

# Applications in Energy and Combustion Science

## BLASTNet: A Call for Community-Involvement Big Data in Combustion Machine Learning

--Manuscript Draft--

Manuscript Number:	
Full Title:	BLASTNet: A Call for Community-Involvement Big Data in Combustion Machine Learning
Article Type:	VSI: ML for Reacting Flows – invited
Section/Category:	Modeling and simulation
Keywords:	Big Data; Deep Learning; Direct Numerical Simulation; BLASTNet
Corresponding Author:	Matthias Ihme  UNITED STATES
Corresponding Author Secondary Information:	
Corresponding Author's Institution:	
Corresponding Author's Secondary Institution:	
First Author:	Wai Tong Chung
First Author Secondary Information:	
Order of Authors:	Wai Tong Chung Ki Sung Jung Jacqueline H. Chen Matthias Ihme
Order of Authors Secondary Information:	
Abstract:	<p>Many state-of-the-art machine learning (ML) fields rely on large datasets and massive deep learning models (with <math>O(10^9)</math> trainable parameters) to predict target variables accurately without overfitting. Within combustion, a wealth of data exists in the form of high-fidelity simulation data and detailed measurements that have been accumulating since the past decade. Yet, this data remains distributed and can be difficult to access. In this work, we present a realistic framework that combines (i) community involvement, (ii) public data repositories, and (iii) lossy compression algorithms for enabling access to high-fidelity data via a network-of-datasets approach. This Bearable Large Accessible Scientific Training Network-of-Datasets (BLASTNet) is consolidated on a community-hosted web-platform (at <a href="https://blastnet.github.io/">https://blastnet.github.io/</a>), and is targeted towards improving accessibility to diverse scientific data for deep learning algorithms. For datasets that exceed the storage limitations in public ML repositories, we propose employing lossy compression algorithms on high-fidelity data, at the cost of introducing controllable amounts of error to the data. This framework leverages the well-known robustness of modern deep learning methods to noisy data, which we demonstrate is applicable in combustion machine learning (CombML) by training deep learning models on lossy direct numerical simulation (DNS) data in two completely different CombML problems – one in combustion regime classification and the other in filtered reaction rate regression. Our results show that combustion DNS data can be compressed by at least 10-fold without affecting deep learning models, and that the resulting lossy errors can even improve their training. We thus call on the research community to contribute to opening a bearable pathway towards accessible big data in combustion.</p>
Suggested Reviewers:	<p>Mathis Bode Research Centre Julich Julich Supercomputing Centre <a href="mailto:m.bode@itv.rwth-aachen.de">m.bode@itv.rwth-aachen.de</a></p> <p>Pinaki Pal Argonne National Laboratory</p>

	pal@anl.gov
	Tarek Echecki NC State University techekk@ncsu.edu
	Evatt Hawkes University of New South Wales evatt.hawkes@unsw.edu.au
<b>Opposed Reviewers:</b>	Michael E Mueller Princeton University muellerm@princeton.edu conflict of interest
	Shashank Yellapantula National Renewable Energy Laboratory Shashank.Yellapantula@nrel.gov conflict of interest
<b>Additional Information:</b>	
<b>Question</b>	<b>Response</b>

**Declaration of interests**

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

☐The authors declare the following financial interests/personal relationships which may be considered as potential competing interests:

# BLASTNet: A Call for Community-Involved Big Data in Combustion Machine Learning

Wai Tong Chung<sup>a,\*</sup>, Ki Sung Jung<sup>b</sup>, Jacqueline H. Chen<sup>b</sup>, Matthias Ihme<sup>a,c</sup>

<sup>a</sup>*Department of Mechanical Engineering, Stanford University, Stanford, CA 94305, USA*

<sup>b</sup>*Combustion Research Facility, Sandia National Laboratories, Livermore, CA 94550, USA*

<sup>c</sup>*SLAC National Accelerator Laboratory, Menlo Park, CA 94025, USA*

---

## Abstract

Many state-of-the-art machine learning (ML) fields rely on large datasets and massive deep learning models (with  $\mathcal{O}(10^9)$  trainable parameters) to predict target variables accurately without overfitting. Within combustion, a wealth of data exists in the form of high-fidelity simulation data and detailed measurements that have been accumulating since the past decade. Yet, this data remains distributed and can be difficult to access. In this work, we present a realistic framework which combines (i) community involvement, (ii) public data repositories, and (iii) lossy compression algorithms for enabling access to high-fidelity data via a network-of-datasets approach. This Bearable Large Accessible Scientific Training Network-of-Datasets (BLASTNet) is consolidated on a community-hosted web-platform (at <https://blastnet.github.io/>), and is targeted towards improving accessibility to diverse scientific data for deep learning algorithms. For datasets that exceed the storage limitations in public ML repositories, we propose employing lossy compression algorithms on high-fidelity data, at the cost of introducing controllable amounts of error to the data. This framework leverages the well-known robustness of modern deep learning methods to noisy data, which we demonstrate is applicable in combustion machine learning (CombML) by training deep learning models on lossy direct numerical simulation (DNS) data in two completely different CombML problems – one in combustion regime classification and the other in filtered reaction rate regression. Our results show that combustion DNS data can be compressed by at least 10-fold without affecting deep learning models, and that the resulting lossy errors can even improve their training. We thus call on the research community to

contribute to opening a bearable pathway towards accessible big data in combustion.

*Keywords:* Big Data, Deep Learning, Direct Numerical Simulation , BLASTNet

---

## Contents

<b>1</b>	<b>Background</b>	<b>3</b>
1.1	Introduction: A big view of machine learning . . . . .	3
1.2	Requirements and pathways towards massive deep learning datasets in CombML	5
1.3	Dimensionality reduction and lossy compression . . . . .	7
1.4	BLASTNet: A big data framework for the combustion community . . . . .	9
1.5	Objectives . . . . .	10
<b>2</b>	<b>DNS Dataset</b>	<b>11</b>
2.1	Classification dataset . . . . .	13
2.2	Regression dataset . . . . .	14
<b>3</b>	<b>Methods</b>	<b>14</b>
3.1	Deep Learning . . . . .	14
3.2	Lossy Compression . . . . .	16
3.3	Evaluation metrics . . . . .	17
<b>4</b>	<b>Results</b>	<b>18</b>
4.1	Effects of Lossy Compression on Data . . . . .	18
4.1.1	Classification . . . . .	20
4.1.2	Regression . . . . .	21
4.2	Deep Learning Predictions . . . . .	22
4.2.1	Classification . . . . .	24
4.2.2	Regression . . . . .	26

---

\*Corresponding author:

*Email address:* wtchung@stanford.edu (Wai Tong Chung)

## 5 Conclusions

28

## Appendix A Training and Validation

36

### 1. Background

#### 1.1. Introduction: A big view of machine learning

Combustion machine learning (CombML) offers numerous opportunities in predictive modeling, scientific discoveries, and intelligent control [1]. One of the most crucial aspects of machine learning (ML) is the availability of data, which in combustion, typically exist in the form of simulation data and experimental measurements. In many ML fields outside of combustion, massive and diverse datasets are the key components in ensuring high predictive accuracy and good generalizability [2].

