

Digital Discovery

Accepted Manuscript

This article can be cited before page numbers have been issued, to do this please use: F. Bolaños-García, J. Commenge and L. Falk, *Digital Discovery*, 2026, DOI: 10.1039/D5DD00526D.



This is an Accepted Manuscript, which has been through the Royal Society of Chemistry peer review process and has been accepted for publication.

Accepted Manuscripts are published online shortly after acceptance, before technical editing, formatting and proof reading. Using this free service, authors can make their results available to the community, in citable form, before we publish the edited article. We will replace this Accepted Manuscript with the edited and formatted Advance Article as soon as it is available.

You can find more information about Accepted Manuscripts in the [Information for Authors](#).

Please note that technical editing may introduce minor changes to the text and/or graphics, which may alter content. The journal's standard [Terms & Conditions](#) and the [Ethical guidelines](#) still apply. In no event shall the Royal Society of Chemistry be held responsible for any errors or omissions in this Accepted Manuscript or any consequences arising from the use of any information it contains.

Gradient-enhanced neural networks for model parameter estimation applied to flow chemistry automated platforms

Francisco Bolaños-García^{1,*}, Jean-Marc Commenge¹, and Laurent Falk¹

¹Université de Lorraine, CNRS, LRGP, Nancy F-54000, France

*Corresponding author: francisco.bolanos-garcia@univ-lorraine.fr

Abstract

The acceleration of chemical process development through flow chemistry depends on obtaining reliable kinetic models. Model-based design of experiments (MBDoE) has been successfully applied with flow chemistry platforms for estimating reaction rate parameters for low-complexity models. However, its use with dynamic experiments involving computationally expensive models remains challenging, particularly when experimental conditions must be suggested in real time. Surrogate models can approximate complex models, but often lack the ability to reproduce the derivatives of the original model, which are essential for MBDoE as a gradient-based approach. This work investigates the use of gradient-enhanced neural networks as surrogate models for parameter estimation within an MBDoE framework. By incorporating gradient information, the surrogate model is able to reproduce the local sensitivity structure of the original first-principles model, ensuring predictive accuracy for both output values and gradients needed for parameter estimation. A case study on competitive-consecutive reactions demonstrates that artificial neural networks (ANNs) that were trained with gradient information can be used for parameter estimation, while substantially reducing computational cost by a factor of approximately 200,000. This enables sequential MBDoE suitable for real-time applications.



1 Introduction

Nowadays, laboratories seek to automatize their workflows¹. High-throughput experimentation (HTE) and flow chemistry platforms have emerged as essential tools for exploring broad chemical spaces and fine tune reaction conditions². The goal is to intensify the data-generation process, which is useful for model development and efficient exploration while reducing the time and resources required. Flow chemistry platforms are well-suited for sequential experimental planning³, and a fully automated system can be achieved by integrating algorithms that perform model fitting using real-time data from process analytical tools (PATs)⁴, enabling a complete closed-loop workflow.

Identifying an appropriate model structure and estimating the values of model parameters are essential steps for the development of process engineering models for different system. In reactor design, for example, kinetic models that describe the chemical reaction should be obtained at an early stage, as they are fundamental for further scale-up, defining control strategies, and optimization. Systematic procedures, known as model-based design of experiments (MBD_{oE}), have been therefore proposed and applied for model discrimination^{5,6} and parameter estimation^{7–9} with the objective of performing more informative experiments. Examples of these strategies in flow chemistry include the estimation of reaction rate parameters under steady-state conditions for single reactions, such as the Diels–Alder reaction¹⁰, and for multistep reactions, such as a nucleophilic aromatic substitution (S_NAr)¹¹. They have also been applied in batch-like flow platforms, where flow ramps are used to extract kinetic information, as demonstrated for the esterification of benzoic acid with ethanol¹². However, there have been no reports of MBD_{oE} applied to platforms using autosamplers and injection loops to intensify the experimental workflow¹³, even though this approach could reduce reagent consumption and lower both cost and time in chemical process development. In these systems, small volumes of reagents are sequentially injected into a reactor whose volume is significantly larger than the injected volume. As a result, the reactor operates under transient conditions.

To deploy a sequential planning on automated platforms, some open-source tools have already been designed. For example EFCOSS¹⁴, that provides a modular environment coupling numerical simulation codes with optimization packages. Pyomo.DOE¹⁵ formulates the MBD_{oE} problem as a stochastic program and uses nonlinear sensitivity analysis. Another is GPdoemd¹⁶, which trains Gaussian process (GP) surrogate models to support model discrimination when candidate models correspond to computationally expensive black-box simulators and gradient information with respect to model parameters is not readily available. In that framework, GP surrogates are constructed from sampled model evaluations and used



54 to evaluate design criteria based on their predictive distributions.

55 In contrast, the focus of the present work is on parameter estimation for first-principles models, where
56 parameter sensitivities can be computed, but the integration of MBDoE with self-guided flow chemistry
57 setups requires fast numerical optimization. During the procedure, a parameter estimation is performed
58 using the available measurements, while a second optimization problem determines the experimental
59 conditions that maximize the information content of subsequent experiments¹⁷. When model evaluations
60 are computationally intensive, these optimization steps become a bottleneck, limiting MBDoE deployment
61 in automated platforms by constraining the speed at which new experiments can be executed. A solution
62 to this has been presented by Friso and Galvanin¹⁸, where an optimization-free procedure is performed in
63 which a relative information value is used to compare possible experimental conditions. However, for the
64 chosen design to be near the unknown true optimal design depends on the pre-defined set of candidate
65 experiments.

66 In process systems engineering, surrogate models have emerged as a tool to tackle complex models
67 while saving computational time and resources¹⁹. Within experimental design frameworks, Artificial Neu-
68 ral Networks (ANNs) have been employed to identify suitable kinetic models from experimental data^{20,21},
69 showing the potential of using surrogate models for model building. However, these applications have
70 typically not extended to parameter estimation, which requires solving an inverse problem to find model
71 parameters from the measurements. Other methodologies have used physics-informed neural networks
72 (PINNs). With them, a simultaneous training of the surrogate model parameters and the kinetic param-
73 eters is performed²², which can be challenging if not enough data is available.

74 One strategy to improve the training of neural networks consists in including gradient information,
75 as matching the Jacobian matrix provides additional information per training sample and can be inter-
76 preted as a form of data augmentation^{23,24}. The idea is to apply surrogate models that are trained not
77 only with the model output, but also using its corresponding derivatives. Named Sobolev-trained neural
78 networks²³, they benefit from current libraries, e.g. Pytorch²⁵, that compute these gradients efficiently
79 through automatic differentiation (AD). Previously, gradient-enhanced ANNs needed the explicit deriva-
80 tion of the gradient^{26,27}, limiting its application. These types of models have already been applied to
81 process engineering for gradient-based optimization²⁸, where they have been proven to be beneficial in
82 scenarios where small data sets are available.

83 Gradients also play a crucial role in MBDoE procedures, so their calculation should be accurate. The
84 sensitivities of the model with respect to its parameters, captured by the gradients, are fundamental for



85 defining the design criteria. In addition, gradients are used within the optimization routine that seeks to
86 optimize these criteria¹⁷.

87 Motivated by this, the present work explores the use of ANNs enhanced by a derivative training for
88 parameter estimation through a sequential planning, with the aim of allowing the deployment of such
89 strategy in automated platforms that must handle computationally expensive model evaluations. To
90 authors knowledge, this work represents the first application of gradient-enhanced surrogate models to
91 sequential MBDoE for parameter estimation with experiments under transient conditions. The following
92 section introduces the background of the MBDoE procedure, followed by a description of the neural
93 network enhancement with gradient information illustrated through an example, and finally a case study
94 where parameters are estimated from in silico flow chemistry experiments under unsteady conditions.

95 2 Background: Model Based Design of Experiments

96 A model built from experimental data is defined as a function that predicts the observed output as a
97 function of the experimental conditions and parameters. Mathematically, this is expressed as

$$\hat{y} = f(\boldsymbol{\xi}, \boldsymbol{\theta}) \quad (1)$$

98 where $\boldsymbol{\xi}$ denotes the vector of operating conditions inside the design space Ξ , defined by equipment
99 constraints or by a region of interest based on prior knowledge, and $\boldsymbol{\theta}$ denotes the vector of model param-
100 eters inside the model parameter space Θ . For instance, in a flow reactor, $\boldsymbol{\xi}$ may include temperature,
101 residence time, and inlet concentrations, while $\boldsymbol{\theta}$ may represent unknown kinetic constants to be esti-
102 mated. In this work, model parameters refer to physically meaningful parameters (e.g., kinetic constants),
103 as opposed to non-physical parameters like the weights and biases of data-driven models. This function
104 can be a first-principles model, denoted by $f()$, or a data-driven model, such as a neural network, de-
105 noted by $f_{NN}()$. Both models take as inputs the operating conditions and physically meaningful model
106 parameters, and return predicted measurable outputs such as concentrations or yields.

107 From a set of preliminary experiments, hypotheses can be made to propose a set of plausible models
108 that describe the observed behaviour. After selection of model structures, it is important to be sure they
109 are identifiable and distinguishable²⁹. If, from the data available, it is not possible to discriminate between
110 the best two models, a sequential procedure to produce more observations to maximize discrepancy
111 between predictions can be performed³⁰.



112 Once we assume the selected model structure is adequate, experiments that increase the information
 113 available for the estimation of parameters must be executed. This part of the MBDoE is the focus of
 114 this work. The definition of the parameter space Θ should be based on some prior knowledge on the
 115 system behaviour from the preliminary experiments. In order to obtain a parameter estimation $\hat{\theta}$ from
 116 experimental data, an estimator of the fit quality of the model with respect to observations should be
 117 defined. The maximum likelihood estimator (MLE) is mostly used due to its consistency, efficiency and
 118 asymptotic normality under modest assumptions³¹. This estimator maximizes the likelihood function
 119 $\ell(\mathbf{y}|\theta)$, defined as the joint probability of observing the experimental data \mathbf{y}_e ($e = 1, \dots, N_e$) given the
 120 N_p model parameters. To properly quantify it, the uncertainty in the measurements must be taken
 121 into account with the variance-covariance matrix of the experimental errors, $\Sigma_{\mathbf{y}}$. Its diagonal elements
 122 represent the variance of the uncertainty associated with each of the N_y measurements within the e^{th}
 123 experiment, and the off-diagonal elements the covariance between pairs of them^{32,33}.

$$\ln \ell(\mathbf{y}|\theta) = -\frac{N_e N_y}{2} \ln 2\pi - \frac{N_e}{2} \ln(\det \Sigma_{\mathbf{y}}) - \frac{1}{2} \sum_{e=1}^{N_e} (\mathbf{y}_e - \hat{\mathbf{y}}_e)^T \Sigma_{\mathbf{y}}^{-1} (\mathbf{y}_e - \hat{\mathbf{y}}_e) \quad (2)$$

$$\hat{\theta} = \arg \max_{\theta} \ln \ell(\mathbf{y}|\theta) \quad (3)$$

$$\text{s.t. } \hat{\theta} \in \Theta$$

124
 125 To improve the estimation with further experiments, it is important to quantitatively assess the
 126 influence of the parameters values on the outputs of the proposed model. This is represented by the
 127 sensitivity matrix $\mathbf{S}[N_y \times N_p]$, which is a local measurement depending on current parameter values and
 128 experimental conditions. For example, it captures how variations in a rate constant affect predicted outlet
 129 concentrations under a given temperature and residence time.

