



**HAL**  
open science

# Clustering-Enhanced Deep Learning Method for Computation of Full Detailed Thermochemical States via Solver-Based Adaptive Sampling

Xi Chen, Cédric Mehl, Thibault Faney, Florent Di Meglio

► **To cite this version:**

Xi Chen, Cédric Mehl, Thibault Faney, Florent Di Meglio. Clustering-Enhanced Deep Learning Method for Computation of Full Detailed Thermochemical States via Solver-Based Adaptive Sampling. *Energy & Fuels*, 2023, 37 (18), pp.14222-14239. hal-04341667

**HAL Id: hal-04341667**

**<https://ifp.hal.science/hal-04341667v1>**

Submitted on 13 Dec 2023

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Clustering-enhanced deep learning method for computation of full detailed thermochemical states via solver-based adaptive sampling

Xi Chen,<sup>\*,†,‡</sup> Cédric Mehl,<sup>‡</sup> Thibault Faney,<sup>‡</sup> and Florent Di Meglio<sup>†</sup>

<sup>†</sup>*Centre Automatique et Systèmes, Mines Paris - PSL, 60 Bd Saint-Michel, 75272, Paris, France*

<sup>‡</sup>*IFPEN, 1 et 4 av. du Bois Preau, 92852, Rueil-Malmaison, France, Institut Carnot IFPEN Transports Energie*

E-mail: xi.chen1@minesparis.psl.eu

## Abstract

Detailed chemistry computations are indispensable in numerous complex simulation tasks, which focus on accurately capturing the ignition process or predicting pollutant levels. Machine learning method is a modern data-driven approach for predicting full detailed thermochemical state-to-state behavior in reacting flow simulations. By combining unsupervised clustering algorithms to subdivide the composition space, the complexity of adaptive regression models for temporal dynamics can be significantly reduced. In this article, a more compact dataset is generated, which is essential for the clustering algorithm, by leveraging the adaptive CVODE solver time steps for data augmentation for stiff reactive states. A learning workflow that utilizes a deep residual network model (ResNet) in conjunction with an adaptive clustering algorithm is proposed. This approach aims to replace the stiff ODE direct integration solver traditionally used for computing thermochemical species' state-to-state temporal evolution

14 for detailed chemistry simulations. The learning models are adaptively trained using the  
15 K-Means clustering algorithm in the nonlinear transformation space for different sub-  
16 spaces of dynamic systems. Three test cases:  $H_2$  (9 species),  $C_2H_4$  (32 species), and  
17  $CH_4$  (53 species), are investigated, each exhibiting varying complexities. The study  
18 demonstrates that the iterative predictions of thermochemical states align well with  
19 the results obtained from direct numerical integration. Additionally, employing multi-  
20 ple adaptive regression models in subdomains yields superior performance compared to  
21 a single regression model prediction case.

## 22 1 Introduction

23 With the growth of computing power, simulating reacting flows in space and time at various  
24 scales has become feasible. However, typical chemical processes entail a multitude of species  
25 and a vast number of chemical reactions, coupled with diffusive and convective transport  
26 phenomena. As a result, computationally demanding tasks, such as simulating industrial  
27 processes like combustion in fuel engines or glass furnaces, and simulation in complex chem-  
28 ical reactors, present significant challenges. Modeling chemical processes involves non-linear  
29 stiff ordinary differential equations (ODEs) and requires appropriate solvers like CVODE.<sup>1</sup>  
30 Integrating the chemical ODEs often becomes the bottleneck in reactive flow simulations.

31 State-of-the-art methods for chemistry accelerations traditionally rely on chemistry re-  
32 ductions and chemistry tabulation methods. In the case of chemistry reductions, only the  
33 major chemical species are considered, while unimportant species and reactions are excluded  
34 from the complex mechanisms. Various approaches, such as the Quasi-Steady State Approx-  
35 imation (QSSA)<sup>2</sup> or Directed Relation Graph (DRG),<sup>3</sup> have been proposed to derive reduced  
36 mechanisms. These reduced mechanisms are then integrated using direct integration solvers,  
37 utilizing available computational resources. However, reduced mechanisms containing a large  
38 number of species are still necessary for studying pollutants like nitrogen compounds. Al-  
39 ternatively, tabulation methods, such as Look-Up Tables (LUT)<sup>4</sup> and In-Situ Adaptive Tab-

40 ulation (ISAT),<sup>5</sup> can be employed to pre-compute chemical solutions based on canonical  
41 problems using detailed chemical mechanisms. The pre-computed terms are stored in tables  
42 with a reduced set of variables and utilized later for chemistry integrations. Nevertheless,  
43 the curse of dimensionality limits their application to simple cases, as the storage memory  
44 exponentially increases with the number of chemical species.

45 In addition to traditional chemistry tabulation, machine learning algorithms using arti-  
46 ficial neural networks offer an alternative method for accelerating detailed chemistry inte-  
47 gration. Machine learning has recently gained significant interest in the applied science and  
48 engineering fields.<sup>6</sup> It has been widely applied to computer vision (CV), natural language  
49 processing (NLP) for several years, and more recently to various scientific computation prob-  
50 lems, such as for reactive flow simulations.<sup>7,8</sup> In theory, deep learning models can serve as  
51 surrogate models capable of fitting any non-linear model<sup>9,10</sup> Previous studies have demon-  
52 strated that neural networks can replace traditional integration methods for chemical accel-  
53 erations.<sup>11,12</sup> One notable advantage of utilizing deep neural networks (DNNs) over chemical  
54 look-up tables is their low memory requirement.<sup>13</sup> The number of parameters needed for a  
55 DNN is much smaller compared to the number of multi-dimensional points that would have  
56 to be stored in a table for a complex chemical reaction. Complex neural network models  
57 with moderate memory storage capacity can handle high-dimensional inputs, making them  
58 suitable for large chemical mechanisms, whereas Look-Up Tables and ISAT store only a  
59 limited number of input variables. Further studies have extended the application of neu-  
60 ral networks to even more complex cases. Sen and Menon<sup>14</sup> employed neural networks for  
61 chemistry integration in turbulent premixed and non-premixed abstract problems, including  
62 syngas combustion with a 14-species mechanism and methane combustion with a 16-species  
63 mechanism. Wan et al.<sup>15</sup> utilized ANN chemistry integration for the direct numerical sim-  
64 ulation of a syngas turbulent oxy-flame with side-wall effects, incorporating an 11-species  
65 reduced mechanism.

66 Most of these research studies demonstrate the efficacy of deep learning models in accel-

erating the computation of source terms compared to direct numerical simulations. Recent investigations have further expanded the scope by focusing on more complex fuels and problems. For instance, Ranade et al.<sup>16</sup> employed DNN models to reconstruct the temporal pyrolysis and oxidation processes of complex fuels based on experimental data. Similarly, Ding et al.<sup>17</sup> applied multiple DNN models to predict the time evolution of each of the 32 states involved in  $CH_4$  combustion, with individual neural networks dedicated to each chemical species. Furthermore, An et al.<sup>18</sup> employed neural networks which replace the direct integration of the reduced skeletal mechanism of kerosene with  $C_{10}H_{22}$  in the context of engine system simulations, where the total dimension number is strongly reduced to 41 with 132 elemental reactions neglecting many minor species. These studies reveal the potential of employing neural networks to accelerate chemical integration for complex fuels while maintaining good prediction accuracy. In addition, embedding augmented physical information into the loss function of neural networks represents a more recent approach to ensure that the prediction results align with the numerical solution of ODEs.<sup>19</sup> However, research on this approach has so far been limited to simple cases, necessitating further extensions and investigations. Moreover, several studies also combine machine learning methods and traditional approach to tackle the complex fuel problems, including using data-driven method to reduce the full chemistry to subspace manifold by linear and nonlinear models,<sup>20-23</sup> and using neural networks to predict the flamelet generated manifolds.<sup>24-27</sup>

Splitting the chemical manifold into several subdomains and training multiple ANNs can improve learning efficiencies by facilitating the training process and reducing local model complexity. Various methods have been employed to efficiently separate subdomains and enhance prediction accuracy. Physical descriptions of the combustion process, such as progress variables<sup>28</sup> or the rate of temperature increase,<sup>29</sup> can effectively partition the subdomains and improve prediction accuracy. Another approach is to utilize data-driven classification algorithms. One such algorithm is the self-organizing map (SOM) with a neural network structure, which preserves the topological structure of the data and produces a low-dimensional

94 representation to partition input data points.<sup>18,30,31</sup> Another option is the K-Means parti-  
95 tioning algorithm. This algorithm partitions the data by calculating the distance between  
96 each data point and a fixed number of centroids, which represent different clusters. In earlier  
97 research, Perini et al.<sup>32</sup> applied an optimal K-Means algorithm to partition the full chem-  
98 istry states into subdomains with similar reactive conditions, successfully accelerating the  
99 simulation time of combustion with high-dimensional full chemistry in engine environments.  
100 Similarly, Barwey et al.<sup>33</sup> used a predefined number of clusters to partition the data and  
101 demonstrated its ability to separate stiff reactive regions from non-stiff regions in combustion  
102 processes. Additionally, Nguyen et al.<sup>34</sup> employed a hierarchical K-Means algorithm with  
103 predefined cluster numbers. These advances underscore the effectiveness of employing non-  
104 supervised partitioning algorithms in complex combustion simulations, enabling improved  
105 prediction accuracy and efficient handling of different combustion regimes.

106 The performance of ANN training results strongly relies on the quality of the dataset used.  
107 Generating an appropriate dataset that includes information from different scales of chemical  
108 species is a challenging task. In early studies of simple chemical mechanisms, researchers  
109 employed random sampling of chemical species within predefined operating ranges<sup>11, 12</sup> Han  
110 et al. applied a fixed small-scale sampling time step to ensure the capture of fast ignition  
111 and heat release processes. Additionally, they introduced input simulation noise during  
112 each resolution time step to generate a more robust dataset. In a study by Zhang et al.,<sup>35</sup>  
113 different strategies for data-driven generative sampling methods were compared. The authors  
114 proposed a new multi-scale sampling method for different scales of chemical species. In this  
115 approach, the species were classified into major and minor groups, and random sampling was  
116 performed at fixed intervals using either normal or logarithmic distributions. Each of these  
117 methods aims to ensure the dataset captures the essential features of the chemical processes  
118 at different scales.

119 The review of these studies highlights the potential of using ANN-based chemistry in-  
120 tegration for more complex fuels. However, combining non-supervised K-Means clustering

121 algorithm with local adaptive temporal full detailed chemistry computation for high dimen-  
122 sion complex fuels has not been evaluated. Additionally, multi-scale data sampling becomes  
123 more challenging for high-dimensional complex fuels with numerous minor species. In this  
124 paper, we propose a workflow for learning 0D chemical kinetics using unsupervised clustering  
125 combined with deep residual networks. The performance of models by varying the number  
126 of clusters and network hyperparameters is systematically analysed. To generate a compact  
127 dataset, the adaptive time marching steps of the CVODE solver is applied to generate the  
128 sampling sequence, which allows to obtain more data points in the fast ignition and reac-  
129 tion regions. Multiple models are trained individually for different composition spaces, with  
130 subdomains of the dynamical system identified using the K-Means clustering algorithm in  
131 logarithmic space. The results, along with the analysis of model training and iterative infer-  
132 ence statistics, demonstrate that more accurate predictions can be achieved by empirically  
133 selecting optimal cluster numbers as hyperparameters for different combustion cases. Fur-  
134 thermore, the deep residual network outperforms traditional multi-layer perceptrons (MLPs)  
135 in terms of training performance.

136 The remaining sections of the paper are organized as follows. Section 2 provides a gen-  
137 eral description of the chemistry evolution equations. Section 3 presents a novel strategy  
138 for generating the training datasets based on a canonical multiple initial conditions numer-  
139 ical experiment. Section 4 introduces the general workflow, including data pre-processing,  
140 clustering algorithm, and the proposed regression model. Section 5 focuses on the evalua-  
141 tion of data clustering, model training, and statistical inference based on actual 0D ignition  
142 simulations.

## 143 2 Physical problem formulation

144 This section details the chemical kinetics problem. The ODEs describing the problem are  
145 detailed in Sec. 2.1 and the numerical solving in Sec. 2.2.

## 146 2.1 Combustion equations

147 Reactive flows involve the coupling between convection, diffusion and chemical reactions. A  
148 wide range of physical scales are present in the flow and the resolution of the fully coupled  
149 system is often out of reach. To address this challenge, a common approach is to employ  
150 operator splitting methods<sup>36, 37</sup> which solves the chemistry and transport phenomena sepa-  
151 rately. In a multi-dimensional simulation, the chemical system in each computational cell is  
152 described by a set of ODEs in time:

$$\begin{aligned}\frac{dY_s}{dt} &= \frac{M_s}{\rho} \dot{\omega}_s \quad \forall s \in \{1, \dots, N_s\} \\ \frac{dT}{dt} &= -\frac{1}{\rho C_p} \sum_{s=1}^{N_s} h_s \dot{\omega}_s\end{aligned}\tag{1}$$

153 where  $N_s$  represents the total number of chemical species. The variables  $Y_s$ ,  $h_s$ ,  $M_s$ , and  
154  $\dot{\omega}_s$  denote the mass fraction, molar enthalpy, molar weight, and chemical reaction rate for  
155 species  $s$ , respectively. The variables  $T$ ,  $C_p$ , and  $\rho$  represent the temperature, specific heat  
156 capacity at constant pressure, and density of the gas mixture, respectively. The reactors are  
157 assumed to be adiabatic and homogeneous. The system is expressed as follows in a generic  
158 form:

$$\begin{aligned}\dot{\mathbf{S}}(t) &= \mathbf{f}(\mathbf{S}(t), q) \quad t \in [0, t_{n+1} - t_n] \\ \mathbf{S}(t = 0) &= [T(t_n), Y_1(t_n), \dots, Y_{N_s}(t_n)]^T\end{aligned}\tag{2}$$

159 where  $q$  is a variable regrouping all thermodynamic constants.

## 160 2.2 Chemical kinetics solver

161 A common solver for solving the system described in 2 is CVODE.<sup>1</sup> This solver is particularly  
162 employed in CANTERA,<sup>38</sup> which is a chemical computation software used for calculating  
163 thermodynamic and chemical species terms in reactive flow systems, including the present  
164 work. CVODE is a multi-step solver with variable order and step sizes, and it utilizes



165 dynamic adaptive time stepping. To approximate the solution  $\mathbf{S}(t_n) = \mathbf{S}^{(n)}$  at time  $t_n$ ,  
166 CVODE solves the following algebraic equation:

$$\sum_{i=0}^{K_1} \alpha_{n,i} \mathbf{S}^{(n-1)} + h_n \sum_{i=0}^{K_2} \beta_{n,i} \dot{\mathbf{S}}^{(n-i)} = 0 \quad (3)$$

167 where  $h_n = t_{n+1} - t_n$  is the time step. To tackle stiff problems, Backward Differentiation  
168 Formulas (BDF) is specifically considered in fixed-leading coefficient (FLC) form,<sup>1</sup> defined  
169 by  $K_1 = q$  and  $K_2 = 0$ . The order  $q$  ranges from 1 and 5. The coefficients are fixed by the  
170 method type, its order, the recent history of the step sizes and the normalization  $\alpha_{n,0} = -1$ .  
171 The standard CVODE chemical solver in CANTERA effectively handles thermochemical  
172 states using implicit and combined stiff/non-stiff solvers.<sup>38</sup> By adaptively resolving the stiff  
173 and non-stiff regions with different time steps, the CVODE solver refines the step sizes in  
174 areas with rapid reaction rates. **Figure 1 illustrates the simulations of temperature with**  
175 **dynamically adjusted time steps during the temporal evolution of  $H_2$ ,  $C_2H_4$ , and  $CH_4$  cases,**  
176 **where the x axis for time evolution is plotted under the logarithmic scale.** The initial  
177 temperature and equivalence ratio are set to  $(T_0, \phi) = (1700K, 1.0)$ . It can be observed  
178 that the CVODE solver refines the time step  $dt_{cvsode}$  in fast-reacting regions while using  
179 larger time steps in the starting and equilibrium regions.

