 3:** An analytical form for the partial derivatives of the DaPC NN representation with respect to each particular weight can be obtained by applying the chain rule of differentiation [35]. Considering the recursive encapsulation in equation (6) and the orthonormal representation in equation (13), we obtain the following analytical form for the partial derivative [4] of the response $\mathcal{R}^{(\mathcal{L}, \mathcal{N})}$ in the node $\mathcal{N}$ of the layer $\mathcal{L}$ with respect to a weight $w_k^{(\mathcal{L}_*, \mathcal{N}_*)}$ ($\forall k \in [1, M^{(\mathcal{L}_*)}]$) in the node $\mathcal{N}_*$ of the layer $\mathcal{L}_*$ when $\mathcal{L} > \mathcal{L}_*$:

$$
\frac{\partial \mathcal{R}^{(\mathcal{L}, \mathcal{N})}}{\partial w_k^{(\mathcal{L}_*, \mathcal{N}_*)}} = \sum_{i=0}^{M^{(\mathcal{L})}} w_i^{(\mathcal{L}, \mathcal{N})} \sum_{j=0}^{N^{(\mathcal{L}-1)}} \left( \frac{\partial \Psi_i^{(\mathcal{L})}}{\partial \mathcal{A}^{(\mathcal{L})}} \cdot \right.
$$
$$
\left. \cdot \frac{\partial \mathcal{A}^{(\mathcal{L})}}{\partial \mathcal{R}^{(\mathcal{L}-1,j)}} \frac{\partial \mathcal{R}^{(\mathcal{L}-1,j)}}{\partial w_k^{(\mathcal{L}_*, \mathcal{N}_*)}} \right)
$$

(18)

and the following simplified analytical form when $\mathcal{L} = \mathcal{L}_*$:

$$
\frac{\partial \mathcal{R}^{(\mathcal{L}_*, \mathcal{N}_*)}}{\partial w_k^{(\mathcal{L}_*, \mathcal{N}_*)}} = \Psi_k^{(\mathcal{L}_*)} \Big[ \mathcal{A}^{(\mathcal{L}_*)}(\mathcal{R}^{(\mathcal{L}_*-1,1)}),
$$
$$
..., \mathcal{A}^{(\mathcal{L}_*)}(\mathcal{R}^{(\mathcal{L}_*-1, N^{\mathcal{L}_*-1})}) \Big],
$$

(19)

where, similar to above, the same analytical relations for partial derivatives hold for the DaPC NN response $\mathcal{R}$ as the response of the last layer $\boldsymbol{\mathcal{R}}^{(L)}$.

**Remark 1:** The data-driven aPC representation $\mathcal{R}_{\mathrm{aPC}}(\boldsymbol{\omega})$ in Section 2.1 is a particular case of a DaPC NN representation $\mathcal{R}_{\mathrm{DaPCNN}}(\boldsymbol{\omega})$ when the number of hidden layers is equal to one:

$$
\mathcal{R}_{\mathrm{aPC}}(\boldsymbol{\omega}) = \mathcal{R}_{\mathrm{DaPCNN}}(\boldsymbol{\omega}), \quad \text{for } \mathcal{L} = 1.
$$

(20)

---

[4] Partial derivatives are functions of the inputs $\boldsymbol{\omega}$

**Remark 2:** The conventional DANN representation $\mathcal{R}_{\text{DANN}}(\boldsymbol{\omega})$ in Section 2.2 is a particular case of the DaPC NN representation $\mathcal{R}_{\text{DaPCNN}}(\boldsymbol{\omega})$ when the degree of expansion in each layer is equal to one and the basis is defined according to equation (11) for a standard Gaussian distribution of neuronal signals in each hidden node:

$$\mathcal{R}_{\text{DANN}}(\boldsymbol{\omega}) = \mathcal{R}_{\text{DaPCNN}}(\boldsymbol{\omega}), \text{ for } d^{\mathcal{L}} = 1, \forall \mathcal{L} = 1, \ldots, L. \tag{21}$$

**Remark 3:** The DaPC NN offers a possibility to reflect non-linearities of neural responses both by introducing the high-order expansion $d^{(\mathcal{L})}$ for the responses $\mathcal{R}^{(\mathcal{L}, \mathcal{N})}$ of the hidden nodes and by introducing non-linearity in the activation functions $\mathcal{A}^{(\mathcal{L})}$. Theoretically, there are no restrictions on the use of activation functions within the current DaPC NN framework according to equation (7). The modeler has the freedom to determine an activation function if it would be advantageous for a specific modeling procedure since equations (18) and (19) provide the analytical form of the partial derivatives that include the activation function. However, simultaneously introducing non-linearities both in expansion degrees and in activation functions could lead to difficulties in network training caused by strong numerical round-off errors. We would rather suggest to employ linear activation functions while shifting the representation of non-linear dependence to the orthonormal decomposition. Especially, any linear normalization (similar to batch normalization [41]) will mitigate potential numerical challenges during the training procedure. Therefore, in the current paper we suggest to employ the analytical relations introduced in *Property 1* to have a handle on total (explained) variances for the overall network. This can be achieved via the following normalized activation function:

$$\mathcal{A}^{(\mathcal{L})}(\mathcal{I}) = \frac{\mathcal{I} - \mu(\mathcal{I})}{\sigma(\mathcal{I})}. \tag{22}$$

## 2.4 Computation of weights

### 2.4.1 Training data

The exact forms of aPC in Section 2.1, DANN in Section 2.2 and DaPC NN in Section 2.3 are determined via corresponding weights that are typically computed using training data in a training procedure. Let us denote the training inputs by $\boldsymbol{\omega}_T = \{\boldsymbol{\omega}_{T_1}, \ldots, \boldsymbol{\omega}_{T_N}\}$ and the corresponding training responses by $\mathbf{R}_{T_N} = \{R_{T_1}, \ldots, R_{T_N}\}$, where $T_N$ is a value greater than zero representing the size of training data sets $\boldsymbol{D}_{T_N} = \{(\boldsymbol{\omega}_{Ti}, R_{Ti}), i = 1, \ldots, T_N\}$. Such data could be obtained from runs of an original physical model for surrogate modelling [75, 15], or directly from a given database for other learning tasks [28, 56, 40]. In order to assure transparency in the current paper, we will consider a fixed training set of size $T_N$, avoiding any data mining procedure such as active learning [78, 93]. Thus, the data set $\boldsymbol{D}_{T_N}$ is the complete training set for the discussed ML approaches.

### 2.4.2 Training procedure

Deep structures like conventional DANNs in Section 2.2 or the newly introduced DaPC NN in Section 2.3 require solving a non-linear system of equations to obtain the unknown weights $w_i^{(\mathcal{L}, \mathcal{N})}$. There are two possible ways to solve the corresponding system of non-linear equations. The first way is based on deterministic approaches, such as gradient-based search [88] or the Levenberg-Marquardt algorithm [61] that have been widely used in deep learning. In the current paper, we will follow such a popular way for training of DaPC NN to ensure that the training process is comparable to the traditional practices for conventional DANN structures. However, for the sake of completeness we would like to mention, that the underlying problem is ill-posed, so that a unique deterministic solution to the training problem may not exist [2, 121]. Therefore, the second way helps to solve the challenges related to ill-posedness of the problem. It is based on stochastic inference [51] by seeing the weights as random variables, and then conditioning the weights $w_i^{(\mathcal{L}, \mathcal{N})}$ on the available training data [83]. However, straightforward stochastic approaches such as Monte Carlo [99] or even Markov chain Monte Carlo [34] are computationally very expensive, and they suffer in cases with high non-linearity and high-dimensionality [81]. Recently, there are some works trying to overcome the computational problem of stochastic inference assuming Gaussianity and mutual independence of weights [17]. Alternatively, combinations of deterministic approaches with a regularization helping to deal with overfitting [35] have been very popular in the last decades for DANNs.

Therefore, in the current paper, we will adopt the widely used approach, and extend the optimization procedure by Tikhonov or ridge regularization [106], to ensure that the DaPC NN training is consistent with the current practice for conventional DANN structures. o ensure that the training of the DaPC NN is consistent with the current methodology used for conventional DANN structures. Thus, we consider the following optimization problem regarding the unknown

weights forming the vector $\boldsymbol{w} = \left\{ w_i^{(\mathcal{L},\mathcal{N})}, i = 1, \ldots, N_w \right\}$ for all neurons distributed within the all layers :

$$
\boldsymbol{w} = \arg\min_{\boldsymbol{w}} \left[ \frac{1}{T_N} \sum_{i=1}^{T_N} \left( \mathcal{R}(\boldsymbol{\omega_{T_i}}) - R_{T_i} \right)^2 \right.
$$
$$
\left. + \frac{1}{N_w} \sum_{i=1}^{N_w} \left( w_i^{(\mathcal{L},\mathcal{N})} \right)^2 \right],
$$

(23)

where the weights should be optimized simultaneously to minimize the so-called loss function [35] in equation (23). The first term in equation (23) reflects the deviation of the response $\mathcal{R}$ from the training data via the mean squared error (MSE). The second term corresponds to a regularization known as mean square weights (MSW) [5]. Other loss functions are also known in the literature, depending on the overall objective such as a Bayesian interpolation [60], physical regularization [47, 84] or classification tasks [13, 56] . Moreover, the scalability to larger datasets relies on the chosen architecture and training methods. The effective model complexity proposed in [70] identifies specific regimes where increasing the number of training samples may negatively impact performance.