For example, in computer vision, a state-of-the-art ML field, massive and diverse datasets such as the ImageNet [3] image recognition dataset (170 GB, 1000 classes, 1.4M labeled images) have enabled ML methods to out-perform human capabilities in image recognition [4, 5]. This achievement was made possible by the co-existence of deep learning architectures, such as the 152-layer deep ResNet [5], and the aforementioned ImageNet dataset, along with its corresponding community-involved image recognition competition [4], where researchers could develop ML methods without the laborious task of data collection, and compare results in a transparent manner via an accessible benchmark dataset.

In contrast, datasets found in flow physics, such as the ( $\sim 500$  TB) Johns Hopkins Turbulence Database [6], are not as diverse (9 flow configurations) but can be much greater in size due to increased degree-of-freedom and resolution requirements when compared to digital images. The fidelity and quality of this type of dataset is highly beneficial for applications in detailed scientific analysis, but its lack of diversity, when compared to other datasets [3, 7] from the broader ML community, can be detrimental for training ML algorithms, especially for predicting in unseen configurations. In order to meet this challenge, the flow physics community has developed knowledge-guided ML [8], where domain knowledge can be leveraged

1  
2  
3 towards augmenting datasets, constraining optimization routines, and customizing model  
4 architectures to learn well from small scientific datasets, *i.e.*, the *small data* regime.

5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65  
Outside of flow physics, ML research tends to focus on *big data*. Many improvements  
(including breakthroughs in model architecture such as residual blocks [5], batch normaliza-  
tion [9], and rectified linear units [10]) in deep learning have been tailored towards developing  
*big models* [11] that gain higher predictive accuracy with growing amounts of data [2]. We  
note that both small and big data paradigms do not necessarily compete, and good results  
have been achieved within CombML by combining ideas from both approaches.

Recent developments in big data ML could inspire potential research directions for  
CombML. In natural language processing (NLP), *foundation models* [12] have led to state-  
of-the-art accuracies in a wide range of language prediction tasks. A foundation model is  
a broadly accessible and big ML model (typically with  $\mathcal{O}(10^9)$  trainable parameters) that  
has been pre-trained on massive and diverse datasets, which can then be fine-tuned at later  
stages, by further training with smaller specific datasets (through transfer learning [13]), for  
application to specific problems. This eliminates the need to build and train a powerful ML  
model from scratch, and reduces the amount of data required to solve a tailored ML problem  
after the foundation model has been pre-trained and shared. With this new paradigm, one  
can envision a future development where only small amounts of additional data is needed to  
fine-tune pre-trained CombML foundation models in order to make accurate and affordable  
predictions of flame physics and chemistry in unseen combustion configurations. However,  
this ML approach is currently largely feasible only in NLP, where low-dimensional readily  
labeled text data can be easily mined. In computer vision, while the practice of transfer  
learning still persists, foundation models are comparatively nascent due to dimensionality of  
images (height, width, and color channels:  $N_H \times N_W \times N_C$ ), and the larger cost of generating  
labels, which typically involves manually annotating images for image recognition or object  
detection.

In CombML, the massive, diverse, and labeled dataset required to eventually develop  
foundation models can certainly exist. A recent review [1] on CombML identified over 200  
direct numerical simulation (DNS) cases, which can potentially serve as the basis of a pub-

lic CombML database. We envision that this database can be further populated with a wide variety of existing experimental measurements and large-eddy simulation (LES) data, as well as future data that is expected to grow in complexity and size with advancements in measurement techniques and computational capabilities. Since simulation and experimental data are readily labeled with high-resolution quantities, CombML does not face challenges tied to laboriously annotating datasets, as seen in computer vision. Instead, this community faces the Herculean challenge of storing and accessing data with much higher degrees-of-freedom ( $N_H \times N_W \times N_L \times N_t \times N_\phi$  with dimensions of length, time, and number of scalars). This becomes especially true when considering the scale of data from peta/exascale simulations [14, 15] and high-speed measurements [16].

In summary, massive, diverse, and public combustion datasets are necessary to advance CombML within the big data paradigm. Specifically, the existence of these datasets would enable CombML researchers:

- To minimize the laborious task of data collection, which enables researchers to focus on advancing CombML techniques.
- To make objective and transparent evaluations of predictive accuracy from different ML approaches on *common* datasets.
- To further leverage existing architectural advances from the big data paradigm, and to foster a CombML paradigm that aligns with the broader ML community.
- To improve accessibility to state-of-the-art transfer learning practices towards eventually building CombML foundation models that can solve a wide range of scientific and engineering problems.

## 1.2. Requirements and pathways towards massive deep learning datasets in CombML

We now discuss a set of requirements for a big CombML dataset, which we note are different to the requirements of centralized high-fidelity databases [17]:



- Massive and diverse:** Large and diverse datasets are crucial for ensuring good accuracy and generalizability in state-of-the-art ML algorithms [2]. For example, super-resolution models [18] in computer vision, which have also been applied towards turbulence modeling [19], are typically trained with  $\mathcal{O}(10^3)$  samples [7] of high-resolution images with great diversity. To establish a similar diverse dataset in CombML, we propose a living dataset that continuously accumulates towards at least a total of 1000 individual snapshots from 100 different configurations. Since this volume of data cannot be easily generated from any individual researcher, a community-involved approach should be considered.
- Accessible and consolidated:** Significant resources will be required to store and share at least 1000 snapshots of high-dimensional data without careful treatment. While services, such as Globus [20], currently enable researchers to access data directly from computing and storage facilities, the private permissions required for this service can hinder accessibility. Public accessibility to scientific data is typically achieved by building a centralized database, such as with the Johns Hopkins Turbulence Database [6] or the Sloan Digital Sky Survey [21]. These centralized scientific public datasets typically require dedicated storage infrastructure which consist of a database cluster, web interface system, and dedicated infrastructure for data analysis. While this approach has lead to reliable sources of scientific data, this can incur significant capital costs, as well as additional costs and human labor for maintaining and updating the centralized database. An alternate approach would be to leverage open-source and free ML repositories such as Kaggle [22], which are currently restricted by a  $\mathcal{O}(100)$  GB limit that may not be sufficient for high-fidelity data, as a single snapshot of petascale DNS data can often exceed this limit.
- Sufficient data quality:** The availability of good quality data is without a doubt important to data-driven methods. However, we must emphasize that this dataset must be sufficiently good for training big supervised ML algorithms. In this context, a recent study [23] demonstrated that ImageNet and other popular benchmark datasets

contain up to 10% label error. Despite these errors in training data, ML continues to transform numerous engineering and scientific endeavors. This is because modern deep learning algorithms are inherently robust to noisy data [24]. In fact, it is well-known that introducing small amounts of noise to a training set can be beneficial for improving the generalization of neural networks [25], and is a common form of *data augmentation* [26]. This has significant implications towards the use of compression and dimensionality reduction algorithms for mitigating storage constraints. However, since some combustion applications involve safety-critical conditions, we note that the use of noisy data with ML under these conditions should be treated with caution and thoroughly investigated prior to deployment.