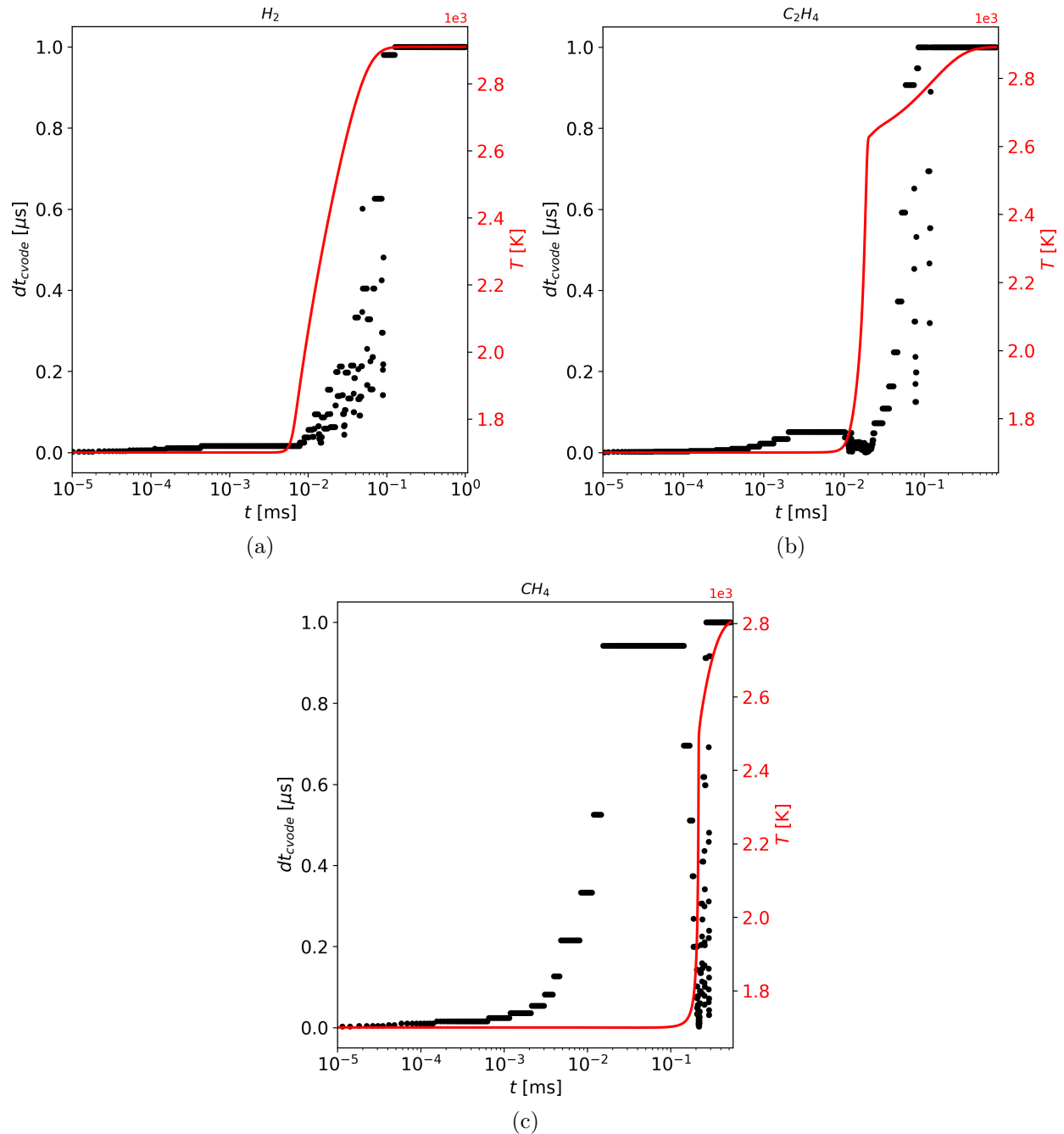


Figure 1: Dynamic adaptive time steps used by CVODE solver with (a)  $H_2/air$ , (b)  $C_2H_4/air$ , and (c)  $CH_4/air$  cases, where the simulations are set with  $T_0 = 1700.0K$  and  $\phi = 1.0$ . The red lines represent the numerical solution of temperature, and the black dot points represent the local adaptive time steps given by CVODE solver.

180 During a 3D simulation using operator splitting, the ODE solver typically performs nu-  
 181 merous adaptive time steps within a single computational fluid dynamics (CFD) time step

182  $dt_{cfd} = t_{n+1} - t_n$ . The goal is to replace these individual time steps with a single function  
 183 call that approximates  $S(t + dt_{cfd})$  based on  $S(t)$ . In this paper,  $dt_{cfd}$  is set to a fixed value  
 184 of  $dt_{cfd} = 10^{-6}$  seconds. In simulations using operator splitting approaches, the time step  
 185 is usually limited by stability limits on convection and diffusion (CFL and Fourier numbers  
 186 are typically defined) and thus dictated by the CFD solver. The extension of the current  
 187 strategy to variable  $dt_{cfd}$  is out-of-the-scope of the present paper and will be tackled in fu-  
 188 ture research. The constant  $dt_{cfd}$  value selected here is a typical value found in large eddy  
 189 simulations of systems such as gas turbines or engines. The following section outlines how  
 190 CVODE is utilized to generate the dataset that will be used for training a model capable of  
 191 achieving this objective.

192 The following section outlines how CVODE is utilized to generate the dataset that will  
 193 be used for training a model capable of achieving this objective.

### 194 3 Dataset generation

195 In this paper, the study of three different fuel combustion scenarios involving the fuels  $H_2$ ,  
 196  $C_2H_4$  and  $CH_4$ , as outlined in table 1. The selection of these cases progressively increases the  
 197 complexity of the thermochemical reaction system.

Table 1: Summary of three different studied fuel cases

| Mechanism              | number of species | number of reactions |
|------------------------|-------------------|---------------------|
| $H_2$ <sup>39</sup>    | 9                 | 19                  |
| $C_2H_4$ <sup>40</sup> | 32                | 206                 |
| $CH_4$ <sup>41</sup>   | 53                | 325                 |

#### 198 3.1 Generation of simulation trajectories

199 The design of experiments to generate the training dataset consists in specifying multiple  
 200 initial conditions for the temperature and species mass fractions for the set of ODE 2.

201 However, not all initial species' mass fractions are physically relevant: in this work, we  
 202 assume that the initial composition is a pure mixture of fuel and air, where the mass fractions  
 203 for all other species are set to zero. A common practice in combustion is to specify the initial  
 204 mass fractions values with a more significant variable called the equivalence ratio  $\phi$ , which is  
 205 a measure of the excess of fuel in the mixture with respect to stoichiometry, and expressed  
 206 as:

$$\phi = \frac{m_{fuel}/m_{ox}}{(m_{fuel}/m_{ox})_{st}} \quad (4)$$

207 where  $m$  are masses and  $n$  are numbers of moles, and the suffix  $st$  refers to stoichiometric  
 208 conditions. The equivalence ratio provides a measure of the excess or deficiency of fuel in  
 209 the mixture. By specifying the equivalence ratio, we can effectively control the initial mass  
 210 fractions of the species in the combustion simulation. The range of initial conditions (IC)  
 211 for the simulations is limited. The pressure is kept constant at the standard atmospheric  
 212 pressure of 1 atm. The temperature ranges from  $1600K$  to  $1800K$ , while the equivalence ratio  
 213 varies from 0.7 to 1.5. To generate the training dataset, a total of 1000 initial conditions  
 214 are randomly selected within these intervals using Latin Hypercube Sampling (LHS), as  
 215 illustrated in Figure 2. Each initial condition is then used to simulate a trajectory that  
 216 represents the evolution of species mass fractions and temperature over time. The resulting  
 217 database consists of pairs  $(S(t), S(t + dt_{cfd}))$  sampled along these trajectories. Since our  
 218 objective is to simulate over multiple time steps starting from a new given initial condition  
 219  $(T_0, \phi)_{new}$ , it is important to split the database based on trajectories rather than individual  
 220 pairs across all trajectories. all pairs of data  $(S(t), S(t + dt_{cfd}))$  from a specific training  
 221 trajectory will be assigned to the training database. The total dataset is splitted with a ratio  
 222 of 75%/15%/10% for train, validation, and test datasets, respectively, based on trajectories.  
 223 A criterion is determined to end a given simulation by defining the following variable:

$$\tau = \left| \frac{T_{eq} - T(t)}{T_{eq}} \right| \quad (5)$$

224 Where  $T(t)$  denotes the temperature of the local time step and  $T_{eq}$  denotes the tem-  
 225 perature at the equilibrium state that is determined for given initial conditions  $(\phi, T_0)$  by  
 226 a simple thermodynamic equilibrium computation. The simulations are terminated when  
 227 the total simulation time reaches  $10^{-3}s$ . At this point, it is assumed that all variables have  
 228 reached their equilibrium states. Alternatively, if the simulation reaches convergence with  
 229 a given tolerance  $\tau$ , the simulation is also considered complete. In this work, the tolerance  
 230  $\tau$  is set to  $10^{-4}$  to ensure that all simulations from different initial conditions reach their  
 231 final equilibrium states. In the next section, a novel strategy for generating the dataset is  
 232 presented, which takes advantage of the time step adaptation feature of the CVODE solver.

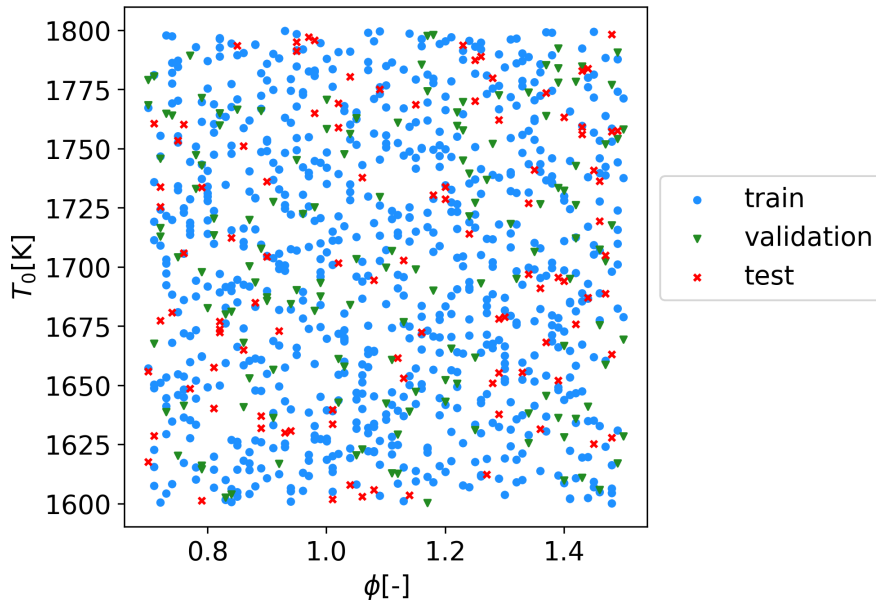


Figure 2: Initial conditions distribution

### 233 3.2 Data acquisition from simulation trajectories

234 The data pairs  $(\mathbf{S}(t), \mathbf{S}(t + dt_{efd}))$  are acquired within the generated **trajectories in the**  
 235 **chemical manifold**. However, the states evolution rates vary with time. Roughly speaking,  
 236 reactions are "slow" close to the equilibrium and "fast" after ignition. A naive strategy

237 would consist in selecting sampling points from the target trajectory during the simulation  
 238 with a regular time step, noted as  $dt_s$ . In other words, the dataset would contain pairs of  
 239 the form  $(\mathbf{S}(kdt_s), \mathbf{S}(kdt_s + dt_{cfd}))$ , where  $dt_s$  denotes the uniform sampling time step. The  
 240 straightforward approach of selecting data points from the target trajectory with a regular  
 241 time step ( $dt_s$ ) may result in an imbalanced dataset. This method would result in fewer data  
 242 points in the ignition region, where chemical reactions occur rapidly. As a consequence, the  
 243 dataset would be imbalanced with a skewed distribution of data points. This imbalance can  
 244 be observed in the temperature and  $CO_2$  mass fraction distributions, as shown in Figure 3  
 245 for the case of  $C_2H_4$ .

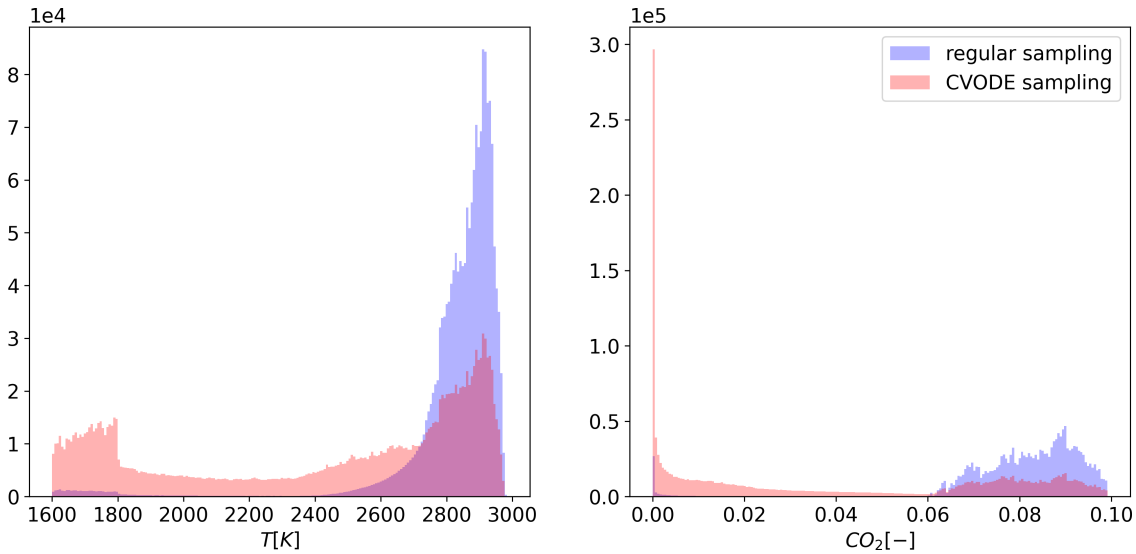


Figure 3: The sampling points distribution of (a) temperature and (b)  $CO_2$  values generated by regular sampling method (blue) and CVODE sampling method (red) for the  $C_2H_4/air$  case. The total size of data points generated by two strategies is around  $1.5 \times 10^6$ .

246 In this research, the dynamic time step adaptation feature of the CVODE solver is utilized  
 247 to generate the sampling sequence. This feature allows us to generate a more balanced  
 248 dataset by leveraging the time step adaptation performed by CVODE during the simulation,  
 249 as presented in section 2.2. To generate the dataset, a format of  $(\mathbf{S}(t_{cnode}), \mathbf{S}(t_{cnode} + dt_{cfd}))$   
 250 is applied, where  $t_{cnode}$  represents the time sequence obtained from the CVODE solver's time  
 251 step adaptation. To achieve this, the simulations by CVODE run twice at each time step:

252 once to obtain the points  $\mathbf{S}(t_{cvode})$ , and a second time to obtain the sequence  $\mathbf{S}(t_{cvode} + dt_{cfd})$ .  
253 By taking adaptive time steps, this approach ensures that there are more sampled points in  
254 fast reaction regions, which is beneficial for training the model. Furthermore, the maximum  
255 CVODE evolution time step is limited to  $dt_{cfd}$  to ensure a sufficient resolution in the near-  
256 equilibrium region. This restriction helps maintain accuracy in those regions. The resulting  
257 dataset is better balanced compared to using a constant time step, as demonstrated in  
258 Figure 3. This balanced dataset is crucial for training neural networks effectively. In the  
259 next section, a detailed description of the training pipeline is provided.

## 260 4 Learning methodology

261 The objective in this study is to replace the direct integration by deep ANN models (DNN)  
262 for each CFD time step resolution, as illustrated in 4. The objective is to have a model that  
263 can predict the thermochemical states at time  $t_0 + k\Delta t$  as  $\mathbf{S}(t_0 + k\Delta t) = \mathbf{f}_{DL}^{(k)}(\mathbf{S}(t_0); \theta)$ , where  
264  $\theta$  represents the parameters of the learning model. **In this workflow  $\mathbf{S}$  represents the states**  
265 **of temperature and each chemical species. Indeed, for such systems the temperature might**  
266 **be directly estimated from species mass fractions and enthalpy. In an attempt to make the**  
267 **model more general and not solely tailored for constant pressure low Mach problems, the**  
268 **temperature is included as an input of the DNN in the context of this research.** In this  
269 section, a detailed explanation of how the dataset is processed to train the neural networks  
270 for predicting the thermochemical states from their values at time  $t$  is introduced.