$$\mathbf{S}_e(\hat{\theta}, \xi_e) = \begin{bmatrix} \frac{\partial \hat{y}_1}{\partial \theta_1} & \cdots & \frac{\partial \hat{y}_1}{\partial \theta_{N_p}} \\ \vdots & \ddots & \vdots \\ \frac{\partial \hat{y}_{N_y}}{\partial \theta_1} & \cdots & \frac{\partial \hat{y}_{N_y}}{\partial \theta_{N_p}} \end{bmatrix} \quad (4)$$

130 From this, a square matrix $N_p \times N_p$ known as the Fisher Information Matrix (FIM), \mathbf{F} , can be
 131 computed with Equation 5. This step is crucial as the variance-covariance matrix of the estimated
 132 parameters $\Sigma_{\hat{\theta}}$ can be approximated by the inverse of the FIM evaluated at $\hat{\theta}, \xi^{33}$, useful to determine
 133 the confidence region of the parameters and the correlation between them. In this context, experiments



134 that produce stronger sensitivity of concentrations to kinetic parameters result in a more informative
 135 FIM. As stated above, the MLE is asymptotically Gaussian and unbiased. This means, being θ^* the
 136 vector of true parameters, the distribution of $\hat{\theta}$ tends to $\mathcal{N}(\theta^*, \mathbf{F}(\theta^*)^{-1})$ as $N_e \rightarrow \infty$ ³¹.

137 In practice, however, only a limited number of experiments is available. As a result, this asymptotic
 138 approximation may not hold exactly, and its accuracy depends on several factors: the proximity of the
 139 estimated parameters to their true values, the degree of model nonlinearity with respect to the parameters,
 140 and the signal-to-noise ratio of the measurements. Its widespread use in MBDoE is motivated by its
 141 computational tractability and its effectiveness for comparing alternative experimental designs, rather
 142 than its ability to provide an exact quantification of uncertainty³³, remaining as a practical and widely
 143 used tool for ranking candidate experiments before performing them.

$$\mathbf{F}_{N_e}(\hat{\theta}, \xi) = \sum_{e=1}^{N_e} \mathbf{S}_e(\hat{\theta}, \xi_e)^T \Sigma_y^{-1} \mathbf{S}_e(\hat{\theta}, \xi_e) \quad (5)$$

$$\Sigma_{\hat{\theta}} \approx \mathbf{F}_{N_e}(\hat{\theta}, \xi)^{-1} \quad (6)$$

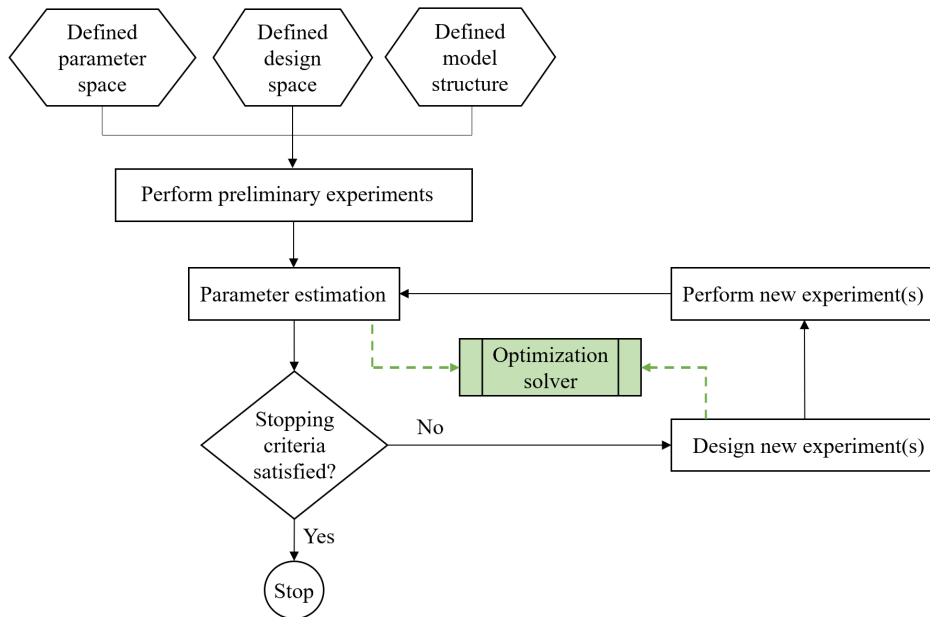


Figure 1: MBDoE procedure focused on parameter estimation. An optimization solver is called at two key stages: to guide the design of new, informative experiments and to perform the parameter estimation.



144 To facilitate sequential experimental design, a scalar function of the FIM is typically introduced
 145 to quantify the information content associated with a given set of operating conditions. This scalar
 146 metric serves as the design criterion and is maximized at each step to inform the selection of subsequent
 147 experiments. This enables ranking candidate operating conditions (e.g., different temperatures or flow
 148 rates) according to their expected contribution to parameter estimation. The optimal design is, therefore,
 149 the one that maximizes the selected criterion $\psi(\cdot)$.

$$\xi_{N_e+1}^{OPT} = \arg \max_{\xi_{N_e+1}} \psi(\mathbf{F}_{N_e} + \mathbf{S}_{N_e+1}(\hat{\boldsymbol{\theta}}, \xi_{N_e+1})^T \Sigma_{\mathbf{y}}^{-1} \mathbf{S}_{N_e+1}(\hat{\boldsymbol{\theta}}, \xi_{N_e+1})) \quad (7)$$

$$\text{s.t. } \xi_{N_e+1}^{OPT} \in \Xi$$

151 Several design criteria exist to decide the best next experiment¹⁷. In this work, the D-optimality
 152 is employed, as it is the most widely used, where the determinant of the FIM is the scalar property
 153 to be maximized. The aim is to find the experimental conditions that minimize the uncertainty of the
 154 parameter estimation, as increasing the determinant of the FIM reduces the volume of the confidence
 155 region of the estimated model parameters. In practice, this corresponds to selecting the next set of
 156 operating conditions that is expected to give the most informative measurements for refining the kinetic
 157 parameters.

$$\psi_D(\mathbf{F}(\hat{\boldsymbol{\theta}}, \boldsymbol{\xi})) := \det(\mathbf{F}(\hat{\boldsymbol{\theta}}, \boldsymbol{\xi})) \quad (8)$$

158 The whole workflow is presented in Figure 1 illustrating the iterative interplay between model evalu-
 159 ation, parameter estimation, and experiment selection. Stopping criteria could be either a fixed number
 160 of experiments based on a pre-defined experimental budget, or a threshold value when assessing the
 161 adequacy and reliability of the estimated parameters by, for example, comparing χ^2 and t-values with
 162 reference values from the corresponding distribution with $N_e \times N_y - N_p$ degrees of freedom¹⁷. The focus
 163 of our work is on parameter estimation and the subsequent design steps, which require efficient compu-
 164 tation of sensitivities, that is, first-order derivatives of the model outputs with respect to the parameters
 165 (Equation 4). Since analytical expressions for these derivatives are rarely available, they are often ob-
 166 tained numerically, a process that can become prohibitively expensive in automated platforms when using
 167 finite-difference schemes.

168 Overall, the workflow takes as inputs a candidate model, experimental data (e.g., concentration mea-
 169 surements), and admissible operating ranges (e.g., temperature and flow limits), and returns updated



170 parameter estimates together with the next optimal experimental conditions to be tested, guiding subse-
171 quent experiments.

172 3 Gradient-enhanced training of neural networks

173 Artificial neural networks (ANNs) are able to approximate any continuous function. Their architecture
174 is based on multiple layers, where each one of them is formed by a number of neurons (also called
175 perceptrons)³⁴. The output of a single layer is expressed as

$$\mathbf{z}_l = \phi_l(\mathbf{W}_l^T \mathbf{z}_{l-1} + \mathbf{b}_l), \quad \forall l = 1, \dots, L \quad (9)$$

176 Where $\mathbf{z}_{l-1}[N_{l-1} \times 1]$ and $\mathbf{z}_l[N_l \times 1]$ denote the input and output vectors, respectively, of the layer
177 l , where N_l represents the number of neurons of this layer. $\mathbf{W}_l[N_{l-1} \times N_l]$ is the matrix of weights,
178 $\mathbf{b}_l[N_l \times 1]$ the vector of biases for the layer l , and $\phi_l()$ is the activation function of layer l . The most
179 common activation functions for the hidden layers are the rectified linear unit (ReLU), sigmoid, and
180 hyperbolic tangent, while a linear function is often used for the output layer in regression tasks. So, when
181 using an ANN as a black box function $\hat{\mathbf{y}} = f_{NN}(\mathbf{x})$, the output of the final layer \mathbf{z}_L corresponds to the
182 prediction value $\hat{\mathbf{y}}$ and the input vector \mathbf{x} is the input \mathbf{z}_0 of the first hidden layer.

183 Different hyperparameters, such as the number of layers (L) and the number of neurons per layer (N_l
184 $\forall l = 1, \dots, L$), must be specified. The surrogate model is then trained on a designated training dataset of
185 size N_s to determine the optimal weights and biases by solving an optimization problem that minimizes
186 a chosen loss function. For regression tasks, commonly used loss functions include the mean absolute
187 error (MAE) and the mean squared error (MSE). During training, a backpropagation is performed where
188 gradients of the loss with respect to all weights and biases are computed using the chain rule via automatic
189 differentiation (AD). This AD framework can also provide the sensitivities of $f_{NN}(\mathbf{x})$ not only during
190 training but also when evaluating the ANN, as illustrated in Figure 2.

$$\text{MSE}(f(\mathbf{x}_n), f_{NN}(\mathbf{x}_n)) := \frac{1}{N_y} \|f(\mathbf{x}_n) - f_{NN}(\mathbf{x}_n)\|_2^2 \quad (10)$$

191 If the loss function is restricted to the MSE (Equation 10) between function outputs, the ANN can
192 predict accurately the function values, but may perform poorly if required to yield the derivatives of the
193 function with respect to the inputs, which are essential in gradient-based applications such as MBDofE for
194 parameter estimation. This limitation is significant since the surrogate model replaces the mechanistic



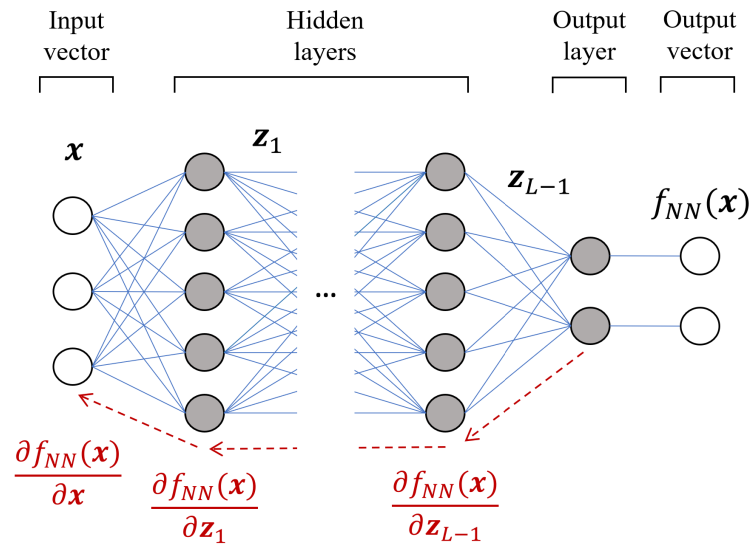


Figure 2: Example of a fully connected artificial neural network with multiple hidden layers. Grey circles represent the neurons. The dashed red arrows illustrate the backpropagation step during training, where the gradients are computed.