### 1.3. Dimensionality reduction and lossy compression

Combustion modeling has embraced dimensionality reduction methods for chemical reduction, resulting in compact chemical models in turbulent reacting flows with an acceptable amount of error. Interpretable data-driven dimensionality techniques such as principal component analysis (PCA) [27] have also been employed to identify optimal low-dimensional manifolds that can be transported through conservation equations [28, 29]. A related practice involves projecting large dimensions onto low-order manifolds by leveraging well-understood physical principles behind representative flame configurations. This approach has resulted in the formulation of models such as the Burke-Schumann solution [30], the flame-prolongation in intrinsic lower-dimensional manifold [31], the flamelet-generated manifold method [32], and the flamelet/progress variable method [33, 34].

Since big ML algorithms are robust to noisy data [24], dimensionality reduction algorithms can be applied towards high-fidelity data that exceed size restrictions before storage in public ML repositories. However, errors obtained during PCA reduction can be difficult to control, which may result in unpredictable behavior if present in an ML dataset. More complex dimensionality reduction methods such as autoencoders [35] have been shown to be more effective (but less interpretable) than PCA at compressing data while avoiding significant information loss [36, 37], but can still be difficult to control and are computationally

expensive.

Recently, lossy compression algorithms [38] have gained popularity in applications with high-fidelity data due to increasing storage and I/O bottlenecks as computational capabilities and high-speed measurements outgrow disk capabilities. Similar to dimensional reduction techniques, these algorithms reduce the size of data, while introducing small errors to the compressed data. This is in contrast to lossless compression algorithms, which preserve all information during compression. As shown in Table 1, lossy compression algorithms can achieve significantly higher compression ratios (defined as the ratio between the sizes of original data and compressed data, respectively) than lossless compression. In addition, many of these lossy compression methods have been tailored towards compressing high-fidelity scientific data at tractable computational costs and include error-boundedness, which enable users to determine and control the desired error/fidelity of the compressed data. Thus, these methods can be employed towards guaranteeing a level of desired quality when compressing ML training data.

Compressor	Type	Compression Ratio
Deduplication [39]	Lossless	1.5 ~ 3
gzip [40]	Lossless	1.5 ~ 2
FPC [41]	Lossless	1.2 ~ $\mathcal{O}(10)$
ISABELA [42]	Lossy	2.1 ~ $\mathcal{O}(100)$
SZ2 [43]	Lossy	3 ~ $\mathcal{O}(100)$
ZFP [44]	Lossy	3 ~ $\mathcal{O}(100)$
TTHRESH [45]	Lossy	5.1 ~ $\mathcal{O}(100)$

Table 1: Comparison of compression ratios achieved by compression algorithms on scientific datasets. Adapted from [38].

Even an  $\mathcal{O}(10)$ -fold compression could turn the storage of high-fidelity combustion simulation data into a bearable task. For instance, a ten-fold compression on petascale DNS data (with 200 GB per snapshot) would result in a few compressed snapshots that can be readily

shared on public ML repositories such as Kaggle [22]. This process could be repeated for multiple DNS configurations, with links to each distributed dataset curated and hosted on a single community-maintained webpage. Employing such an approach, which we detail in Section 1.4, would eliminate the time and labor required to build and maintain a centralized database by making use of the open-source nature of the broader ML community.

#### 1.4. BLASTNet: A big data framework for the combustion community

In this work, we propose an affordable weakly centralized framework combining the use of lossy compression algorithms with public open-source data repositories and community involvement for sharing massive and diverse deep learning training data for combustion. In particular, this framework is targeted towards improving the diversity of accessible scientific training data, and thus serves a distinct purpose when compared to existing high-fidelity databases [6].

Figure 1 summarizes our proposed framework, *Bearable Large Accessible Scientific Training Network-of-Datasets (BLASTNet)*. BLASTNet is aimed at providing accessibility to raw simulation and measurement data (from a diverse range of configurations), which can be employed for solving a wide range of deep learning problems. This data is shared through Kaggle [22], which has an interface amenable for scientific clusters and also provides the ability to register digital-object-identifiers (DOI) for each dataset. In cases where a single sample of data exceeds ( $\mathcal{O}(100)$  GB limit) storage limits in Kaggle, this data is compressed at a desired level of error, with an error-bounded lossy compression algorithm. Here, we recommend the use of a consistent compression algorithm, SZ2 [43], so that all lossy compressed datasets can be shared in consistent data formats.

The link to, description of, and all other metadata (boundary conditions, initial conditions, fuel composition, DOI) from the dataset can then be shared onto a community-hosted webpage [46], at <https://blastnet.github.io/>, which curates all existing distributed ML datasets and provides a *centralized search interface* to enable convenient public access. In addition, this webpage provides *tutorials* for compressing, decompressing, sharing, and accessing the lossy data. BLASTNet also sets *standards* (further detailed in Section 5), and

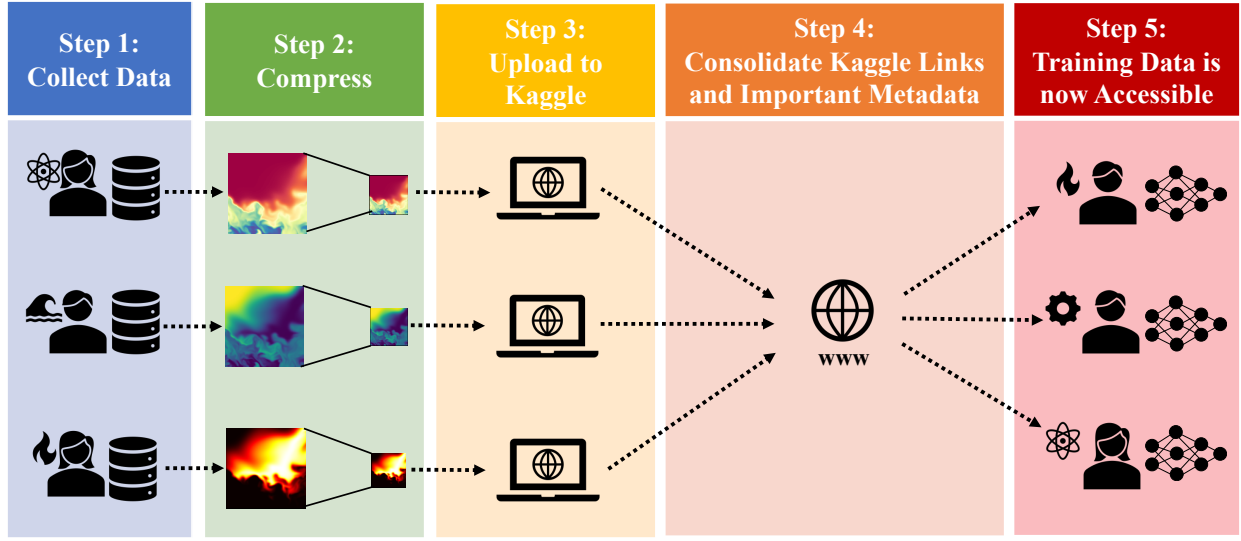


Figure 1: BLASTNet: A community-involved pathway to big combustion data at <https://blastnet.github.io/>

screens the data to ensure that these standards are met. A community *discussion forum* is also hosted on BLASTNet in order to receive continuous feedback from users and to provide a platform for additional support to users. Importantly, to ensure that fair attribution is provided in this open-source project, a version update will be applied to BLASTNet each time a new contribution is provided by the research community to include each individual contributor into BLASTNet’s list of authors, which is a common practice in open-source software [47].