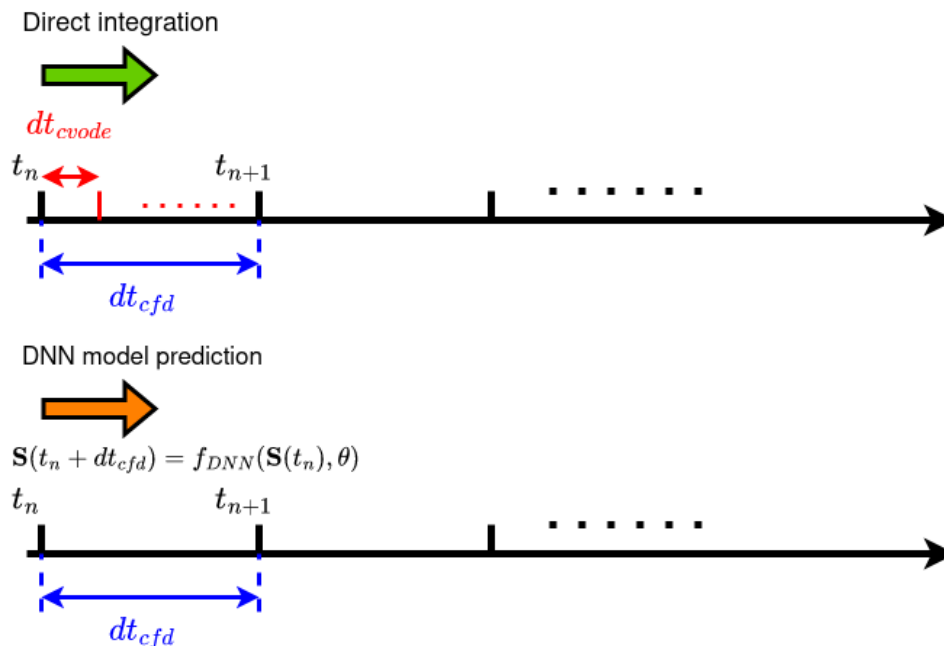


Figure 4: The workflow for the states prediction with time evolution

## 271 4.1 Data Pre-processing

272 To address the issue of chemical species values being in different scales and having ex-  
 273 tremely small values skewed towards zero, a pre-treatment of the chemical species values  
 274 is performed using a logarithmic transformation. The transformed state vector is denoted  
 275 as  $\mathbf{S}_t = [T, \ln(Y_j)]$ , where  $j = 1, \dots, N_s$ . This transformation ensures that the variables are  
 276 scaled to similar ranges, which is particularly beneficial for improving the prediction accu-  
 277 racy of minor chemical species. However, before applying the logarithmic transformation, it  
 278 is necessary to clip the zero values by a threshold. This is because the logarithmic transfor-  
 279 mation is only applicable to positive values. Empirical thresholds are determined for each  
 280 dataset, where the thresholds used for the  $H_2$ ,  $C_2H_4$ , and  $CH_4$  cases are  $10^{-10}$ ,  $10^{-12}$ , and  
 281  $10^{-28}$ , respectively. These thresholds have been chosen to preserve the accuracy of reactor  
 282 computations. **Two aspects seem to have an important impact on the threshold choice:**

- 283 • The smallest species mass fraction magnitude differ from one chemical mechanism to
- 284 another. For instance here, the maximum magnitude is around  $10^{-8}$  and the maximum



285 magnitude is around  $10^{-40}$  for  $C_2H_4$  and  $CH_4$ . The integration needs to take these  
286 species into account to recover the proper chemical evolution and this sets a higher  
287 bound for the threshold.

- 288 • The thresholds need to be high enough to suppress numerical artifacts which arise from  
289 the CVODE numerical integration for some chemical species with small mass fractions.

290 A quantitative comparison of training results with different thresholds for  $C_2H_4$  case are  
291 provided in to supporting information as an example.

292 In addition to the logarithmic transformation, a scaling of the data is applied to facilitate  
293 the training of the DNN models. A standard normalization technique is employed, which  
294 sets the mean of each feature to zero and the variance to one. This scaling helps in ensuring  
295 that all features contribute effectively to the learning process.

## 296 4.2 Data clustering

297 A strategy to simplify the learning process is to separate the composition space into sev-  
298 eral sub-domains. This division is achieved using the K-Means method, which is a non-  
299 supervised clustering algorithm. The K-Means algorithm partitions the sampling points  
300 into sub-domains by minimizing the average squared distance between  $K$  centroids and the  
301 sampling points. The K-Means algorithm considers the mean distance in the logarithmic  
302 scale to maintain consistency with the data pre-processing used for DNN regression. Each  
303 sub-domain corresponds to a separate region of the composition space. In this approach, a  
304 distinct DNN model is trained for each sub-domain. This allows for more specialized and  
305 focused modeling within each region of the composition space. In this work, the K-Means al-  
306 gorithm is implemented based on K-Means++, which is an improved version of the standard  
307 K-Means algorithm. K-Means++ enhances the clustering quality compared to the standard  
308 algorithm. More information about this algorithm can be found in appendix A.

### 309 4.3 Neural network model

310 Neural networks, particularly multiple-layer perceptrons (MLP), have been extensively used  
311 as surrogates for combustion kinetics solvers. MLPs are popular due to their simplicity and  
312 computational efficiency. Numerous studies have demonstrated the accuracy of MLPs in  
313 combustion problems<sup>7,8</sup> However, when dealing with large chemical mechanisms, accurate  
314 regression using MLPs may require a significant number of parameters. This can lead to  
315 optimization challenges, such as the vanishing gradient problem commonly encountered in  
316 machine learning applications.<sup>6</sup> The vanishing gradient problem refers to the issue where the  
317 gradients during backpropagation diminish as they propagate through many layers, making  
318 it difficult for the network to learn effectively. To address this challenge, an alternative  
319 approach used in this work is the application of residual layer structures with shortcut  
320 connections, known as ResNet.<sup>42</sup> Similar topology of ResNet design has been utilized in  
321 various physico-chemistry computations, such as fluid flash computations<sup>43</sup> and flamelet  
322 progress variables tabulation.<sup>26</sup> The ResNet architecture has been shown to mitigate the  
323 vanishing gradient problem and enable efficient learning for complex problems. Figure 5  
324 provides an illustration of the ResNet architecture. The ResNet architecture in our research  
325 consists of basic units called resblocks, which include two standard hidden layers. The key  
326 principle of ResNet is the inclusion of shortcut connections that bypass one or more layers,  
327 allowing each layer to predict an increment rather than a direct value. This enables the  
328 network to learn residual information, which facilitates the training process. If  $x$  is the input  
329 of a resblock, the output  $y$  is given by:

$$\mathbf{y} = \sigma(\mathcal{F}(\mathbf{x}) + \mathbf{x}) \quad (6)$$

330 where  $\mathcal{F}(\mathbf{x})$  is the feed-forward neural network composed of the two hidden layers and  $\sigma$   
331 represents the non-linear activation function. The number of residual blocks is denoted as  
332  $n_r$ , and the neuron number in each hidden layer is represented as  $n_e$ . The swish activation

333 function<sup>44</sup> will be used throughout this work:

$$\text{swish}(x) = \frac{x}{1 + e^{-\beta x}} \quad (7)$$

334 In this equation,  $\beta$  is set by 1 by default as indicated in<sup>44</sup>. In practical applications, par-  
 335 ticularly for regression problems, the use of the *swish* activation function has been found to  
 336 be more efficient for the optimization process compared to commonly used activation func-  
 337 tions such as *ReLU*, *sigmoid*, or *tanh*.<sup>29,43</sup> This observation has been verified through various  
 338 experiments conducted in this study. A summary of the DNN-based model architecture is  
 339 provided in Fig. 6.

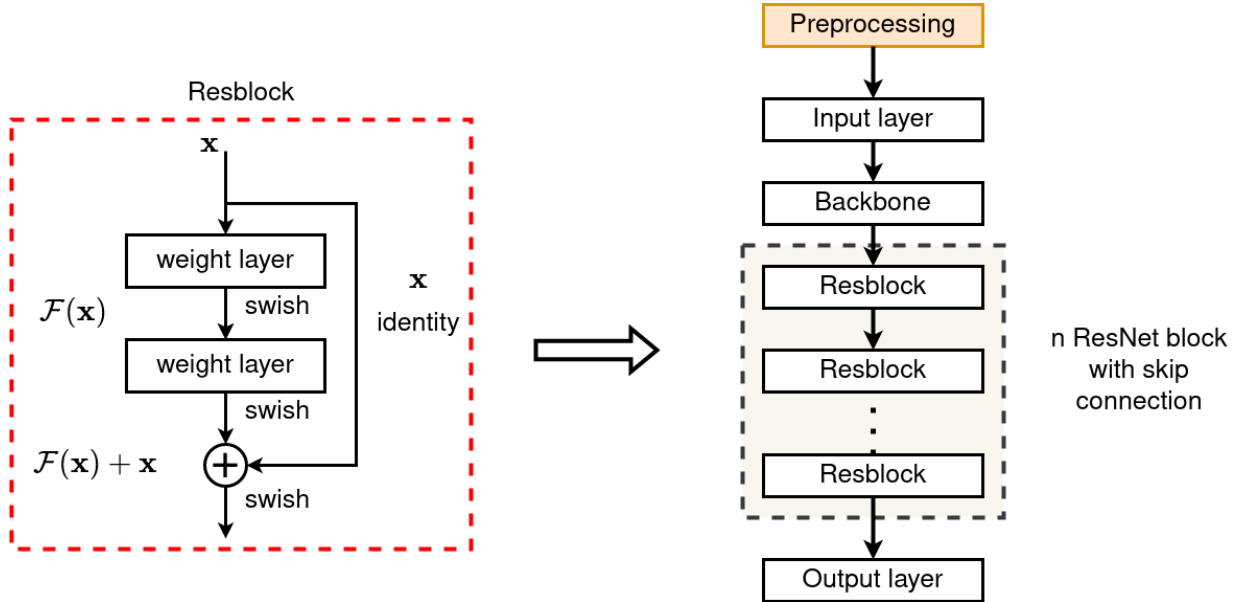


Figure 5: Neural network learning model structure design

#### 340 4.4 Model training

341 In the present work, the Mean Squared Error (MSE) is used as the loss function for model  
 342 training. The DNN models are trained using the Tensorflow 2.10.1 framework.<sup>45</sup> For opti-  
 343 mization, the Adam algorithm is chosen.<sup>46</sup> The model parameters are initialized using the  
 344 Glorot Uniform method, and an exponential decay strategy is applied to the initial learning

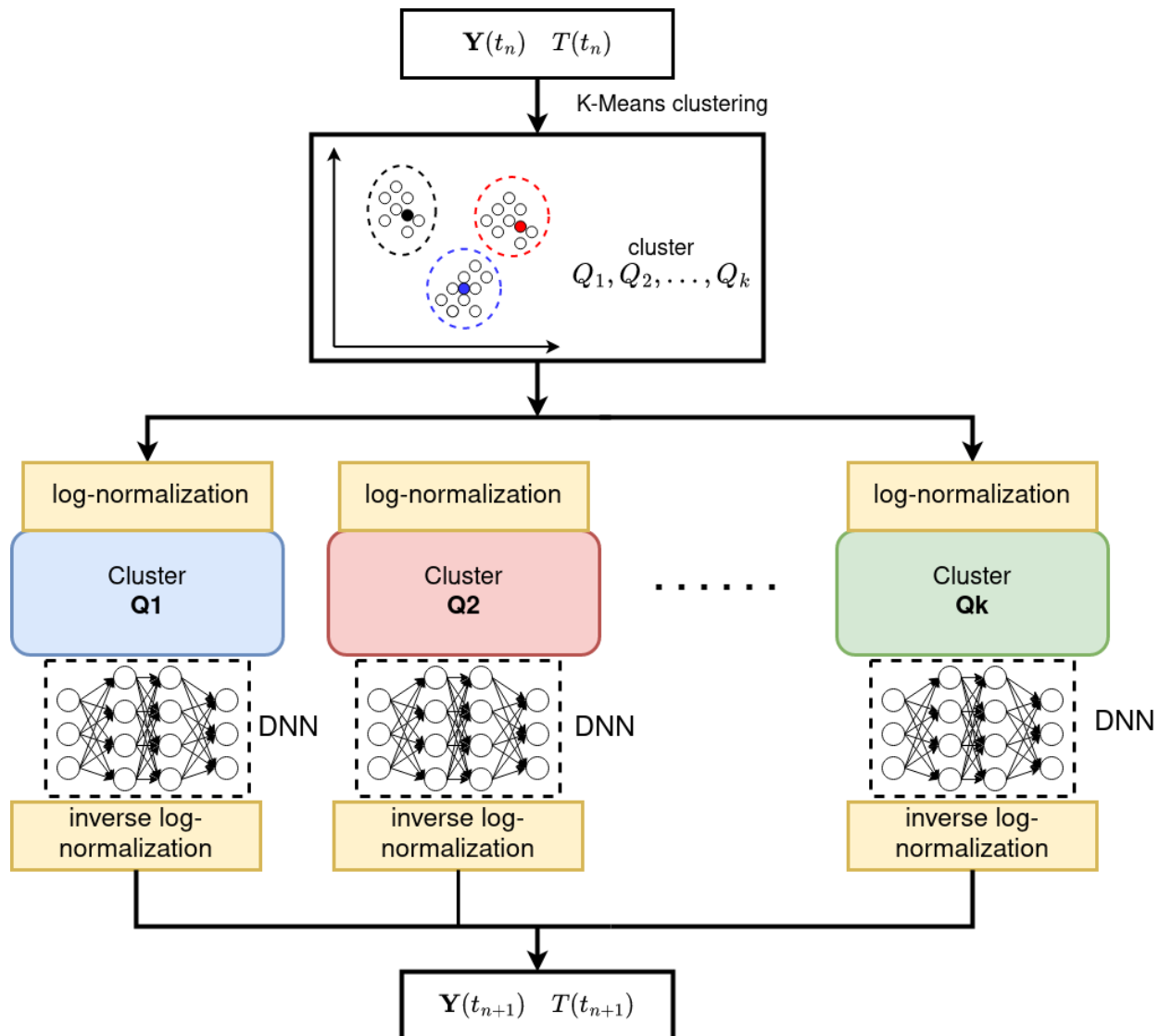


Figure 6: The global learning algorithm for the states prediction with the separation of composition space

345 rate during the training epochs. In this research, the initial learning rate for all three cases  
 346 is set to 0.005. A learning decay rate is employed to aid the optimization process,<sup>47</sup> and it  
 347 is defined as  $r(k) = r_0 \times \eta^{\frac{k}{N}}$ , where  $r_0$ ,  $\eta$ ,  $k$ , and  $N$  represent the initial learning rate, the  
 348 decay rate, the  $k^{th}$  step, and the total number of decay steps, respectively. In this study, the  
 349 decay rate is empirically set to 0.92 and the number of decay steps to 650.

## 350 4.5 Model performance evaluation

351 The DNN model serves as an operator for numerical time-stepping, denoted as  $\mathbf{f}_{DL}$ . By  
 352 utilizing this learned operator, we can simulate the trajectory of a chemical reactor with  
 353 a fixed initial condition by iteratively computing  $\mathbf{S}(t_0 + k\Delta t) = \mathbf{f}_{DL}^{(k)}(\mathbf{S}(t_0 + (k - 1)\Delta t); \theta)$ .  
 354 During the time evolution, the thermochemical states are computed iteratively in each time  
 355 step resolution until the end of the simulation. Therefore, it is insufficient to evaluate  
 356 the prediction performance for a single input-output pair of the models. The error may  
 357 accumulate and lead to divergence after multiple iterations. Hence, it is crucial to assess the  
 358 overall inference performance over multiple iterations for the entire simulation.

