

1 manuscript.tex

2 **Deep Learning Dynamical Latencies for the Analysis and Reduction of**
3 **Combustion Chemistry Kinetics**

4 Luisa Castellanos^{1,2,*}, Rodolfo S. M. Freitas², Alessandro Parente², Francesco Contino¹

5 ¹*Université Catholique de Louvain, École polytechnique de Louvain,*
6 *Institute of Mechanics, Materials and Civil Engineering,*
7 *Ottignies-Louvain-la-Neuve, Belgium and*

8 ²*Université Libre de Bruxelles, École polytechnique de Bruxelles,*
9 *Aero-Thermo-Mechanics Laboratory, Brussels, Belgium*

10 (*luisa.castellanos@uclouvain.be)

11 (Dated: September 8, 2023)

12 **Abstract**

13 The modeling of chemical kinetics holds many challenges, as well as a necessity for more efficient mod-
14 eling techniques, together with dimensionality reduction techniques. This work studies the application of
15 time-lag auto-encoders (TAEs) for the analysis of combustion chemistry kinetics. Such a technique allows a
16 better reconstruction of the thermochemical temporal advancement in relation to traditional reduction tech-
17 niques (Principal Component Analysis or PCA) while applying a potential denoising operation. Moreover,
18 the reduced manifolds or latencies are provided with physical meaning, which further analysis, gives insight
19 into key chemical reactions and interactions between chemical species, allowing for a deeper understanding
20 of the chemical mechanism itself.

This is the author's peer reviewed, accepted manuscript. However, the online version of record will be different from this version once it has been copyedited and typeset.

PLEASE CITE THIS ARTICLE AS DOI: 10.1063/5.0167110

Accepted to *Phys. Fluids* 10.1063/5.0167110

21 **NOMENCLATURE**

22 **Abbreviations**

23 AE Autoencoder

24 CFD Computational Fluid Dynamics

25 CombML Combustion Machine Learning

26 D Decoder

27 DMD Dynamic Mode Decomposition

28 E Encoder

29 HODMD High-Order Dynamic Mode Decomposition

30 ISAT In Situ Adaptive Tabulation

31 LHS Latin Hypercube Sampling

32 LPCA Local Principal Component Analysis

33 MAE Mean Absolute Error

34 NNs Neural Networks

35 PCA Principal Component Analysis

36 SPARC Sample-Partitioning Adaptive Reduced Chemistry

37 TAE Time-lag Autoencoder

38 TDAC Tabulation of Dynamic Adaptive Chemistry

39 **Greek letters**

40 β Regularization constant

41 β_1 Adam's optimizer first constant

42 β_2 Adam's optimizer second constant

This is the author's peer reviewed, accepted manuscript. However, the online version of record will be different from this version once it has been copyedited and typeset.

PLEASE CITE THIS ARTICLE AS DOI: 10.1063/5.0167110

Accepted to Phys. Fluids 10.1063/5.0167110

43	ϕ	Equivalence ratio
44	ψ	Encoder Neural Network vectors
45	τ	Time interval
46	ξ	Decoder Neural Network vectors
47	Latin letters	
48	\mathbb{R}	Dimensionality
49	\mathcal{L}_{TAE}	Time-lagged reconstruction loss
50	i_j	Coefficient of importance for the j-th variable
51	K	Number of time series
52	M	Number of time samples
53	N	Number of variables of high-dimensional space
54	n	Number of variables of low-dimensional space
55	Q	Number of permutation repetitions
56	R^2	Coefficient of determination
57	s	Linear regression coefficient of determination
58	$s_{q,j}$	Linear regression coefficient of determination for the j-th variable in q permutation repetition
59		
60	T	Temperature
61	T_0	Initial temperature
62	t_0	First time step
63	t_{nt}	Last time step
64	X	Thermochemical state vector

- 65 X_0 Initial thermochemical state
- 66 Y_i Mass fraction of i -th species
- 67 Z Reduced manifold size

68 I. INTRODUCTION

69 Reacting flows are a common topic when it comes to industrial applications; this also means
 70 that efficient techniques that foresee their behavior are a necessity. Simulation techniques are a
 71 common technique used for the chemistry reactions' approximations, more precisely, CFD tech-
 72 niques have become part of common knowledge applications. Unfortunately, many of the current
 73 simulation techniques are unfeasible for their application in real industrial cases, mainly due to
 74 the computational overhead that they suppose. In brief, the number of equations increases linearly
 75 in the measure that more chemical species are added to the model, which becomes necessary to
 76 accurately describe the chemical evolution of the reacting system [1], at the same time, the com-
 77 putational cost for the solution increases exponentially in accordance to the number of chemical
 78 species. Additionally, there are many other physical phenomena to be considered, such as tur-
 79 bulence models, thermal irradiation models, acoustics, etc., from which, each one considers the
 80 addition of phenomenological closure models. For high-fidelity detailed chemistry mechanisms,
 81 indeed, the chemical reaction mechanism is frequently the most computationally demanding of
 82 a computational combustion model, as a reference, above 50 species, it takes about 90% of the
 83 computational effort [2].

84 As an answer to this issue, many efforts have been invested in the development of reduced
 85 efficient chemistry representation. Among these approaches, a very popular technique is the In
 86 Situ Adaptative Tabulation (ISAT) [3], which consists in storing thermochemical states through
 87 binary trees, however, the principal drawback of this technique is its high memory requirements,
 88 which makes the technique unviable for cases considering an elevated number of species, which
 89 are most of the real-cases scenarios [4].

90 Going further, many other techniques for reduced representations have been applied, such as
 91 the application of limited chemical steps [5] in which many reactions that lead to minor species are
 92 neglected; additionally, we can mention the development of skeleton mechanisms [6] that neglect
 93 many of the intermediate species while considering just major species and final products; it is to

This is the author's peer reviewed, accepted manuscript. However, the online version of record will be different from this version once it has been copyedited and typeset.

PLEASE CITE THIS ARTICLE AS DOI: 10.1063/5.0167110

Accepted to Phys. Fluids 10.1063/5.0167110

94 mention the Tabulation of Dynamic Adaptive Chemistry (TDAC) methodology, in which there is
 95 a coupling of ISAT tabulation of previous results, plus a generation of skeleton mechanisms on
 96 the fly [7], however, the calculations of the reduced mechanisms can be demanding. Many other
 97 approaches can be mentioned, each one of them presenting its own advantages and disadvantages,
 98 however, none of the developed techniques represents a general solution to all the possible cases.
 99 With the current development of numerous Machine Learning techniques, it is possible to find
 100 a variety of works in which they are applied to combustion problems, giving place to the area of
 101 Combustion Machine Learning (CombML [8]). From these applications, it is important to mention
 102 the application of Principal Component Analysis (PCA) [9], a feature reduction technique, and the
 103 application of Neural Networks (NNs) as a substitute for the ODE time integrator [10], [11].

104 When it comes to PCA applications, its use as a feature reduction technique is limited since the
 105 resulting components are usually a linear combination of chemical species, which supposes some
 106 difficulties for the treatment of source terms ([12]). This technique has found a more efficient
 107 use as the application of Local PCA (LPCA) for clustering, as it is done in SPARC (Sample-
 108 Partitioning Adaptive Reduced Chemistry) chemistry models [13]. Another popular technique for
 109 feature reduction is the Dynamic Mode Decomposition (DMD), which has been proven useful for
 110 the identification of principal modes or frequency patterns in flows through matrix decompositions
 111 [14]. For the case of complex dynamics, High-Order Dynamic Mode Decomposition (HODMD)
 112 has proven effective [15] and introduces a windowing concept to apply DMD locally [16, 17].

113 On the other hand, NNs can be used as universal approximators, which means they are suitable
 114 for the representation of chemistry non-linearities, some works that present this kind of applica-
 115 tions are such as [18, 19]; additionally, it is to mention hybrid applications, in which techniques
 116 such as PCA or HODMD decompositions are used for a dimensionality reduction of datasets, the
 117 same data that is used afterward for the training of neural networks (NNs), which integrate the
 118 variables in time [20]; the application of these methodologies reduces significantly the computa-
 119 tional cost of the hyperparameters optimization. However, as it is observed in [21], a drawback
 120 of their implementation is the huge dimensions of the training datasets, which might require the
 121 application of generative models. Moreover, the application of neural networks and CombML
 122 techniques in general, address the problem mostly from a statistical perspective in most cases and
 123 give back very little understanding of principal species that play a role in the mechanism, as well
 124 as the principal reactions involved.