### 1.5. Objectives

The objectives of this work can thus be summarized as follows:

- To advocate the benefits of a massive, diverse, and distributed CombML datasets for deep learning.
- To introduce a platform, at <https://blastnet.github.io/>, for a community-involved network-of-datasets (BLASTNet).

- To demonstrate lossy compression as an affordable and expedited pathway for storing and sharing state-of-the-art high-fidelity data.
- To quantify the compression gained from lossy algorithms and to demonstrate the robustness and limitations of deep learning algorithms to the resulting lossy errors.
- To call on the combustion community to contribute data to BLASTNet.

We note that a key component of BLASTNet operates under the assumption that deep learning methods are robust to the controllable amounts of noise introduced during lossy compression. To investigate the applicability of this assumption within combustion, we apply a lossy compression algorithm (SZ2 [43]) to DNS data of a turbulent lifted hydrogen jet flame in heated co-flow [48], and study the effects of lossy data on training deep learning models in two completely different ML problems namely, combustion regime classification and filtered reaction rate regression. The investigated DNS dataset is described further in Section 2, while the chosen lossy compression algorithm and deep learning architecture are detailed in Section 3. Results from this investigation are presented in Section 4, before concluding in Section 5.

## 2. DNS Dataset

A three-dimensional DNS dataset from a previous study [48] of a turbulent lifted hydrogen jet flame in heated co-flow air is used to demonstrate the robustness of deep learning models to lossy errors. Figure 2 shows the schematic of the DNS configuration. A diluted fuel mixture (65% H<sub>2</sub> and 35% N<sub>2</sub> by volume) is issued from the central slot at an inlet temperature of 400 K. This central jet is surrounded on either side by co-flowing heated air streams with an inlet temperature of 850 K, at atmospheric pressure. The mean inlet axial velocity  $U_{in}$  is given by:

$$U_{in} = U_c + \frac{U_{jet} - U_c}{2} \left( \tanh\left(\frac{y+H/2}{0.1H}\right) - \tanh\left(\frac{y-H/2}{0.1H}\right) \right), \quad (1)$$

where  $U_{jet} = 240 \text{ ms}^{-1}$  represents the mean inlet jet velocity,  $U_c = 2 \text{ ms}^{-1}$  the mean inlet co-flow velocity, and  $H = 2 \text{ mm}$  the jet width at the inlet, respectively. Velocity fluctuations,

obtained by generating an auxiliary homogeneous isotropic turbulence field, are fed from the inlet using the Taylor hypothesis.

This  $2000 \times 1600 \times 400$  computational domain is  $15H \times 20H \times 3H$  in the streamwise  $x$ -, transverse  $y$ -, and spanwise  $z$ -directions, respectively, resulting in a total of 1.28 billion cells. A uniform grid size of  $15 \mu\text{m}$  is placed in  $x$ - and  $z$ -directions, while the  $y$ -directional grid is algebraically stretched outside the flame and shear zones. Improved non-reflecting boundary conditions [49, 50] are adopted in the  $x$ - and  $y$ -directions and periodic boundary conditions are applied in the  $z$ -direction.

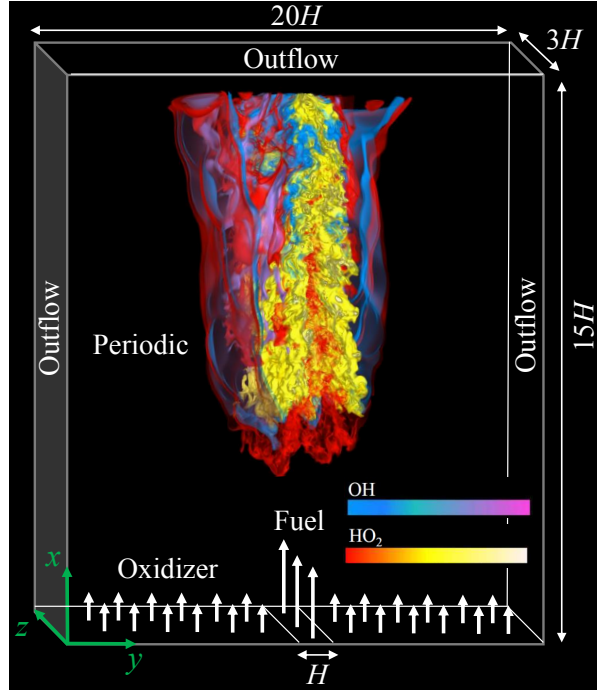


Figure 2:  $\text{H}_2$ -air direct numerical simulation [48] data used in this study.

The Sandia DNS code, S3D [51] was employed for solving the compressible Navier-Stokes, species continuity, and total energy equations. The employed detailed  $\text{H}_2$ -air chemical mechanism composed of 9 species ( $\text{H}_2$ ,  $\text{O}_2$ ,  $\text{H}_2\text{O}$ ,  $\text{O}$ ,  $\text{H}$ ,  $\text{OH}$ ,  $\text{HO}_2$ ,  $\text{H}_2\text{O}_2$ , and  $\text{N}_2$ ) and 21 elementary reaction steps, was developed by Li et al. [17]. In the present study, a  $1200 \times 300 \times 200$  sub-region of the DNS field (*i.e.*, a left half branch of the lifted jet flame) is sampled from the a single 124 GB DNS snapshot, in order to reduce computational costs during training

and analysis while maintaining the fidelity of the flame structure. We employ this 70M-cell subvolume to demonstrate the robustness of deep learning models to noise from lossy compression algorithms in both classification and regression problems, as specified in Section 2.1 and Section 2.2, respectively.