359 To compare the results with reference numerical simulations, a global accumulative log-  
 360 arithmic mean average percentage error  $\mathcal{M}_i$  is used to evaluate the performance of chemical  
 361 species trajectories under different initial conditions. Additionally, the normal mean average  
 362 percentage error  $\mathcal{M}_0$  is employed to evaluate the temperature prediction. It is important to  
 363 note that since the models are trained using data in logarithmic space for chemical species,  
 364 the logarithmic error metric is used for consistency. These introduced errors are denoted  
 365 mathematically by:

$$\begin{aligned} \mathcal{M}_0(T_0, p_0, \phi) &= \frac{1}{N_{iter}} \sum_{k=1}^{N_{iter}} \left| \frac{T^{pred}(kdt) - T(kdt)}{T(kdt)} \right| \\ \mathcal{M}_i(T_0, p_0, \phi) &= \frac{1}{N_{iter}} \sum_{k=1}^{N_{iter}} \left| \frac{\ln(Y_i^{pred}(kdt)) - \ln(Y_i(kdt))}{\ln(Y_i(kdt))} \right| \quad i = 1, \dots, N_s \end{aligned} \quad (8)$$

where  $i$  represents each state component (temperature and chemical species), and  $N_{iter}$  is the number of total iterations until the final time step prediction which is specific for each simulation with different initial conditions. Each  $\mathcal{M}_i$  is computed for one trajectory with one initial condition, for one state component. The overall average errors  $\mathcal{M}$  are computed by averaging values of all dimensions, which is written as:

$$\mathcal{M}(T_0, p_0, \phi) = \frac{\mathcal{M}_0(T_0, p_0, \phi) + \sum_{i=1}^{N_s} \mathcal{M}_i(T_0, p_0, \phi)}{N_s + 1} \quad i = 1, \dots, N_s \quad (9)$$

366 The global errors  $\mathcal{M}$  and errors for each state component  $\mathcal{M}_i$  are statistically evaluated  
 367 for all 100 initial conditions in the test set. This metric accounts for the potential error  
 368 accumulation that may arise when repeatedly using a neural network to predict the evolution  
 369 of chemical states. By taking into account the complete range of initial conditions, the overall  
 370 performance and reliability of the neural network model can be evaluated in accurately  
 371 capturing the dynamics of the chemical reactions.

## 372 5 Results

373 This section presents the results of the DNN model training performance and evaluates the  
 374 simulation results achieved using the trained models. Primarily, a parametric assessment of  
 375 the K-means clustering is conducted in Section 5.1. Then, the investigation of the perfor-  
 376 mance of the DNN surrogate model on 0-D reactor simulations is carried out in Section 5.2.  
 377 Finally, a particular study about the influence of the database generation method, comparing  
 378 regular and CVODE-based sampling is introduced in Section 5.3.

### 379 5.1 Analysis of data clustering

The number of K-means clusters needs to be selected beforehand as there is no method for  
 systematic selection. A practical approach to evaluate the clustering efficiency is to compute  
 the distortion  $\mathcal{D}$ , which corresponds to the squared Euclidean distance between the data  
 points and the centroids:

$$\mathcal{D} = \sum_{i=1}^K \sum_{\mathbf{S} \in \Omega_i} \left\| \hat{\mathbf{S}}_i - \mathbf{d}_i \right\|^2 \quad (10)$$

380 where  $\mathbf{S}_i$  is the logarithm of the state vector and  $\hat{\mathbf{S}}_i$  its normalized counter-part.  $\mathbf{d}_i$  represent  
 381 the coordinates of the  $K$  centroids. The evolution of the distortion with the number of  
 382 clusters is shown in Figure 7 for  $H_2$ ,  $C_2H_4$  and  $CH_4$ . It can be observed that as the number  
 383 of clusters increases, the distortion decreases monotonically, with a steeper slope initially and

384 a smoother evolution thereafter. In this study, the optimal number of clusters for each fuel  
385 is determined empirically within a predefined range by evaluating inference simulations. It  
386 should be noted that while the data is separated into subdomains with a converged distortion  
387 value, the local data distribution may not be optimal for model training, which can lead to  
388 inference instabilities during iterative predictions. Models are evaluated with up to six  
389 clusters, as the distortion value does not significantly decrease for larger numbers of clusters,  
390 indicating that all sampling points have already been optimally partitioned around local  
391 centroids. Moreover, with an increasing number of clusters after six, some subdomains  
392 may have insufficient data for effective local model training. After partitioning the dataset,  
393 individual learning models are trained and evaluated for each subdomain.

## 394 5.2 Model training evaluations

395 Several model training strategies are crucial for achieving optimal performance in this re-  
396 search. Firstly, trajectory learning under the logarithmic scale is necessary to address the  
397 prediction of extremely small values. After the data preprocessing, the rescaled data for each  
398 subdomain model in different cases are shown in Figure 8. The logarithmic transformation  
399 helps to transform the original data distribution into a more homogeneous distribution and  
400 removes extreme values skewed towards zero.<sup>48</sup> Data standardization, on the other hand,  
401 rescales the original data to have a mean of zero and a unit variance.

402 Training with a deeper neural network structure, including residual shortcuts, is more  
403 efficient compared to traditional multiple-layer perceptrons (MLP). This is demonstrated by  
404 comparing a standard MLP with the same number of hidden layers and neurons in each  
405 layer to our optimal trained ResNet models. Additionally, the activation function *swish*  
406 is expected to improve the optimization process compared to the commonly used *ReLU*  
407 activation function in previous research papers<sup>15, 34</sup>. The evolution of the loss function over  
408 training epochs for models in a cluster containing strongly reactive states is depicted in  
409 Figure 9. It can be observed that the *swish* activation function significantly promotes

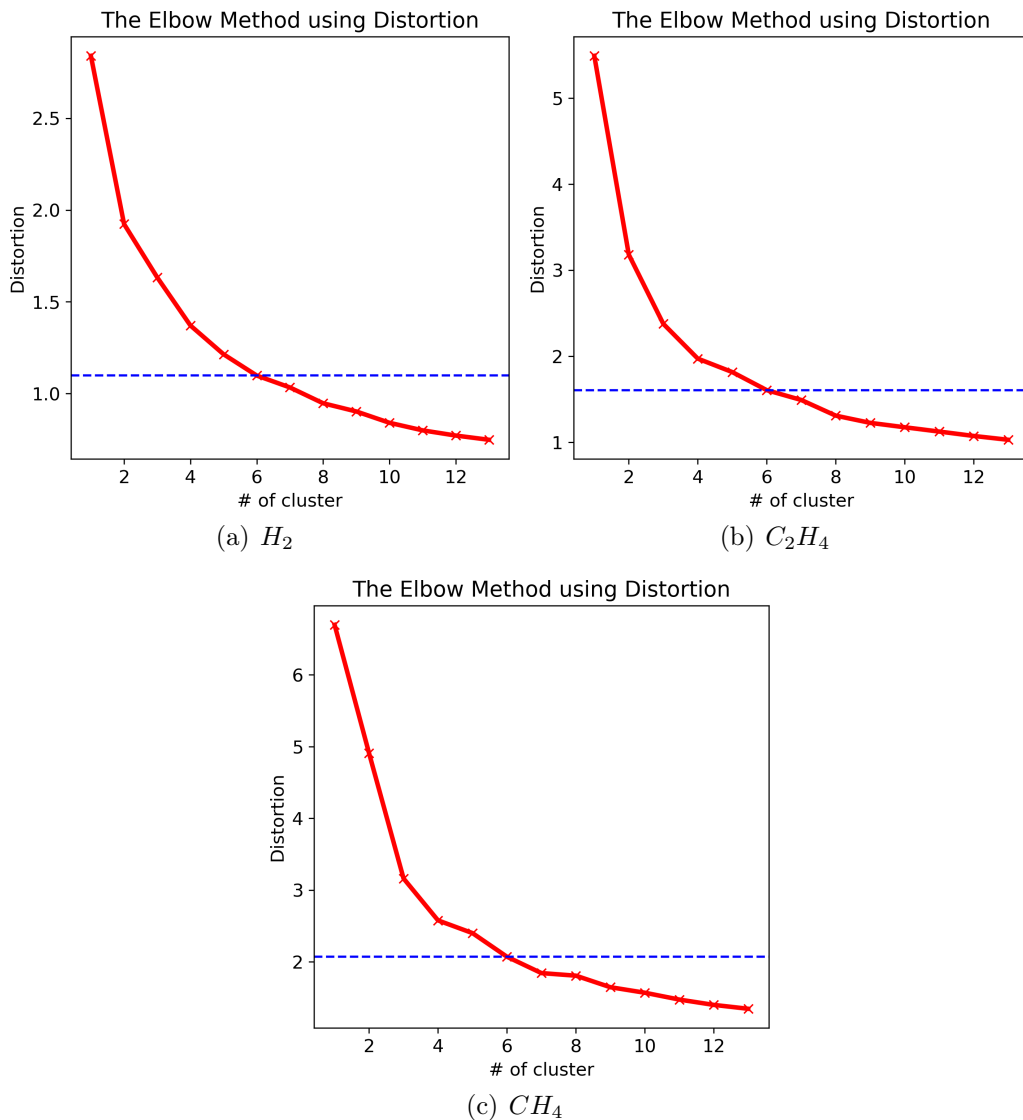
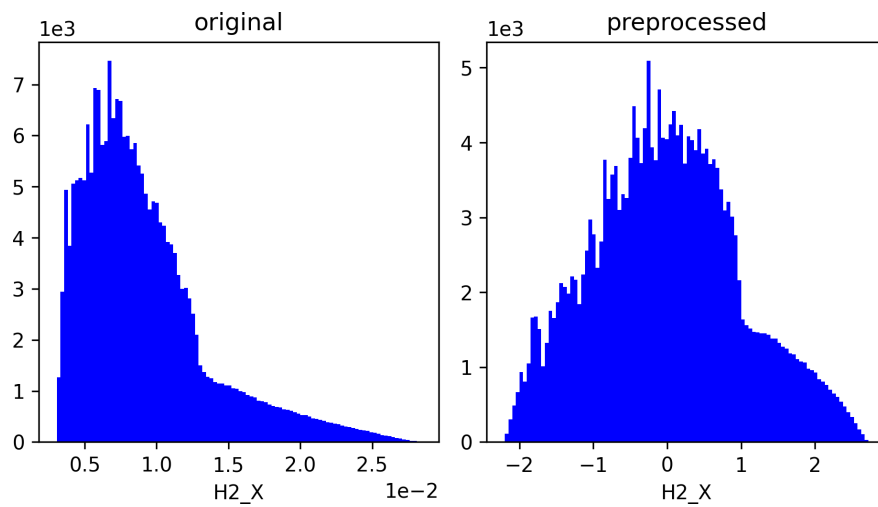


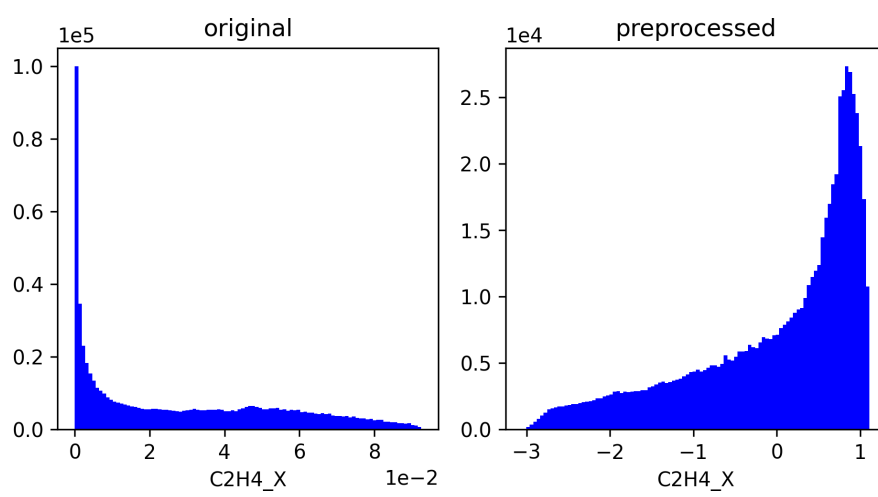
Figure 7: distortion values over cluster numbers for (a)  $H_2/air$ , (b)  $C_2H_4/air$ , and (c)  $CH_4/air$  cases

410 the optimization process, leading to smaller loss function values compared to the *ReLU*  
 411 activation function. Furthermore, the residual networks with skip-connection structures also  
 412 improve the optimization results for reactive zones. By employing the selected strategy in this  
 413 work, which is a residual neural network model with the *swish* activation function, we achieve  
 414 the lowest loss values during the optimization compared to other strategies. Additionally,  
 415 there is no significant difference between training and validation losses, indicating that there  
 416 is no overfitting during the training process. Similar conclusions regarding the absence of

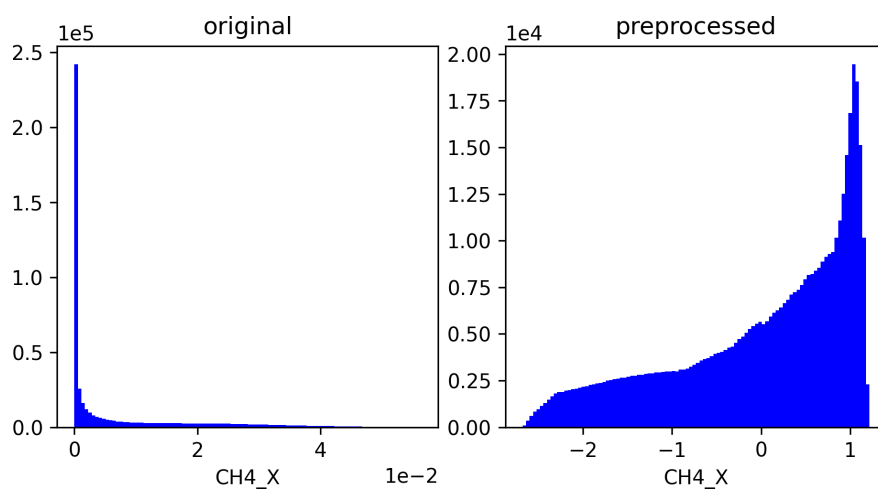




(a)



(b)



(c)

Figure 8: Distribution of training data from reactive subdomains before and after the preprocessing with nonlinear logarithmic transformation: (a)  $H_2$ , (b)  $C_2H_4$ , and (c)  $CH_4$ .

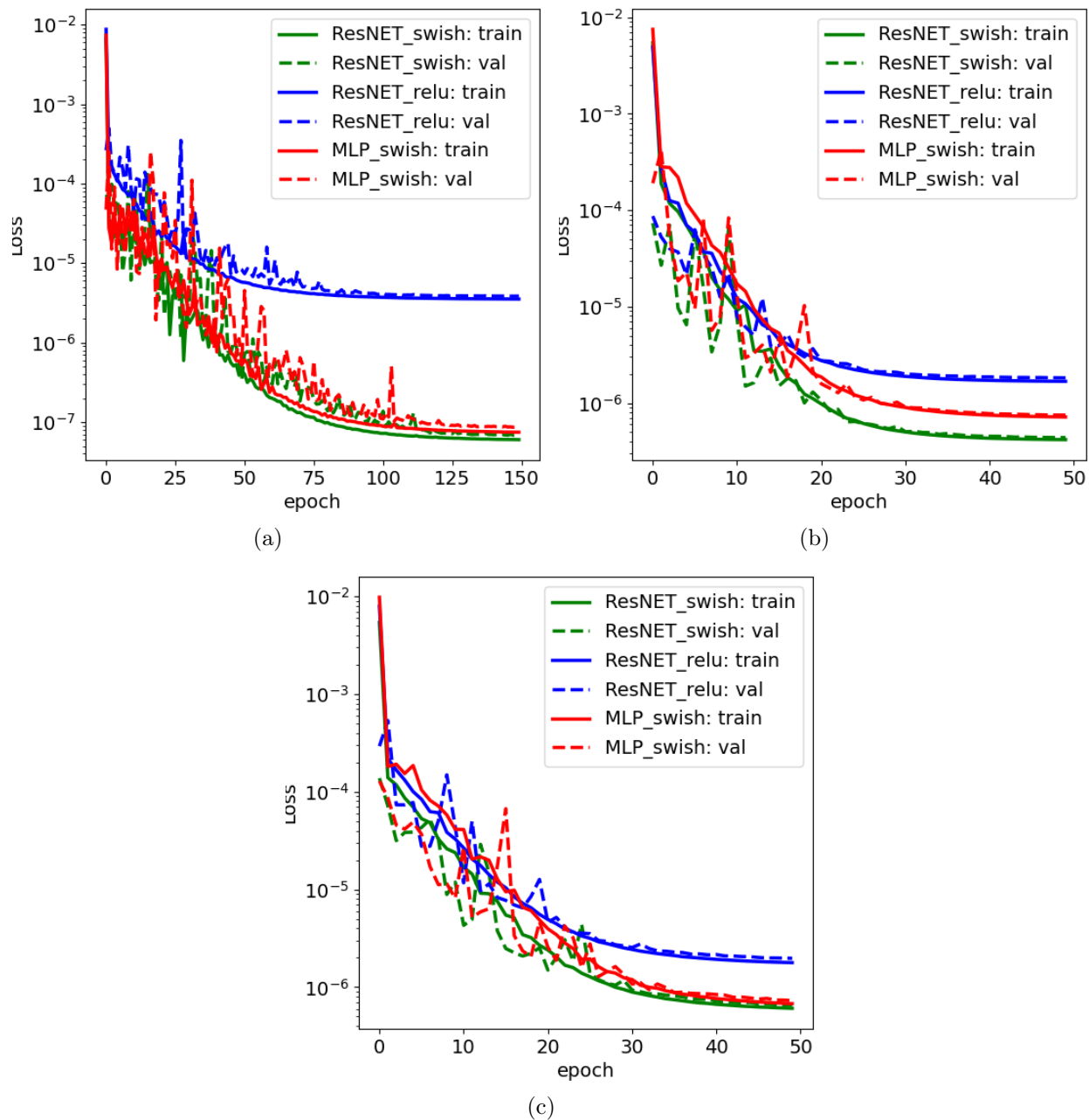


Figure 9: Training and validation loss function values(MSE) over epoch number for models of reactive subdomains for 3 cases (a)  $H_2/air$ , (b)  $C_2H_4/air$ , and (c)  $CH_4/air$ .