In order to directly compare the DaPC NN structure with the conventional DANN under identical conditions, we choose the Levenberg-Marquardt algorithm [61] to solve [6] equation (23) for all structures due to its robustness. We also use the exact same training procedure for all structures using the exactly same loss-function in equation (23) and the exactly same training data sets $\boldsymbol{D}_{T_N}$. To provide gradients to the Levenberg-Marquardt algorithm, we use the analytical derivatives for DaPC NN introduced in *Property 3* of Section 2.3.2, whereas we use typical backpropagation to provide gradients. Thus, the orthogonal bases in DaPC NNs are updated simultaneously with the weights during the training procedure (see also Figure 2). Moreover, as pointed out in Section 2.3.1, the weights determine uniquely the corresponding data-driven orthonormal bases in equations (2) and (2) for each layer of DaPC NNs by solving the linear system of equations (5) introduced in Section 2.1. Consequently, the DaPC NN can be seen as a black box, where only the unknown weights need to be determined using one or another training procedure. This makes it easy for users to operate the DaPC NN in a similar way to the conventional DANN, without requiring them to have any special knowledge of its internal workings.

The training procedure of ML approaches with single-layer ML representations [29] like the PCE usually determines the unknown weights by solving a linear system of equations. Therefore, the unknown weights $w_i$ of the aPC representation in Section 2.1 are directly obtained by solving a linear system of equations (5) (see details in [73, 12, 11]). For numerical robustness, we use Moore-Penrose inversion [68, 82, 14] (also known as pseudoinverse) to solve the regularized [106] least-squares problem of equation (23) as in the original aPC work [75] that is available in Matlab file exchange [76].

In the current paper, we focus on fundamental issues of deep neuronal network structures exploring the concept of orthonormal decomposition and the possibility to introduce high-order neural interaction, rather than delving into specifications of the loss function or training algorithm. Therefore, the reader is invited to adopt the introduced DaPC NN structure for own needs, introducing any own specific loss function or training procedure. The DaPC NN is available online for the reader through Matlab file exchange [77].

### 2.4.3   Technical remark on DaPC NN bases

Technically, the training procedure specified in equation (23) remains the same for the DaPC NN as well as for the DANN and the aPC, and the authors of the current paper do not aim at any novelty here. However, we would like to underline that the DaPC NN constructs its data-driven orthonormal bases implicitly during the training procedure, based on the input distribution and upon the responses of hidden nodes that react to adjustments of the weights during the training procedure. This means, the orthonormal bases on all layers are adaptively recomputed during the iterative training procedure while searching for the optimal weights. It could be seen as re-adjusting of orthonormal bases simultaneously within the iteration procedure over the weights.

## 3   Illustration of Performance and Comparison

In the previous Section 2, we have introduced the DaPC NN using the data-driven orthonormal decomposition from aPC theory. The current section will illustrate the performance of the suggested DaPC NN, comparing it with the

---

[5]See also Bayesian regularized neural networks in [71])

[6]We are aiming to not impose any preferences on the training algorithms in the current paper and we encourage the developers of training algorithms to incorporate their approaches to the DaPC NN framework.

performance of a conventional DANN structure and also with an aPC expansion. To do so, we will employ exactly the same training data set $\boldsymbol{D}_{T_N} = \left\{ (\boldsymbol{\omega}_{Ti}, R_{Ti}), i = 1, \ldots, T_N \right\}$ of the size $T_N$ for all three ML approaches. The quality of the training itself (e.g. MSE) will not be strongly discussed as the training procedure provides only a small discrepancy between the ML model response $\mathcal{R}(\boldsymbol{\omega}_T)$ and the training data $\mathbf{R}_{T_N}$. [7] Such a discrepancy partially indicates whether the chosen architecture, loss function and training algorithm are appropriate for analysed problem. Instead, we rather focus on the prediction ability for a separate data set that has not been used during the training procedure. Therefore, we will employ a validation data set [44] $\boldsymbol{D}_{V_N} = \left\{ (\boldsymbol{\omega}_{Vi}, R_{Vi}), i = 1, \ldots, V_N \right\}$ of the size $V_N$ with the validation inputs $\boldsymbol{\omega}_{V_N} = \{ \boldsymbol{\omega}_{V_1}, \ldots, \boldsymbol{\omega}_{V_N} \}$ and the corresponding validation responses $\mathbf{R}_{V_N} = \{ R_{V_1}, \ldots, R_{V_N} \}$.

Because the quality of the prediction is strongly dependent on the size $T_N$ of the available training data set [75, 39], we will assess the prediction quality for validation data sets $\boldsymbol{D}_{V_N}$ under various sizes of the training data set $T_N$. To measure the prediction ability, we will assess the mean square error, which is usually used in ML, as well as a weighted error that will be defined below in eq. (24),(25),(26), which is used in the uncertainty quantification community. The mean square error could be also normalized by the variance, converting it to the coefficient of determination if desired. The corresponding estimates of the mean square error $MSE_{T_N} [\mathcal{R}(\boldsymbol{\omega}_{V_N})]$, the relative error of mean $E_{T_N}^{\mu} [\mathcal{R}(\boldsymbol{\omega}_{V_N})]$ and the relative error of standard deviation $E_{T_N}^{\sigma} [\mathcal{R}(\boldsymbol{\omega}_{V_N})]$ are defined as:

$$\text{MSE}_{T_N} [\mathcal{R}(\boldsymbol{\omega}_{V_N})] = \frac{1}{V_N} \sum_{i=1}^{T_N} \left( \mathcal{R}(\boldsymbol{\omega}_{V_i}) - R_{V_i} \right)^2 , \tag{24}$$

$$E_{T_N}^{\mu} [\mathcal{R}(\boldsymbol{\omega}_{V_N})] = \frac{\mu [\mathcal{R}(\boldsymbol{\omega}_{V_N})] - \mu [\mathbf{R}_{V_N}]}{\mu [\mathbf{R}_{V_N}]} , \tag{25}$$

$$E_{T_N}^{\sigma} [\mathcal{R}(\boldsymbol{\omega}_{V_N})] = \frac{\sigma [\mathcal{R}(\boldsymbol{\omega}_{V_N})] - \sigma [\mathbf{R}_{V_N}]}{\sigma [\mathbf{R}_{V_N}]} , \tag{26}$$

where $\mu [\cdot]$ is the expected value and $\sigma [\cdot]$ is the standard deviation over the validation data set. The expected value $\mu [\cdot]$ and standard deviation $\sigma [\cdot]$ are obtained via the empirical mean and standard deviation over the validation set $\boldsymbol{D}_{V_N}$.

Because we are interested in the overall reliability of the ML response and to guarantee fairness in comparison, we omit the powerful property of aPC in computing the mean and variance without additional evaluations of the response. Instead, we will estimate them numerically using the final aPC response constructed on the validation data set, similarly to DANN and DaPC NN.

As test cases, we use examples of functional approximation as in surrogate modeling. To test different aspects, we consider an Ishigami problem with three inputs [42] in Section 3.2, an ON-10 problem with ten inputs [80] in Section 3.3 and also a carbon dioxide benchmark problem with non-linear shock propagation [52] in Section 3.4.

### 3.1 Architecture and technical specification of ML models

Architectures of the considered ML models used in the current and upcoming test cases are specified via the total number of layers $L$, corresponding numbers of nodes $N^{(\mathcal{L})}$ for each layer $\mathcal{L}$ as $[N^{(1)}, \ldots, N^{(L)}]$, degrees of non-linearity $d^{(\mathcal{L})}$ for on each layer $\mathcal{L}$ as $[d^{(1)}, \ldots, d^{(L)}]$, activation functions $\mathcal{A}^{(\mathcal{L})}$ (same for all layers $\mathcal{L}$) and loss function ($\mathcal{LF}$). For the sake of transparency, we will also provide the total number of unknown weights or coefficients as $N_w$ in upcoming Sections. This is particularly useful for aPC and DaPC NN, as the number of unknown weights may not be as intuitive as for conventional DANN. We will adopt the architectures of the analyzed ML models taking into consideration the size $T_N$ of the training data sets. Nevertheless, we would like to emphasize that the setup of the corresponding architectures are the responsibility of the ML modeler. For our study, we selected the architectures with equal care and effort for all three compared approaches, trying to achieve the most reliable results. Nevertheless, the degrees of freedom in choosing the number of layers, the number of nodes per layer, the degree of non-linearity (if allowed) and the activation function pose strong challenges onto the modeling procedure, requiring a deep understanding of each underlying ML approach. For instance, in addition to the standard settings for the conventional DANN architecture, the DaPC NN architecture incorporates additional degrees of freedom by choosing the non-linearity degree on specific layers, which introduces unknown weights associated with the number of input neurons from the previous layer, as discussed in Section 2.2. Exploring all possible configurations of deep architectures using a trial-and-error approach is

---
[7]

not feasible due to the vast number of combinations. Consequently, developing an adaptive strategy for setting up the architecture could significantly improve the performance of deep multi-layer representations. This approach would require further research to minimize the potential for subjectivity in the modeling procedure. The authors of the current paper recognize the potential of the Bayesian framework for optimizing hyper-parameters such as the number of layers, number of nodes, and degree of non-linearity. However, direct Bayesian analysis seems to be computationally intensive, and it may be necessary to rely on approximate indicators that involve certain assumptions [80]. However, this topic goes beyond the scope of the current paper but could be the subject of future research.