125 Nonetheless, this neural network function approximation can be used for dimensionality reduc-

This is the author's peer reviewed, accepted manuscript. However, the online version of record will be different from this version once it has been copyedited and typeset.

PLEASE CITE THIS ARTICLE AS DOI: 10.1063/5.0167110

Accepted to Phys. Fluids 10.1063/5.0167110

126 tion, as it is pictured in autoencoder (AE) applications, such as [22] and [23]. As expected, even
 127 if an AE provides with a better representation of chemistry non-linearities, the resulting reduced
 128 representations are difficult to interpret, leading to model closure issues. In the present work, a
 129 time shift will be added to an AE network architecture, which gives place to a time-lag autoen-
 130 coder (TAE) [24]. The application of such a network allows to assess an AE topology for the task
 131 of time integration itself, as well as it adds a temporal characterization to the reduced components,
 132 which becomes a dynamical characterization of the chemistry. Additionally, this kind of network
 133 allows for the development of reduced dynamic models [25], which are expressed in terms of the
 134 network's reduced components. Furthermore, it should be assessed if the reduced components
 135 enable the recognition of principal chemical features, yielding better insights about the chemistry
 136 mechanism behavior from a dynamic perspective.

137 In the next sections, a conceptualization of the TAE approach is provided with its applicability
 138 as a feature extractor of chemistry kinetics (Sec. II A). Additionally, some notions regarding PCA
 139 as a reduction technique are given, principally when used in the form of neural networks (subsec-
 140 tion II B). A proof-of-concept study for hydrogen-air homogeneous autoignition combustion and
 141 a discussion of the capability of TAE's for chemistry estimation, such an experiment intends to
 142 provide some general guidelines that will help for the technique implementation in other reactive
 143 scenarios. The paper ends with a summary of our main findings and future perspectives to leverage
 144 the application of the proposed methodology in combustion science.

145 II. THEORETICAL BACKGROUND

146 A. Time-lagged Autoencoders

147 Auto-encoders (AE) are a type of neural network that is used for finding reduced representa-
 148 tions of a given input vector [22]. In its original form, the network hyperparameters search to
 149 minimize the error of a loss function while reconstructing the same input parameters as the net-
 150 work's outputs, through the application of two main operations, which are named the encoder (E)
 151 and decoder (D) networks. These networks are separated by a bottleneck layer, from which it is
 152 possible to obtain this reduced representation. Generally speaking, the encoder network is the one
 153 that takes the input parameters and transforms them into the reduced representation, while the de-
 154 coder network takes this reduced or latent representation, and reconstructs the original parameters.

155 As suggested by the technique's name, a TAE follows the same concept and architecture as an
 156 AE, with the difference that a time shift is applied between the network's inputs and outputs. In
 157 particular, the inputs and outputs to be applied are the respective thermochemical states $X \in \mathbb{R}^N$
 158 of N state variables (i.e. species mass fractions and temperature), in the time interval $\tau \in t \subset \mathbb{R}^+$.
 159 More specifically, it is assumed that the time interval τ is discretized by M time steps with a
 160 time-step size Δt . Thus, the modeling of a dynamical system is posed as a time-series problem
 161 propagating from the initial state X_0 .

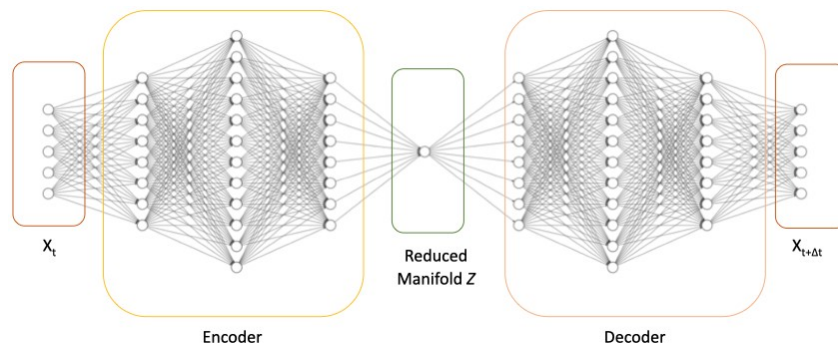


FIG. 1: Schematic that describes the composition of a TAE network, such network architecture is to be applied for a dynamical analysis of chemistry kinetics.

162 Mathematically, the encoder network can be understood as a sequence of nonlinear transfor-
 163 mations that brings the vector X from dimensionality $\mathbb{R}^N \rightarrow \mathbb{R}^n$,

$$164 \quad Z_t = E(X_t). \quad (1)$$

165 Here, Z represents the low-dimensional latent space $Z \in \mathbb{R}^n$ with $n \ll N$; for the case of the de-
 166 coder, the transformation is such that enables the reconstruction of $\mathbb{R}^n \rightarrow \mathbb{R}^N$, but with a temporal
 167 advancement, therefore, it is possible to write:

$$168 \quad \tilde{X}_{t+\Delta t} = D(Z_t) = D(E(X_t)). \quad (2)$$

169 In the present study, fully connected neural networks are employed to define the encoding and
 170 decoding models, as shown in Fig. 1. Thus, the final result is a deep learning model $\tilde{X}_{t+\Delta t} =$
 171 $D_{\xi}(E_{\psi}(X_t))$ parameterized by vectors ξ and ψ , like trained deep neural networks. Thus, the final
 172 goal is to find an encoding and decoding which minimizes the time-lagged reconstruction loss
 173 \mathcal{L}_{TAE} ,

$$174 \quad \mathcal{L}_{TAE} = \sum_0^{M-1} \|X_{t+\Delta t} - \tilde{X}_{t+\Delta t}\|^2, \quad (3)$$

175 for a family of functional mappings E and D . Additionally, neural networks are trained through the
 176 minimization of a loss function, which ideally, should reflect physical laws constraints. Following
 177 this aim, we proposed a physics-aware loss function to ensure mass conservation at each time-step:

$$178 \quad \mathcal{L}_{TAE} = \frac{1}{M-1} \sum_0^{M-1} \|X_{t+\Delta t} - \tilde{X}_{t+\Delta t}\|^2 + \frac{\beta}{M-1} \sum_0^{M-1} \|\sum Y_{t+\Delta t} - \sum \tilde{Y}_{t+\Delta t}\|^2. \quad (4)$$

179 Variables with $(\tilde{\cdot})$ refer to the TAE's output, and Y stands for the species mass fractions, which
 180 should sum up to one. β is the regularization constant, that is to say, how much the violation of
 181 mass conservation will be penalized; lastly, M stands for the number of time samples.

182 B. Comparison with other typical reduction methods

183 Considering that TAE can be used as a feature extraction technique, it is important to understand
 184 the benefits that a TAE implementation holds in relation to traditional decompose techniques. Even
 185 though it is true that PCA gained so much popularity due to its low computational cost and the
 186 fact that almost any data matrix can be decomposed by the technique [26], [27]. Moreover, PCA
 187 is also a linear technique, which means that it assumes the existence of a linear map that can
 188 represent the dataset, which is not necessarily true in many real-world cases [28]. Regardless of
 189 this intrinsic characteristic, even if the method can be applied successfully for the feature extraction
 190 of the thermochemical state [29], the resulting variables might not be of good significance for the
 191 application of predictive methods for the chemistry temporal advancement, which can be rather
 192 complex.

193 Since autoencoders hold a non-linear character, therefore, the question to answer is how much
 194 the reconstruction cost is altered by the application of a TAE while keeping a constant number of
 195 features for both methods. Additionally, it should be emphasized that TAE gives direct information

196 about the thermochemical state evolution, while PCA techniques just provide information about a
 197 present state. Because it is intended to analyze both techniques under equality of conditions, a one-
 198 layer autoencoder is proposed to mimic the PCA technique, called a PCA-like autoencoder. Since
 199 PCA is a linear technique, the PCA-like autoencoder is constructed by applying a linear activation
 200 function in the hidden layer. Thus, the feature variables are obtained from a linear transformation
 201 of the original thermochemical state. A diagram of this type of autoencoder is visualized in Figure
 202 2.