### 2.1. Classification dataset

Within CombML, classification can be useful for optimizing numerical computations [52], detecting catastrophic events [53], and identifying combustion regimes [54]. As detailed in Table 2, we generate five classes of labels for the present classification problem, with the use of the flame index FI [55], progress variable  $C = Y_{\text{H}_2\text{O}}$ , and mixture fraction  $Z$ , as defined by Bilger [56]. The flame index is defined by:

$$\text{FI} = \frac{\nabla Y_{\text{H}_2} \cdot \nabla Y_{\text{O}_2}}{\|\nabla Y_{\text{H}_2}\| \cdot \|\nabla Y_{\text{O}_2}\|}. \quad (2)$$

These five labels were chosen (i) to account for a well-balanced proportion of classes, (ii) to investigate the effects of lossy compression on fine thresholds, and (iii) to investigate the effects of the gradient operator in Equation (2) in magnifying lossy errors. For each label, we extract four flow features  $\{Z, C, Y_{\text{H}_2}, Y_{\text{O}_2}\}$ , and then divide the data into 268 subvolumes, each with  $256 \times 256 \times 3$  cells. Note that 3 cells in the  $z$ -axis is sufficient for preserving spatial information in these samples, since this configuration is homogeneous in the spanwise direction.

Label	Definition
Premixed Flame	$(C > 0.01)$ and $(\text{FI} > 0)$ for all $Z$
Non-premixed Flame	$(C > 0.01)$ and $(\text{FI} \leq 0)$ for all $Z$
Air	$(C \leq 0.01)$ and $(Z \leq 0.01)$
Fuel	$(C \leq 0.01)$ and $(Z > 0.90)$
Fuel-air Mixture	$(C \leq 0.01)$ and $(0.01 < Z \leq 0.90)$

Table 2: Classification labels generated with flame index FI, progress variable  $C$ , and mixture fraction  $Z$ .



## 2.2. Regression dataset

Within CombML, regression is particularly popular for constructing turbulence closure [57], modeling thermodynamics and chemistry [58], and parameterizing combustion manifolds [29]. Here, we generate our regression label by filtering and down-sampling the DNS data to evaluate the Favre-filtered progress variable reaction rate  $\tilde{\omega}_C$ :

$$\tilde{\omega}_C = \frac{\overline{\rho \dot{\omega}_C}}{\bar{\rho}}, \quad (3a)$$

with

$$\bar{\omega}_C(\mathbf{x}) = \int_V \dot{\omega}_C^{DNS}(\mathbf{y}) G(\mathbf{x} - \mathbf{y}, \Delta_F) d\mathbf{y}, \quad (3b)$$

$$G(\mathbf{x} - \mathbf{y}, \Delta_F) = \left( \frac{6}{\pi \Delta_F^2} \right)^{3/2} \exp \left[ \frac{-6(\mathbf{x} - \mathbf{y})^2}{\Delta_F^2} \right], \quad (3c)$$

where  $\bar{\cdot}$  denotes a filtered quantity,  $\tilde{\cdot}$  is a Favre-filtered quantity,  $G$  is a Gaussian filter, and  $\Delta_F = 8\Delta$  is the filter width, which is prescribed to be 8-times larger than the DNS cell width  $\Delta$ . This filter width corresponds to 3 cells (in a corresponding LES) for sufficiently resolving a laminar flame thickness of 0.3 mm, which is evaluated through a stoichiometric 1D premixed flame calculation. The quantity  $\tilde{\omega}_C$  from turbulence-chemistry interaction is of interest within CombML, as shown in other studies [59, 60], and is a suitable quantity to test the robustness of ML models to lossy errors, due to the presence of the exponential operator in the Arrhenius term, which can significantly magnify lossy errors. For each label, we extract two flow features  $\{Z, C\}$  from this dataset, and then divide the data into 177 sub-volumes (each with  $32 \times 32 \times 3$  cells) that encompass the flame region.

## 3. Methods

### 3.1. Deep Learning

Figure 3 shows the deep learning architecture used in both classification and regression problems. This 3-D convolutional neural network (CNN) architecture is based on the autoencoder architecture by Glaws et al. [61], with the input channel  $N_F^{in}$  of the model

modified to suit the present classification ( $N_F^{in} = 4$ ) and regression ( $N_F^{in} = 2$ ) datasets and the filter width reduced to 3. The present network contains 93 layers and approximately 1M trainable parameters, with weights initialized via Xavier initialization [62], and contains 12 residual blocks [5] near the input and output, for improving training and avoiding vanishing gradients during back-propagation. A key component of this architecture is its autoencoder structure. Autoencoder networks can be thought of as a non-linear PCA [35], where raw features are automatically processed by the encoder into an embedded form which can then be forward-propagated by the decoder to generate complex predictions.

For the classification problem, a softmax output activation with five filters  $N_F^{out} = 5$  (for the five classes) is used together with a categorical cross-entropy loss function, while a linear output activation with a single filter  $N_F^{out} = 1$  is used for the regression problem with a mean-absolute-error (MAE) loss function. Train and validation procedures are further detailed in Appendix A.

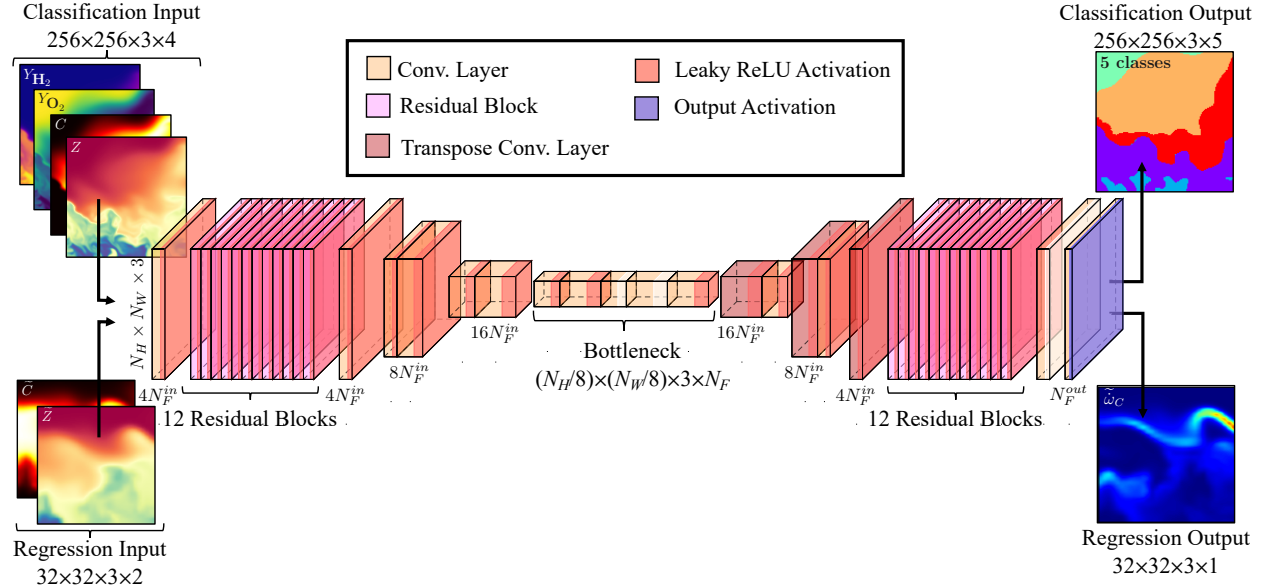


Figure 3: Present 3-D CNN architecture. The number of filters in each layer  $N_F$  is represented in terms of the number of input channels  $N_F^{in}$ .