We encourage the reader to test the suggested DaPC NN approach and also known aPC and DANN approaches for own needs. The conventional DANN approach used here can be found under Matlab software [62] using the fitnet functionality. The aPC and DaPC NN approaches are available online for the reader through Matlab file exchange [76] and [77], correspondingly.

### 3.2 Ishigami test case

Table 1: Architectures of aPC, DANN and DaPC NN for Ishigami test case.

| ML model | $T_N$ | $L$ | $N^{(\mathcal{L})}$ | $d^{(\mathcal{L})}$ | $\mathcal{A}^{(\mathcal{L})}$ | $N_w$ | $\mathcal{LF}$ |
|---|---|---|---|---|---|---|---|
| aPC | 10 | 1 | 1 | 2 | none | 10 | MSE+MSW |
| | 100 | 1 | 1 | 4 | none | 35 | MSE+MSW |
| | 1000 | 1 | 1 | 6 | none | 84 | MSE+MSW |
| DANN | 10 | 3 | [6,6,1] | [1,1,1] | eq. (9) | 78 | MSE+MSW |
| | 100 | 4 | [9,6,3,1] | [1,1,1,1] | eq. (9) | 127 | MSE+MSW |
| | 1000 | 4 | [10,8,6,1] | [1,1,1,1] | eq. (9) | 189 | MSE+MSW |
| DaPC NN | 10 | 2 | [3,1] | [2,2] | eq. (22) | 40 | MSE+MSW |
| | 100 | 3 | [3,3,1] | [2,2,2] | eq. (22) | 70 | MSE+MSW |
| | 1000 | 3 | [3,3,1] | [3,3,2] | eq. (22) | 130 | MSE+MSW |

#### 3.2.1 Problem set up

As the first test case, we will employ the widely used Ishigami function [42], as it shows strong non-linearity accompanied with non-monotonicity:

$$\mathcal{M}(\boldsymbol{\omega}) = \sin(\omega_1) + a\sin^2(\omega_2) + b\omega_3^4 \sin(\omega_1), \tag{27}$$

where we will use the particular case with $a = 7$ and $b = 0.1$. The distribution of the three input random variables $\boldsymbol{\omega}$ is given by mutually independent uniform distributions with $\omega_i \sim \mathcal{U}(-\pi, \pi)$.



Figure 3: Prediction of aPC, DANN and DaPC NN for Training (left) and Validation (right) data: Ishigami test case with three inputs and the size of training data set equal to 10.

To obtain a training data set, we use Sobol sequences [100] for the underlying distribution of inputs. To analyze how the quality of the prediction depends on the size of the training data set, we use the training data sets $\boldsymbol{D}_{T_N}$ of size $T_N$ equal to 10, 100 and 1000. To assess the quality of prediction, we generate a validation data set $\boldsymbol{D}_{V_N}$ of the size $V_N = 10^3$, via Monte Carlo sampling [99].

Figure 4: Prediction of aPC, DANN and DaPC NN for Training (left) and Validation (right) data: Ishigami test case with three inputs and the size of training data set equal to 100.



Figure 5: Prediction of aPC, DANN and DaPC NN for Training (left) and Validation (right) data: Ishigami test case with three inputs and the size of training data set equal to 1000.

### 3.2.2 Training and Validation

Figures 3-5 demonstrate how predictions made by aPC, conventional DANN and DaPC NN match the training and validation data sets for the size $T_N$ of training data equals to 10, 100 and 1000. In principle, all considered ML approaches show the necessary flexibility to approximate the training data of different sizes, but the DaPC NN shows a superior performance with the smallest scatter in the left plots of Figures 3-5. It is remarkable here that the DaPC NN is equipped with fewer degrees of freedom (number of weights) in comparison to the DANN . Table 1 show the technical specifications for the Ishigami test case. Additionally, Table 1 indicated the total number of unknown weights/coefficients by $N_w$ as for aPC and DaPC NN it could be less intuitive as for conventional DANN. Overall, the selection of the ML model architectures should take into account the size of the training data sets $T_N$ to provide certain flexibility during the training procedure.

However, flexibility alone is not sufficient to make reliable prediction outside of the training data sample and only shows quality of the training approach . Therefore, we pay stronger attention to the right plots in Figures 3-5, reflecting the validation performance. All three approaches are powerless in case of the very small data set used for training (right plot in Figure 3). As expected, increasing the size of training data set to 100 (right plot in Figure 4) or even 1000 (right plot in Figure 5) helps to overcome these issues during the validation phase. In particular, the DaPC NN demonstrates the best performance in the validation phase as well, even with a moderate sample size of training data. Figures 3-5 indicate that the flexibility of DANN tends to have over-fitting issues regarding the regularisation as specified in equation (23). Opposite to theDANN, the DaPC NN shows less issues with over-fitting, inhering the orthonormality from the aPC representation and also the required flexibility from DANN. Due to its definition, the aPC itself seems to have not enough of flexibility to capture the validation data set as well as the DaPC NN. We also have observed in the current and upcoming test cases, that omitting regularization in loss function only slightly degrades the results of DaPC NN and aPC, whereas the conventional DANN shows extremely poor prediction. Hence, pure comparison focusing on least-square solution without regularization seems to be not attractive.

Figure 6: Performance of aPC, DANN and DaPC NN for Ishigami test case with three inputs: Convergence of Mean Square Error, Relative Error of Mean Value and Relative Error of Standard Deviation relatively the reference validation data set.

### 3.2.3   Evidence of Convergence

It Figures 3-5, it appears that DANNs have a stronger spreading of the point clouds and also higher density of spreading in comparison to DaPC NN and aPC models. However, determining the density of spreading based on visual inspection alone can be challenging. Therefore, we will pay closer attention the metrics introduced in the beginning of Section, which provides more clarity on this matter. Figure 6 shows the convergence in terms of mean square error, relative error of the mean value and relative error of the standard deviation as a function of the training data size $T_N$, applied to the validation data set $\boldsymbol{D}_{V_N}$ for aPC expansion, conventional DANN and DaPC NN. The figure reveals that a faster convergence has been reached for DaPC NN in terms of all investigated convergence metrics, summarising the observations made in the previous section. The relative errors of mean and the standard deviation are smaller with the aPC than with the conventional DANN.

### 3.3   ON-10 test case

Table 2: Architectures of aPC, DANN and DaPC NN for ON-10 test case.

| ML model | $T_N$ | $L$ | $N^{(\mathcal{L})}$ | $d^{(\mathcal{L})}$ | $\mathcal{A}^{(\mathcal{L})}$ | $N_w$ | $\mathcal{LF}$ |
|---|---|---|---|---|---|---|---|
| aPC | 100 | 1 | 1 | 2 | none | 66 | MSE+MSW |
| | 500 | 1 | 1 | 3 | none | 286 | MSE+MSW |
| | 1000 | 1 | 1 | 4 | none | 1001 | MSE+MSW |
| DANN | 100 | 4 | [9,6,3,1] | [1,1,1,1] | eq. (9) | 175 | MSE+MSW |
| | 500 | 4 | [15,10,5,1] | [1,1,1,1] | eq. (9) | 386 | MSE+MSW |
| | 1000 | 4 | [15,10,5,1] | [1,1,1,1] | eq. (9) | 386 | MSE+MSW |
| DaPC NN | 100 | 2 | [5,1] | [2,3] | eq. (22) | 386 | MSE+MSW |
| | 500 | 2 | [10,1] | [2,3] | eq. (22) | 946 | MSE+MSW |
| | 1000 | 2 | [10,1] | [2,3] | eq. (22) | 946 | MSE+MSW |

### 3.3.1   Problem set up

As second test case, we will consider a non-linear analytical function $\mathcal{M}(\boldsymbol{\omega})$ of ten ($n = 10$) uncertain inputs $\boldsymbol{\omega} = \{\omega_1, \ldots, \omega_n\}$ from the paper [80]:

$$
\begin{aligned}
\mathcal{M}(\boldsymbol{\omega}) =& (\omega_1^2 + \omega_2 - 1)^2 + \omega_1^2 + 0.1\omega_1 \exp(\omega_2) \\
& + 1 + \sum_{i=3}^{n} \frac{\omega_i^3}{i},
\end{aligned}
\tag{28}
$$

where the inputs $\boldsymbol{\omega}$ in equation (28) are considered to be independent uniformly distributed with $\omega_i \sim \mathcal{U}(-5, 5)$ for $i = 1 \ldots 10$.