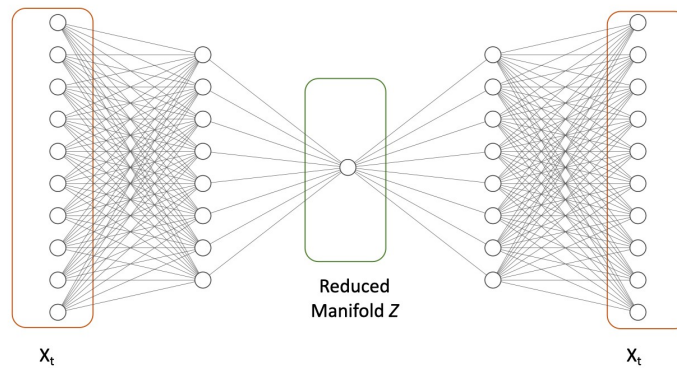


FIG. 2: Schematic of the parts of a PCA-like autoencoder, it is to notice the difference in complexity regarding the displayer TAE architecture.

203 III. COMPUTATIONAL EXPERIMENT.

204 In the present work, we test the capability of TAE as a feature extraction technique, while it
 205 performs the time integration of the thermochemical variables, this can be traduced into a time
 206 integration that eliminates unnecessary inputs; in this case, both capabilities will be applied to
 207 hydrogen combustion. The dataset for training the machine learning models is obtained with
 208 isobaric batch reactors simulations, developed with the Cantera software [30]. The simulations
 209 were performed using a constant time step, which should be small enough to ensure a sufficient
 210 number of time samples for different ignition conditions, additionally, it should provide a good
 211 description of the chemical species' mass fractions time evolution, this can be ensured if the time
 212 step is in agreement with the inverse reaction rate chemical time scale [31], which is in order of

213 10^{-4} , therefore, in order to capture all the chemical time scales with detail, the applied time step
 214 is $10^{-7}s$ for a time length of $0.5ms$. The reduced version of the University of San Diego chemistry
 215 mechanism for hydrogen combustion was used [32], which consists of 9 species and 21 reactions.
 216 Since non-reacting species are neglected, the resulting thermochemical state holds dimensionality
 217 $X \in \mathbb{R}^9$, which considers 8 reacting chemical species, plus temperature. The application of a
 218 small mechanism helps to understand the principles behind this application without any loss of
 219 generality.

220 Furthermore, correct modeling of chemical kinetics is characterized by multi-physics and
 221 multi-scale phenomena, in which the properties might vary by several ranges of magnitudes. As
 222 the main goal here is to build predictive models capable of producing reliable predictions for dif-
 223 ferent conditions, the reconstruction model must be able to bear a wide range of scenarios defined
 224 by state variables, also called control variables, such as temperature, pressure, and equivalence
 225 ratio in various flame configurations. So, the training dataset is composed of K time-series of an
 226 isobaric zero-dimensional reactor using different ignition conditions.

227 For the different ignition conditions, different values of equivalence ratio (ϕ) and initial temper-
 228 ature are considered. A Latin Hypercube Sampling (LHS) was implemented [33] with the Python
 229 PyDOE library [34]. A uniform distribution is used to obtain 100 samples in a region delimited
 230 by $\phi \in [0.9, 1.2]$ and $T \in [1100, 1220]$. The number of samples was chosen through a try-error
 231 approach in which the number of ignition cases was progressively increased from 50 to 75, and
 232 finally to 100, where a satisfactory accuracy was achieved. Here, it is worth highlighting that for
 233 cases with higher dimensionality such as methane combustion the number of ignition cases needed
 234 to train the TAE model might be larger. In such cases, generative models can be used to build large
 235 training datasets [21]. The resulting sampling is visible in Figure 3.

236 The thermochemical state vectors are defined using mass fractions values and temperature val-
 237 ues, the dataset normalization is done with a general maximum value per quantity, meaning that
 238 each quantity is divided by its maximum occurrence among all the datasets. This allows all the
 239 datasets to lie in the same manifold. The applied time shift equals a single time step.

240 At the same time, different sizes of bottleneck layers were explored; the latent spaces range is
 241 described by $Z \in [1, 4]$, therefore, four different networks are meant to be studied. The networks
 242 are built using the Python Tensorflow package[35], each of the networks' hyperparameters has
 243 been tuned using the optimization functions available in the package Keras Tuner [36], and this
 244 optimization process provides a set of different networks that differ between them just in the num-

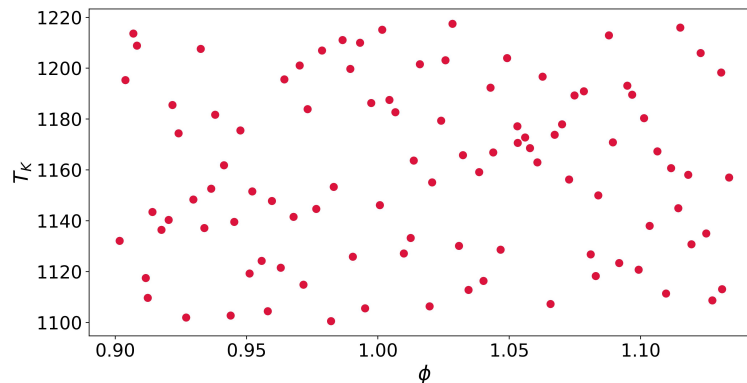


FIG. 3: LHS sampling of ignition conditions considered for training data

245 ber of neurons per layer, which can be observed in Table I, these values are specifically obtained
 246 making use of a Random Search algorithm, while the rest parameters are obtained making use of
 247 grid search.

248 The Leaky ReLU functions used, hold an activation constant of 1×10^{-3} , which allows a better
 249 representation of combustion non-linearities. Moreover, a sigmoid activation function is used in
 250 the latent layer to constrain the latent space in $Z \in [0, 1]$ [24]. The neural network parameters are
 251 initialized using a Glorot normal [37]. The physics-aware loss function available in Equation (4)
 252 is employed to train the TAE model. The Adam optimizer [38] is employed with a learning rate
 253 of 1×10^{-3} , $\beta_1 = 0.9$, and $\beta_2 = 0.999$, for a total of 2000 epochs. Early stopping is used with
 254 a minimum change of 1×10^{-4} for five epochs. Due to the computational overhead that implies
 255 using the whole dataset, the optimization process involves the utilization of 75% of the whole
 256 dataset, while posterior cross-validation is implemented using the whole dataset. In each phase,
 257 80% of the data is used as a training dataset while the remaining 20% is used as validation data.

258 Moreover, the resulting latent space will be used for the determination of latent chemical
 259 species; with latent species, it is intended chemical species that hold a strong correlation with
 260 the latent variables. The aim is to assess whether these chemical carriers may bear the potential
 261 to act as input features for the development of physically explainable reduced-order models. The
 262 identification is carried out through a correlation analysis using Kendall's Tau B correlation in-
 263 dex, due to its boundedness and resistance to outliers [39]. This feature identification is achieved

TABLE I: TAE Network's Architectures

Operation	Number Layer	Number of Neurons				Activation Function
		Z=1	Z=2	Z=3	Z=4	
Input Layer	1	9	9	9	9	Linear
Encoder	2	161	245	318	210	Leaky ReLu
	3	339	319	227	186	Leaky ReLu
	4	46	79	96	76	Leaky ReLu
	5	322	336	285	324	Leaky ReLu
Code Layer	6	1	2	3	4	Sigmoid
Decoder	7	322	336	285	324	Leaky ReLu
	8	46	79	96	76	Leaky ReLu
	9	339	319	227	186	Leaky ReLu
	10	161	245	318	210	Leaky ReLu
Output Layer	11	9	9	9	9	Linear

264 through the usage of the Scipy package [40].