195 model that provides a detailed but computationally expensive representation of the experimental setup.
 196 As shown by Czarnecki²³, Sobolev-trained neural networks minimize the distance between the true func-
 197 tion $f(\mathbf{x})$ and the surrogate model $f_{NN}(\mathbf{x})$ at the Sobolev-space by adding a loss term for each of the j
 198 elements in the set of available and relevant derivative orders \mathcal{J} .

$$\mathcal{L} = \frac{1}{N_s} \sum_{n=1}^{N_s} \text{MSE}(f(\mathbf{x}_n), f_{NN}(\mathbf{x}_n)) + \frac{1}{N_s} \sum_{j \in \mathcal{J}} \sum_{n=1}^{N_s} \lambda_j \text{MSE}(D_{\mathbf{x}}^j f(\mathbf{x}_n), D_{\mathbf{x}}^j f_{NN}(\mathbf{x}_n)) \quad (11)$$

199 To ensure all terms in the loss function are of the same order of magnitude, Tsay²⁸ proposed a simple
 200 method to scale the derivative terms and select a proper value for each weight λ_j , avoiding the challenge
 201 of relying on heuristic tuning. Since it is common practice to normalize the input and output vectors
 202 before training the neural network, the derivatives used for training must be scaled consistently with
 203 Equation 12. We denote a scaled variable with an overbar symbol (e.g., \bar{x}). If resulting derivatives are
 204 also normalized using the min-max scaler, then $\bar{D}_{\bar{\mathbf{x}}}^j \bar{f}(\bar{\mathbf{x}}_n) \in [0, 1]$, and both terms in the loss function will
 205 have the same order of magnitude.

$$D_{\bar{\mathbf{x}}}^j \bar{f}(\bar{\mathbf{x}}_n) = \left. \frac{\partial^j \bar{f}}{\partial \bar{x}^j} \right|_{\bar{\mathbf{x}}_n} = \frac{\text{range}(x_n)}{\text{range}(f(x_n))} \left. \frac{\partial^j f}{\partial x^j} \right|_{x_n} \quad (12)$$



$$\mathcal{L} = \frac{1}{N_s} \sum_{n=1}^{N_s} \text{MSE}(\bar{f}(\mathbf{x}_n), f_{NN}(\bar{\mathbf{x}}_n)) + \frac{1}{N_s} \sum_{j \in \mathcal{J}} \sum_{n=1}^{N_s} \text{MSE}(\bar{D}_{\bar{\mathbf{x}}}^j \bar{f}(\bar{\mathbf{x}}_n), \bar{D}_{\bar{\mathbf{x}}}^j f_{NN}(\bar{\mathbf{x}}_n)) \quad (13)$$

206 Where the range is defined as the maximum minus the minimum value. Within the proposed frame-
 207 work of integrating a neural network into the standard MBDofE for parameter estimation, the input vector
 208 \mathbf{x}_n consists of the experimental conditions and the parameter vector $[\xi_n, \theta_n]$. The whole set of gradients
 209 $D_{\bar{\mathbf{x}}} \bar{f}$ is not needed and just the first derivative with respect to θ_n will be computed. This avoids the
 210 need to construct complete gradient information over all inputs, providing a distinct structural advantage
 211 over other surrogate models, i.e., gradient-enhanced Gaussian process (GEM) approaches, which generally
 212 require complete gradient information across all inputs to formulate stable joint covariance matrices³⁵.

$$D_{\bar{\theta}} \bar{f}(\xi_n, \bar{\theta}_n) = \left. \frac{\partial \bar{f}}{\partial \theta} \right|_{\xi_n, \bar{\theta}_n} = \frac{\text{range}(\theta_n)}{\text{range}(f(\xi_n, \theta_n))} \left. \frac{\partial f}{\partial \theta} \right|_{\xi_n, \theta_n} \quad (14)$$

$$\mathcal{L} = \frac{1}{N_s} \sum_{n=1}^{N_s} \text{MSE}(\bar{f}(\mathbf{x}_n), f_{NN}(\bar{\mathbf{x}}_n)) + \frac{1}{N_s} \sum_{n=1}^{N_s} \text{MSE}(\bar{D}_{\bar{\theta}} \bar{f}(\xi_n, \bar{\theta}_n), \bar{D}_{\bar{\theta}} f_{NN}(\bar{\xi}_n, \bar{\theta}_n)) \quad (15)$$

213 As discussed previously, the gradients of the neural network output with respect to its inputs can be
 214 obtained by AD, already available in libraries such as Pytorch (version 2.8.0). For the first-principles
 215 model, several strategies can be employed to obtain these derivatives. While analytical expressions
 216 offer exact results, they are often impractical for complex systems. Numerical finite differences provide
 217 a straightforward alternative but are computationally costly and prone to numerical error. A more
 218 systematic alternative is AD, which provides exact derivatives up to machine precision. Depending on
 219 the problem size, AD can be applied in forward or backward mode. In the context of models governed
 220 by ordinary differential equations (ODEs), the corresponding sensitivity equations (forward mode) or
 221 adjoint equations (backward mode) can be solved alongside the system dynamics to efficiently compute
 222 parameter sensitivities^{36,37}.

223 4 Example using a gradient-enhanced neural network

224 To demonstrate the capabilities of a gradient-enhanced neural network, the following example considers a
 225 first-order reaction $A \rightarrow B$ with kinetic parameter k . The reaction occurs under steady-state conditions
 226 in a tubular microreactor with axial dispersion. Calculations were done assuming a tube of 5 mL with an
 227 internal diameter of 0.75 mm, a system that can be perfectly used in flow chemistry applications. In non



228 ideal tubular reactors, the Peclet (Pe) number must be determined to estimate the axial dispersion. For
 229 this, assuming laminar flow, the dispersion coefficient \mathbf{D} was calculated using the Aris-Taylor equation
 230 (Equation 16). Then, for first-order reactions, Equation 18 is used to determine the conversion X_A at
 231 the outlet of the tubular reactor³⁸.

$$\mathbf{D} = \mathcal{D} + \frac{u^2 d_t^2}{192\mathcal{D}} \quad (16)$$

$$\text{Pe} = \frac{uL}{\mathbf{D}} \quad (17)$$

$$X_A = 1 - \frac{C_A}{C_{A,0}} = 1 - \frac{4a \exp(\text{Pe}/2)}{(1+a)^2 \exp(a \text{Pe}/2) - (1-a)^2 \exp(-a \text{Pe}/2)} \quad (18)$$

$$a = \sqrt{1 + \frac{4k\tau}{\text{Pe}}} \quad (19)$$

232 where u denotes the linear velocity, d_t the internal diameter of the reactor, \mathcal{D} the molecular diffusion,
 233 L the length of the reactor, τ the space time, and k the kinetic constant. It was assumed that $\mathcal{D} =$
 234 $1 \times 10^{-9} \text{ m}^2/\text{s}$, a typical value for liquids. This first-principles model will serve as the reference to be
 235 replaced by a surrogate model.

236 A dataset $\{(\tau_n, k_n, X_{A_n}, dX_{A_n}/dk_n)\}_{n=1}^{100}$ was created by a standard Latin Hypercube Sampling (LHS)
 237 over the domain of the experimental conditions, i.e., the space time ($\tau \in [10 \text{ s}, 90 \text{ s}]$), and the domain of
 238 model parameters, i.e., the kinetic constant ($k \in [0.001 \text{ s}^{-1}, 0.1 \text{ s}^{-1}]$). This dataset was then randomly
 239 partitioned into three subsets: a training set (70%), a validation set (15%), and an independent test set
 240 (15%). The architecture of the surrogate model is a fully connected neural network with two inputs,
 241 τ and k , one hidden layer of 10 neurons, and a single output neuron for X_A . The number of neurons
 242 in the hidden layer (10) was determined through preliminary testing, as increasing this number did not
 243 result in any further improvement in the surrogate model's predictive accuracy. For scaling the input
 244 vector, the boundaries of the design space were used as min-max values, and not the min-max values
 245 from the sampled dataset. Two surrogate models were trained by minimizing the loss function using the
 246 Adam optimizer: one using only the function values X_A , and another using both the function values and
 247 the derivatives dX_A/dk . Training was performed for a maximum of 20,000 epochs with early stopping
 248 based on the validation loss, using a patience of 500 epochs. The trained network weights and biases
 249 corresponding to the lowest validation loss were retained, and final performance metrics were computed



250 on the independent test set.

251 A comparative analysis was performed to determine the most effective activation function for the ANN
252 architecture. The evaluation involved training 10 networks, each with a random weight initialization. The
253 Root Mean Squared Error (RMSE) of the model output and of the gradient obtained through AD of the
254 ANN was recorded. The median RMSE across the 10 trained networks was used as a performance
255 metric. This comparison was carried out for both types of surrogate models, the standard trained and
256 the enhanced with gradient information.

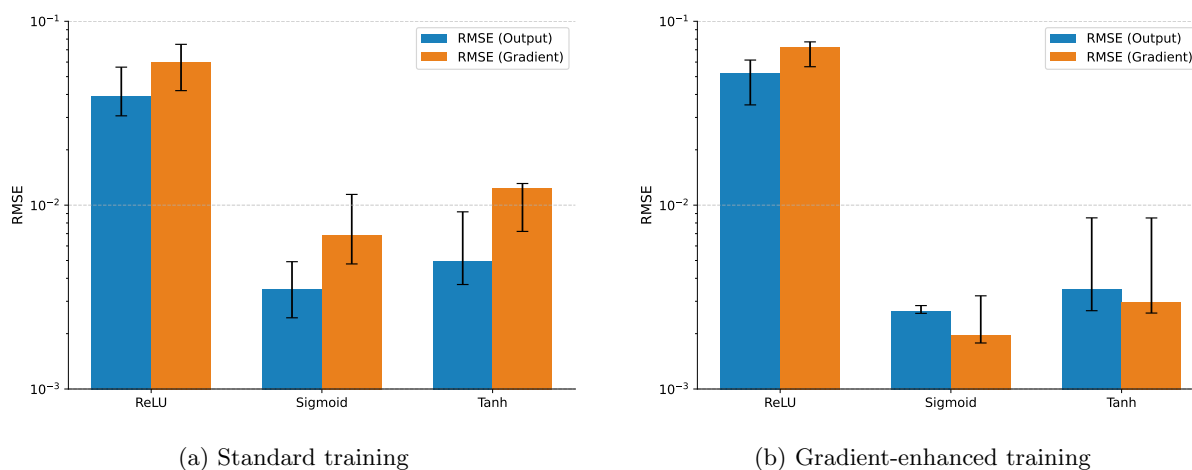


Figure 3: Median RMSE obtained from 10 independent neural network trainings for each activation function. Error bars indicate the computed interquartile range.