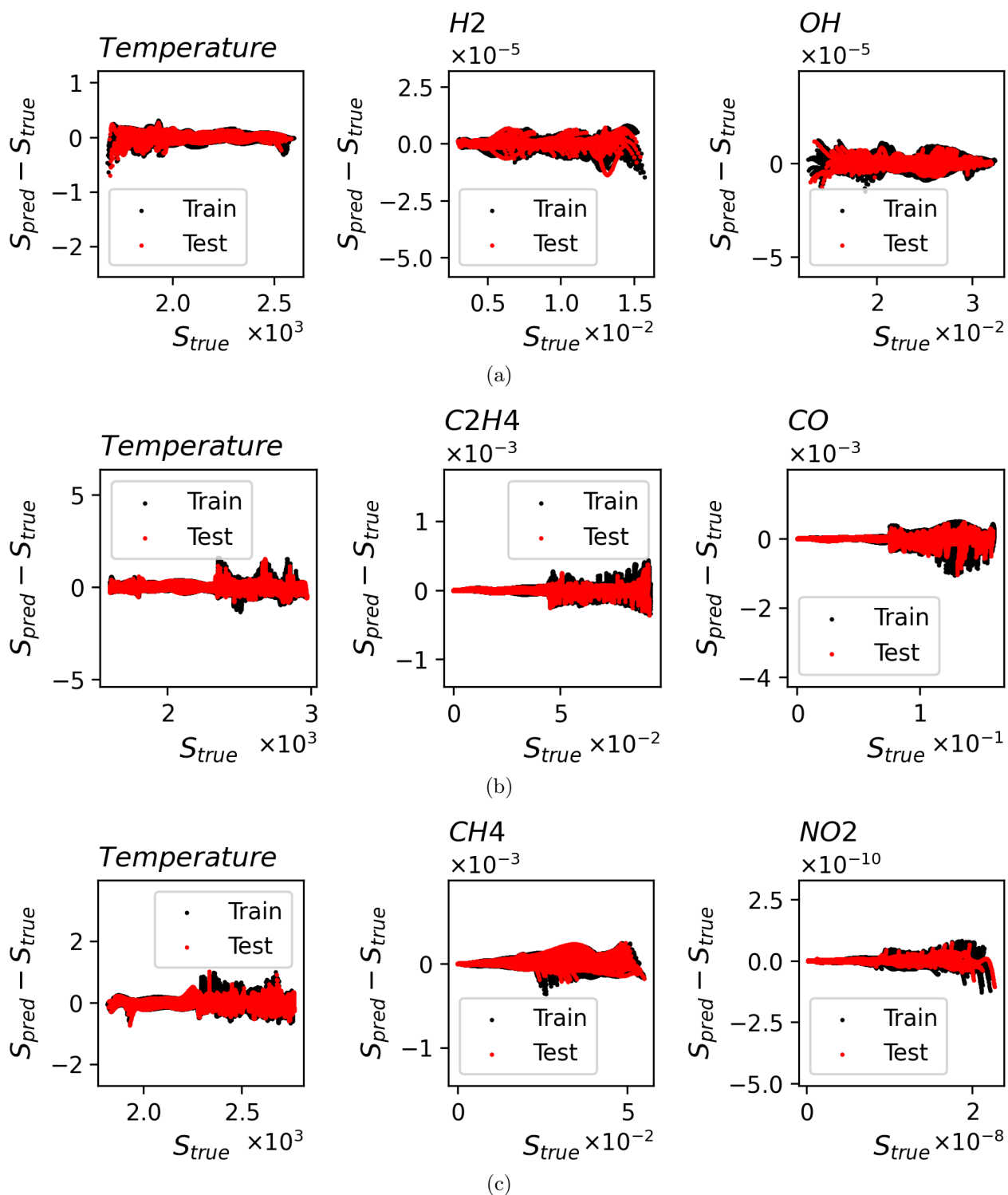


Figure 10: Parity plots of true output values and predicted errors for (a)  $H_2/air$ , (b)  $C_2H_4/air$ , and (c)  $CH_4/air$ . The chosen plot elements here are temperature, the fuel and a radical minor species, which are predicted by models of a reactive subdomain.

417 overfitting can also be drawn from Figure 10, which presents the parity plots of the true  
418 output values  $\mathbf{S}_{true}$  and the errors between the true and predicted values  $\mathbf{S}_{pred} - \mathbf{S}_{true}$ . Most  
419 of error values between the true and predicted states are much smaller than the state values  
420 under physical scales for each dimension.

421 Tables 2, 3, and 4 provide statistics of  $\mathcal{M}$ , including the average, minimum, and maximum  
422 values, for inference simulations of 100 initial conditions in the test set. The total sizes of  
423 the trained models are 330 kB, 3.2 MB, and 3.8 MB for the  $H_2$ ,  $C_2H_4$ , and  $CH_4$  cases,  
424 respectively. Box plots in Figures 11, 12, and 13 illustrate the statistical distribution of  
425 temperature and species mass fractions. The boxplots represent the distribution of the  
426 errors  $\mathcal{M}_i$  for each chemical species and of  $\mathcal{M}_0$  for temperature. The errors are computed  
427 for 100 test simulations and the y axis denotes the error values. The inner box extends  
428 from the first quartile  $Q1$  to the third quartile  $Q3$  of the distribution, while the whiskers  
429 extend from  $Q1 - 1.5 * IQR$  to  $Q3 + 1.5 * IQR$  where the interquartile range  $IQR$  equals  
430  $Q3 - Q1$ . The black dots are determined to be outliers which represent the highest and  
431 lowest values of the error distribution. From the overall statistical results, it is observed  
432 that methods with multiple models are more efficient compared to the baseline method with  
433 only one cluster. However, in some cases where the mean error  $\mathcal{M}$  is low, there are still a  
434 few extreme error values. This observation highlights the issue of trajectory divergence from  
435 the correct path during the simulation, which can be attributed to the limited robustness  
436 of the local models. The optimal number of clusters leading to the best performance, with  
437 the lowest  $\mathcal{M}$ , in the predefined workflow is found to be 6, 4, and 5 for the  $H_2$ ,  $C_2H_4$ , and  
438  $CH_4$  cases, respectively. The data is divided into different subdomains, including preheated  
439 mixing zones with extreme values skewed towards zero, stiff reactive zones, and burn-up  
440 zones. The best models for each fuel will be considered in the subsequent analysis. Figure  
441 14 demonstrates the partitioned clusters represented by different colors on the manifold of  
442 mass fractions of  $O_2$  and  $H_2O$  in logarithmic scale over the progress variable defined by  
443  $c = \frac{T - T_0}{T_{eq} - T_0}$ . It is evident that the total manifolds are piecewise separated into subdomains.

444 Each subdomain exhibits reduced complexity, enabling the construction of learning models  
 445 with simpler structures and fewer parameters for all subdomains.

Table 2: Statistics for experiments of  $H_2/air$  case

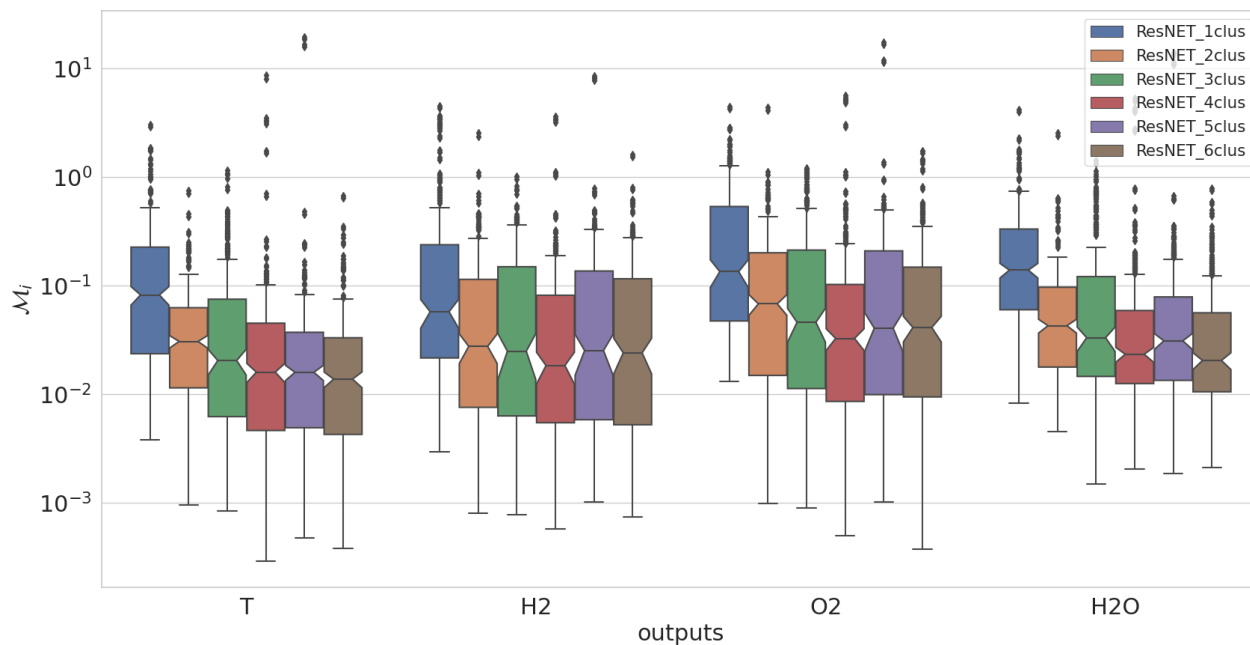
| Cluster number | $n_r$ | $n_e$ | mean $\mathcal{M}(\%)$ | minimum $\mathcal{M}(\%)$ | maximum $\mathcal{M}(\%)$ |
|----------------|-------|-------|------------------------|---------------------------|---------------------------|
| 1              | 1     | 120   | 0.301                  | 0.043                     | 2.913                     |
| 2              | 1     | 85    | 0.097                  | 0.016                     | 1.446                     |
| 3              | 1     | 70    | 0.117                  | 0.020                     | 1.056                     |
| 4              | 1     | 60    | 0.068                  | 0.010                     | 3.571                     |
| 5              | 1     | 54    | 0.081                  | 0.023                     | 11.87                     |
| 6              | 1     | 49    | 0.068                  | 0.012                     | 0.643                     |

Table 3: Statistics for experiments of  $C_2H_4/air$  case

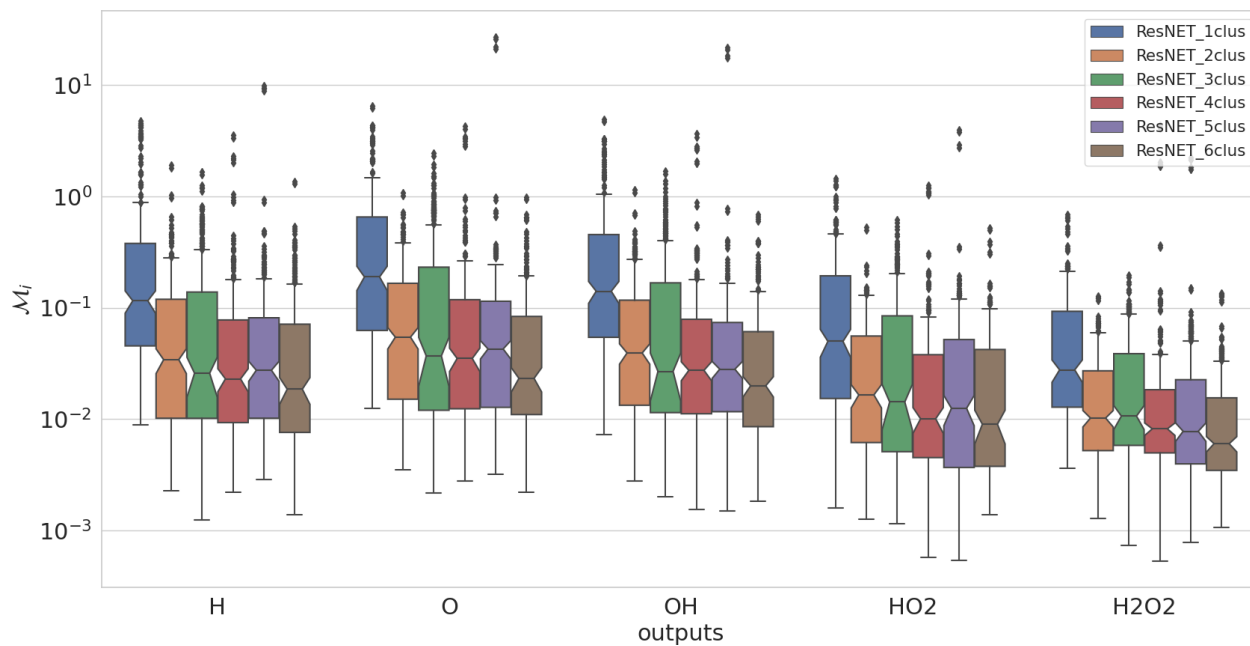
| Cluster number | $n_r$ | $n_e$ | mean $\mathcal{M}(\%)$ | minimum $\mathcal{M}(\%)$ | maximum $\mathcal{M}(\%)$ |
|----------------|-------|-------|------------------------|---------------------------|---------------------------|
| 1              | 2     | 300   | 0.971                  | 0.047                     | 45.30                     |
| 2              | 2     | 212   | 0.366                  | 0.070                     | 3.571                     |
| 3              | 2     | 174   | 0.039                  | 0.009                     | 187.94                    |
| 4              | 2     | 150   | 0.033                  | 0.010                     | 0.249                     |
| 5              | 2     | 135   | 0.040                  | 0.010                     | 0.312                     |
| 6              | 2     | 123   | 0.043                  | 0.010                     | 0.179                     |

Table 4: Statistics for experiments of  $CH_4/air$  case

| Cluster number | $n_r$ | $n_e$ | mean $\mathcal{M}(\%)$ | minimum $\mathcal{M}(\%)$ | maximum $\mathcal{M}(\%)$ |
|----------------|-------|-------|------------------------|---------------------------|---------------------------|
| 1              | 2     | 350   | 0.947                  | 0.199                     | 6.907                     |
| 2              | 2     | 248   | 0.882                  | 0.138                     | 9.391                     |
| 3              | 2     | 202   | 0.332                  | 0.036                     | 4.421                     |
| 4              | 2     | 175   | 0.153                  | 0.042                     | 2.159                     |
| 5              | 2     | 157   | 0.152                  | 0.033                     | 1.102                     |
| 6              | 2     | 143   | 0.173                  | 0.030                     | 1.583                     |