265 IV. RESULTS.

266 As a first step, it is important to assess the quality of reconstruction that is achieved using the
 267 TAE networks in training data. For this, the training ignition cases will be reconstructed, and
 268 the quality of these reconstructions will be assessed with the analysis of the R^2 scores [41], which
 269 evaluates the quality of the model's output, scoring it within a constrained range of $[0, 1]$, where 0 is
 270 given to a poorly performing model, and 1 to a perfect prediction. A general R^2 is to be calculated
 271 for each ignition case, this means calculating the score for each one of the components of the
 272 thermochemical variables and averaging these results. This process provides a mean R^2 score
 273 of 0.991, with a standard deviation of 4.6×10^{-3} , such results are an indicator of a satisfactory
 274 reconstruction. The histogram with the distribution of the respective R^2 scores is available in
 275 Figure 4, this for a TAE model $z = 1$, in which a lower accuracy is to be expected.

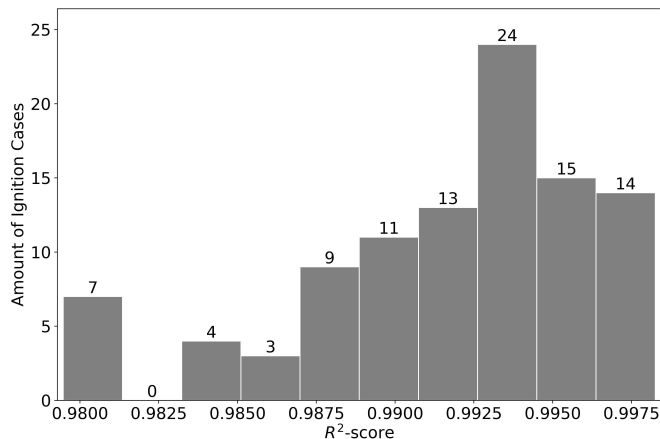


FIG. 4: Distribution of R² scores for training ignition cases tested in TAE model $z = 1$, the y-axis describes the count of ignition cases for each range of values, while the x-axis presents the ranges of R² scores.

276 Since the studied models give a good reconstruction in training data, the outcomes of this
 277 study will be organized as follows: first, an analysis of the TAE-reduced manifolds will be pre-
 278 sented. For this test, ignition conditions that can be seen as interpolation cases are considered;
 279 interpolation makes reference to conditions contained inside the limits of training data. Lastly, a
 280 comparison of TAE as a reduction technique will be performed regarding PCA capabilities for a
 281 range of latent space dimensions or a number of principal components, adding in this case, some
 282 extrapolation conditions, that is to say, conditions outside of the training limits. The interpolation
 283 and extrapolation conditions considered are available in Tables II and III respectively.

TABLE II: Interpolation Ignition cases

Case	$T[K]$	ϕ
1	1160	0.93
2	1200	1.0
3	1130	1.10

TABLE III: Extrapolation Ignition cases

Case	$T[K]$	ϕ
1	1160	0.85
2	1300	1.0

284 A. Time-lag Auto-encoder Reduced Manifolds.

285 For obtaining these manifolds, the TAE network is fed with K time series that describe the
 286 hydrogen-air homogeneous autoignition problem in different ignition scenarios at the time interval
 287 $[t_0, t_{n-1}]$, where n_t stands for the number of time steps. The expected output is the future values
 288 of the time series, it is to say, the state vectors for the time interval $[t_{0+\Delta t}, t_{n_t}]$.

289 Going further, we evaluate the ability of the TAE model to generalize to unseen ignition sce-
 290 narios. Specifically, we assess the ability of the TAE integrator in interpolation and extrapolation
 291 scenarios. The interpolation scenario represents the homogeneous ignition cases described in Ta-
 292 ble II. The R^2 score will be used again to assess the quality of reconstruction; this assessment is
 293 the first step in the analysis since the validity of the latent manifold is dependant on the outputs
 294 reconstruction accuracy. The R^2 score values are available in Figure 5.

295 After observation, it is visible that can see that the models retain a good representation of
 296 the dynamic behavior, returning R^2 scores larger than 0.96. Moreover, we note that in the latent
 297 space with a dimension equal to 1, the model accuracy decreases. That might be explained by the
 298 fact that combustion is a highly non-linear phenomenon, and latent spaces with larger dimensions
 299 are required to describe accurately the whole thermochemical states, mainly for more complex
 300 chemical mechanisms.

301 Once the accuracy of the reconstruction is assessed, the next objective is to observe the latent
 302 manifolds for different bottleneck sizes, in order to assess the minimum number of latent variables
 303 for a good reconstruction. Such manifolds will be associated with chemical species, in order to
 304 obtain what can be considered as *chemical carriers*. To perform such an association, a correlation
 305 analysis is performed using the Kendall Tau B correlation index due to its robustness and resis-
 306 tance to outliers [39]; in other words, there is a comparison of the temporal behavior between the
 307 thermochemical variables, and the reduced latencies from the same ignition case. In general, when
 308 there is a similarity between two objects, they are likely to hold a high positive correlation index

This is the author's peer reviewed, accepted manuscript. However, the online version of record will be different from this version once it has been copyedited and typeset.

PLEASE CITE THIS ARTICLE AS DOI: 10.1063/5.0167110

Accepted to Phys. Fluids 10.1063/5.0167110

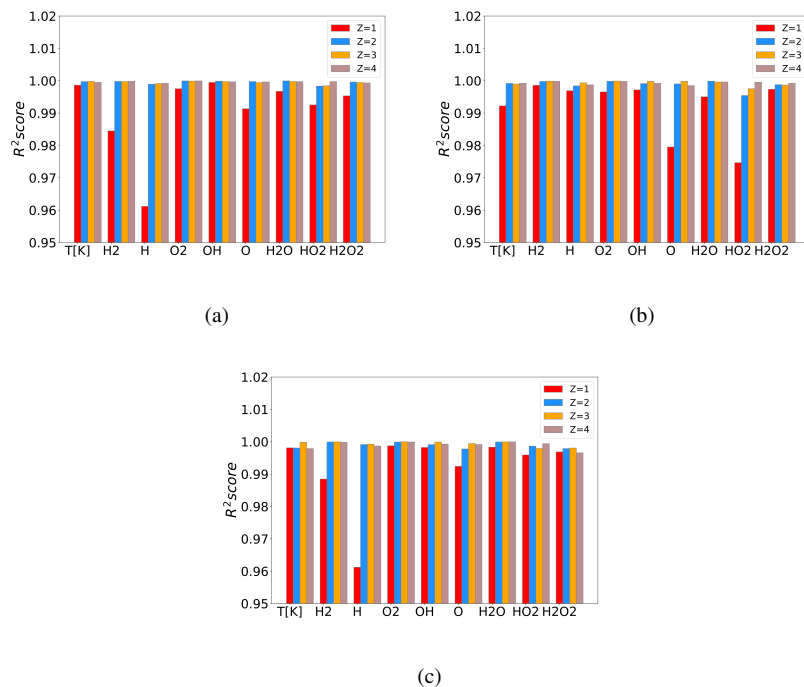


FIG. 5: Interpolation cases R2 scores for all reduced manifolds sizes (Z), y-axes describes the number of latent variables considered by the model, (a) ignition case ($\phi = 0.93, T_0 = 1160K$), (b) ignition case ($\phi = 1.0, T_0 = 1200K$), (c) ignition case ($\phi = 1.10, T_0 = 1130K$).

309 [42]. The resulting associated species are available in Table IV.

310 The chemical carriers presented in Table IV are identifiable due to the dataset preprocessing,
 311 since the easiness of pattern or key features recognition is related to the representation of datasets
 312 under study [42]. It is important to highlight the presence of just one set of chemical carriers, which
 313 is due to the repetitiveness of the same chemical carriers for all the interpolation ignition cases con-
 314 sidered, for the cases in extrapolation, the majority of chemical carriers are retained, however, one
 315 chemical specie might change in each case, i.e. for extrapolation case 1 ($\phi = 0.85, T_0 = 1160K$) the
 316 combination T, O is favored in $z = 2$, while in extrapolation case 2 ($\phi = 1.0, T_0 = 1300K$) the com-
 317 bination T, H is preferred, nevertheless, this change in preference is due to a very small variation
 318 in correlation indexes, around 0.05 maximum, holding always a high similarity with the chemical

TABLE IV: Chemical Carriers Identification

Manifold Size (Z)	Chemical Carriers
1	T
2	T, O
3	H ₂ O ₂ , T, OH
4	O ₂ , O, T, T

319 carriers obtained in interpolation; this can be interpreted as a potential change in the chemical
 320 carriers in agreement to the analyzed combustion regime when applied to a wider range of ignition
 321 cases. However, it is seen that temperature is a key variable appearing as a chemical carrier for all
 322 models. Moreover, in all cases we can see that for the latent space with dimensionality 4, there is a
 323 repetitiveness of the temperature. That may suggest that a latent space with dimension 3 is enough
 324 to describe accurately the whole thermochemical states, and for larger dimensions, the model can
 325 be over-parametrized; such over-parametrization is likely a product of high-dimensionality in the
 326 latencies manifold size, making it difficult to recognize patterns properly [43], producing a non-
 327 optimum identification; it is to keep in mind that the chemical carriers are obtained through an
 328 interpretation of the patterns given by the neural network's signals. From this analysis, it is also
 329 possible to relate the TAE decomposition with the DMD-like techniques, considering that both
 330 obtain a dynamical characterization, that is however, from different perspectives, while DMD-like
 331 techniques use a regression approach [16], while TAE uses a variational approach [24].