For the current test case, we will consider training data sets $\boldsymbol{D}_{T_N}$ of the sizes $T_N$ equal to 100, 500 and 1000. Similar to the previous test case, we will generate our training data according to the Sobol sequence [100] based on the

underlying distributions of inputs $\boldsymbol{\omega}$. As before, we employ a validation data set $\{(\boldsymbol{\omega}_{Vi}, R_{Vi}), i = 1, \ldots, V_N\}$ of the size $V_N = 10^3$ generated by Monte Carlo sampling [99] from the distribution of inputs $\boldsymbol{\omega}$.

### 3.3.2   Prediction and Validation

Figures 7-9 show the predictions made by the aPC, DANN, and DaPC NN compared with the corresponding data during the training and validation phases. Architectures of the considered ML models are presented in Table 2 of Appendix A. With a low number of training data (Figure 7), the DANN and DaPC NN a show better ability to capture the training data compared to the aPC expansion. Here, the aPC demonstrates a lack of flexibility during training. During validation, however, all methods fail to produce consistent results when the small set is used as shown in Figure 7. When trained with more data (Figure 8 and Figure 9), DaPC NN produces predictions with relatively high accuracy compared to aPC and DANN even during validation; where aPC and DANN do not benefit significantly from the increase of training data. Moreover, the orthonormal structure of the DaPC NN representation handles the overfitting issues extremely well. Table 2 presents the technical specifications for the ON-10 test case, including the total number of unknown weights/coefficients denoted by $N_w$ for all ML models. Again, when selecting an appropriate ML model architecture, it is important to consider the size of the training data set $T_N$ that can influence the flexibility of the training process and the predictive power of the model. For example, the number of unknown weights in the DaPC NN can be two times higher than the size of the training data.

It seems that the high non-linearity of the considered test case strongly limits the aPC to obtain the necessary polynomial representation. A high degree of freedom is very helpful in this case. Opposite to that, the DANN seems to suffer significantly from over-fitting, regardless of the regularization in that scenario. Overall, the DaPC NN predictions are more consistent, with lower scatter compared to the other methods for the considered test case.



Figure 7: Prediction of aPC, DANN and DaPC NN for Training (left) and Validation (right) data: ON-10 test case with ten inputs and the size of training data set equal to 100.



Figure 8: Prediction of aPC, DANN and DaPC NN for Training (left) and Validation (right) data: ON-10 test case with ten inputs and the size of training data set equal to 500.

Figure 9: Prediction of aPC, DANN and DaPC NN for Training (left) and Validation (right) data: ON-10 test case with ten inputs and the size of training data set equal to 1000

### 3.3.3 Evidence of convergence

Figure 10 shows convergence in terms of mean square error, relative error of mean value and relative error of standard deviation over the training data size $T_N$. Both the aPC and DANN converge to a similar point, meaning that the performance of both methods is comparable. Figure 10 also indicates that increasing the number of training points does not improve the aPC and DANN performance significantly, again confirming our observation in the previous Section 3.3.2. The DaPCNN, on the other hand, distinctly outperforms the aPC and DANN, with approximately two orders of magnitudes lower prediction errors in validation. The effect of increasing the number of training points can be seen clearly in the DaPCNN performance, showing better learning capacity compared to aPC and DANN.



Figure 10: Performance of aPC, DANN and DaPC NN for ON-10 test case with ten inputs: Convergence of Mean Square Error, Relative Error of Mean Value and Relative Error of Standard Deviation relatively the reference validation data set.

### 3.4 CO$_2$ Benchmark problem



Figure 11: Distributions of injection rate $(m^3/s)$, relative permeability degree (-), and reservoir porosity (-)

19

Table 3: Architectures of aPC, DANN and DaPC NN for $CO_2$ benchmark problem employing Sobol sequences.

| ML model | $T_N$ | $L$ | $N^{(\mathcal{L})}$ | $d^{(\mathcal{L})}$ | $\mathcal{A}^{(\mathcal{L})}$ | $N_w$ | $\mathcal{LF}$ |
|---|---|---|---|---|---|---|---|
| aPC | 100 | 1 | 1 | 3 | none | 20 | MSE+MSW |
| | 500 | 1 | 1 | 5 | none | 56 | MSE+MSW |
| | 1000 | 1 | 1 | 10 | none | 286 | MSE+MSW |
| DANN | 100 | 3 | [8,6,1] | [1,1,1,1] | eq. (9) | 93 | MSE+MSW |
| | 500 | 3 | [8,6,1] | [1,1,1,1] | eq. (9) | 93 | MSE+MSW |
| | 1000 | 4 | [10,8,6,1] | [1,1,1,1] | eq. (9) | 189 | MSE+MSW |
| DaPC NN | 100 | 2 | [3,1] | [2,2] | eq. (22) | 40 | MSE+MSW |
| | 500 | 2 | [3,3] | [2,2] | eq. (22) | 80 | MSE+MSW |
| | 1000 | 2 | [3,1] | [3,3] | eq. (22) | 80 | MSE+MSW |

### 3.4.1 Benchmark set up

As third test case, we will consider a carbon dioxide ($CO_2$) benchmark model that has already been used to compare different ML approaches in the earlier paper [52]. The particularity of that problem consists in a strong shock propagation, where various ML approaches relying on smooth functions could suffer a lot. The problem refers to a multi-phase flow in porous media, where $CO_2$ is injected into a deep aquifer and then spreads in a geological formation. This yields a pressure build-up and a plume evolution. $CO_2$ injection into the subsurface could be a possible practice to mitigate the $CO_2$ emission into the atmosphere. The $CO_2$ benchmark model proposed by Köppel et al. [54] is a reduced version of the model in a benchmark problem defined in the paper [25]. This reduction consists of a radial flow in the vicinity of the injection well, and is made primarily due to the high computational demand of the original $CO_2$ model. It is assumed that the fluid properties such as the density and the viscosity are constant, and all processes are isothermal. The $CO_2$ and the brine build two separate and immiscible phases, and mutual dissolution is neglected. Additionally, the formation is isotropically rigid and chemically inert, and capillary pressure is negligible. We consider the $CO_2$ saturation to be the quantity of interest as a model response that is a function of the coordinates for space $x$ and time $t$ as introduced in [54]. Overall, the considered $CO_2$ benchmark problem is strongly non-linear because the $CO_2$ saturation spreads as a strongly non-linear front that could be challenging to capture via surrogates. For detailed information on the governing equations, the modeling assumptions and numerical approaches, the reader is referred to the original publication [54].

Following the comparison study [54], we consider the combined effects of three sources of uncertainty. We take into account the uncertainty of boundary conditions via the injection rate, the uncertainty of constitutive relation introduced via in the relative permeability definitions and the uncertainty of material properties represented by the porosity of the geological formation. Figure 11 shows the distribution of these three inputs taken from the original data set [53].

Similar to the previous test cases, we will train the aPC, DANN and DaPC NN employing the exactly same training data sets $\boldsymbol{D}_{T_N}$ while varying the size $T_N$. However, the final trained ML representation could depend not only on the size $T_N$ of the training data, but also on how the training data have been constructed. Indeed, at least for the polynomial representation, the training points must be selected with dedicated strategies in order to avoid additional oscillations known as the Runge phenomenon [89]. In particular, wrong training points could lead to a very strong oscillations near the discontinuities (Gibbs effect) of the considered $CO_2$ benchmark. Therefore, it will be also interesting to see how the Runge phenomenon may affect the quality of DANN and DaPC NN training.

In order to address that question, we will consider two training strategies. In the first training strategy, we will use the Sobol sequence [100] to construct the training sets of sizes $T_N$ equals to 100, 500 and 1000 similarly to the previous test cases. In the second training strategy, we will employ Gaussian integration points [110, 32] to generate the training data sets of sizes $T_N$ equal to 27, 216 and 1331 according to the available distribution [75] of inputs displayed in Figure 11. In short, the optimal choice [110] of training points (i.e. Gaussian quadrature rule) corresponds to the roots of the polynomial of one degree higher than the order used in the polynomial representation. The Gaussian integration points form a full tensor (FT) grid in the space of inputs. This is what leads to sizes of the training data sets equal to 27, 216 and 1331 related to the polynomial representation of $3^{rd}$, $5^{th}$ and $10^{th}$ degrees correspondingly. Strictly, this training rule can be fully satisfied for aPC representation following the Gaussian quadrature rule [110], but there is no proof for multi-layered structures such as DANN or DaPC NN indicating that the corresponding Gaussian training strategy could mitigate the Runge phenomenon. To assess the quality of prediction, we will employ the validation data set $\boldsymbol{D}_{V_N}$ of the size $V_N = 10^4$ [53] generated according to the variability of the inputs (Monte Carlo approach) shown in Figure 11.