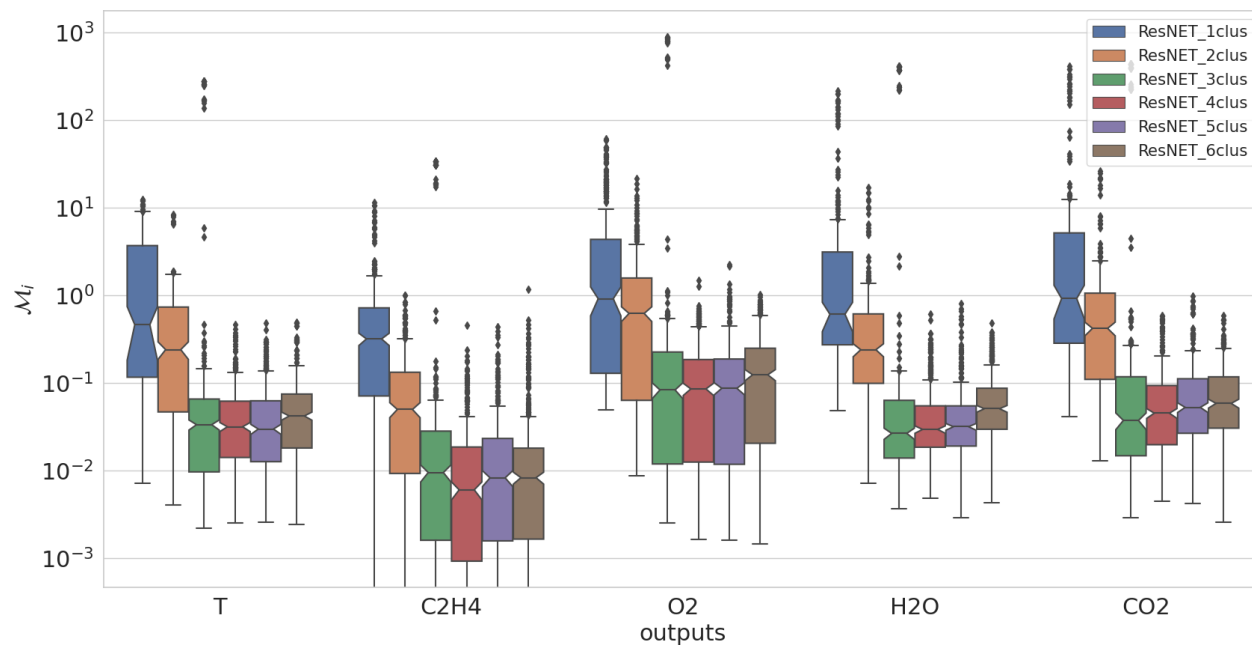
(a)



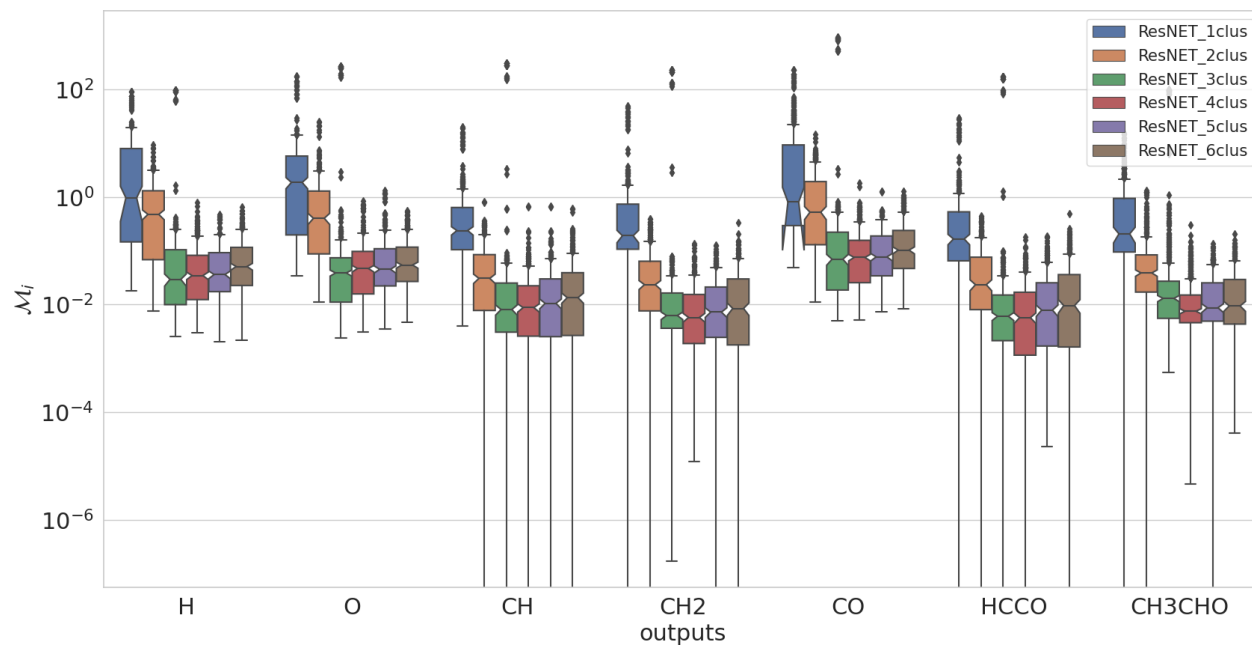
(b)

Figure 11: Box plot of statistical logMAPE errors for 100 a posteriori test simulations of  $H_2/air$  case with (a) temperature and major chemical species and (b) several minor chemical species

446 Furthermore, additional evaluations of iterative predictions are conducted by increasing  
 447 the number of simulations within the predefined range of initial conditions. In this case,



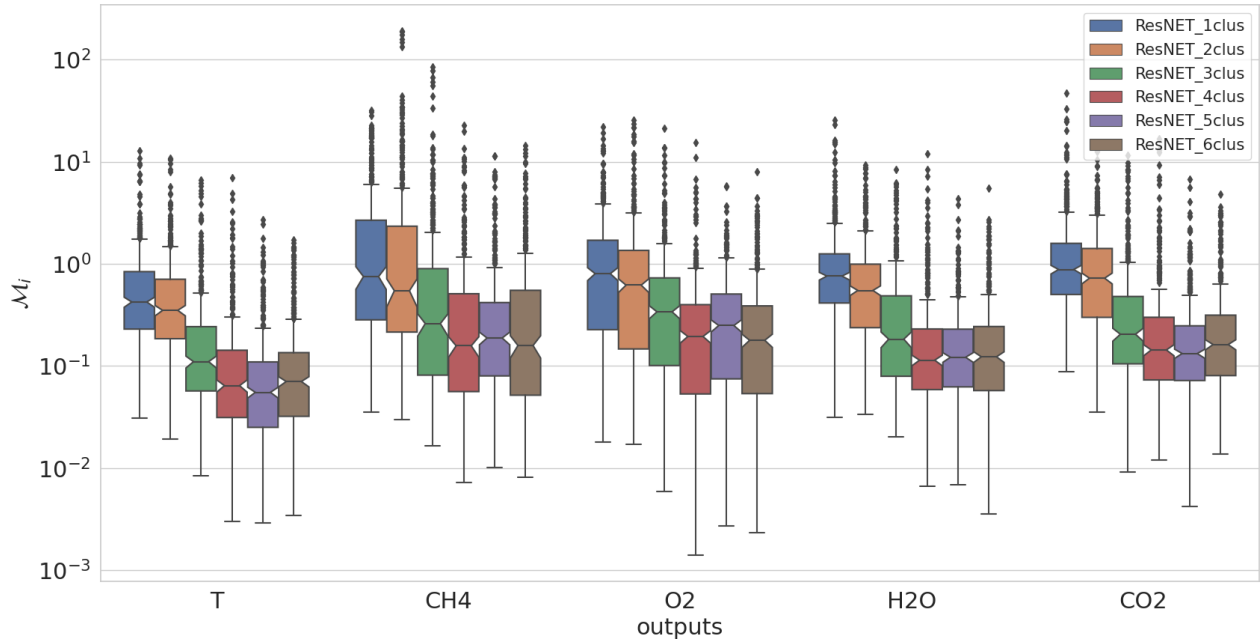
(a)



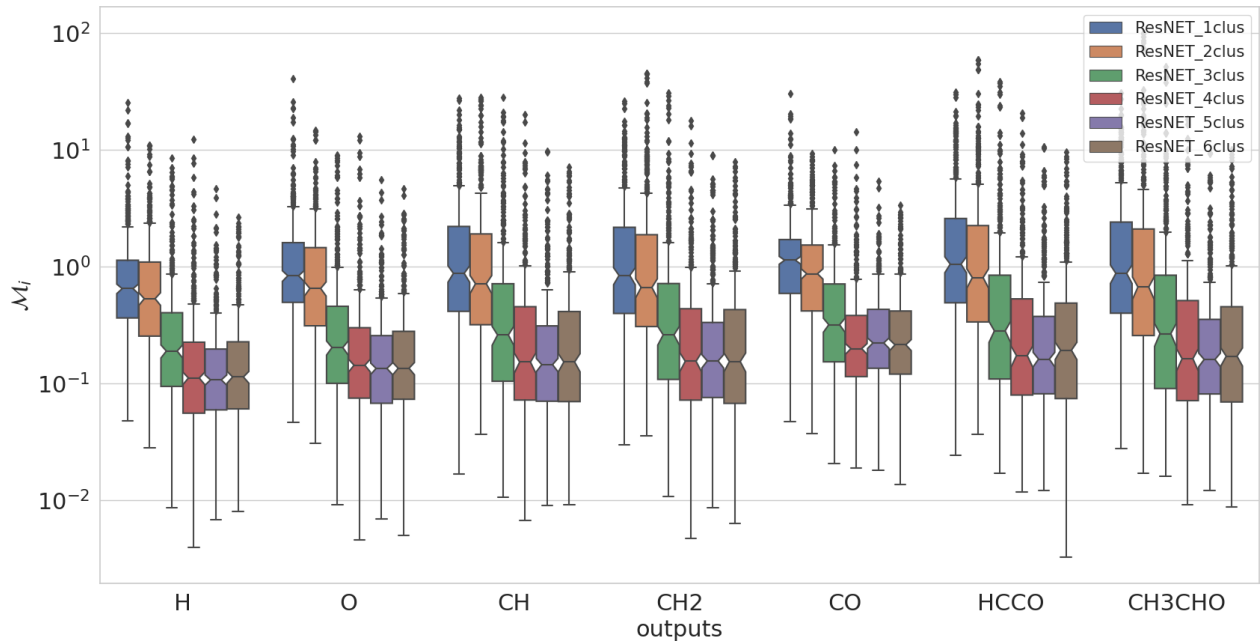
(b)

Figure 12: Box plot of statistical logMAPE errors for 100 a posteriori test simulations of  $C_2H_4/air$  case with (a) temperature and major chemical species and (b) several minor chemical species

448 1680 simulations based on model predictions are generated, and the overall accumulative  
 449 logarithmic Mean Average Percentage Error (logMAPE) (9) is computed for each simulation.



(a)



(b)

Figure 13: Box plot of statistical logMAPE errors for 100 a posteriori test simulations of  $CH_4/air$  case with (a) temperature and major chemical species and (b) several minor chemical species

450 Cubic interpolation is utilized to generalize the 2D error distribution functions. The design of  
 451 the experiment with a large number of initial conditions is fixed within the same predefined



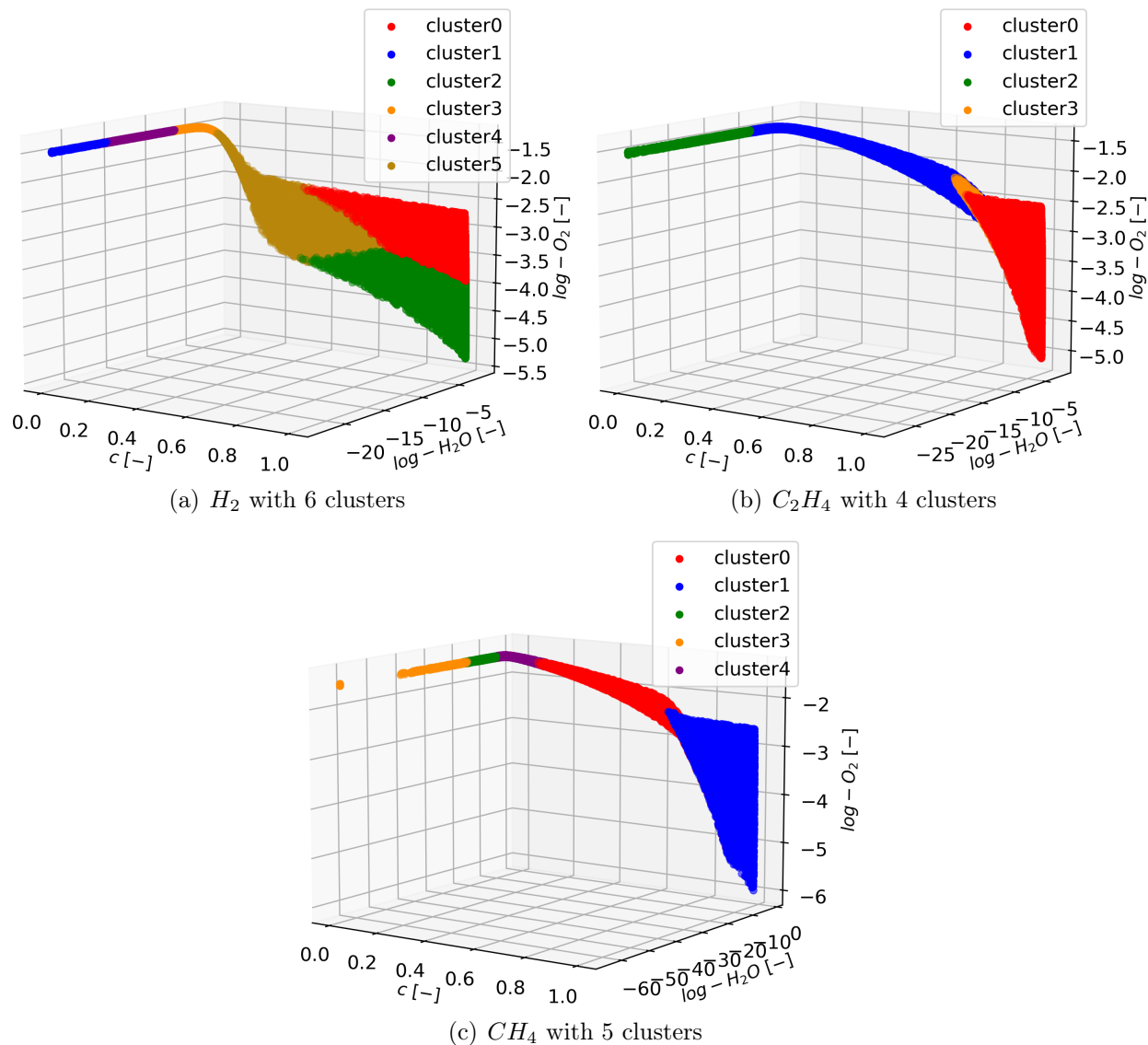


Figure 14: 3D clustering plots of manifolds for mass fractions of  $O_2$  and  $H_2O$  over the evolution of progress variable: (a)  $H_2/air$ , (b)  $C_2H_4/air$ , and (c)  $CH_4/air$ .

452 range that was used for training the models, with regular samplings of  $T_0$  and  $\phi$  at intervals  
 453 of  $\Delta T_0 = 5K$  and  $\Delta\phi = 0.01$ . The distribution of overall average MAPE errors is depicted  
 454 in Figure 15. It can be observed that simulations based on the best models exhibit overall  
 455 average MAPE errors lower than 1%, with a slight decrease in accuracy near the boundaries of  
 456 the domain. Due to its higher complexity, the  $CH_4/air$  case exhibits larger overall logMAPE  
 457 errors, but they are still within an acceptable range.