332 Aiming to show the degree of correlation between the *chemical carriers* in the interpolation
 333 cases, a visual comparison between the chemical carrier and the respective latent variable is dis-
 334 played in Figure 6, while an example of the obtained correlation indexes is available in Figure 7.
 335 This is done for the ignition case ($\phi = 1.0, T_0 = 1200K$), for a latent space size $z = 2$.

336 On the other hand, Figure 8 shows the latent space behavior of $Z = 2$ for the three ignition cases
 337 under study. It is noted an identical behavior, which just differs in relation to the magnitude of the
 338 latent variables at each time step for the different ignition scenarios. Furthermore, there is a clear
 339 time shift for each ignition case due to the different ignition delay times, which are accurately
 340 represented.

341 The repetition of *chemical carriers* for different ignition conditions, suggests that there might

This is the author's peer reviewed, accepted manuscript. However, the online version of record will be different from this version once it has been copyedited and typeset.

PLEASE CITE THIS ARTICLE AS DOI: 10.1063/5.0167110

Accepted to Phys. Fluids 10.1063/5.0167110

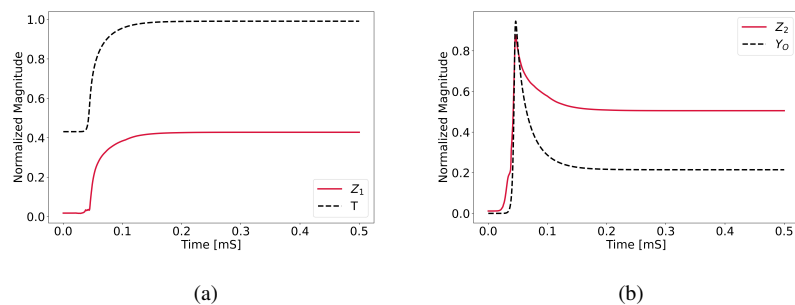


FIG. 6: Latent variables associations for the model manifold size $z = 2$, ignition case ($\phi = 1.0, T_0 = 1200K$), (a) presents the behavior of the first latent variables with its associated thermochemical variable, while (b) presents the same items for the second latent variables; it is to mention that the same association is repeated for all the tested ignition cases.

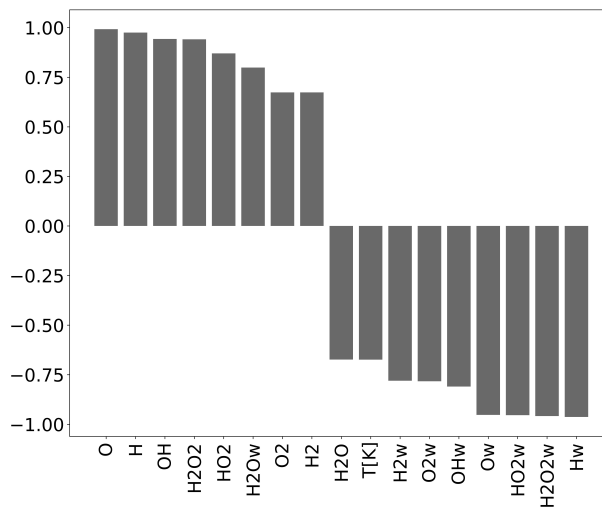


FIG. 7: Kendall Tau B correlation indexes for the second latent variable, manifold size $z = 2$, ignition case ($\phi = 1.0, T_0 = 1200K$), the tendency to extreme correlation indexes is repeated in all cases.

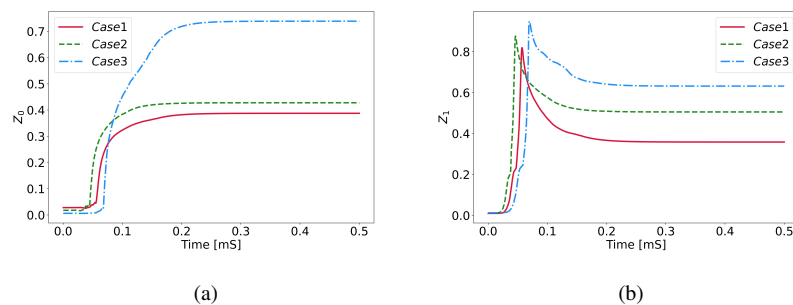


FIG. 8: Latent space visualization for manifold size $z = 2$, (a) shows the first latent variable, and (b) the second latent variable. The Case 1 curve stands for the ignition case ($\phi = 0.93, T_0 = 1160K$), Case 2 for ignition case ($\phi = 1.0, T_0 = 1200K$), and Case 3 for ignition case ($\phi = 1.10, T_0 = 1130K$)

342 be key chemical species that are capable of explaining the variance of a chemical reaction. Aiming
 343 to verify this fact, a permutation importance test is proposed [44]; conceptually, the test consists
 344 of the elaboration of a model that relates a set of inputs and output variables, once the model is
 345 established, each input variable will be permuted individually. If the permuted variable is mean-
 346 ingful for the description of the output, the error in predictions will increase, and therefore a high
 347 permutation weight will be assigned.

348 For testing purposes, a linear regression model will relate the latent chemical species with the
 349 future values of the thermochemical state vector, the intentionality of this is to relate the direct
 350 statistical relationships between the variables, which can be very key for finding chemical species
 351 interactions or dependencies. For the scoring process, an R^2 score is used again. Additionally,
 352 the Sci-Kit Learn built-in function for the score of determination, which means that it allows also
 353 negative values, a situation that reflects the model performing worst than a constant model with
 354 the mean value as an output; this consideration enables the R^2 score to be in the range of $[-\infty, 1]$.

355 Moreover, the Sci-Kit Learn permutation importance scoring follows the Equation:

$$356 \quad i_j = s - \frac{1}{Q} \sum_{n=1}^Q s_{q,j}. \quad (5)$$

357 where i_j stands for the coefficient of importance for the j -th variable, s is the R^2 score from the
358 linear regression without permuting the input variables, and $s_{q,j}$ is the resulting R^2 with a single
359 scrambled input variable. Q refers to the number of repetitions in which the permutation is
360 performed and the determination score is re-calculated [45].

361 However, Equation (5) also provides a warning, and it is that from case to case, the magnitude
362 of importance weights can vary greatly, not only depending on the number of latent variables
363 considered but also in relation to the scrambled variable. Therefore, for each output variable, the
364 resulting importance weights will be added, and this total will be used to normalize the actual
365 importance weights; such a procedure searches to allocate all the weights in the same numerical
366 scale for better interpretation. For the case of the latent space $z = 1$, normalization is applied with
367 the consideration of the maximum importance weight among all the output variables. Figure 9
368 shows the results of the permutation tests. Here, the importance of the chemical carriers for each
369 thermochemical variable is displayed in proportion to the height of the bar.