257 The median RMSE reported in Figure 3 exhibits the performance for three different activation func-
258 tions: ReLU, sigmoid and tanh. Under standard training, the sigmoid and tanh functions yielded similar
259 performance, although the tanh activation showed greater variability. A common characteristic across all
260 three activation functions was a greater error in predicting the gradient than in predicting the function
261 value itself. When incorporating gradient information into the surrogate model training, the performance
262 for the sigmoid and tanh functions improved. In contrast, the ReLU function consistently produced the
263 highest RMSE and showed no significant improvement with gradient-informed training, an expected out-
264 come given its non-differentiability at zero. Regarding computational cost, the inclusion of the gradient
265 loss term did not lead to an increase in training time. On the contrary, the gradient-enhanced training
266 exhibited a slightly lower average training time (10 s) compared to standard training (12 s).

267 Because the sigmoid activation achieved the lowest RMSE for both the model output X_A and its
268 gradient with respect to k , it was selected as the preferred choice for further analysis. An ablation study



269 on the gradient loss weight was conducted to assess its influence on model performance. In addition
 270 to the baseline formulation (Equation 15), weights of 0.5 and 2.0 for the gradient loss were evaluated.
 271 Both alternatives resulted in higher average total MSE values, 1.78×10^{-4} and 1.02×10^{-4} , respectively,
 272 compared to 3.91×10^{-5} obtained with the baseline weighting. These results indicate that the chosen
 273 weighting provides a near-optimal balance between fitting the function value and its gradient.

274 Figure 4 compares the true function and gradient values with the predictions from the standard trained
 275 ANNs with the sigmoid activation function. The comparison is performed at $\tau = \{10s, 50s, 90s\}$, which
 276 correspond to the lower boundary, centre, and upper boundary of the operating range used to generate
 277 the training dataset. The prediction of the conversion X_A is accurate for all three space time values and
 278 over the whole range of k values. This is not the case for the prediction of dX_A/dk in the lower region
 279 of k values of the training set, especially at both boundary values of τ , where the predictions are less
 280 accurate. This is a common behaviour of ANNs and why they tend to perform badly when extrapolating.
 281 Figure 5 shows the improvement on both the conversion and its first-order derivative prediction when
 282 adding the gradient information during training, particularly increasing accuracy near the lower values of
 283 k . This means that using a gradient-enhanced training yielded a surrogate model that will behave more
 284 alike the first-principles model, a relevant advantage in gradient-based applications.

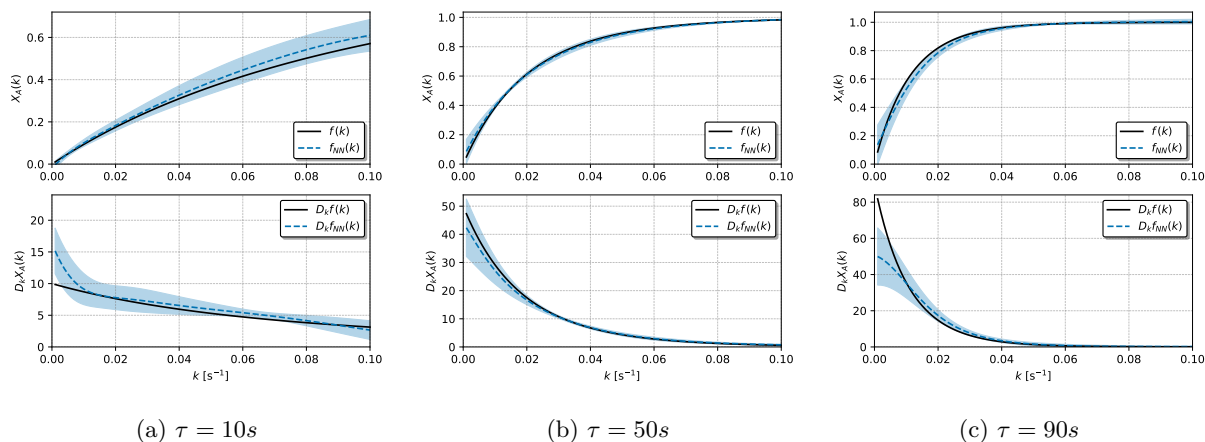


Figure 4: Predictive accuracy for the conversion $X_A(k)$ and its sensitivity with respect to the kinetic constant $D_k X_A(k)$ obtained using ANNs trained without gradient information. The solid black line represents the true model output, the blue dashed line indicates the mean prediction of 10 ANNs, and the shaded area shows the mean \pm one standard deviation.

285 Then, if the surrogate model is used by a MBDofE framework to determine the most informative
 286 experimental condition for the estimation of k , the sensitivities are computed. As just one parameter is



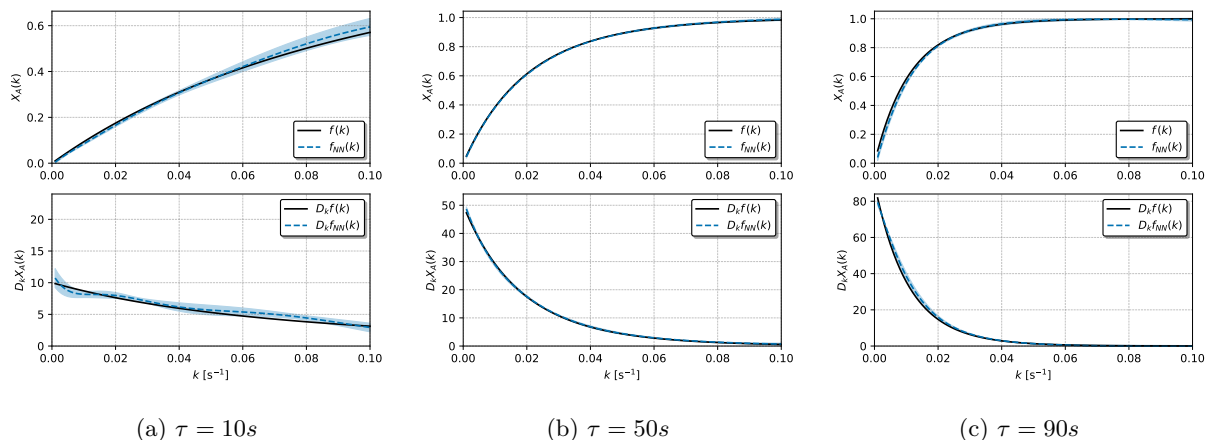


Figure 5: Predictive accuracy for the conversion $X_A(k)$ and its sensitivity with respect to the kinetic constant $D_k X_A(k)$ obtained using ANNs trained with gradient information. The solid black line represents the true model output, the blue dashed line indicates the mean prediction of 10 ANNs, and the shaded area shows the mean \pm one standard deviation.

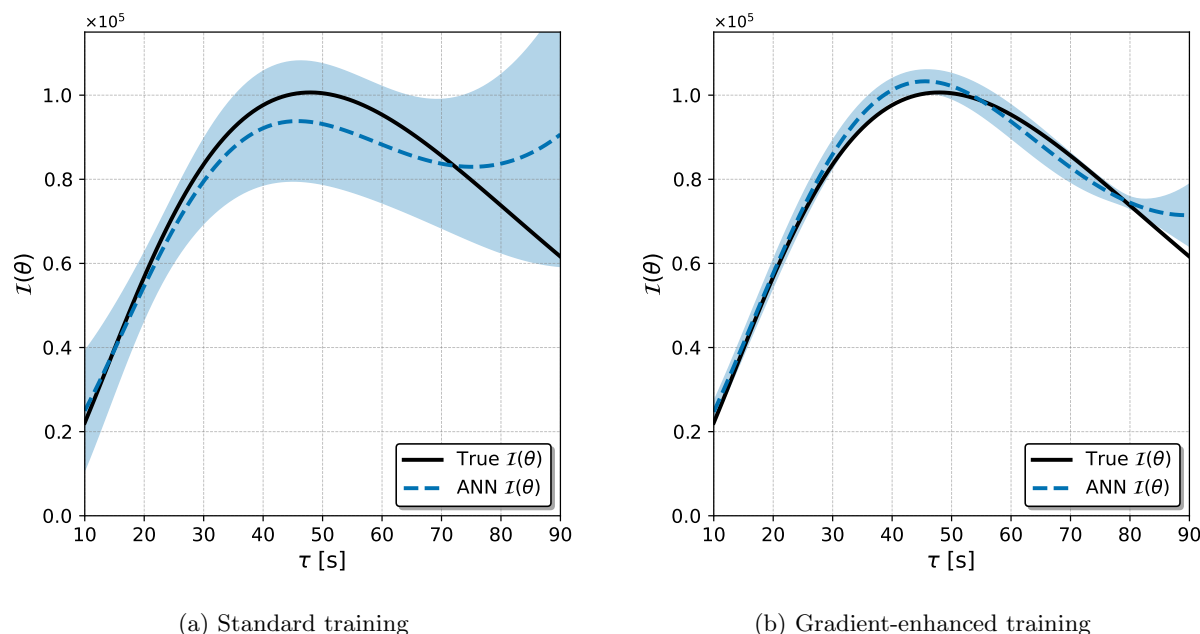


Figure 6: Comparison between the true Fisher information (solid black line) and the prediction made by 10 randomly initialized neural networks (blue dashed line). The shaded area represents the mean \pm one standard deviation.

287 being estimated, the FIM becomes a scalar value $\mathcal{I}(\theta)$. This value was determined assuming a measure-
 288 ment error in the observed conversion, $\varepsilon \sim \mathcal{N}(0, \sigma_e^2)$, with $\sigma_e = 0.05$. Figure 6 presents the predicted



289 and true Fisher information for all possible values of the space time with $k = 0.022 \text{ s}^{-1}$, a value for
290 which $dX_A/dk > 0$ over the whole design space so the conversion is sensitive to variations in the kinetic
291 constant. When comparing between the standard and the gradient-enhanced training, the latter is more
292 accurate in the prediction of the Fisher information. This is relevant when performing an experimental
293 planning to find the optimal value of τ that reduces the uncertainty on our estimation of k . Based on
294 the first-principles model (full black line), an optimal space time for an experiment is close to 47 s, and
295 the same decision would be taken if the gradient-enhanced ANN is used. In contrast, an ANN with a
296 standard training would hesitate between the values of 47 s and 90 s, which could misguide a sequential
297 experimental strategy toward less informative experiments.