Table 4: Architectures of aPC, DANN and DaPC NN for $CO_2$ benchmark problem employing Gaussian integration points.

| ML model | $T_N$ | $L$ | $N^{(\mathcal{L})}$ | $d^{(\mathcal{L})}$ | $\mathcal{A}^{(\mathcal{L})}$ | $N_w$ | $\mathcal{LF}$ |
|----------|-------|-----|---------------------|---------------------|-------------------------------|-------|----------------|
| aPC | 27 | 1 | 1 | 2 | none | 10 | MSE+MSW |
|  | 216 | 1 | 1 | 5 | none | 56 | MSE+MSW |
|  | 1331 | 1 | 1 | 10 | none | 286 | MSE+MSW |
| DANN | 27 | 3 | [8,6,1] | [1,1,1,1] | eq. (9) | 93 | MSE+MSW |
|  | 216 | 3 | [8,6,1] | [1,1,1,1] | eq. (9) | 93 | MSE+MSW |
|  | 1331 | 4 | [10,8,6,1] | [1,1,1,1] | eq. (9) | 189 | MSE+MSW |
| DaPC NN | 27 | 2 | [3,1] | [2,2] | eq. (22) | 40 | MSE+MSW |
|  | 216 | 2 | [3,1] | [2,2] | eq. (22) | 40 | MSE+MSW |
|  | 1331 | 2 | [3,1] | [3,3] | eq. (22) | 80 | MSE+MSW |

### 3.4.2 Prediction and Validation

We will reproduce the $CO_2$ saturation along the radial distance from the injection well for the fixed time instance of 100 days in accordance with the $CO_2$ benchmark scenario [54]. To do so, we will train aPC expansion, conventional DANN and DaPC NN for each numerical discretion cell (i.e. 250 times). In this sense, we will construct 250 networks/expansions that seek to capture features of $CO_2$ displacement in the subsurface. The related ML architectures are presented in Table 3 . Figure 12 demonstrates the performance of the considered ML models during validation, using the Sobol sequence training data with $T_N = 1000$. All three approaches are far away in capturing the reference values of mean and standard deviation for $CO_2$ saturation even considering a sufficiently large size of the training data. All three ML models suffer strongly from the Gibbs effect (i.e oscillations) caused by the strong non-linearity of shock propagation (see examples in [54]). The regularization in equation (23) seems to show the most significant effect for the aPC. Opposite to that, using Gaussian integration points as training data helps to mitigate the Runge phenomenon and to reduce the Gibbs effect substantially not only for the aPC as expected, but also for the DaPC NN. Corresponding numbers of training points and technical specifications of the considered ML models are presented in Table 4 . Figures 13-15 demonstrate the predictions made by aPC expansion, conventional DANN and DaPC NN. It seems that the DaPC NN inherits partially the aPC properties that help overcome the Runge phenomenon, and the optimal training strategy relying on the distribution of the inputs helps to construct acceptable ML model at low computational costs (e.g. Figures 13 or Figures 14). Unfortunately, the training strategy based on Gaussian integration points does not improve the performance of the conventional DANN. Indeed, as we have clarified in Section 2.2.2, the conventional DANN structure relies on Gaussian distribution of inputs with zero mean and unit variance, which is not fulfilled in the considered $CO_2$ benchmark case.



Figure 12: Prediction of mean (left) and standard deviation (right) using aPC, DANN and DaPC NN for Validation data: $CO_2$ benchmark problem using 1000 Sobol sequences as training data set.

### 3.4.3 Evidence of convergence

We will assess the convergence in terms of mean square error, mean value and standard deviation of $CO_2$ saturation over the training data size $T_N$ for both previously mentioned training strategies. In order to make the quantities of interest in correspondence with the original benchmark study [52] and its follow-up [87], we will consider the

Figure 13: Prediction of mean (left) and standard deviation (right) using aPC, DANN and DaPC NN for Validation data: $CO_2$ benchmark problem using 27 Gaussian integration points as training data set.



Figure 14: Prediction of mean (left) and standard deviation (right) using aPC, DANN and DaPC NN for Validation data: $CO_2$ benchmark problem using 216 Gaussian integration points as training data set.



Figure 15: Prediction of mean (left) and standard deviation (right) using aPC, DANN and DaPC NN for Validation data: $CO_2$ benchmark problem using 1331 Gaussian integration points as training data set.

following space-averaged quantities:

$$\text{MSE}_{T_N} = \frac{1}{250} \cdot \frac{1}{V_N} \sum_{i=1}^{V_N} \|\mathcal{R}(\boldsymbol{\omega}_{V_i}) - R_{V_i}\|_{L^2}^2 \,, \tag{29}$$

$$E_{T_N}^{\mu} = \frac{1}{250} \left\| \mu[\mathcal{R}(\boldsymbol{\omega}_{V_N})] - \mu[\mathbf{R}_{V_N}] \right\|_{L^2} \,, \tag{30}$$

$$E_{T_N}^{\sigma} = \frac{1}{250} \left\| \sigma[\mathcal{R}(\boldsymbol{\omega}_{V_N})] - \sigma[\mathbf{R}_{V_N}] \right\|_{L^2} \,. \tag{31}$$

Figure 16 shows convergence in terms of the absolute averaged mean square error in equation (29), absolute averaged error of mean value in equation (30) and absolute averaged error of standard deviation in equation (31) over the training data size $T_N$ for aPC, conventional DANN and DaPC NN. Figure 16 displays the convergence of the analyzed ML approaches trained on the Sobol sequences and as well on the Gaussian integration points (marked by superscript FT).

We observe that the performance of DANN does not profit from an increasing training data size. Also, it does not benefit from the considered choices of training points and seems to be affected by oscillations caused by the strong non-linearity of the underlying problem. Additionally, increasing the training data size does not help the aPC expansion if the training points are not optimally distributed (see more details in [75]), but it is clearly visible how the use of Gaussian integration points help overcome the Runge phenomenon. Moreover, the Gaussian integration points strongly help in DaPC NN training and also reduce the Gibbs effect, rendering the results more reliable. We also would like to remark, that turning the DaPC NN architecture into a one-layer structure reproduces the aPC results under the same technical specification. However, we have selected the multi-layered DaPC NN architecture as specified in Table 4 for illustration of performance in order to provide the broader picture. Overall, we conclude that all considered ML approaches suffer a lot in capturing the strong shock propagation in the $CO_2$ benchmark model, as all three rely on smooth functions.



Figure 16: Performance of aPC, DANN and DaPC NN trained on Sobol sequences and Gaussian integration points for $CO_2$ benchmark problem: Convergence of Absolute Averaged Mean Square Error, Absolute Averaged Error of Mean Value and Absolute Averaged Error of Standard Deviation relatively the reference validation data set.

### 3.5 Final Remarks

In the current paper, we offer a view on neural signal processing in deep artificial neural networks from the PCE perspective, introducing in that way the DaPC NN. By employing PCE fundamentals, we demonstrate that orthonormality conditions in each node of the conventional DANN structure are fulfilled if the neural signal propagating through the multi-layer architecture was normally distributed with zero mean and unit variance. This situation is not necessarily satisfied for the majority of data-driven applications and, hence, could lead to redundant representation, where one neural signal could contain partial information that is also coming from other neurons. Moreover, introducing the PCE into the DANN structure provides an opportunity to go beyond the linear weighted superposition of single univariate neurons on each node.

From the modelling perspective, the user is prompted to specify the DaPC NN architecture through the number of layers, number of nodes per layer, activation function and loss function similar as in a conventional DANNs. Additionally, the novel DaPC NN requires specification of the desired degree of non-linearity for each hidden layer. The latter aspect leads to high-order weighted superposition on each node of the network, so that non-linear activation functions become optional. This reduces the potential for subjectivity in the modeling procedure. However, it also introduces a new responsibility of choosing the degree of non-linearity, which in turn raises new research questions.

Technically, the DaPC NN requires a similar training procedure as any conventional DANN, and all trained weights determine automatically the corresponding multi-variate data-driven orthonormal bases for all layers of DaPC NN. The DaPC NN Matlab Toolbox is available online and users are invited to adopt it for own needs [77].

To illustrate the performance of the DaPC NN, we have investigated three test cases, comparing the DaPC NN with the conventional DANN and also with aPC expansion. To do so, we have employed identical training procedures for DANN, DaPC NN and aPC, minimizing the mean squared error, regularized via ridge regression. In the training procedure, we employ exactly the same training and validation data sets for all three ML approaches. Additionally, we

vary the size of the training data set and assess the convergence of the different ML approaches on a validation data set by evaluating the mean square error, error of mean and error of standard deviation.

In the first and second test cases, we observe that the DaPC NN reaches superior results in comparison to the aPC and DANN, even considering that we use only a moderate training data size. In the third $CO_2$ benchmark test case, we observe that all considered ML approaches suffer a lot in capturing the strong $CO_2$ shock propagation, as they all rely on smooth functions. However, using Gaussian integration points as training dataset helps to mitigate the Runge phenomenon and reduce the Gibbs effect substantially  not only for the aPC but also for the DaPC NN. This emphasizes that the DaPC NN inherits partially the aPC properties that help to overcome the Runge phenomenon. Therefore, an optimal training strategy relying on the distribution of the inputs helps to construct an acceptable ML model at low computational costs. Unfortunately, the performance of the DANN does not profit from the considered choices of training points. Increasing the training data size does not seem to be helpful against oscillations caused by the strong non-linearity of the underlying problem. The latter point can be explained by the derivations in the current paper, where the conventional DANN structure is shown to be optimal for Gaussian distribution of inputs (and also Gaussian signal processing on each node) that is hardly fulfilled in the considered $CO_2$ benchmark.