370 It is observed that just in the case of latent space $z = 1$ the input variables do not hold influence
371 in every output of the mechanism, this finding is also representative of the lower R^2 scores that
372 come from the network's reconstruction in Figure 5. Additionally, from this permutation impor-
373 tance test it is to mention the Figure 9 (c), which the analysis is made for the manifold size $z = 3$.
374 In general, the most important radicals for the ignition development in hydrogen combustion are
375 OH and HO2 [46], however, in this case, the network has chosen as the third chemical carrier the
376 H2O2 chemical species, this is sustained since the TAE network searches for statistical patterns,
377 and therefore, HO2 is not considered since H2O2 holds part of its statistical variance, which is
378 only surpassed by the influence of temperature, therefore, the network judges HO2 as temperature
379 dependent, same that holds significant importance for other chemical species; lastly, in Figure
380 9 (a) it is observed that temperature does not hold a direct significance for H2O2, justifying why
381 this variable is chosen. Another similar analysis is possible from Figure 9 (b), in which the strong
382 relationship between OH, HO2, and Temperature is displayed, at the same time, the interactions
383 between the H and O species is highlighted; the strong bond between T and OH does not appear
384 in Figure 9 (c) since OH is considered as a chemical carrier. This analysis reinforces the fact that
385 TAE highlights the chemical species that are likely to preserve the variance of the mechanism,
386 moreover, such a result suggests that TAE's temporal integration is likely to eliminate unnecessary
387 variables for the description of future thermochemical spaces.

This is the author's peer reviewed, accepted manuscript. However, the online version of record will be different from this version once it has been copyedited and typeset.

PLEASE CITE THIS ARTICLE AS DOI: 10.1063/5.0167110

Accepted to Phys. Fluids 10.1063/5.0167110

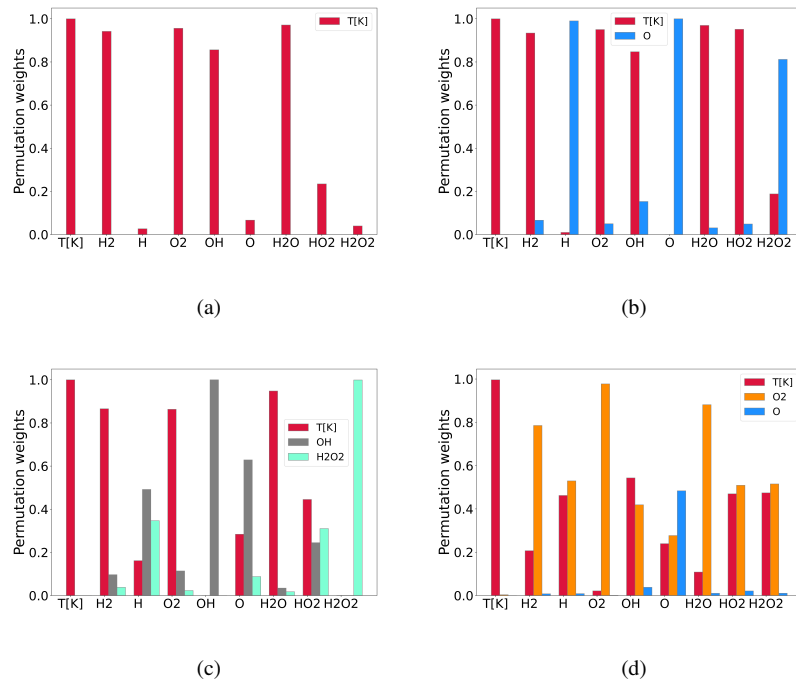


FIG. 9: Permutation importance test weights for chemical carriers for ignition case ($\phi = 1.10, T_0 = 1130K$); (a) latent space $z = 1$, (b) latent space $z = 2$, (c) latent space $z = 3$, and (d) latent space $z = 4$. In each picture, the bars' height is proportional to the importance of each chemical carrier for each of the considered thermochemical variables, which are displayed on the x-axis.

388 B. Comparison with Principal Component Analysis.

389 The most important point of comparison regarding other reduction techniques is the error they
 390 introduce at the moment of reconstruction. To avoid circumstances that could benefit one model
 391 over another, the PCA-Autoencoder is trained with the same normalization and datasets used for
 392 the TAE training, afterward, their performance is evaluated in interpolation cases.

393 For this case, the chosen scoring technique is the Mean Absolute Error (MAE) [47], which gives
 394 an estimation of the average deviation between the predicted values and the real ones, the inclusion

395 of this scoring technique ensures a more comprehensive evaluation of the model's performance.
396 Additionally, since the chemical variables are constrained between values of zero and one, the
397 MAE also gives information on the average percentual error in the predictions. The results for
398 TAE and PCA reconstructions are available in Figure 10, it is to notice the significant difference in
399 the reconstruction errors, which are around four times higher in the PCA technique than in TAE.

400 From Figure 10, it is evident that TAE presents better reconstruction capabilities when con-
401 sidering each chemical species separately; it is to highlight the order difference in the error that
402 is visible in the Figures comparison; since the error is significantly lower for the TAE recon-
403 structions, it is possible to state this technique as more accurate. Moreover, PCA is not able to
404 reproduce the fast phenomena in the chemical reaction, fact that is visible in Figure 11, in which
405 different reconstruction curves are presented for a latent space size $z = 2$, it is observable in the
406 OH and HO2 curves that PCA fails to represent any behavior that takes place for small periods of
407 time, in other words, it is unable to reproduce small chemical time scales.

408 As expected, the reconstruction error is higher in extrapolation cases, as it is possible to observe
409 in Figure 12. However, it is to remark that reconstruction behaviors maintain a similar behavior to
410 the ones in interpolation conditions, this indicates that models actually associate the extrapolation
411 case with the closest known conditions. For the TAE case, the temperature deviations might be
412 also a product of the MSE reconstruction that is applied, since the physically aware term from
413 Equation 4 regularizes only chemical mass fractions.

414 In agreement with the provided results, although PCA provides a much faster method to reduce
415 the dimensionality of combustion systems via matrix decomposition [48], it is observed that TAE
416 reconstruction is more accurate than PCA. Going further, TAE provides interpretable latent spaces
417 allowing an understanding of the underlying physics of the chemical kinetics that may be useful to
418 construct reduced chemical mechanisms. However, when it comes to defining a minimum amount
419 of reduced variables, PCA brings results more straightforwardly, considering the eigenvalues of the
420 matrix decomposition, while TAE requires space exploration for the determination of optimums.

This is the author's peer reviewed, accepted manuscript. However, the online version of record will be different from this version once it has been copyedited and typeset.

PLEASE CITE THIS ARTICLE AS DOI: 10.1063/5.0167110

Accepted to Phys. Fluids 10.1063/5.0167110

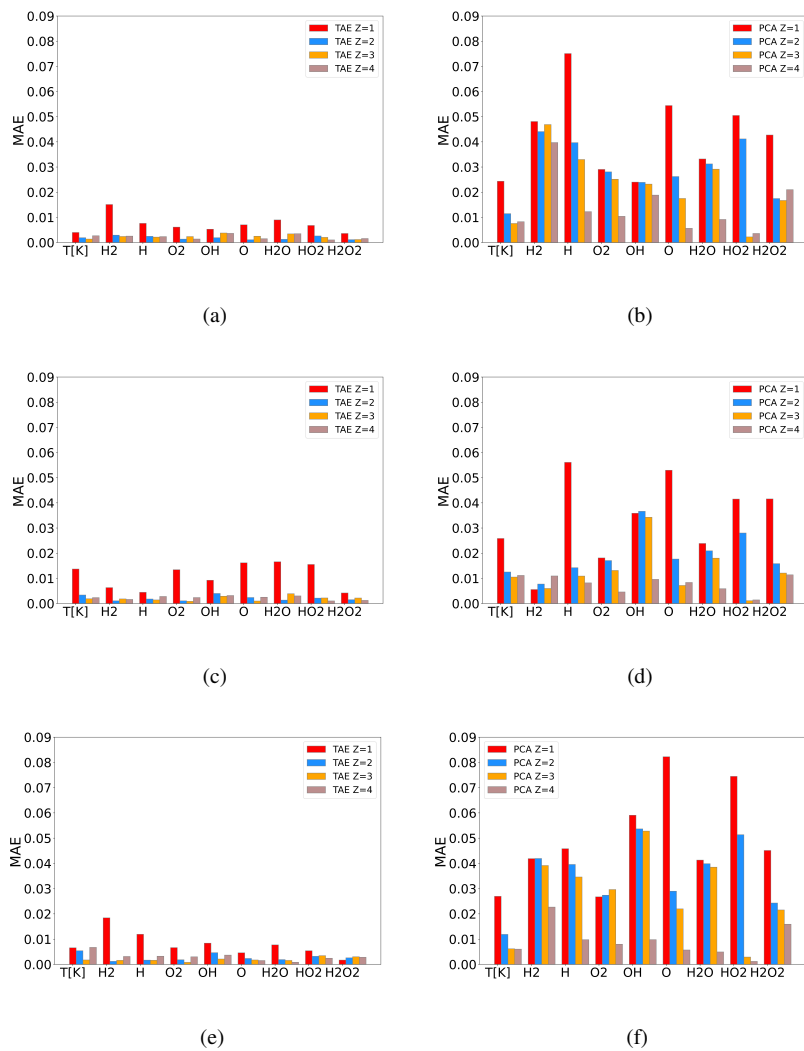


FIG. 10: TAE and PCA MAE scores for interpolation cases and all manifold sizes (Z), figures at the left present TAE's values, while figures at the right present PCA's ones; (a) and (b) present ignition case ($\phi = 0.93, T_0 = 1160K$), (c) and (d) the ignition case ($\phi = 1.0, T_0 = 1200K$), and, (e) and (f) ignition case ($\phi = 1.10, T_0 = 1130K$). All the graphs hold the same y-axis limits for a better comparison of magnitudes between techniques.