298 5 Case study: Competitive consecutive reactions in a discrete 299 injection flow platform

300 In chemical development, reducing reagent consumption is desirable to lower both cost and time required.
301 Slug flow platforms address this by isolating small reactive volumes (slugs) between gas bubbles, limiting
302 axial dispersion and enhancing radial mixing by recirculation within the slug³⁹. However, their automa-
303 tion across a broad design space remains challenging. An alternative approach is to inject a discrete
304 reactive volume directly into a continuous solvent stream (a rectangular pulse of reactants concentra-
305 tion), without any gas separator¹³. The resulting axial dispersion produces a concentration gradient
306 along the microreactor.

307 To gain knowledge on model parameters from this kind of systems is not trivial, as neglecting the
308 effect of dispersion will lead to an erroneous estimation of the kinetic constants. A detailed model is
309 then needed to account for the real physical behaviour and precisely estimate the intrinsic values. The
310 proposed case study consists of two competitive-consecutive reactions, the iodination of L-tyrosine. In
311 this scheme, L-tyrosine (A) reacts with iodine in water (B) to form 3-iodotyrosine (R), and then R can
312 further react with B to give 3,5-diiodotyrosine (S). Both reactions are second order⁸, with reaction rates
313 denoted by r and kinetic constants by k .



$$r_1 = k_1 C_A C_B \quad (20)$$



$$r_2 = k_2 C_R C_B \quad (21)$$

5.1 First-principles model

A first-principles model describing this system was built, and simulations were performed to generate training data for the surrogate neural network. The model inputs consist of the experimental conditions ξ , corresponding to the space time τ and the inlet concentration ratio of tyrosine to total concentration, and the model parameters θ , i.e., both kinetic constants k_1 and k_2 . The mass balances over the four species along the tubular reactor with axial dispersion yields a set of partial differential equations:

$$\frac{\partial C_i(z, t)}{\partial t} = \mathbf{D} \frac{\partial^2 C_i(z, t)}{\partial z^2} - u \frac{\partial C_i(z, t)}{\partial z} + \nu_{i,1} r_1 + \nu_{i,2} r_2 \quad \forall i \in \{A, B, R, S\} \quad (22)$$

With the following boundary conditions:

$$C_i(z = 0, t) = \begin{cases} C_{i,0} & 0 \leq t \leq t_{inj} \\ 0 & t > t_{inj} \end{cases} \quad (23)$$

$$\left. \frac{\partial C_i(z, t)}{\partial z} \right|_{z=L} = 0 \quad (24)$$

where C_i is the concentration and ν_i is the stoichiometric coefficient of the i^{th} chemical species, t is the time after injection of the reactive volume started, t_{inj} is the injection duration, z the spatial coordinate, and u the linear velocity.

The model output $f(\xi, \theta)$ corresponds to a unique experimental measurement, the area under the curve of iodine concentration at the reactor outlet as a function of time (Equation 25). Iodine concentration was selected as the measurable variable because it can be directly monitored in practice using an inline UV spectrometer.

$$f(\xi, \theta) = \int_0^{3\tau} C_B(z = L, t) dt \quad (25)$$

To solve the set of PDEs numerically, we used the method of lines. The axial coordinate z was discretized into 1000 uniformly spaced nodes, and both first- and second-order spatial derivatives were approximated using centred finite differences. The discretization converted the PDE into a system of



333 ODEs, solved with a Runge-Kutta method from the SciPy library (version 1.16.2) in Python (version
334 3.11.7)⁴⁰.

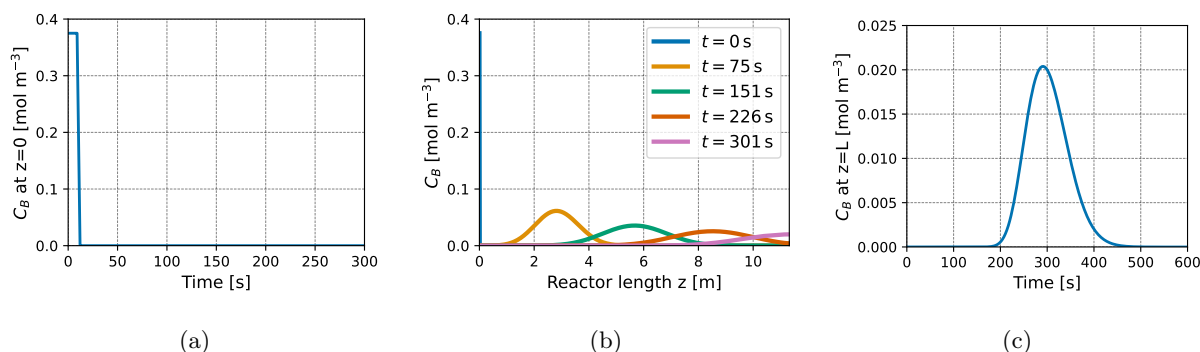


Figure 7: Transient behaviour of species B in a tubular reactor with axial dispersion. (a) Rectangular-wave inlet concentration of B. (b) Concentration of B along the reactor length at selected times. (c) Concentration of B at the reactor outlet.

335 For the simulations, a 5 mL tubular reactor with an internal diameter of 0.75 mm and a length of 11.32
336 m was considered. A molecular diffusion of 1×10^{-9} m²/s was assumed, and the dispersion coefficient
337 D was determined as described in the previous section. The total concentration, $(C_{A,0} + C_{B,0})$, was
338 held constant at 0.75 mol/m³, which defined the height of the rectangular injection profile. The ratio
339 t_{inj}/τ was fixed at 0.04, corresponding to a 200 μ L injection volume. As an example, Figure 7 shows
340 the simulated behaviour of species B for a space time of 5 min, an injection duration of 12 seconds, an
341 injection concentration of $C_B = 0.375$ mol/m³, $k_1 = 0.05$ m³ s⁻¹ mol⁻¹, and $k_2 = 0.02$ m³ s⁻¹ mol⁻¹.
342 The simulation starts with the discrete injection of the reagents (Figure 7a). The initial concentration
343 of B decreases along the reactor length due to the combined effects of reaction and dispersion (Figure
344 7b). The outlet concentration of B is tracked as a function of time (Figure 7c), mimicking the response
345 that could be measured with an inline UV detector. The area under this curve is then used as the single
346 output of our first-principles model.

347 5.2 ANNs training and testing

348 A Latin Hypercube Sampling was performed to generate 100 full model simulations to train the ANNs.
349 The dataset was generated by varying $\tau \in [1 \text{ min}, 10 \text{ min}]$, the initial concentration $C_{A,0} \in [0.0375$
350 mol/m³, 0.7125 mol/m³], and both reaction rate constants, k_1 and k_2 , within the range $[0.01 \text{ m}^3 \text{ s}^{-1}$
351 mol⁻¹, 0.1 m³ s⁻¹ mol⁻¹]. Simulations were run to estimate the value of the peak area for B at the outlet



352 of the reactor (Equation 25), and the values of $\partial f(\boldsymbol{\xi}, \boldsymbol{\theta})/\partial \theta_p$ were determined by the central difference
 353 method. The resulting peak areas ranged from 0.01 to 15.02 mol s m⁻³. The dataset was randomly
 354 divided into three subsets: 70% for training, 15% for validation, and 15% reserved as an independent
 355 test set. The data used to train the ANN should not be confused with the experimental data set used
 356 for parameter estimation.

357 The architecture of the neural network consisted on 4 inputs (τ , $C_{A,0}/0.75$, k_1 , k_2), 1 hidden layer of
 358 16 neurons, and one output layer corresponding to the area under the curve of B. The hidden layer size
 359 of 16 neurons was selected based on preliminary tests, as increasing the number of neurons beyond this
 360 value did not lead to any improvement in the surrogate model's accuracy. A training for 10 000 epochs
 361 with early stopping (a patience parameter of 500) was performed. Both sigmoid and tanh functions were
 362 tested by training 10 neural networks with a randomly initialized set of weights.

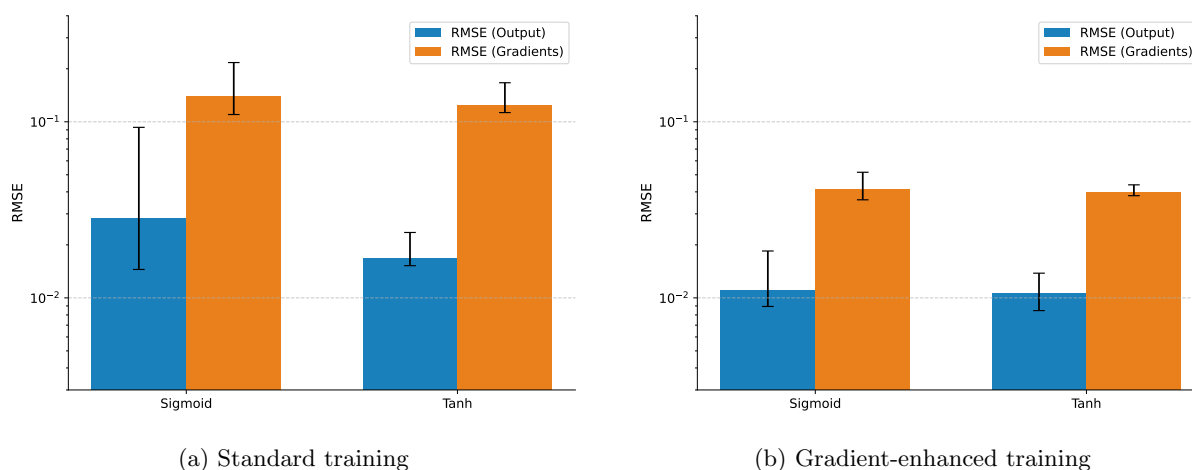


Figure 8: Median RMSE obtained from 10 independent neural network trainings for each activation function. Error bars indicate the computed interquartile range.