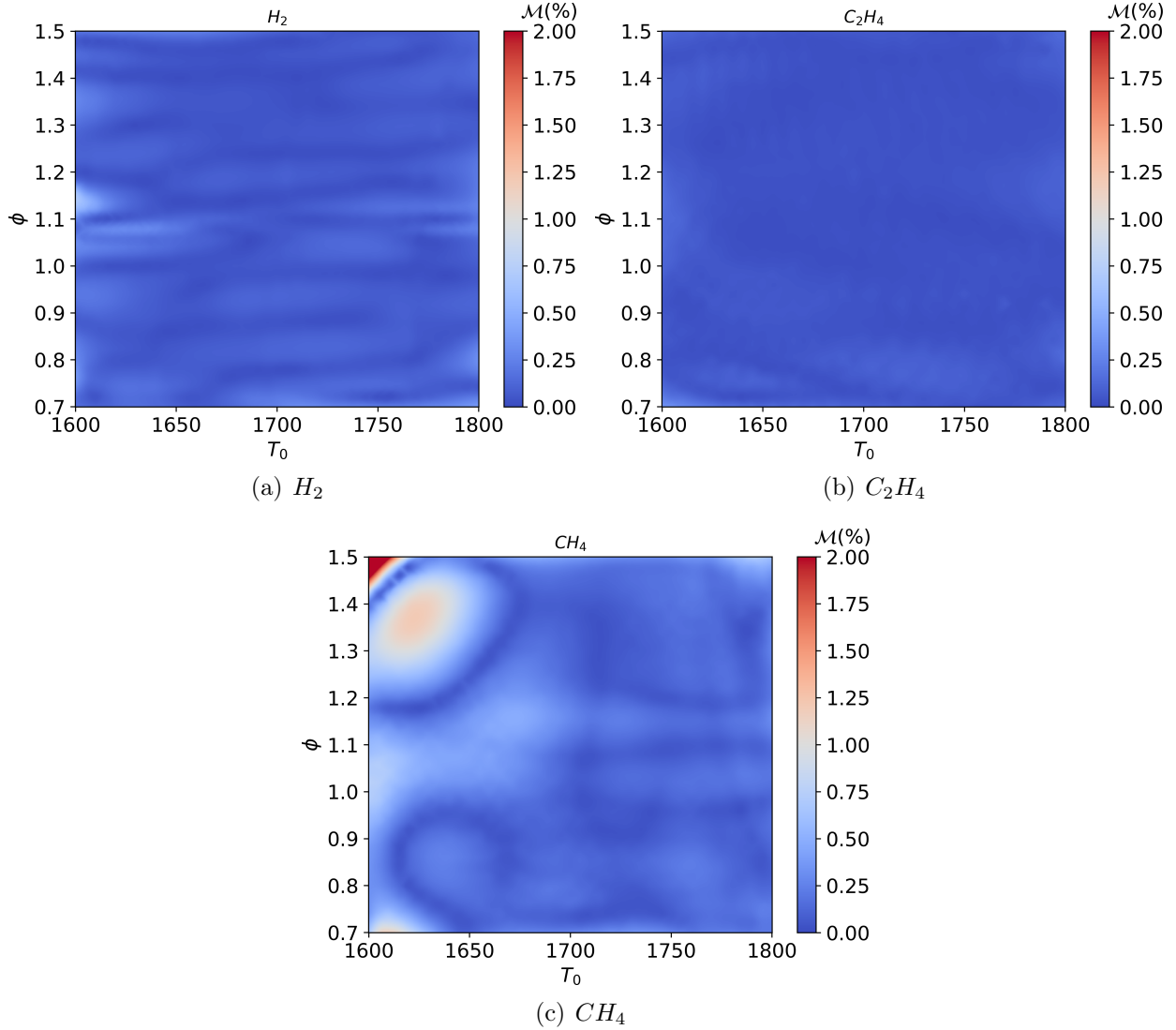


Figure 15: 2D distribution of mean relative errors of all dimensions for (a)  $H_2/air$ , (b)  $C_2H_4/air$ , and (c)  $CH_4/air$  cases within designed initial condition zone, using refined initial conditions sampling and cubic interpolation

### 458 5.3 Comparison of sampling strategies

459 The results of comparing predictions using regular sampling and the CVODE-based sampling  
 460 method are shown in Table 5. The chosen clusters are based on the optimal clusters for  
 461 each case, which is 6, 4 and 5 for  $H_2$ ,  $C_2H_4$  and  $CH_4$ . The results demonstrate that the  
 462 CVODE-based sampling method improves the overall prediction performance, as indicated  
 463 by smaller mean, minimum, and maximum errors ( $\mathcal{M}$ ) for all three cases. Regular sampling

Table 5: Statistics for experiments of  $H_2/air$ ,  $C_2H_4/air$  and  $CH_4/air$  case using regular sampling method and CVODE sampling method

|                           | $H_2$  | $C_2H_4$ | $CH_4$ |
|---------------------------|--------|----------|--------|
| <b>Regular sampling</b>   |        |          |        |
| mean $\mathcal{M}(\%)$    | 0.170  | 0.051    | 0.332  |
| maximum $\mathcal{M}(\%)$ | 1103.2 | 0.434    | 1.880  |
| minimum $\mathcal{M}(\%)$ | 0.018  | 0.019    | 0.063  |
| <b>CVODE sampling</b>     |        |          |        |
| mean $\mathcal{M}(\%)$    | 0.069  | 0.033    | 0.152  |
| maximum $\mathcal{M}(\%)$ | 0.643  | 0.249    | 1.102  |
| minimum $\mathcal{M}(\%)$ | 0.012  | 0.010    | 0.034  |

464 can lead to unbalanced clustering and uneven data distribution, resulting in unstable iterative  
 465 predictions, as observed in Table 5 for the  $H_2$  case. By employing adaptive time steps, the  
 466 CVODE-based method generates more data points from fast reaction regions, leading to a  
 467 more robust unsupervised clustering model with a balanced distribution of data points.

## 468 5.4 Simulation of 0D reactors using DNN

469 To illustrate the prediction of specific trajectories, two simulations with different initial con-  
 470 ditions near the boundaries of the dataset were performed. One simulation was conducted  
 471 with a low temperature of  $1620.0K$  and an equivalence ratio of 0.75, while the other simu-  
 472 lation had a high temperature of  $1790.0K$  and an equivalence ratio of 1.45. The results for  
 473 temperature are shown in Figure 16 for the three fuels. It can be observed that the predic-  
 474 tions generated by the models closely match the results of the direct numerical simulations.  
 475 **Additional inference results which show the predictions for a subset of the chemical species,**  
 476 **and results which display the predictions in logarithmic predictive space, with clustering**  
 477 **zones marked by different colors can be found in supplementary materials.** The DNN results  
 478 exhibit good agreement with the direct numerical simulations in both the auto-ignition zones  
 479 and equilibrium zones, compared to the exact numerical simulations. The cluster index of  
 480 states is predicted after unsupervised classification by the clustering algorithm, resulting in

481 states belonging to different subdomains at different time steps. The predictions in loga-  
 482 rithmic space are accurate within each partitioned subdomain of states, and the state values  
 483 skewed toward zero are accurately predicted at the start of the simulations.

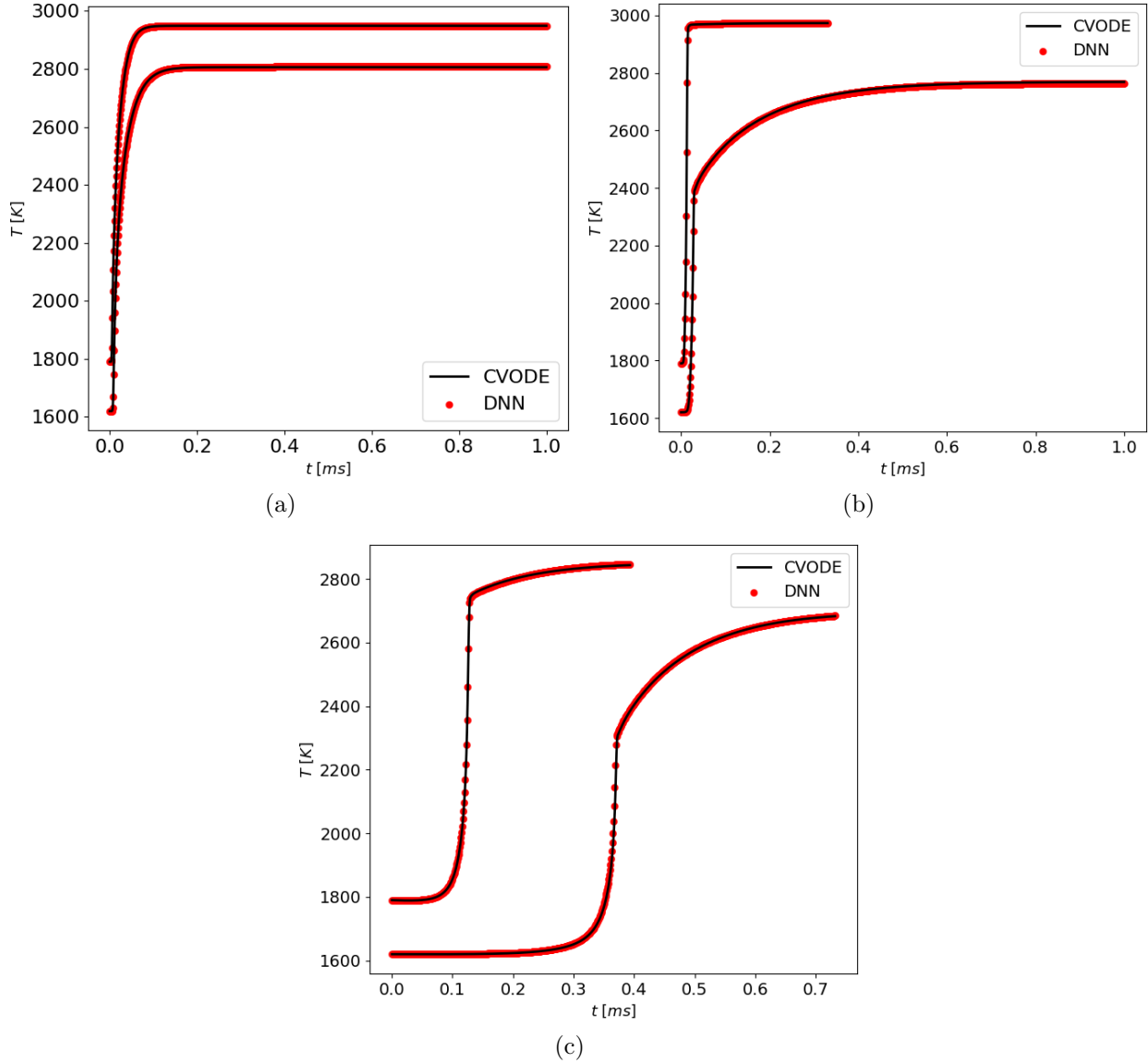


Figure 16: Continuous inference simulation of (a)  $H_2/air$ , (b)  $C_2H_4/air$ , and (c)  $CH_4/air$  cases for temperature with 2 trajectories:  $T_{01} = 1620.0K, \phi_1 = 0.75$  and  $T_{02} = 1790.0K, \phi_2 = 1.45$ .

## 484 5.5 Computing performance

485 The performance of the DNNs in terms of computational cost is assessed in this section.  
486 A single 0-D reactor is computed and the cost of the simulation is analyzed. The Python  
487 framework used in the above DNN study is not adapted to performance analysis as the Ten-  
488 sorflow DNN inference involves overhead costs. Tests on computational speed are therefore  
489 performed using the NNICE<sup>49</sup> library, which is an in-house C++ code providing lightweight  
490 DNN inference capabilities without relying on Machine Learning libraries C++ APIs. The  
491 NNICE library is easily coupled to any CFD code and is relevant for estimating gains which  
492 can be expected by replacing a direct integrator with DNNs.

493 The CVODE and DNN integrations are timed at each time step. We focus here on the  
494  $C_2H_4$  case to demonstrate the decrease in computational cost using DNNS. The selected  
495 initial condition is  $T_0 = 1700.0K$  and  $\phi = 1.0$ . We use here 4 clusters, as it has been shown  
496 to lead to the best results for  $C_2H_4$ . Several DNN architectures, summarized in Table 6, are  
497 tested to illustrate the impact of DNN size on acceleration. We verified that all networks  
498 lead to satisfactory results. Figure 17 shows the acceleration ratio  $t_{cvode}/t_{DNN}$  for each  
499 iteration during the simulation, where  $t_{cvode}$  and  $t_{DNN}$  are the time spent in CVODE and  
500 DNN integration, respectively.

501 The results clearly demonstrate that the chemistry integration performed by the DNN is  
502 5 to 30 times faster than CVODE. More specifically, DNN predictions for both stiff (from 0.0s  
503 to approx.  $4 \times 10^{-5}s$ ) and non-stiff regions exhibit the same computing efficiency, whereas  
504 the CVODE solver incurs higher costs due to the use of implicit schemes involving Jacobian  
505 matrix computations. This leads to a significant gain in computational cost for the most  
506 reactive states, as seen in Figure 17 in the first instants of the simulation, corresponding to  
507 the ignition phase. Moreover, the computation acceleration is enhanced when the size of the  
508 network decreases. This suggests that for an optimal computing efficiency, a hyper-parameter  
509 search on the network size could bring benefits.

Table 6: Models information for DNN performance testing

| DNN model | $n_r$ | $n_e$ | model size for each cluster |
|-----------|-------|-------|-----------------------------|
| $DNN_1$   | 1     | 150   | 462.7 kB                    |
| $DNN_2$   | 2     | 120   | 554.7 kB                    |
| $DNN_3$   | 2     | 150   | 830.9 kB                    |

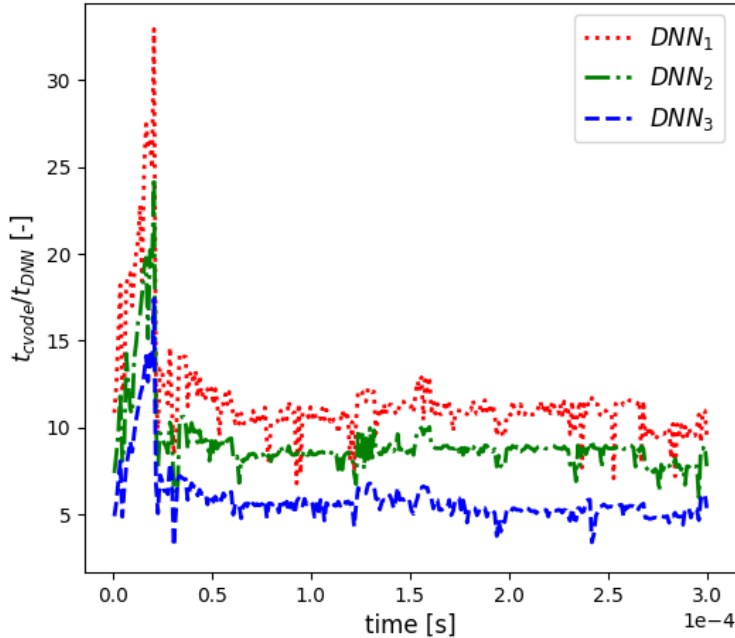


Figure 17: The acceleration ratio  $t_{cvode}/t_{DNN}$  between CVODE resolution and DNN prediction resolution for each time step, curves for three models with different sizes.

## 510 6 Conclusion and perspective

511 In this research, a deep learning surrogate model has been successfully trained and applied  
512 to the predictions of 0D combustion simulations using different fuels of increasing complex-  
513 ity. Several key strategies have been employed, including a new sampling method based  
514 on adaptive time steps of the CVODE solver to achieve a balanced data distribution, the  
515 application of a non-supervised K-Means algorithm to separate the dataset into subdomains  
516 in logarithmic space, and the use of residual networks to improve optimization during the  
517 training process. The analysis of the inference results demonstrates that by partitioning the  
518 sampling points into subdomains, the overall complexity of the model prediction process

519 can be reduced while achieving optimal results. Proper selection of hyperparameters plays  
520 a crucial role in obtaining accurate iterative predictions.

521 As an extension of this work, the proposed sampling method based on implicit numerical  
522 stiff solvers can be further applied to sample thermochemical data points in the context of  
523 2D/3D combustion processes involving convection and diffusion terms. This would allow the  
524 learning workflow to be extended to more complex problems that are coupled with convection  
525 and diffusion in computational fluid dynamics.