This is the author's peer reviewed, accepted manuscript. However, the online version of record will be different from this version once it has been copyedited and typeset.

PLEASE CITE THIS ARTICLE AS DOI: 10.1063/5.0167110

Accepted to Phys. Fluids 10.1063/5.0167110

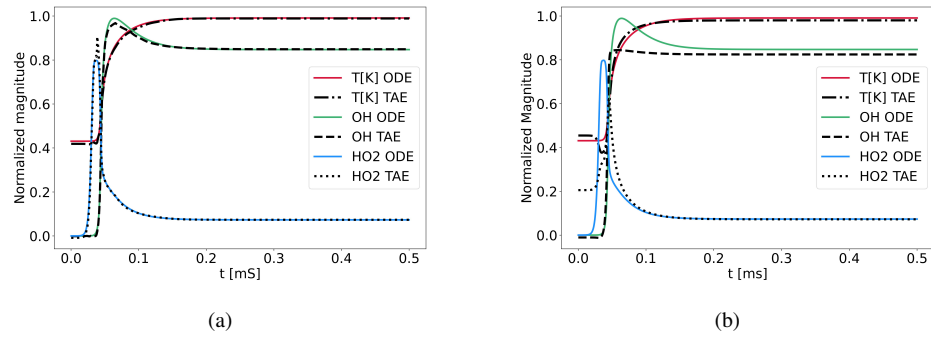


FIG. 11: Normalized chemical species reconstruction for a latent space size $z = 2$ in interpolation, at the left the TAE reconstructions are available, while at the right the PCA reconstructions are displayed for an ignition case $\phi = 1.0, T_0 = 1200K$.

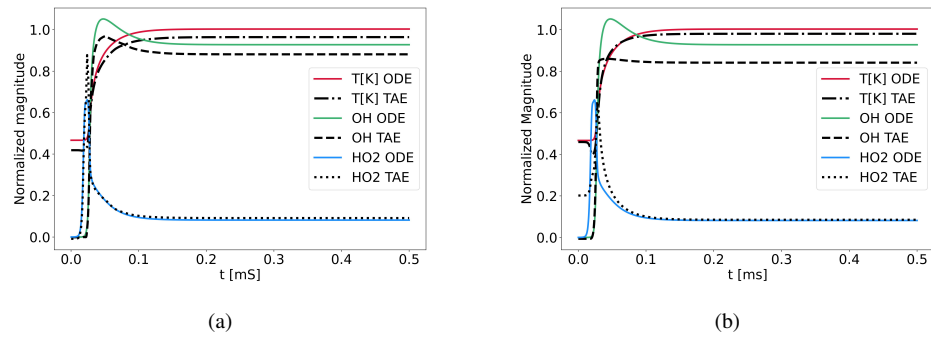


FIG. 12: Normalized chemical species reconstruction for a latent space size $z = 2$ in extrapolation, at the left the TAE reconstructions are available, while at the right the PCA reconstructions are displayed for an ignition case $\phi = 1.0, T_0 = 1300K$.

421 **V. CONCLUSION.**

422 The modeling of combustion chemistry kinetics is a challenging problem for which, no abso-
423 lute solutions have been found. Therefore, exploration of new perspectives and methodologies
424 is a must. Time-lag auto-encoders (TAE) are an interesting ML model, which provides a time
425 integration scheme while allowing the identification of reduced representations. This paper cen-
426 ters on TAE applications for dimensionality reduction, which holds a dynamical characterization
427 of chemistry kinetics. As a dimensionality reduction technique, TAE provides a higher qual-
428 ity of reconstruction in relation to standard techniques (principally PCA), while providing at the
429 same time, a physical interpretation of the reduced manifolds (*chemical carriers* identification).
430 Considering that the *chemical carriers* are the thermochemical variables that hold all the statis-
431 tical variance of the mechanism, the set of network transformations can be easily interpreted as
432 a denoising operation, in which the influence of statistically meaningless variables is canceled.
433 Moreover, with the appropriate treatment, the *chemical carriers* could lead to the study and better
434 interpretation of key chemical pathways, improving the knowledge of chemical mechanisms.

435 Further work must be developed for the assessment of TAE's application as a time integrator,
436 as well as to assess its benefits regarding simpler neural network architectures. Additionally, better
437 model training techniques should be tested for the improvement of autoregressive models, since
438 most of the literature cases consider just the mass conservation principle leaving the temperature
439 variable without any sort of regularization or constraint, allowing for the arising of outliers [49].
440 Lastly, possible applications for the identified latencies for the development of alternative surrogate
441 models should be evaluated.

442 **ACKNOWLEDGMENTS**

443 This research is funded by the Fonds de la Recherche Scientifique de Belgique (F.R.S.-F.N.R.S)
444 project METRIC, under grant No.T.0175.21.
445 Computational resources have been provided by the Consortium des Équipements de Calcul Inten-
446 sif (CÉCI), funded by the Fonds de la Recherche Scientifique de Belgique (F.R.S.-FNRS) under
447 Grant No. 2.5020.11 and by the Walloon Region. Also, this work has received funding from the
448 European Union's Horizon 2020 research and innovation programme under Marie Skłodowska-
449 Curie grant agreement No 801505.

450 **DATA AVAILABILITY STATEMENT**

451 The data that support the findings of this study are openly available in Time-lag-AE-Chemistry-
452 Analysis-Reduction at <https://sandbox.zenodo.org/record/1238601> reference number [1238601].

-
- 453 [1] D. Veynante and L. Vervisch, *Progress in Energy and Combustion Science* **28**, 193 (2002).
454 [2] A. Avdić, G. Kuenne, F. di Mare, and J. Janicka, *Combustion and Flame* **175**, 201 (2017), special
455 Issue in Honor of Norbert Peters.
456 [3] S. B. Pope, *Combustion Theory and Modelling* **1**, 41 (1997), publisher: Taylor & Francis _eprint:
457 <https://doi.org/10.1080/713665229>.
458 [4] M. Singer, S. Pope, and H. Najm, *Combustion and Flame - COMBUST FLAME* **147**, 150 (2006).
459 [5] P. Boivin, C. Jiménez, A. L. Sánchez, and F. A. Williams, *Proceedings of the Combustion Institute*
460 **33**, 517 (2011).
461 [6] P. Pepiot and H. Pitsch, 4th Jt Meet US Sect Combust Inst (2005).
462 [7] F. Contino, H. Jeanmart, T. Lucchini, and G. D'Errico, *Proceedings of the Combustion Institute* **33**,
463 3057 (2011).
464 [8] M. Ihme, W. T. Chung, and A. A. Mishra, *Progress in Energy and Combustion Science* **91**, 101010
465 (2022).
466 [9] A. Parente, J. Sutherland, L. Tognotti, and P. Smith, *Proceedings of the Combustion Institute* **32**, 1579
467 (2009).
468 [10] O. Owoyele and P. Pal, *Energy and AI* **7** (2021), 10.1016/j.egyai.2021.100118.
469 [11] R. Malpica Galassi, P. P. Ciottoli, M. Valorani, and H. G. Im, *Journal of Computational Physics* **451**,
470 110875 (2022).
471 [12] J. C. Sutherland and A. Parente, *Proceedings of the Combustion Institute* **32**, 1563 (2009).
472 [13] G. D'Alessio, A. Cuoci, G. Aversano, M. Bracconi, A. Stagni, and A. Parente, *Energies* **13**, 2567
473 (2020).
474 [14] P. J. Schmid, *Annual Review of Fluid Mechanics* **54**, 225 (2022), [https://doi.org/10.1146/annurev-
475 fluid-030121-015835](https://doi.org/10.1146/annurev-fluid-030121-015835).
476 [15] S. Le Clainche and J. M. Vega, *Complexity* **2018**, 6920783 (2018).