### 3.2. Lossy Compression

In this work, we employ the SZ2 compressor [43], which combines curve-fitting, the Lorenzo predictor, and data quantization for compressing scientific data. In principle, SZ2 (i) partitions field variables into clusters, (ii) iteratively searches for regression functions that can approximate each cluster with a guaranteed error-bound, and (iii) stores the quantized regression coefficients of the function and indices of the field variables for recomputing the original data during decompression. Data can be compressed effectively since the quantized coefficients and indices are much smaller than the original variables. Compression and decompression of the 12 quantities in the thermo-chemical state-space for the present 72M subvolume requires a total of approximately 35 seconds wall-clock-time on a single CPU. We note that SZ2 has been reported to be at least 2-times faster than the other lossy compressors listed in Table 1 [43, 45].

Thus, SZ2 meets the criteria described in Section 1.3 for compressing high-fidelity data for a large public training database: (i) capable of high compression ratios, (ii) fast, and (iii) allows for bounded error control. While a global error bound is typically used for controlling errors in other compressors [42, 45], SZ2 allows for control via both global error bound, which guarantees that the lossy error in all cells do not exceed a single user-defined value, as well as the point-wise relative error bound [63]  $b_p$ , which guarantees that the lossy error in each cell does not exceed a user-defined percentage of the compressed value. Figure 4 demonstrates the range of lossy data obtained via point-wise relative error control and a corresponding global relative error control, on a curve obtained from the maximum conditional progress variable  $\max(C|Z)$ . The use of point-wise error control is seen to ensure that values near zero are preserved, with positive values guaranteed to stay positive after lossy compression. Point-wise error control also preserves steep gradients more effectively than global error control between  $Z = 0$  and  $Z = 0.3$ . Both these properties are important for preserving the fidelity of steep gradients and small values of chemical species seen commonly within combustion.

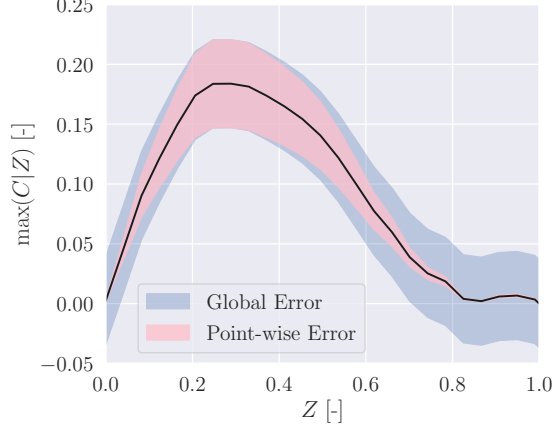


Figure 4: Difference between global and point-wise error bounded control modes in SZ2, illustrated on the maximum conditioned progress variable  $(C|Z)_{max}$ .

### 3.3. Evaluation metrics

For the results discussed in Section 4, we employ several statistical metrics of accuracy to evaluate the effects of lossy compression on data fidelity, and to assess the predictive accuracy of deep learning models. We quantify the quality of lossy compressed data through two popular image quality metrics, *i.e.*, the structural similarity index measure [64] (SSIM) and peak-signal-to-noise-ratio [65] (PSNR). SSIM is evaluated by passing a sliding window  $\Omega$  ( $\Delta_\Omega = 6\Delta$ ) across two scalars  $\phi$  and  $\psi$ , and volume-averaging their statistical quantities:

$$\begin{aligned} \text{SSIM}(\phi, \psi) &= \langle l(\phi, \psi) s(\phi, \psi) r(\phi, \psi) \rangle, \\ &= \left\langle \left( \frac{2\mu_\phi\mu_\psi + c_1}{\mu_\phi^2 + \mu_\psi^2 + c_1} \right) \left( \frac{2\sigma_\phi\sigma_\psi + c_1}{\sigma_\phi^2 + \sigma_\psi^2 + c_1} \right) \left( \frac{\sigma_{\phi\psi} + c_3}{\sigma_\phi\sigma_\psi + c_3} \right) \right\rangle, \end{aligned} \quad (4a)$$

where mean  $\mu_\phi$  and variance  $\sigma_\phi^2$  of the sliding window are:

$$\mu_\phi = \frac{1}{N_\Omega} \int_\Omega \phi d\Omega, \quad (4b)$$

$$\sigma_\phi^2 = \frac{1}{N_\Omega} \int_\Omega (\phi - \mu_\phi)^2 d\Omega, \quad (4c)$$

while  $l$  measures the similarity of  $\mu_{\phi,\psi}$ ,  $s$  measures the similarity of  $\sigma_{\phi,\psi}^2$ ,  $r$  measures correlation of the  $\{\phi, \psi\}$ , and constants  $c_{\{1,2,3\}}$  ensure numerical stability.

PSNR is related to mean-squared-error (MSE):

$$\text{PSNR}(\phi, \psi) = -10 \log_{10} \left( \frac{\max(\phi, \psi)^2}{\text{MSE}(\phi, \psi)} \right). \quad (5)$$

Note that for both metrics, higher values are indicative of higher post-compression quality, with SSIM bounded between -1 and 1, while the highest possible value for PSNR is restricted by the maximum value of a data type, *i.e.*, 48 dB for 8-bit images.

For the classification problem, we evaluate the class accuracy score via a one-versus-all approach, *i.e.*, the number of sample points that have been predicted correctly for a given class divided by the total number of sample points. In the regression problem, SSIM is employed to compare the similarity between filtered progress variable reaction rates from the DNS and from the deep learning models. The normalized mean-squared-error (MSE) for  $N$  number of cells is also employed to measure the difference between the ground truth  $\phi$  and model predictions  $\psi$ :

$$\text{Norm. MSE}(\phi, \psi) = \frac{\sum_{i=1}^N (\phi - \psi)^2}{\sum_{i=1}^N \phi^2}. \quad (6)$$

## 4. Results

### 4.1. Effects of Lossy Compression on Data

This section describes the effects of lossy compression on the training data, while Section 4.2 discusses the effects of training deep learning models with this lossy data.