## 526 A K-Means algorithm

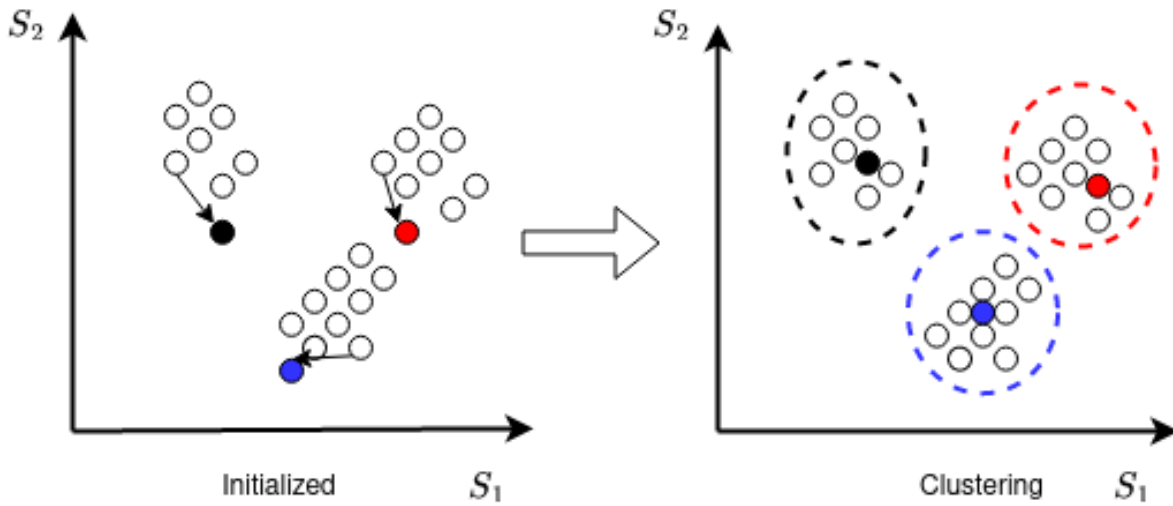


Figure 18: Clustering by K-Means algorithm

The K-Means algorithm, as shown in Figure 18, is a clustering algorithm that partitions a set of  $N$  sampling points based on their similarities, aiming to minimize the average squared distance between  $K$  centroids and the sampling points. Given a set of  $N$  observed sampling points  $\mathbf{S}_1, \mathbf{S}_2, \dots, \mathbf{S}_N$ , the algorithm iteratively minimizes the within-cluster sum of squares (WCSS) as the loss function. The WCSS is defined as the sum of the squared Euclidean

distances between each sampling point and its assigned centroid, which is denoted as:

$$\mathcal{L}_{kmeans} = \sum_{i=1}^K \sum_{\mathbf{S} \in \Omega_i} \left\| \hat{\mathbf{S}}_l - \mathbf{d}_i \right\|_2 \quad (11)$$

527 where the  $\mathbf{d}_i$  denotes the centroid of each cluster,  $\hat{\mathbf{S}}$  represents the data normalization. Notice  
528 that the clustering step is performed in the transformed coordinates after the normalized  
529 logarithmic transformation for chemical species  $Y_j$  with  $\mathbf{S}_l = [T, \ln(Y_j)]$ ,  $j = 1, \dots, N_s$ .

530 The K-Means++ algorithm<sup>50</sup> is an improved version of the standard K-Means algorithm  
531 that addresses the issue of poor clusterings that can occur with the standard approach. It  
532 introduces a more effective initialization step for selecting the initial centroids. The steps of  
533 the centroids initialization by the K-Means++ algorithm are as follows:

- 534 • Initialize the first centroid by randomly selecting one data point from the dataset using  
535 Latin Hypercube Sampling(LHS).
- 536 • For each remaining data point, compute its distance to the nearest centroid.
- 537 • Select the next centroid from the remaining data points with a probability proportional  
538 to the square of the distance to the nearest centroid. This ensures that points further  
539 away from existing centroids are more likely to be selected as new centroids.
- 540 • Repeat steps 2 and 3 until K centroids are selected.

541 By using the K-Means++ algorithm for initialization, the K-Means clustering process  
542 starts with more representative initial centroids, leading to better overall clusterings. This  
543 helps to avoid situations where the algorithm gets stuck in sub-optimal solutions or produces  
544 unbalanced clusters.



## 545 Supporting Information

546 Detailed quantitative analysis of thresholds and additional 0D simulation results for several  
547 chemical species based on optimal DNN models (PDF)

## 548 References

- 549 (1) Cohen, S. D.; Hindmarsh, A. C.; Dubois, P. F. CVODE, a stiff/nonstiff ODE solver in  
550 *C. Computers in physics* **1996**, *10*, 138–143.
- 551 (2) Peters, N.; Rogg, B. *Reduced kinetic mechanisms for applications in combustion sys-*  
552 *tems*; Springer Science & Business Media, 2008; Vol. 15.
- 553 (3) Lu, T.; Law, C. K. A directed relation graph method for mechanism reduction. *Pro-*  
554 *ceedings of the Combustion Institute* **2005**, *30*, 1333–1341.
- 555 (4) Chen, J.-Y.; Kollmann, W.; Dibble, R. Pdf modeling of turbulent nonpremixed methane  
556 jet flames. *Combustion Science and Technology* **1989**, *64*, 315–346.
- 557 (5) Pope, S. Computationally efficient implementation of combustion chemistry using in  
558 situ adaptive tabulation. *Combustion Theory and Modelling* **1997**, *1*, 41–63.
- 559 (6) Goodfellow, I.; Bengio, Y.; Courville, A. *Deep learning*; MIT press, 2016.
- 560 (7) Ihme, M.; Chung, W. T.; Mishra, A. A. Combustion machine learning: Principles,  
561 progress and prospects. *Progress in Energy and Combustion Science* **2022**, *91*, 101010.
- 562 (8) Zhou, L.; Song, Y.; Ji, W.; Wei, H. Machine learning for combustion. *Energy and AI*  
563 **2022**, *7*, 100128.
- 564 (9) Cybenko, G. Approximation by superpositions of a sigmoidal function. *Mathematics of*  
565 *control, signals and systems* **1989**, *2*, 303–314.

- 566 (10) Hornik, K.; Stinchcombe, M.; White, H. Multilayer feedforward networks are universal  
567 approximators. *Neural networks* **1989**, *2*, 359–366.
- 568 (11) Christo, F.; Masri, A.; Nebot, E. Artificial neural network implementation of chemistry  
569 with PDF simulation of H<sub>2</sub>/CO<sub>2</sub> flames. *Combustion and Flame* **1996**, *106*, 406–427.
- 570 (12) Blasco, J.; Fueyo, N.; Dopazo, C.; Ballester, J. Modelling the temporal evolution of a  
571 reduced combustion chemical system with an artificial neural network. *Combustion and  
572 Flame* **1998**, *113*, 38–52.
- 573 (13) Ihme, M.; Schmitt, C.; Pitsch, H. Optimal artificial neural networks and tabulation  
574 methods for chemistry representation in LES of a bluff-body swirl-stabilized flame.  
575 *Proceedings of the Combustion Institute* **2009**, *32*, 1527–1535.
- 576 (14) Sen, B. A.; Hawkes, E. R.; Menon, S. Large eddy simulation of extinction and reignition  
577 with artificial neural networks based chemical kinetics. *Combustion and Flame* **2010**,  
578 *157*, 566–578.
- 579 (15) Wan, K.; Barnaud, C.; Vervisch, L.; Domingo, P. Chemistry reduction using machine  
580 learning trained from non-premixed micro-mixing modeling: Application to DNS of a  
581 syngas turbulent oxy-flame with side-wall effects. *Combustion and Flame* **2020**, *220*,  
582 119–129.
- 583 (16) Ranade, R.; Alqahtani, S.; Farooq, A.; Echehki, T. An ANN based hybrid chemistry  
584 framework for complex fuels. *Fuel* **2019**, *241*, 625–636.
- 585 (17) Ding, T.; Readshaw, T.; Rigopoulos, S.; Jones, W. Machine learning tabulation of ther-  
586 mochemistry in turbulent combustion: An approach based on hybrid flamelet/random  
587 data and multiple multilayer perceptrons. *Combustion and Flame* **2021**, *231*, 111493.
- 588 (18) An, J.; He, G.; Luo, K.; Qin, F.; Liu, B. Artificial neural network based chemical mech-

- 589 anisms for computationally efficient modeling of hydrogen/carbon monoxide/kerosene  
590 combustion. *International Journal of Hydrogen Energy* **2020**, *45*, 29594–29605.
- 591 (19) Ji, W.; Qiu, W.; Shi, Z.; Pan, S.; Deng, S. Stiff-pinn: Physics-informed neural network  
592 for stiff chemical kinetics. *The Journal of Physical Chemistry A* **2021**, *125*, 8098–8106.
- 593 (20) Sutherland, J. C.; Parente, A. Combustion modeling using principal component analy-  
594 sis. *Proceedings of the Combustion Institute* **2009**, *32*, 1563–1570.
- 595 (21) Isaac, B. J.; Coussement, A.; Gicquel, O.; Smith, P. J.; Parente, A. Reduced-order PCA  
596 models for chemical reacting flows. *Combustion and flame* **2014**, *161*, 2785–2800.
- 597 (22) Malik, M. R.; Isaac, B. J.; Coussement, A.; Smith, P. J.; Parente, A. Principal com-  
598 ponent analysis coupled with nonlinear regression for chemistry reduction. *Combustion  
599 and Flame* **2018**, *187*, 30–41.
- 600 (23) Zdybał, K.; Sutherland, J. C.; Parente, A. Manifold-informed state vector subset for  
601 reduced-order modeling. *Proceedings of the Combustion Institute* **2023**, *39*, 5145–5154.
- 602 (24) Zhang, Y.; Xu, S.; Zhong, S.; Bai, X.-S.; Wang, H.; Yao, M. Large eddy simulation  
603 of spray combustion using flamelet generated manifolds combined with artificial neural  
604 networks. *Energy and AI* **2020**, *2*, 100021.
- 605 (25) Prieler, R.; Moser, M.; Eckart, S.; Krause, H.; Hochenauer, C. Machine learning tech-  
606 niques to predict the flame state, temperature and species concentrations in counter-  
607 flow diffusion flames operated with CH<sub>4</sub>/CO/H<sub>2</sub>-air mixtures. *Fuel* **2022**, *326*, 124915.
- 608 (26) Hansinger, M.; Ge, Y.; Pfitzner, M. Deep residual networks for flamelet/progress vari-  
609 able tabulation with application to a piloted flame with inhomogeneous inlet. *Combus-  
610 tion Science and Technology* **2022**, *194*, 1587–1613.
- 611 (27) Li, K.; Rahnama, P.; Novella, R.; Somers, B. Combining flamelet-generated manifold

- 612 and machine learning models in simulation of a non-premixed diffusion flame. *Energy*  
613 *and AI* **2023**, *14*, 100266.
- 614 (28) Chi, C.; Janiga, G.; Thévenin, D. On-the-fly artificial neural network for chemical  
615 kinetics in direct numerical simulations of premixed combustion. *Combustion and Flame*  
616 **2021**, *226*, 467–477.
- 617 (29) Han, X.; Jia, M.; Chang, Y.; Li, Y. An improved approach towards more robust deep  
618 learning models for chemical kinetics. *Combustion and Flame* **2022**, *238*, 111934.
- 619 (30) Blasco, J. A.; Fueyo, N.; Dopazo, C.; Chen, J. A self-organizing-map approach to  
620 chemistry representation in combustion applications. *Combustion Theory and Modelling*  
621 **2000**, *4*, 61–76.
- 622 (31) Franke, L. L.; Chatzopoulos, A. K.; Rigopoulos, S. Tabulation of combustion chemistry  
623 via Artificial Neural Networks (ANNs): Methodology and application to LES-PDF  
624 simulation of Sydney flame L. *Combustion and Flame* **2017**, *185*, 245–260.
- 625 (32) Perini, F. High-dimensional, unsupervised cell clustering for computationally efficient  
626 engine simulations with detailed combustion chemistry. *Fuel* **2013**, *106*, 344–356.
- 627 (33) Barwey, S.; Prakash, S.; Hassanaly, M.; Raman, V. Data-driven classification and mod-  
628 eling of combustion regimes in detonation waves. *Flow, Turbulence and Combustion*  
629 **2021**, *106*, 1065–1089.
- 630 (34) Nguyen, H.-T.; Domingo, P.; Vervisch, L.; Nguyen, P.-D. Machine learning for integrat-  
631 ing combustion chemistry in numerical simulations. *Energy and AI* **2021**, *5*, 100082.
- 632 (35) Zhang, T.; Yi, Y.; Xu, Y.; Chen, Z. X.; Zhang, Y.; Weinan, E.; Xu, Z.-Q. J. A multi-  
633 scale sampling method for accurate and robust deep neural network to predict combus-  
634 tion chemical kinetics. *Combustion and Flame* **2022**, *245*, 112319.

- 635 (36) Lanser, D.; Verwer, J. G. Analysis of operator splitting for advection–diffusion–reaction  
636 problems from air pollution modelling. *Journal of computational and applied mathemat-*  
637 *ics* **1999**, *111*, 201–216.
- 638 (37) MacNamara, S.; Strang, G. *Splitting methods in communication, imaging, science, and*  
639 *engineering*; Springer, 2016; pp 95–114.
- 640 (38) Goodwin, D. G. Cantera c++ user’s guide. *California Institute of Technology* **2002**,
- 641 (39) Ó Conaire, M.; Curran, H. J.; Simmie, J. M.; Pitz, W. J.; Westbrook, C. K. A com-  
642 prehensive modeling study of hydrogen oxidation. *International journal of chemical*  
643 *kinetics* **2004**, *36*, 603–622.
- 644 (40) Luo, Z.; Yoo, C. S.; Richardson, E. S.; Chen, J. H.; Law, C. K.; Lu, T. Chemical  
645 explosive mode analysis for a turbulent lifted ethylene jet flame in highly-heated coflow.  
646 *Combustion and Flame* **2012**, *159*, 265–274.
- 647 (41) Smith, G. GRI-Mech.-An Optimized Detailed Chemical Reaction Mechanism for  
648 Methane Combustion. [http://www. me. berkeley. edu/gri\\_ mech](http://www.me.berkeley.edu/gri_mech) **1999**,
- 649 (42) He, K.; Zhang, X.; Ren, S.; Sun, J. Deep residual learning for image recognition. Pro-  
650 ceedings of the IEEE conference on computer vision and pattern recognition. 2016; pp  
651 770–778.
- 652 (43) Qu, J.; Faney, T.; de Hemptinne, J.-C.; Yousef, S.; Gallinari, P. PTFlash : A vectorized  
653 and parallel deep learning framework for two-phase flash calculation. *Fuel* **2023**, *331*,  
654 125603.
- 655 (44) Ramachandran, P.; Zoph, B.; Le, Q. V. Searching for activation functions. *arXiv*  
656 *preprint arXiv:1710.05941* **2017**,
- 657 (45) Abadi, M. et al. TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems.  
658 2015; <https://www.tensorflow.org/>, Software available from tensorflow.org.

- 659 (46) Kingma, D. P.; Ba, J. Adam: A method for stochastic optimization. *arXiv preprint*  
660 *arXiv:1412.6980* **2014**,
- 661 (47) You, K.; Long, M.; Wang, J.; Jordan, M. I. How does learning rate decay help modern  
662 neural networks? *arXiv preprint arXiv:1908.01878* **2019**,
- 663 (48) Sharma, A. J.; Johnson, R. F.; Kessler, D. A.; Moses, A. Deep learning for scalable  
664 chemical kinetics. AIAA scitech 2020 forum. 2020; p 0181.
- 665 (49) Aubagnac-Karkar, D.; Mehl, C. NNICE: Neural Network Inference in C made Easy.  
666 2023; <https://doi.org/10.5281/zenodo.7645515>.
- 667 (50) Arthur, D.; Vassilvitskii, S. K-means++ the advantages of careful seeding. Proceedings  
668 of the eighteenth annual ACM-SIAM symposium on Discrete algorithms. 2007; pp  
669 1027–1035.