## Summary and Conclusions

The current paper analyzes the neural signal processing in deep artificial neural networks from the point of view of polynomial chaos (PCE) theory. The response on each node of a deep network has been seen from the PCE perspective, making use of orthonormal polynomial basis functions. The proposed generalization of the conventional structure of DANNs towards the DaPC NN decomposes the neural signals, employing an adaptive construction of data-driven multi-variate orthonormal bases on each node in the multi-layer structure. Moreover, by introducing PCE theory to represent the response of each node, it offers an additional opportunity to go beyond the linear weighted superposition of single univariate neurons as in conventional DANN structures. In that sense, the introduced DaPC NN structure can be seen as a generalization of the conventional DANN with incorporation of aPC theory. The newly introduced DaPC NN assures orthonormal decomposition on each node and also offers an additional possibility to account for high-order neural effects. Doing so, the weights gain a clear meaning according to global sensitivity analysis and they reflect the partial contribution of each single neuron (linear univariate terms) or simultaneous combination of neurons (non-linear multivariate terms) to the total variance of the response on each node.

Concluding the analysis of the current paper, we anticipate that avoiding the redundant representation and accounting for high-order neural effects could increase the performance of the neural network. Following the observation in the current paper, we stress that the high non-linearity of the underlying problems could limit the ability of the plain aPC with only one-layer polynomial representation. As opposed to that, the encoded flexibility of DANNs seems to suffer significantly from over-fitting, regardless of the regularization of weights. Overall, the DaPC NN shows the ability to predict quantities of interest more consistently with lower variance in comparison to DANNs. Furthermore, we also observe that omitting the regularization in the loss function only slightly degrades the results of DaPC NN and aPC, whereas the conventional DANN is significantly more affected by it.

Remarking that aspect, we accentuate that joining the fundamentals of homogeneous chaos theory with the deep representation of neural networks requires additional research, where various architectures accompanied by specific loss functions and training algorithms could be investigated.  In particular, concepts employed in recurrent neural networks, neural ordinary differential equations, physical regularization or Bayesian regularization could be adopted directly to the DaPC NN structure. Moreover, the introduced DaPC NN structure opens a pathway to use analytical properties quantifying the importance of neural signal on each node . In particular, implementing an explicit analytical form of data-driven orthonormal polynomial bases in equations (17)-(23) from [75] could significantly accelerate the training process of DaPC NN. Additionally, since the weights of DaPC NN reflect the meaning according to global sensitivity analysis, they could be partially omitted [20] to offer a sparse DaPC NN representation. The latter is highly beneficial in mitigating issues related to overfitting. Additionally, state-of-the-art findings in the PCE community can further be included, such as various sparse learning techniques, multi-element decomposition, Bayesian learning, etc.

From a technical perspective, we would like to clarify that the primary goal of the current paper is not to achieve high computational efficiency in the training procedure. Instead, we offer a unique perspective on the DANN structure. Indeed, the computational time of the current version of DaPC NN is significantly higher than that of DANN, since it requires adaptive computation of the orthonormal basis during the training process. The procedure could be accelerated using explicit analytic relations for a moderate degree of expansion from [75]. Moreover, further speed up could be achieved when employing parallel and GPU-based computing.

## Acknowledgments

## References

[1] M. Abramowitz and A. Stegun, I. Handbook of Mathematical Functions with Formulas, Graphs, and Mathematical Tables. Dover Publications, Inc., New York, 1965.

[2] Jonas Adler and Ozan Öktem. Solving ill-posed inverse problems using iterative deep neural networks. Inverse Problems, 33(12):124007, Nov 2017.

[3] C.C. Aggarwal. Neural Networks and Deep Learning: A Textbook. Springer, 1 edition, 9 2018.

[4] R Ahlfeld, B Belkouchi, and Francesco Montomoli. SAMBA: sparse approximation of moment-based arbitrary polynomial chaos. J. Comput. Phys., 320:1–16, 2016.

[5] NI Akhiezer. The classical moment problem. Hafner Publ. Co., New York, 2, 1965.

[6] Osama Alkhateeb and Nathan Ida. Data-Driven Multi-Element Arbitrary Polynomial Chaos for Uncertainty Quantification in Sensors. IEEE Transactions on Magnetics, 54(3), 2017.

[7] M. Anthony and P.L. Bartlett. Neural Network Learning: Theoretical Foundations. Cambridge University Press, 1999.

[8] Sercan Ö Arık, Mike Chrzanowski, Adam Coates, Gregory Diamos, Andrew Gibiansky, Yongguo Kang, Xian Li, John Miller, Andrew Ng, Jonathan Raiman, et al. Deep voice: Real-time neural text-to-speech. In International Conference on Machine Learning, pages 195–204. PMLR, 2017.

[9] Martin Arjovsky, Amar Shah, and Yoshua Bengio. Unitary evolution recurrent neural networks. In International conference on machine learning, pages 1120–1128. PMLR, 2016.

[10] R. Askey and J.A. Wilson. Some Basic Hypergeometric Orthogonal Polynomials that Generalize Jacobi Polynomials. Number no. 319 in American Mathematical Society: Memoirs of the American Mathematical Society. American Mathematical Society, 1985.

[11] Anthony Curtis Atkinson and Alexander N Donev. Optimum experimental designs, volume 5. Clarendon Press, 1992.

[12] F. Augustin, A. Gilg, M. Paffrath, P. Rentrop, and U. Wever. Polynomial chaos for the approximation of uncertainties: Chances and limits. European Journal of Applied Mathematics, 19(2):149–190, 2008.

[13] Dana Harry Ballard and Christopher M. Brown. Computer Vision. Prentice Hall Professional Technical Reference, 1st edition, 1982.

[14] João Carlos Alves Barata and Mahir Saleh Hussein. The Moore–Penrose pseudoinverse: A tutorial review of the theory. Brazilian Journal of Physics, 42(1):146–165, 2012.

[15] Felix Beckers, Andrés Heredia, Markus Noack, Wolfgang Nowak, Silke Wieprecht, and Sergey Oladyshkin. Bayesian calibration and validation of a large-scale and time-demanding sediment transport model. Water Resources Research, 56(7):e2019WR026966, 2020.

[16] Géraud Blatman and Bruno Sudret. Sparse polynomial chaos expansions and adaptive stochastic finite elements using a regression approach. C. R. Mécanique, 336(6):518–523, 2008.

[17] Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty in neural network. In Francis Bach and David Blei, editors, Proceedings of the 32nd International Conference on Machine Learning, volume 37 of Proceedings of Machine Learning Research, pages 1613–1622, Lille, France, 07–09 Jul 2015. PMLR.

[18] Thierry Bouwmans, Sajid Javed, Maryam Sultana, and Soon Ki Jung. Deep neural network concepts for background subtraction: A systematic review and comparative evaluation. Neural Networks, 117:8–66, 2019.

[19] Mateusz Buda, Atsuto Maki, and Maciej A Mazurowski. A systematic study of the class imbalance problem in convolutional neural networks. Neural Networks, 106:249–259, 2018.

[20] Paul-Christian Bürkner, Ilja Kröker, Sergey Oladyshkin, and Wolfgang Nowak. The sparse polynomial chaos expansion: a fully bayesian approach with joint priors on the coefficients and global selection of terms. Journal of Computational Physics, 488:112210, 2023.

[21] R. H. Cameron and W. T. Martin. The orthogonal development of non-linear functionals in series of Fourier-Hermite functionals. Ann. of Math. (2), 48:385–392, 1947.

[22] K. Cheng, Lu Z, Xia S., Oladyshkin S., and Nowak W. Surrogate models for uncertainty quantification: A unified bayesian framework. Computer Methods in Applied Mechanics and Engineering, page Under Review, 2021.

[23] Grigorios G Chrysos, Stylianos Moschoglou, Giorgos Bouritsas, Yannis Panagakis, Jiankang Deng, and Stefanos Zafeiriou. P-nets: Deep polynomial neural networks. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 7325–7335, 2020.

[24] Gioele Ciaparrone, Francisco Luque Sánchez, Siham Tabik, Luigi Troiano, Roberto Tagliaferri, and Francisco Herrera. Deep learning in video multi-object tracking: A survey. Neurocomputing, 381:61–88, 2020.

[25] Holger Class, Anozie Ebigbo, Rainer Helmig, Helge K Dahle, Jan M Nordbotten, Michael A Celia, Pascal Audigane, Melanie Darcis, Jonathan Ennis-King, Yaqing Fan, et al. A benchmark study on problems related to $CO_2$ storage in geologic formations. Computational Geosciences, 13(4):409, 2009.

[26] Corinna Cortes and Vladimir Vapnik. Support-vector networks. Machine learning, 20(3):273–297, 1995.

[27] Noel AC Cressie. Spatial Prediction and Kriging. Statistics for Spatial Data (Cressie NAC, ed). New York: John Wiley & Sons, pages 105–209, 1993.