We first compress flowfields required to solve both regression and classification problems with SZ2. Figure 5 demonstrates that the total compression ratio, from 1% to 50% point-wise error bound  $b_p$ , ranges from 7- to 20-fold compression. Even if we consider only the lowest compression ratio seen in compressing  $\text{H}_2\text{O}_2$  mass fraction, a 4-fold compression of the 124 GB DNS solution file, would enable at least 3 snapshots of this data to be shared as a single dataset on Kaggle. Data compression could be repeated on other flow configurations of a similar scale, and shared via the framework presented in Figure 1 for building a diverse network-of-datasets. Figure 5 also shows that greater compression ratios are seen in state

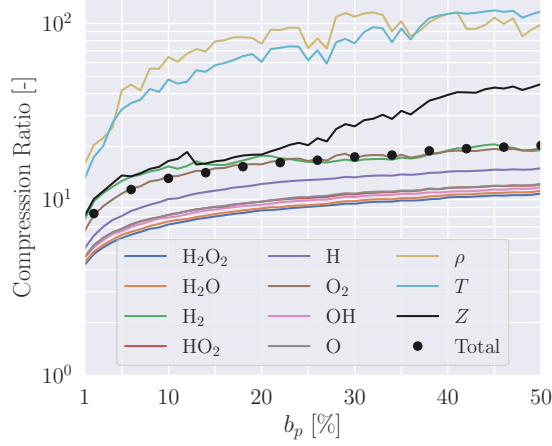


Figure 5: SZ2 compression ratio of different scalar flowfields at varying point-wise error bounds  $b_p$ . Species quantities are expressed in mass fractions.

variables such as density  $\rho$  and temperature  $T$  than in mixture fraction  $Z$  and chemical species.

Visualization of the uncompressed and lossy-compressed flowfields in Figure 6 provides better insight into the different compression ratios exhibited by different quantities. In mixture fraction (left), temperature  $T$  (center), and OH mass fraction  $Y_{\text{OH}}$  (right), PSNR is shown to decrease with increasing point-wise error bound  $b_p$  when compared to the uncompressed flowfields in Figure 6a. SSIM also decreases for  $Z$  and  $T$  with increasing  $b_p$ , but is preserved for  $Y_{\text{OH}}$ . At large settings of  $b_p$ , distortions first appear in regions with large magnitudes and small gradients, as is expected from the point-wise error control, as discussed with Figure 4. As such, large field distortions are first clearly observable in temperature and mixture fraction  $Z$  at  $b_p = 20\%$  and  $b_p = 40\%$  in Figures 6b and 6c, respectively, with no temperature fluctuations visible at  $b_p = 40\%$  in Figure 6c. In these cases, compression should be performed with smaller values of  $b_p$  to preserve flow structure. These results also demonstrate that the point-wise error bound is suited for preserving the large gradients and small magnitudes as seen in  $Z$  at  $b_p = 10\%$  in Figure 6b, and in all  $b_p$  for OH mass fraction  $Y_{\text{OH}}$  (Figure 6a,b,c). This property is useful for preserving the flowfields of many scalar quantities encountered in reacting flow configurations.

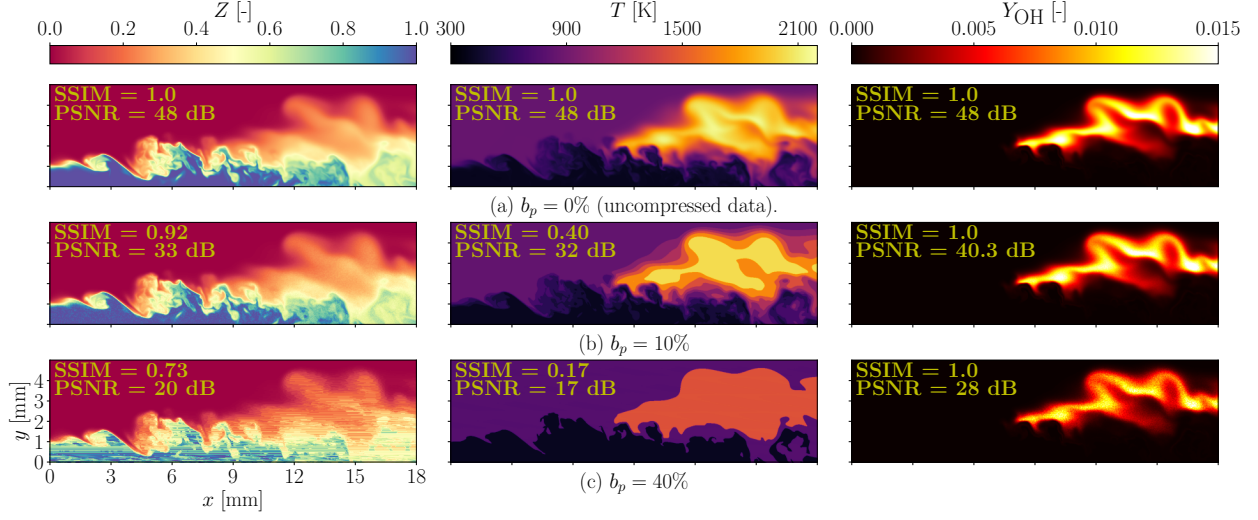


Figure 6: Mixture fraction  $Z$  (left), temperature  $T$  (center), and OH mass fraction  $Y_{\text{OH}}$  (right) from the train set at different levels of maximum point-wise error specified during compression. Quality metrics such as PSNR and SSIM are included in-panel.

#### 4.1.1. Classification

When developing an ML dataset, one can either (i) identify and target a specific supervised learning problem or (ii) share raw data that can then be processed for a target ML problem just before training. Thus, in the present classification problem, we investigate two corresponding scenarios: (i) training with lossy features and clean labels, and (ii) training with lossy features and *post-processed labels* generated from lossy data.

Figure 7 compares the labels used to train the deep learning models at different levels of point-wise error bound  $b_p$ , with Figure 7a showing original uncompressed labels. Significant noise is seen at  $b_p = 10\%$  (Figure 7b), especially in the premixed and non-premixed flame regions, with a 9.3% total label error introduced to the data. This noise is present because lossy errors are magnified by the cell width  $\Delta$  when evaluating scalar gradients used to determine the flame index (Equation (2)). We demonstrate this on a central-differencing

scheme:

$$f(X + \epsilon^l(X)) = \frac{(X_{i+1} + \epsilon_{i+1}^l) - (X_{i-1} + \epsilon_{i-1}^l)}{2\Delta} \quad (7a)$$

$$= \frac{X_{i+1} - X_{i-1}}{2\Delta} + \frac{\epsilon_{i+1}^l - \epsilon_{i-1}^l}{2\Delta}. \quad (7b)$$

Note that in this text, we use the superscript  $\cdot^l$  to denote lossy terms. In the worst case, where lossy errors  $\epsilon_{i+1}^l = b_p X_{i+1}$  and  $\epsilon_{i-1}^l = -b_p X_{i-1}$ :

$$f(X + \epsilon^l(X)) = f(X) + b_p \frac{(X_{i+1} + X_{i-1})}{2\Delta}, \quad (8)$$

which could be significantly larger than:

$$f(X) + \epsilon^l(f) = f(X) + b_p \frac{(X_{i+1} - X_{i-1})}{2\Delta}. \quad (9)$$

Figure 7b shows that the fuel labels also become distorted at  $b_p = 40\%$  as the lossy errors obfuscate the threshold ( $Z \leq 0.01$ ) in generating the labels, as discussed with Table 2, resulting in a total label error of 19.9%.