This is the author's peer reviewed, accepted manuscript. However, the online version of record will be different from this version once it has been copyedited and typeset.

PLEASE CITE THIS ARTICLE AS DOI: 10.1063/5.0167110

Accepted to *Phys. Fluids* 10.1063/5.0167110

- 477 [16] A. Corrochano, G. D'Alessio, A. Parente, and S. L. Clainche, "Hierarchical higher-order dynamic
478 mode decomposition for clustering and feature selection," (2023), arXiv:2301.07976 [physics.flu-
479 dyn].
- 480 [17] A. Corrochano, G. D'Alessio, A. Parente, and S. Le Clainche, *International Journal of Mechanical*
481 *Sciences* **249**, 108219 (2023).
- 482 [18] W. Ji, W. Qiu, Z. Shi, S. Pan, and S. Deng, *The Journal of Physical Chemistry A* **125**, 8098 (2021),
483 pMID: 34463510, <https://doi.org/10.1021/acs.jpca.1c05102>.
- 484 [19] A. Sharma, R. Johnson, D. Kessler, and A. Moses (2020).
- 485 [20] A. Hetherington, A. Corrochano, R. Abadía-Heredia, E. Lazpita, E. Muñoz, P. Díaz, E. Moira,
486 M. López-Martín, and S. L. Clainche, "Modelflows-app: data-driven post-processing and reduced
487 order modelling tools," (2023), arXiv:2305.17150 [cs.CE].
- 488 [21] T. Zhang, Y. Yi, Y. Xu, Z. X. Chen, Y. Zhang, W. E, and Z.-Q. J. Xu, *Combustion and Flame* **245**,
489 112319 (2022).
- 490 [22] G. Hinton and R. Salakhutdinov, *Science (New York, N.Y.)* **313**, 504 (2006).
- 491 [23] P. Zhang and R. Sankaran, *Journal of Machine Learning for Modeling and Computing* **3**, 1 (2022).
- 492 [24] C. Wehmeyer and F. Noé, *The Journal of Chemical Physics* **148** (2017), 10.1063/1.5011399.
- 493 [25] K. Champion, B. Lusch, J. N. Kutz, and S. L. Brunton, *Proceedings of the National Academy of*
494 *Sciences* **116**, 22445 (2019).
- 495 [26] K. P. F.R.S, *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science* **2**,
496 559 (1901), publisher: Taylor & Francis _eprint: <https://doi.org/10.1080/14786440109462720>.
- 497 [27] S. Wold, K. Esbensen, and P. Geladi, *Chemometrics and Intelligent Laboratory Systems* **2**, 37 (1987).
- 498 [28] L. van der Maaten, E. Postma, and H. Herik, *Journal of Machine Learning Research - JMLR* **10**
499 (2007).
- 500 [29] A. Parente and J. Sutherland, *Combustion and Flame* **160**, 340 (2013).
- 501 [30] D. G. Goodwin, R. L. Speth, H. K. Moffat, and B. W. Weber, "Cantera: An Object-oriented Software
502 Toolkit for Chemical Kinetics, Thermodynamics, and Transport Processes," (2021).
- 503 [31] E.-M. Wartha, M. Bösenhofer, and M. Harasek, *Combustion Science and Technology* **193**, 2807
504 (2021), <https://doi.org/10.1080/00102202.2020.1760257>.
- 505 [32] P. Saxena and F. A. Williams, *Combustion and Flame* **145**, 316 (2006).
- 506 [33] M. Cavazzuti, "Design of experiments," in *Optimization Methods: From Theory to Design Scientific*
507 *and Technological Aspects in Mechanics* (Springer Berlin Heidelberg, Berlin, Heidelberg, 2013) pp.

This is the author's peer reviewed, accepted manuscript. However, the online version of record will be different from this version once it has been copyedited and typeset.

PLEASE CITE THIS ARTICLE AS DOI: 10.1063/5.0167110

Accepted to Phys. Fluids 10.1063/5.0167110

- 508 13–42.
- 509 [34] M. Baudin, (2015).
- 510 [35] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S.
511 Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew
512 Harp, Geoffrey Irving, Michael Isard, Y. Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur,
513 Josh Levenberg, Dandelion Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike
514 Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Van-
515 houcke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin
516 Wicke, Yuan Yu, and Xiaoqiang Zheng, “TensorFlow: Large-Scale Machine Learning on Heteroge-
517 neous Systems,” (2015).
- 518 [36] T. O’Malley, E. Bursztein, J. Long, F. Chollet, H. Jin, L. Invernizzi, and others, “KerasTuner,” (2019).
- 519 [37] X. Glorot and Y. Bengio, in *AISTATS* (2010).
- 520 [38] D. Kingma and J. Ba, International Conference on Learning Representations (2014).
- 521 [39] C. Croux and C. Dehon, Tilburg University, Center for Economic Research, Discussion Paper **19**
522 (2010), 10.1007/s10260-010-0142-z.
- 523 [40] P. Virtanen, R. Gommers, T. E. Oliphant, M. Haberland, T. Reddy, D. Cournapeau, E. Burovski, P. Pe-
524 terson, W. Weckesser, J. Bright, S. J. van der Walt, M. Brett, J. Wilson, K. J. Millman, N. Mayorov,
525 A. R. J. Nelson, E. Jones, R. Kern, E. Larson, C. J. Carey, Polat, Y. Feng, E. W. Moore, J. VanderPlas,
526 D. Laxalde, J. Perktold, R. Cimrman, I. Henriksen, E. A. Quintero, C. R. Harris, A. M. Archibald,
527 A. H. Ribeiro, F. Pedregosa, P. van Mulbregt, and SciPy 1.0 Contributors, *Nature Methods* **17**, 261
528 (2020).
- 529 [41] R. B. Darlington and A. F. Hayes, New York, NY: Guilford , 603 (2017).
- 530 [42] R. P. W. Duin and E. Pełkalska, “Object representation, sample size, and data set complexity,” in *Data*
531 *Complexity in Pattern Recognition*, edited by M. Basu and T. K. Ho (Springer London, London, 2006)
532 pp. 25–58.
- 533 [43] C. M. Bishop, *Pattern Recognition and Machine Learning (Information Science and Statistics)*
534 (Springer-Verlag, Berlin, Heidelberg, 2006).
- 535 [44] P. BIECEK, *Explanatory model analysis: Explore, explain, and examine predictive models* (CHAP-
536 MAN amp; HALL CRC, 2022).
- 537 [45] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Pretten-
538 hofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and

This is the author's peer reviewed, accepted manuscript. However, the online version of record will be different from this version once it has been copyedited and typeset.

PLEASE CITE THIS ARTICLE AS DOI: 10.1063/5.0167110

Accepted to *Phys. Fluids* 10.1063/5.0167110

- 539 E. Duchesnay, *Journal of Machine Learning Research* **12**, 2825 (2011).
- 540 [46] M. Blocquet, C. Schoemaeker, D. Amedro, O. Herbinet, F. Battin-Leclerc, and C. Fittschen, *Proc*
541 *Natl Acad Sci U S A* **110**, 20014 (2013).
- 542 [47] C. Sammut and G. I. Webb, eds., “Mean absolute error,” in *Encyclopedia of Machine Learning*
543 (Springer US, Boston, MA, 2010) pp. 652–652.
- 544 [48] I. Jolliffe, *Principal component analysis* (Springer Verlag, New York, 2002).
- 545 [49] D. M. Hawkins, “Introduction,” in *Identification of Outliers* (Springer Netherlands, Dordrecht, 1980)
546 pp. 1–12.