[28] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In 2009 IEEE Conference on Computer Vision and Pattern Recognition, pages 248–255, 2009.

[29] Li Deng. An overview of deep-structured learning for information processing. In Proc. Asian-Pacific Signal & Information Proc. Annual Summit & Conference (APSIPA-ASC), pages 1–14, October 2011.

[30] Oliver G. Ernst, Antje Mugler, Hans-Jörg Starkloff, and Elisabeth Ullmann. On the convergence of generalized polynomial chaos expansions. ESAIM Math. Model. Numer. Anal., 46(2):317–339, 2012.

[31] J. Foo and G.E. Karniadakis. Multi-element probabilistic collocation method in high dimensions. J. Comput. Phys., 229(5):1536–1557, 2010.

[32] Walter Gautschi. Orthogonal polynomials: computation and approximation. Numerical Mathematics and Scientific Computation. Oxford University Press, New York, 2004. Oxford Science Publications.

[33] R. G. Ghanem and P. D. Spanos. Stochastic finite elements: A spectral approach. Springer-Verlag, New York, 1991.

[34] W. Gilks, S. Richardson, and D. Spiegelhalter. Markov chain Monte Carlo in practice. Chapmann & Hall, Boca Raton, 1996.

[35] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. Deep Learning. MIT Press, 2016. http://www.deeplearningbook.org.

[36] Daniel Graupe. Principles of artificial neural networks, volume 7 of Advanced series in circuits and systems. World Scientific Publishing Company, Singapore, 2013.

[37] M.H. Hassoun. Fundamentals of Artificial Neural Networks. A Bradford book. MIT Press, Cambridge, 1995.

[38] S. Hochreiter and J. Schmidhuber. Long short-term memory. Neural computation, 9(8):1735–1780, 1997.

[39] Kurt Hornik, Maxwell Stinchcombe, and Halbert White. Multilayer Feedforward Networks are Universal Approximators. Neural Networks, 2(5):359–366, 1989.

[40] Gao Huang, Zhuang Liu, Laurens van der Maaten, and Kilian Q. Weinberger. Densely connected convolutional networks, 2018.

[41] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift, 2015.

[42] T. Ishigami and T. Homma. An importance quantification technique in uncertainty analysis for computer models. In [1990] Proceedings. First International Symposium on Uncertainty Modeling and Analysis, pages 398–403, 1990.

[43] A.G. Ivakhnenko and V.G. Lapa. Cybernetics and forecasting techniques. 1967.

[44] Gareth James, Daniela Witten, Trevor Hastie, and Robert Tibshirani. An Introduction to Statistical Learning: With Applications in R. Springer Publishing Company, Incorporated, New York, 2014.

[45] A.S. Witteveen Jeroen, Sunetra Sarkar, and Bijl Hester. Modeling physical uncertainties in dynamic stall induced fluid–structure interaction of turbine blades using arbitrary polynomial chaos. Computers and Structures, 85(11-14):866–878, 2007.

[46] Kui Jia, Shuai Li, Yuxin Wen, Tongliang Liu, and Dacheng Tao. Orthogonal deep neural networks. IEEE Transactions on Pattern Analysis and Machine Intelligence, 43, 2019.

[47] Xiaowei Jia, Jared Willard, Anuj Karpatne, Jordan S. Read, Jacob A. Zwart, Michael Steinbach, and Vipin Kumar. Physics-guided machine learning for scientific discovery: An application in simulating lake temperature profiles. ACM/IMS Trans. Data Sci., 2(3), May 2021.

[48] Fazle Karim, Somshubra Majumdar, Houshang Darabi, and Samuel Harford. Multivariate LSTM-FCNs for time series classification. Neural Networks, 116:237–245, 2019.

[49] S. Karlin. Total Positivity, volume I. Stanford University Press, Stanford, 1968.

[50] A. Keese and H. G. Matthies. Sparse quadrature as an alternative to Monte Carlo for stochastic finite element techniques. Proc. Appl. Math. Mech., 3:493–494, 2003.

[51] Andrei Nikolaevich Kolmogorov and Albert T Bharucha-Reid. Foundations of the theory of probability: Second English Edition. Dover Publications, Inc., New York, 2018.

[52] M. Köppel, I. Kröker, and C. Rohde. Intrusive uncertainty quantification for hyperbolic-elliptic systems governing two-phase flow in heterogeneous porous media. Comput. Geosci., 21(4):807–832, 2017.

[53] Markus Köppel, Fabian Franzelin, Ilja Kröker, Sergey Oladyshkin, Gabriele Santin, Dominik Wittwar, Andrea Barth, Bernard Haasdonk, Wolfgang Nowak, Dirk Pflüger, and Christian Rohde. Datasets and executables of data-driven uncertainty quantification benchmark in carbon dioxide storage.

[54] Markus Köppel, Fabian Franzelin, Ilja Kröker, Sergey Oladyshkin, Gabriele Santin, Dominik Wittwar, Andrea Barth, Bernard Haasdonk, Wolfgang Nowak, Dirk Pflüger, and Christian Rohde. Comparison of data-driven uncertainty quantification methods for a carbon dioxide storage benchmark scenario. Computational Geosciences, Nov 2019.

[55] Daniel G Krige. A statistical approach to some basic mine valuation problems on the witwatersrand. Journal of the Southern African Institute of Mining and Metallurgy, 52(6):119–139, 1951.

[56] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. ImageNet Classification with Deep Convolutional Neural Networks. Commun. ACM, 60(6):84–90, May 2017.

[57] Ilja Kröker, Wolfgang Nowak, and Christian Rohde. A stochastically and spatially adaptive parallel scheme for uncertain and nonlinear two-phase flow problems. Computational Geosciences, 19(2):269–284, 2015.

[58] Heng Li and Dongxiao Zhang. Probabilistic collocation method for flow in porous media: Comparisons with other stochastic methods. Water Resources Research, 43(9):1–13, 2007.

[59] G. Lin and A.M. Tartakovsky. An efficient, high-order probabilistic collocation method on sparse grids for three-dimensional flow and solute transport in randomly heterogeneous porous media. Adv. Water Res., 32(5):712–722, 2009.

[60] David JC MacKay. Bayesian interpolation. Neural computation, 4(3):415–447, 1992.

[61] Donald W. Marquardt. An algorithm for least-squares estimation of nonlinear parameters. Journal of the Society for Industrial and Applied Mathematics, 11(2):431–441, 1963.

[62] MATLAB. version 9.7.0.1216025 (r2019b). https://www.mathworks.com/help/stats/fitrgp.html, 2019.

[63] John McCarthy. Review of the question of artificial intelligence. Annals of the History of Computing, 10(3):224–229, 1988.

[64] Warren S McCulloch and Walter Pitts. A logical calculus of the ideas immanent in nervous activity. The bulletin of mathematical biophysics, 5(4):115–133, 1943.

[65] Zakaria Mhammedi, Andrew Hellicar, Ashfaqur Rahman, and James Bailey. Efficient orthogonal parametrisation of recurrent neural networks using householder reflections. In International Conference on Machine Learning, pages 2401–2409. PMLR, 2017.

[66] Hrushikesh Narhar Mhaskar and Charles A Micchelli. How to choose an activation function. In Advances in Neural Information Processing Systems, pages 319–326, Denver, 1994.

[67] Risto Miikkulainen, Jason Liang, Elliot Meyerson, Aditya Rawal, Daniel Fink, Olivier Francon, Bala Raju, Hormoz Shahrzad, Arshak Navruzyan, and Nigel Duffy. Evolving deep neural networks. In Artificial intelligence in the age of neural networks and brain computing, pages 293–312. Elsevier, 2019.

[68] Eliakim H Moore. On the reciprocal of the general algebraic matrix. Bull. Am. Math. Soc., 26:394–395, 1920.

[69] Maryam M Najafabadi, Flavio Villanustre, Taghi M Khoshgoftaar, Naeem Seliya, Randall Wald, and Edin Muharemagic. Deep learning applications and challenges in big data analytics. Journal of big data, 2(1):1–21, 2015.

[70] Preetum Nakkiran, Gal Kaplun, Yamini Bansal, Tristan Yang, Boaz Barak, and Ilya Sutskever. Deep double descent: Where bigger models and more data hurt. Journal of Statistical Mechanics: Theory and Experiment, 2021(12):124003, 2021.

[71] Hayrettin Okut. Bayesian regularized neural networks for small n big p data. Artificial neural networks-models and applications, 2016.

[72] S. Oladyshkin, H. Class, R. Helmig, and W. Nowak. A concept for data-driven uncertainty quantification and its application to carbon dioxide storage in geological formations. Adv. Water Res., 34:1508–1518, 2011.

[73] S. Oladyshkin, H. Class, R. Helmig, and W. Nowak. An integrative approach to robust design and probabilistic risk assessment for $CO_2$ storage in geological formations. Comput. Geosci., 15(3):565–577, 2011.

[74] S Oladyshkin, FPJ De Barros, and W Nowak. Global sensitivity analysis: a flexible and efficient framework with an example from stochastic hydrogeology. Advances in Water Resources, 37:10–22, 2012.