363 Figure 8 shows the RMSE obtained with both activation functions. Under standard training, the
 364 tanh activation achieved lower RMSE values than the sigmoid function for both output and gradient
 365 predictions. The higher accuracy of tanh in this case may be attributed to its resemblance to the error
 366 function (erf), which is commonly used to describe cumulative distributions such as the area under the
 367 curve after a pulse injection. Incorporating gradient information into the training reduced RMSEs in
 368 all cases relative to standard training, except for the output predictions with tanh, leading to similar
 369 performance between the two activation functions. The advantage provided by the resemblance between
 370 tanh and erf during standard training appears to diminish once gradient terms are included in the loss



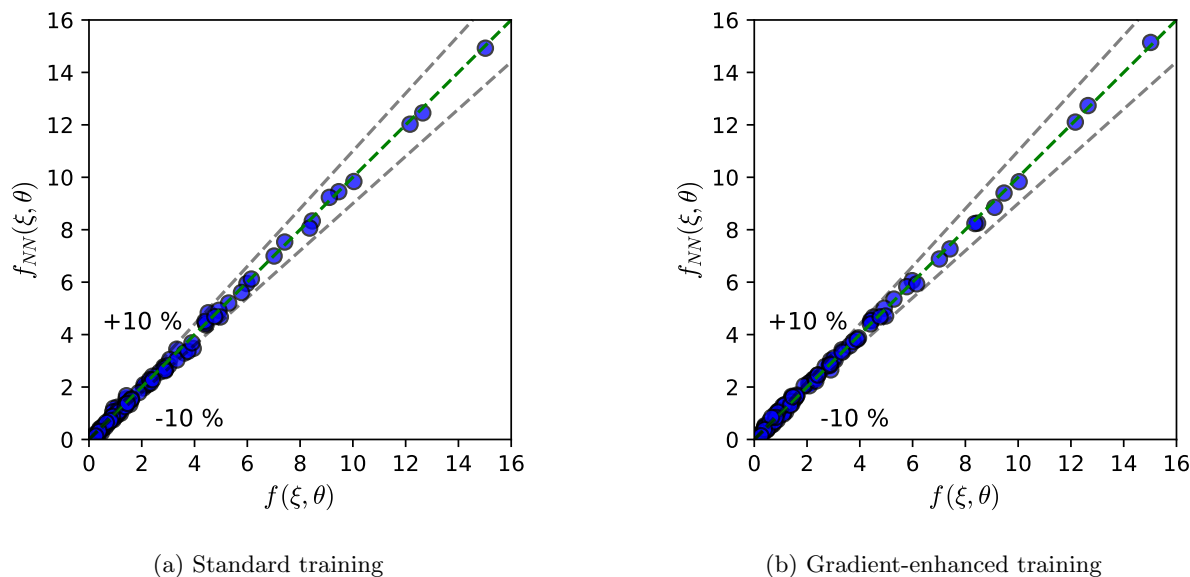


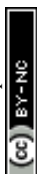
Figure 9: Parity plots comparing the ANN predicted output $f_{NN}(\xi, \theta)$ against the true values obtained from the full mechanistic model $f(\xi, \theta)$.

371 function. In addition, consistent with observation at Section 4, the inclusion of gradient information did
 372 not introduce additional computational cost. The gradient-enhanced models converged with an average
 373 training time of 4.2 s compared to 5 s for the standard approach.

374 Since the tanh outperformed the sigmoid activation in approximating the true function values under
 375 standard training, it was selected for further comparison with the gradient-enhanced neural network.
 376 Alternative gradient loss weights of 0.5 and 2.0 were tested, yielding higher average total MSE values
 377 (2.51×10^{-3} and 3.01×10^{-3} , respectively) than the baseline value of 2.06×10^{-3} obtained with weight of
 378 1.0.

379 From each of the two sets of 10 ANNs, one with standard training and one with gradient-enhanced
 380 training, the model with the lowest total RMSE (sum of output and gradient errors) was selected. These
 381 two surrogate models were then compared to evaluate their predictive ability and performance when
 382 applied to a MBDofE framework for parameter estimation. Both the best standard and gradient-enhanced
 383 networks accurately predicted the peak area of B , as shown in the parity plots of Figure 9, indicating
 384 that the inclusion of gradient information does not compromise the accuracy of output predictions.

385 In contrast, differences arise when comparing gradient predictions. The parity plots from Figures 10
 386 and 11 demonstrate that the gradient-enhanced network provides an improved accuracy for the derivatives
 387 w.r.t both kinetic parameters k_1 and k_2 . In addition, the sign of the gradients is distinguished with



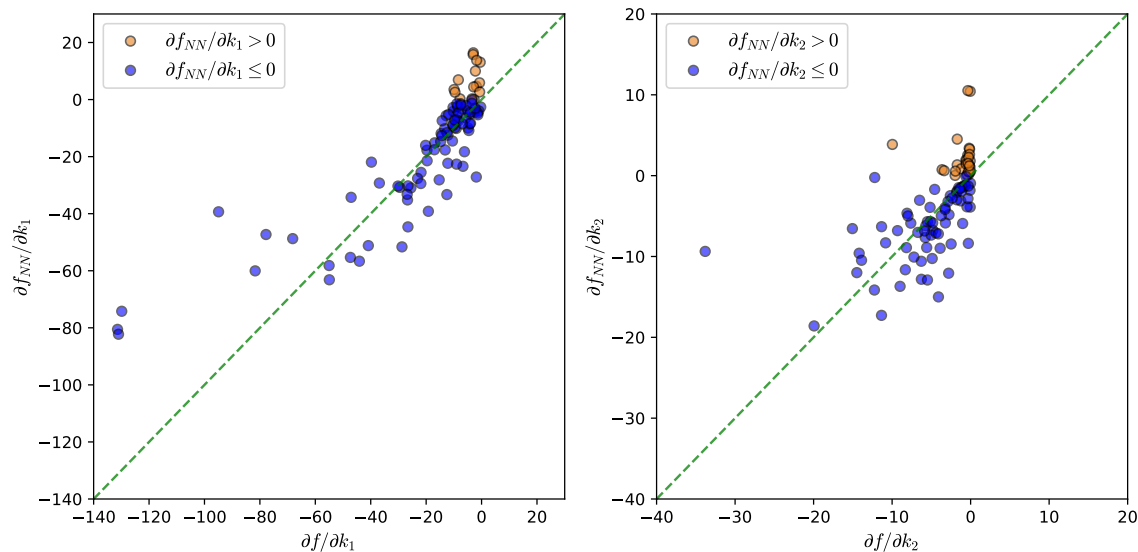


Figure 10: Parity plots comparing the ANN predicted gradients with respect to k_1 and k_2 under standard training against the true gradients obtained from the mechanistic model.

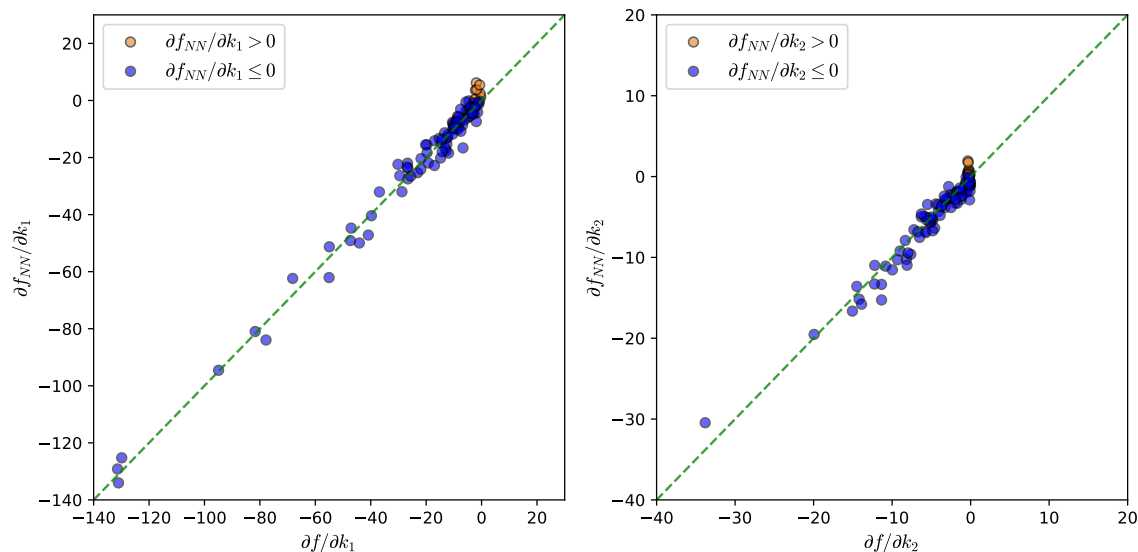


Figure 11: Parity plots for the ANN predicted gradients with respect to k_1 and k_2 under a gradient-enhanced training against the true gradients obtained from the mechanistic model.

388 colours, as all derivatives should be negative for physical consistency: indeed, increasing either kinetic
 389 constant while holding the other fixed increases the consumption of B and reduces the peak area at the
 390 reactor outlet, implying a negative derivative. In the standard training case (Figure 10), several predicted



391 gradients were positive, which is physically unrealistic. Incorporating gradient information (Figure 11)
392 reduced the occurrence of such predictions, and the few remaining positive values were close to zero,
393 despite no hard constraints being imposed during training. Therefore, gradient-enhanced training not
394 only improves the accuracy of gradient predictions but also mitigates the occurrence of gradients that
395 lack physical meaning. As an additional baseline, a standard GP surrogate trained on the same dataset
396 yielded a total MSE (1.89×10^{-2}) and parity plots comparable to those of the standard ANN (Figure
397 S.1), indicating similar limitations in reproducing parameter sensitivities.

398 5.3 ANNs integrated into an MBDoE sequential framework.

399 As a final evaluation, the surrogate models were integrated within a MBDoE framework. To initialize
400 the in-silico experimental planning, the first-principles model was simulated using the kinetic parameters
401 $k_1 = 0.0562 \text{ m}^3 \text{ mol}^{-1} \text{ s}^{-1}$ and $k_2 = 0.0193 \text{ m}^3 \text{ mol}^{-1} \text{ s}^{-1}$, previously reported⁸. A random experimental
402 error, $\varepsilon \sim \mathcal{N}(0, \sigma_e^2)$, was added to the computed area under the curve for species B, assuming a constant
403 σ_e across all measurements. The impact of the noise magnitude is tested below.

404 The campaign began with five preliminary runs selected via LHS. Because the initial design may
405 influence the outcome, different LHS designs comprising five runs were evaluated. Subsequently, the
406 sixth and all following experimental conditions were chosen iteratively according to the D-optimality
407 criterion. Optimization of ψ_D was performed using a differential evolution algorithm, as implemented in
408 SciPy⁴⁰.

409 As discussed in previous sections, a primary motivation for employing surrogate models within an
410 MBDoE framework is the substantial reduction in the computational burden associated with the opti-
411 mization routines used to design each experiment. Computing the sensitivity matrix \mathbf{S} using the full
412 first-principles model with a central difference scheme, for instance, requires two model evaluations for
413 each parameter θ_p at every experimental condition ξ_e . To assess computational cost, we measured the
414 time required to compute the $[5 \times 2]$ sensitivity matrix for the five preliminary experiments on a work-
415 station with 16 cores (Intel Xeon w5-3433, 2.00 GHz) and 256 GB RAM. Using finite differences with the
416 first-principles model, which required 20 full model evaluations, the computation time was 970 ± 143 s.
417 In contrast, the gradient-enhanced ANN required only 5 surrogate model evaluations and completed the
418 same task in 0.005 ± 0.003 s, corresponding to a reduction in computational cost by a factor of approx-
419 imately 200,000. For this comparison, finite differences was preferred over an adjoint approach. Given
420 that only two parameters are estimated, the benefits of employing an adjoint formulation are limited.