[75] S. Oladyshkin and W. Nowak. Data-driven uncertainty quantification using the arbitrary polynomial chaos expansion. Reliab. Eng. Syst. Safe., 106:179–190, 2012.

[76] Sergey Oladyshkin. aPC matlab toolbox: Data-driven arbitrary polynomial chaos. https://www.mathworks.com/matlabcentral/fileexchange/ 72014-apc-matlab-toolbox-data-driven-arbitrary-polynomial-chaos, 2022.

[77] Sergey Oladyshkin. DaPC NN: Deep arbitrary polynomial chaos neural network. https://www.mathworks.com/matlabcentral/fileexchange/ 112110-dapc-nn-deep-arbitrary-polynomial-chaos-neural-network, 2022.

[78] Sergey Oladyshkin, Farid Mohammadi, Ilja Kroeker, and Wolfgang Nowak. Bayesian[3] Active Learning for the Gaussian Process Emulator Using Information Theory. Entropy, 22(8):890, 2020.

[79] Sergey Oladyshkin and Wolfgang Nowak. Incomplete statistical information limits the utility of high-order polynomial chaos expansions. Reliability Engineering & System Safety, 169:137–148, 2018.

[80] Sergey Oladyshkin and Wolfgang Nowak. The Connection between Bayesian Inference and Information Theory for Model Selection, Information Gain and Experimental Design. Entropy, 21(11):1081, 2019.

[81] Theodore Papamarkou, Jacob Hinkle, M. Todd Young, and David Womble. Challenges in Markov chain Monte Carlo for Bayesian neural networks, 2021.

[82] Roger Penrose. On best approximate solutions of linear matrix equations. In Mathematical Proceedings of the Cambridge Philosophical Society, volume 52, pages 17–19. Cambridge University Press, 1956.

[83] Timothy Praditia, Matthias Karlbauer, Sebastian Otte, Sergey Oladyshkin, Martin V Butz, and Wolfgang Nowak. Learning groundwater contaminant diffusion-sorption processes with a finite volume neural network. Water Resources Research, page e2022WR033149, 2022.

[84] Timothy Praditia, Thilo Walser, Sergey Oladyshkin, and Wolfgang Nowak. Improving thermochemical energy storage dynamics forecast with physics-inspired neural network architecture. Energies, 13(15):3873, 2020.

[85] Waseem Rawat and Zenghui Wang. Deep convolutional neural networks for image classification: A comprehensive review. Neural computation, 29(9):2352–2449, 2017.

[86] JR Red-Horse and AS Benjamin. A probabilistic approach to uncertainty quantification with limited information. Reliability Engineering & System Safety, 85(1):183–190, 2004.

[87] Michael F Rehme, Fabian Franzelin, and Dirk Pflüger. B-splines on sparse grids for surrogates in uncertainty quantification. Reliability Engineering & System Safety, 209:107430, 2021.

[88] Sebastian Ruder. An overview of gradient descent optimization algorithms, 2017.

[89] Carl Runge. Über empirische Funktionen und die Interpolation zwischen äquidistanten Ordinaten. Zeitschrift für Mathematik und Physik, 46(224-243):20, 1901.

[90] Arthur L Samuel. Some studies in machine learning using the game of checkers. IBM Journal of research and development, 3(3):210–229, 1959.

[91] Juergen Schmidhuber. Annotated history of modern ai and deep learning, 2022. Technical Report IDSIA-22-22.

[92] Jürgen Schmidhuber. Deep learning in neural networks: An overview. Neural networks, 61:85–117, 2015.

[93] Burr Settles. Active learning literature survey. Computer Sciences Technical Report 1648, University of Wisconsin–Madison, 2009.

[94] Sagar Sharma and Simone Sharma. Activation functions in neural networks. Towards Data Science, 6(12):310–316, 2017.

[95] Baoguang Shi, Xiang Bai, and Cong Yao. An end-to-end trainable neural network for image-based sequence recognition and its application to scene text recognition. IEEE transactions on pattern analysis and machine intelligence, 39(11):2298–2304, 2016.

[96] JA Shohat and JD Tamarkin. The problem of moments, mathematical surveys no. 1. American Mathematical Society, New York, 1950, 1943.

[97] Paz Fink Shustin, Shashanka Ubaru, Vasileios Kalantzis, Lior Horesh, and Haim Avron. Pcenet: High dimensional surrogate modeling for learning uncertainty, 2022.

[98] W. McC. Siebert. On the determinants of moment matrices. Ann. Statist., 17(2):711–721, 1989.

[99] Adrian FM Smith and Alan E Gelfand. Bayesian statistics without tears: a sampling–resampling perspective. The American Statistician, 46(2):84–88, 1992.

[100] Ilya M Sobol', Danil Asotsky, Alexander Kreinin, and Sergei Kucherenko. Construction and comparison of high-dimensional Sobol' generators. Wilmott, 2011(56):64–79, 2011.

[101] Il'ya Meerovich Sobol'. On sensitivity estimation for nonlinear mathematical models. Matematicheskoe modelirovanie, 2(1):112–118, 1990.

[102] J Stieltjes, T. Quelques recherches sur la théorie des quadratures dites méchaniques. Oeuvres I, pages 377–396, 1884.

[103] Bruno Sudret. Global sensitivity analysis using polynomial chaos expansions. Reliability Engineering & System Safety, 93(7):964–979, 2008. Bayesian Networks in Dependability.

[104] T.J. Sullivan. Introduction to Uncertainty Quantification. Texts in Applied Mathematics. Springer International Publishing, Cham, 2015.

[105] Chunwei Tian, Yong Xu, and Wangmeng Zuo. Image denoising using deep CNN with batch renormalization. Neural Networks, 121:461–473, 2020.

[106] Andrey Nikolaevich Tikhonov, Vasiliy Iakovlevich Arsenin, VY Arsenin, et al. Solutions of ill-posed problems. Vh Winston, 1977.

[107] Michael E Tipping. The relevance vector machine. In Advances in neural information processing systems, pages 652–658, 2000.

[108] Vladimir Vapnik and Alexey Chervonenkis. Theory of pattern recognition. Nauka, Moscow, 1974.

[109] John Villadsen and Michael L Michelsen. Solution of differential equation models by polynomial approximation. Prentice-Hall, New Jersey, 1978.

[110] John Villadsen and Michael L Michelsen. Solution of differential equation models by polynomial approximation(book). Englewood Cliffs, N. J., Prentice-Hall, Inc., 1978. 460 p, 1978.

[111] Eugene Vorontsov, Chiheb Trabelsi, Samuel Kadoury, and Chris Pal. On orthogonality and learning recurrent networks with long term dependencies. In International Conference on Machine Learning, pages 3570–3578. PMLR, 2017.

[112] Jiayun Wang, Yubei Chen, Rudrasis Chakraborty, and Stella X Yu. Orthogonal convolutional neural networks. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 11505–11515, 2020.

[113] Norbert Wiener. The homogeneous chaos. American Journal of Mathematics, 60(4):897–936, 1938.

[114] Norbert Wiener. Cybernetics, or Control and Communication in the Animal and the Machine. Actualités Scientifiques et Industrielles [Current Scientific and Industrial Topics], No. 1053. Hermann et Cie., Paris; The Technology Press, Cambridge, Mass.; John Wiley & Sons, Inc., New York, 1948.

[115] Christopher K Williams and Carl Edward Rasmussen. Gaussian processes for machine learning, volume 2. MIT press Cambridge, MA, 2006.

[116] Scott Wisdom, Thomas Powers, John Hershey, Jonathan Le Roux, and Les Atlas. Full-capacity unitary recurrent neural networks. Advances in neural information processing systems, 29, 2016.

[117] Lin Xiao, Bolin Liao, Shuai Li, and Ke Chen. Nonlinear recurrent neural networks for finite-time solution of general time-varying linear matrix equations. Neural Networks, 98:102–113, 2018.

[118] D. Xiu and G. E. Karniadakis. Modeling uncertainty in flow simulations via generalized polynomial chaos. Journal of Computational Physics, 187:137–167, 2003.

[119] D. Xiu and G.E. Karniadakis. The Wiener-Askey Polynomial Chaos for stochastic differential equations. SIAM Journal of Scientific Computing, 24(2):619–644, 2002.

[120] Dongbin Xiu and George Em Karniadakis. The Wiener–Askey polynomial chaos for stochastic differential equations. SIAM Journal on Scientific Computing, 24(2):619–644, 2002.

[121] P. Yee and S. Haykin. Pattern classification as an ill-posed, inverse problem: a regularization approach. In 1993 IEEE International Conference on Acoustics, Speech, and Signal Processing, volume 1, pages 597–600 vol.1, 1993.

[122] Yingqi Zhang, Yaning Liu, George Pau, Sergey Oladyshkin, and Stefan Finsterle. Evaluation of multiple reduced-order models to enhance confidence in global sensitivity analyses. Int. J. Greenh. Gas Control, 49:217–226, 2016.

[123] Xiaohu Zheng, Jun Zhang, Ning Wang, Guijian Tang, and Wen Yao. Mini-data-driven Deep Arbitrary Polynomial Chaos Expansion for Uncertainty Quantification, 2021.