421 Additionally, due to the stiffness of the governing equations, the evaluation of sensitivities via the adjoint
422 method proved to be more computationally demanding (4571 ± 1278 s).

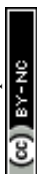
423 To quantify the time needed for a full MBDoE cycle, the number of function evaluations was assessed.
424 On average, 640 model evaluations are required for parameter estimation, along with 339 sensitivity
425 computations to determine the D-optimal experimental conditions. Consequently, identifying the optimal
426 experimental design using the first-principles model requires approximately 3–6 days, compared to only
427 1–5 s when using the surrogate model. This acceleration is critical for an online MBDoE strategy, which
428 relies on a fast iterative calculation of the sensitivity matrix to inform the next optimal experiment in
429 real-time applications.

430 Effectiveness was assessed using two complementary metrics: (i) the distance between estimated and
431 true parameter values, and (ii) the probability that the true parameters lie within the estimated confidence
432 region at the end of the workflow. Results are shown for a total budget of ten experiments, including
433 five preliminary runs, in Figures 12 for two noise levels ($\sigma_e = 0.2$ and 0.4 mol s/m^3 , corresponding to
434 approximately 7.5% and 15% relative error relative to the mean output).

435 To account for variability due to the initial design, the procedure was repeated using ten different
436 LHS initializations, and the median performance after each experiment was reported. In addition, the
437 coverage probability was estimated from 100 independent runs performed under the same experimental
438 budget. The resulting coverage probabilities, along with the median volume ratio and correlation error,
439 are summarized in Tables 1 and 2.

440 For comparison, three surrogate models were evaluated: the standard ANN with the lowest $\text{RMSE}_{\text{total}}$
441 ($\text{RMSE}_{\text{total}} = \text{RMSE}_{\text{output}} + \text{RMSE}_{\text{grad}}$), the standard ANN with the lowest $\text{RMSE}_{\text{output}}$, and the
442 gradient-enhanced ANN with the lowest $\text{RMSE}_{\text{total}}$. The results in Figure 12 reveal that the standard-
443 trained ANNs failed to converge to the true parameter values, even under the lower noise scenario. By
444 contrast, the gradient-enhanced ANN converged to the correct parameters. This behaviour arises because
445 the D-optimal designs often involve experimental conditions located at the boundaries of the design space
446 (Figures S.2-S.4). As demonstrated in the previous section, standard ANNs exhibit reduced accuracy
447 near boundary regions, particularly under limited training data. Consequently, the MBDoE procedure
448 systematically selects experiments in regions where the surrogate is least reliable, creating a feedback
449 loop that reinforces model bias and prevents convergence to the true parameters.

450 As expected in sequential MBDoE, all surrogates guide the algorithm toward experiments that reduce



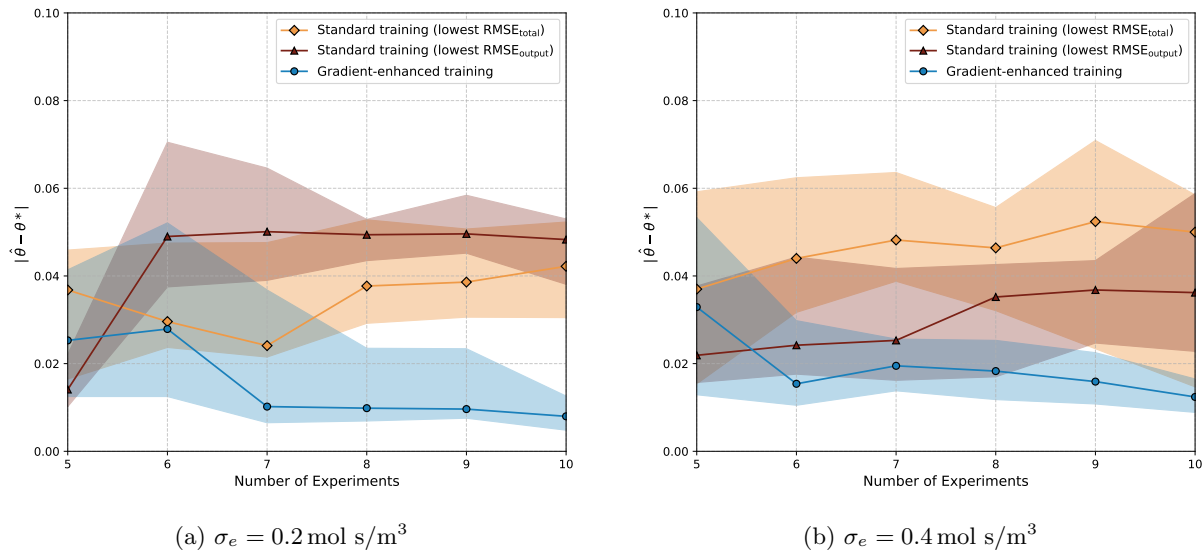


Figure 12: Median distance between the estimated parameters $\hat{\theta}$ and the true parameter values θ^* after each experiment. The shaded area shows the interquartile range of the distribution of values obtained for 10 different initial LHS designs.

parameter uncertainty, as evidenced by the steady decrease in the median determinant of the FIM (Figure S.5). However, reducing uncertainty alone is insufficient; a reliable surrogate must also ensure that the resulting confidence region is centered on the true parameter values. This aspect is captured by the coverage probability.

Table 1: Confidence region accuracy and coverage probabilities for sequential experimental design under lower noise conditions ($\sigma_e = 0.2 \text{ mol s/m}^3$). Volume ratio and correlation error are reported as the median [Interquartile Range] across 100 independent tests.

ANN	Volume ratio	Correlation error	Coverage probability
Standard (Output-Optimized)	0.30 [0.28, 0.34]	0.12 [0.10, 0.15]	4%
Standard (Total-Optimized)	0.26 [0.22, 0.46]	0.09 [0.04, 0.13]	7%
Gradient-enhanced	0.44 [0.39, 0.49]	0.15 [0.12, 0.21]	60%

Under lower noise conditions, the standard ANNs exhibit low coverage probabilities of 4% and 7%. As shown in Figure 12a, this failure is primarily due to their inability to approach the true parameter values. Interestingly, these models still show reasonable alignment in terms of correlation error (Table 1), indicating that their FIM-based covariance matrices adequately describe the spread of their own parameter estimates. This can be visually assessed in the supplementary material (Figures S.6-S.8), which illustrates that while the shape and size of the covariance ellipses generally match the parameter



Table 2: Confidence region accuracy and coverage probabilities for sequential experimental design under higher noise conditions ($\sigma_e = 0.4 \text{ mol s/m}^3$). Volume ratio and correlation error are reported as the median [Interquartile Range] across 100 independent tests.

ANN	Volume ratio	Correlation error	Coverage probability
Standard (Output-Optimized)	0.42 [0.33, 0.55]	0.20 [0.18, 0.23]	46%
Standard (Total-Optimized)	0.25 [0.21, 0.37]	0.07 [0.05, 0.16]	36%
Gradient-enhanced	0.62 [0.52, 0.76]	0.22 [0.13, 0.29]	68%

distributions, they are centered far from the true values due to the flawed local sensitivity structure they provide to the experimental design algorithm. In contrast, the gradient-enhanced ANN minimizes the distance to the true parameters, yielding an improved coverage probability of 60%. This value remains below the nominal 95% and may be attributed to limitations of FIM-based uncertainty quantification at this experimental budget ($N_e = 10$), as the FIM provides a local approximation of the parameter covariance matrix that may not fully capture nonlinear effects. As discussed in Section 2, these factors can lead to an underestimation of uncertainty, even when the surrogate accurately represents local model sensitivities.

Comparing to the higher noise scenarios reveals an important characteristic of the FIM-based uncertainty metrics. Under higher noise, the coverage probabilities for the standard models improve, rising to 46% and 36% (Table 2). This increase, however, does not reflect an improvement in parameter estimation accuracy; the distance trajectories (Figure 12b) show the estimates remain far from the true values. Rather, the higher experimental variance inflates the volume of the resulting confidence ellipses. The coverage probability increases primarily because these wider confidence bounds contain the true parameters more frequently, despite the central estimates remaining off-target. The gradient-enhanced ANN scales its coverage to 68% under these wider bounds.

As an example of a sequential experimental planning for parameter estimation, Figure 13 shows the evolution of the estimated parameters and the associated confidence ellipse. The results are displayed after the five preliminary experiments (LHS), after the sixth optimal experiment, and at subsequent steps, for both the ANN with the lowest $\text{RMSE}_{\text{total}}$ and the gradient-enhanced ANN. This trajectory is representative of the behaviour discussed above. In both cases, the sequential addition of experiments leads to a reduction of the confidence regions, indicating improved parameter precision. However, the gradient-enhanced ANN gives estimates closer to the true parameter values. This improved accuracy and reliability make the gradient-enhanced ANN a promising candidate to replace first-principles models, enabling its use in real-time, automated experimental platforms for flow chemistry.



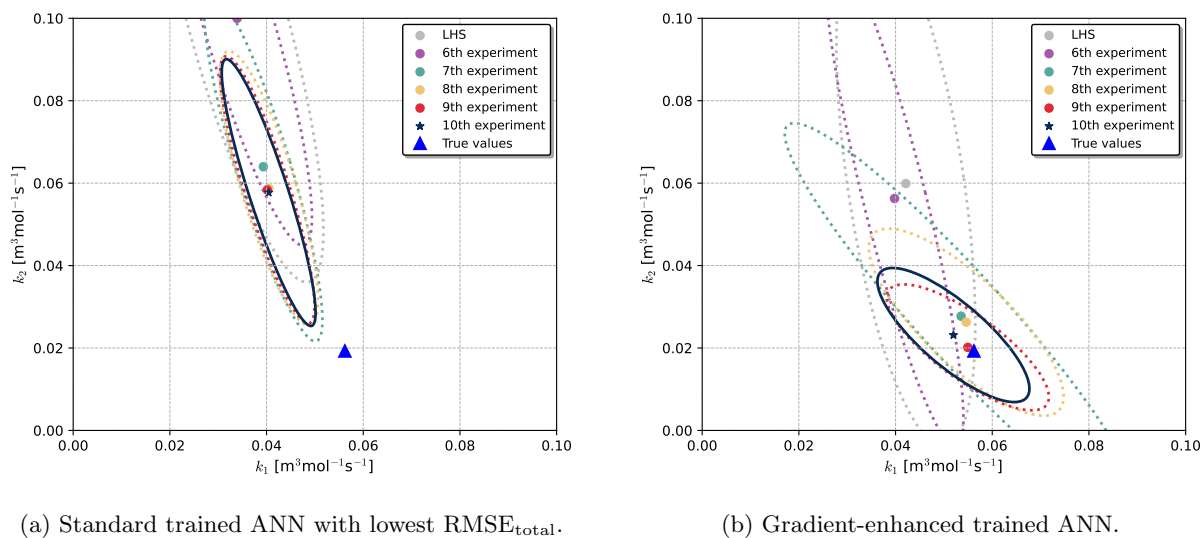


Figure 13: Confidence regions of the estimated kinetic parameters k_1 and k_2 during an MBDoE campaign.

6 Conclusions

One of the main barriers to applying MBDoE in automated flow chemistry platforms is the computational burden associated with complex models, particularly when dynamic experiments must be described. This work addressed this challenge by implementing within the MBDoE a gradient-enhanced neural network as a surrogate model to approximate the behaviour and sensitivities of first-principles models. Because automatic differentiation is available in most neural network implementation software, it is possible to compute derivatives with respect to model parameters during training. By incorporating gradient information, the local sensitivity structure is preserved, so the ANN is able to reproduce both model outputs and parameter sensitivities relevant for parameter estimation.

The case study highlighted the effectiveness of the gradient-enhanced neural network through *in silico* synthetic validation. The use of gradient information during training improved predictive accuracy of model derivatives, mitigated the occurrence of physically implausible gradients, and delivered reliable predictions near boundary values that are typically targeted by MBDoE as highly informative regions. The gradient-enhanced approach demonstrated higher robustness in parameter estimation, converging toward true parameter values in a higher percentage of cases than the standard ANN, supported by an increase in coverage probability. Standard ANN surrogates, however, exhibited poor coverage probability in scenarios with low noise and systematic bias in parameter estimates, indicating that accurate reproduction of sensitivities is key to ensuring reliable parameter estimation within MBDoE. In addition,



504 assessments of optimization timing confirmed a substantial reduction in the computational cost required
505 for iterative model evaluations. As a proof-of-concept, these results establish gradient-enhanced neural
506 networks as promising candidates for substituting first-principles models, providing a framework for their
507 future application in self-driven flow chemistry platforms where experimentation, parameter estimation,
508 and optimization can be achieved in a fully automated closed-loop.

509 Beyond flow chemistry, this framework could be applicable to other domains requiring accelerated
510 parameter estimation in real-time applications. However, the present study relies on synthetic validation
511 and specific noise model assumptions, which may not fully capture experimental variability or potential
512 mismatch between the model and the experimental platform encountered in practice. Validation with
513 real experimental workflows on automated platforms is therefore needed to assess the robustness of the
514 proposed approach. Additionally, limitations arise in systems with a large number of parameters, where
515 surrogate training becomes increasingly challenging and ill-conditioning of the FIM may occur. Future
516 work with larger models could focus on strategies before surrogate training, such as identifying a subset
517 of the most estimable parameters. This reduction in dimensionality would improve numerical stability
518 during FIM inversion and ease the training process by concentrating on the most influential parameters
519 under the given experimental conditions.

520 Author contributions

521 Francisco Bolaños-García: conceptualization, methodology, investigation, software, formal analysis, visu-
522 alization, writing - original draft. Jean-Marc Commenge: conceptualization, formal analysis, methodol-
523 ogy, validation, supervision, writing - review & editing. Laurent Falk: validation, supervision, funding
524 acquisition, writing - review & editing.

525 Conflicts of interest

526 There are no conflicts to declare.

527 Data availability

528 The datasets and source code used for the first-principles model and for the training of standard and
529 gradient-enhanced surrogate models are available in the GitHub repository at <https://github.com/FranBol/ANN->



530 MBDoe. The repository also includes the weights and biases for all neural networks. An archived release of
531 the repository has been deposited on Zenodo and is accessible via DOI: <https://doi.org/10.5281/zenodo.19818341>.
532 Supplementary Information (SI) regarding the experimental conditions proposed for improving parameter
533 estimation and visual representation of the estimated confidence regions is available.

534 Acknowledgements

535 The authors acknowledge the support of the French Agence Nationale de la Recherche (ANR), under
536 grant ANR-22-CE51-0025 (project OPTIFLEX).

537 References

- 538 (1) Delgado-Licona, F.; Addington, D.; Alsaiani, A.; Abolhasani, M. *Nature Chemical Engineering*
539 **2025**, *2*, 277–280.
- 540 (2) Lyall-Brookes, G.; Padgham, A. C.; Slater, A. G. *Digital Discovery* **2025**, *4*, 2364–2400.
- 541 (3) Slattery, A.; Wen, Z.; Tenblad, P.; Sanjosé-Orduna, J.; Pintossi, D.; den Hartog, T.; Noël, T. *Science*
542 **2024**, *383*, Publisher: American Association for the Advancement of Science, eadj1817.
- 543 (4) Rodriguez-Zubiri, M.; Felpin, F.-X. *Org. Process Res. Dev.* **2022**, *26*, Publisher: American Chemical
544 Society, 1766–1793.
- 545 (5) Buzzi Ferraris, G.; Forzatti, P.; Emig, G.; Hofmann, H. *Chemical Engineering Science* **1984**, *39*,
546 81–85.
- 547 (6) Buzzi-Ferraris, G.; Manenti, F. *Chemical Engineering Science* **2009**, *64*, 1061–1074.
- 548 (7) Geremia, M.; Macchietto, S.; Bezzo, F. *Chemical Engineering Science* **2026**, *319*, 122347.
- 549 (8) Mathieu, F.; Commenge, J.-M.; Falk, L.; Lomel, S. *Chemical Engineering Science* **2013**, *104*, 829–
550 838.
- 551 (9) Jiang, Z.; Portha, J.-F.; Commenge, J.-M.; Falk, L. *Chemical Engineering Research and Design*
552 **2019**, *146*, 290–310.
- 553 (10) McMullen, J. P.; Jensen, K. F. *Organic Process Research & Development* **2011**, *15*, 398–407.
- 554 (11) Reizman, B. J.; Jensen, K. F. *Organic Process Research & Development* **2012**, *16*, 1770–1782.
- 555 (12) Waldron, C.; Pankajakshan, A.; Quaglio, M.; Cao, E.; Galvanin, F.; Gavriilidis, A. *React. Chem.*
556 *Eng.* **2020**, *5*, 112–123.



- 557 (13) Senthil Vel, A.; Konan, K. E.; Cortés-Borda, D.; Felpin, F.-X. *Org. Process Res. Dev.* **2023**,
558 acs.oprd.3c00238.
- 559 (14) Rasch, A.; Bücker, H. M. *ACM Trans. Math. Softw.* **2010**, *37*.
- 560 (15) Wang, J.; Dowling, A. W. *AIChE Journal* **2022**, *68*, e17813.
- 561 (16) Olofsson, S.; Hebing, L.; Niedenführ, S.; Deisenroth, M. P.; Misener, R. *Computers & Chemical*
562 *Engineering* **2019**, *125*, 54–70.
- 563 (17) Franceschini, G.; Macchietto, S. *Chemical Engineering Science* **2008**, *63*, Model-Based Experimen-
564 tal Analysis, 4846–4872.
- 565 (18) Friso, A.; Galvanin, F. *Computers & Chemical Engineering* **2024**, *187*, 108724.
- 566 (19) Bhosekar, A.; Ierapetritou, M. *Computers & Chemical Engineering* **2018**, *108*, 250–267.
- 567 (20) Quaglio, M.; Roberts, L.; Bin Jaapar, M. S.; Fraga, E. S.; Dua, V.; Galvanin, F. *Computers &*
568 *Chemical Engineering* **2020**, *135*, 106759.
- 569 (21) Sangoi, E.; Quaglio, M.; Bezzo, F.; Galvanin, F. *Computers & Chemical Engineering* **2024**, *187*,
570 108752.
- 571 (22) Gusmão, G. S.; Retnanto, A. P.; da Cunha, S. C.; Medford, A. J. *Catalysis Today* **2023**, *417*,
572 Transient Kinetics Seminar, 113701.
- 573 (23) Czarnecki, W. M.; Osindero, S.; Jaderberg, M.; Świrszcz, G.; Pascanu, R. Sobolev Training for
574 Neural Networks, 2017.
- 575 (24) Srinivas, S.; Fleuret, F. Knowledge Transfer with Jacobian Matching, 2018.
- 576 (25) Paszke, A. et al. In *Proceedings of the 33rd International Conference on Neural Information Pro-*
577 *cessing Systems*; Curran Associates Inc.: 2019.
- 578 (26) Sellar, R.; Batill, S. In *6th Symposium on Multidisciplinary Analysis and Optimization*, 1996,
579 p 4019.
- 580 (27) Liu, W.; Batill, S. In *8th Symposium on Multidisciplinary Analysis and Optimization*, 2000, p 4923.
- 581 (28) Tsay, C. *Computers & Chemical Engineering* **2021**, *153*, 107419.
- 582 (29) Asprey, S.; Macchietto, S. *Computers & Chemical Engineering* **2000**, *24*, 1261–1267.
- 583 (30) Hunter, W. G.; and, A. M. R. *Technometrics* **1965**, *7*, 307–323.
- 584 (31) Pronzato, L.; Pazman, A., *Design of experiments in nonlinear models*, 2013th ed.; Lecture notes
585 in statistics; Springer: New York, NY, 2013.



- 586 (32) Bard, Y., *Nonlinear parameter estimation*; Academic press New York: 1974; Vol. 1209.
- 587 (33) Walter, É.; Pronzato, L., *Identification de modèles paramétriques à partir de données expérimentales*;
588 Masson: 1994.
- 589 (34) Himmelblau, D. M. *Korean Journal of Chemical Engineering* **2000**, *17*, 373–392.
- 590 (35) Ulaganathan, S.; Couckuyt, I.; Dhaene, T.; Degroote, J.; Laermans, E. *Engineering with Computers*
591 **2015**, *32*, 15–34.
- 592 (36) Sapienza, F.; Bolibar, J.; Schäfer, F.; Groenke, B.; Pal, A.; Boussange, V.; Heimbach, P.; Hooker,
593 G.; Pérez, F.; Persson, P.-O.; Rackauckas, C. Differentiable Programming for Differential Equations:
594 A Review, 2024.
- 595 (37) Chen, R. T. Q.; Rubanova, Y.; Bettencourt, J.; Duvenaud, D. *Neural Ordinary Differential Equations*, 2019.
596
- 597 (38) Levenspiel, O., *Chemical Reaction Engineering*; Wiley: 1999.
- 598 (39) Wagner, F.; Sagmeister, P.; Jusner, C. E.; Tampone, T. G.; Manee, V.; Buono, F. G.; Williams,
599 J. D.; Kappe, C. O. *Advanced Science*, *11*, 2308034.
- 600 (40) Virtanen, P. et al. *Nature Methods* **2020**, *17*, 261–272.



Data availability

The datasets and source code used for the first-principles model and for the training of standard and gradient-enhanced surrogate models are available in the GitHub repository at <https://github.com/FranBol/ANN-MBDoE>. The repository also includes the weights and biases for all neural networks. An archived release of the repository has been deposited on Zenodo and is accessible via DOI: <https://doi.org/10.5281/zenodo.19818341>. Supplementary Information (SI) regarding the experimental conditions proposed for improving parameter estimation and visual representation of the estimated confidence regions is available .