#### 4.1.2. Regression

We consider the same two scenarios from Section 4.1.1: (i) targeting a specific supervised learning problem or (ii) generating labels from shared lossy simulation data. Specifically, we explore scenarios where (i) pre-processed filtered progress variable reaction rate  $\tilde{\omega}_C^l$  are compressed and shared, and where (ii) post-processed filtered progress variable reaction rate  $\tilde{\omega}_C(T^l, p^l, Y_k^l)$  are generated directly from shared lossy data. The pre-processed label  $\tilde{\omega}_C^l$  is generated by (i) evaluating  $\dot{\omega}_C^{DNS}$  through inputting the thermo-chemical vector  $[T, p, Y_k]^T$  from each cell into the chemical mechanism, (ii) applying Favre-filtering (Equation (3a)) to form  $\tilde{\omega}_C^{DNS}$ , and applying lossy compression to form  $\tilde{\omega}_C^l$ . In contrast, the post-processed label  $\tilde{\omega}_C(T^l, p^l, Y_k^l)$  is generated by (i) applying lossy compression on thermo-chemical vector to form  $[T^l, p^l, Y_k^l]^T$ , (ii) evaluating  $\dot{\omega}_C(T^l, p^l, Y_k^l)$  using the chemical mechanism, and (iii) applying Favre-filtering to form  $\tilde{\omega}_C(T^l, p^l, Y_k^l)$ .

Figure 8 presents percentage of lossy errors from pre-processed  $\tilde{\omega}_C^l$  and post-processed  $\tilde{\omega}_C(T^l, p^l, Y_k^l)$  labels, with the uncompressed filtered progress variable reaction rate  $\tilde{\omega}_C$  (Figure 8a). Figure 8b shows that the normalized lossy error from the pre-processed  $\tilde{\omega}_C^l$  never



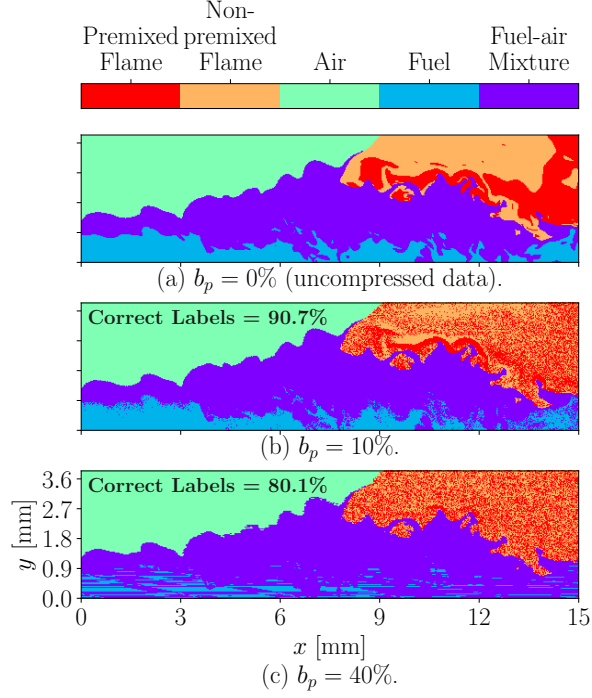


Figure 7: Classification labels generated from lossy data at different levels of point-wise error bounds  $b_p$ .

exceeds the point-wise error bound  $b_p = 40\%$ , while evaluating  $\tilde{\omega}_C(T^l, p^l, Y_k^l)$  from lossy data can result in lossy errors that significantly exceed  $b_p = 4\%$  and  $b_p = 8\%$ , as shown in Figure 8c and Figure 8d, respectively. This is because exponential operators in the Arrhenius term can magnify the lossy errors, which is also seen with gradient operators in Equation (8).

#### 4.2. Deep Learning Predictions

We now explore the effects of lossy data on deep learning. In general, validation and test data do not necessarily match the distribution of the training data, and are usually sampled to represent data encountered after deployment. For instance, when building a data-driven turbulence model in a numerical solver, training data can be extracted from as many different sources as possible to improve generalizability, while validation and test data should match the flow conditions simulated by in the numerical solver [1]. Thus, in the big data framework proposed in Figure 1, we envision a scenario where large quantities of

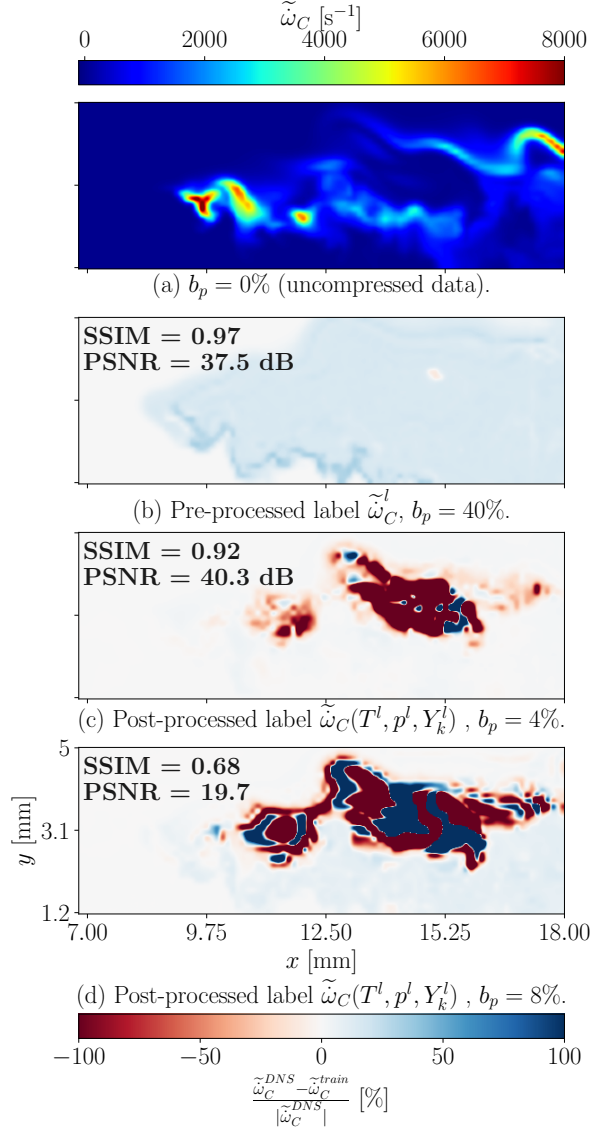


Figure 8: Filtered progress variable reaction rate  $\tilde{\omega}_C$  with percentage errors from lossy compression.

lossy compressed training data can be easily obtained from public repositories, with small quantities of clean test and validation data sampled personally by a user. As such, for the classification and regression problems in Section 4.2.1 and Section 4.2.2, only the training data are lossy-compressed data, while validation and test sets are uncompressed.

#### 4.2.1. Classification

Figure 9a compares class accuracy scores for different levels of point-wise error bounds  $b_p$ , for ML models trained on lossy features and clean labels. A mean class accuracy score of 87% is seen in the baseline case of  $b_p = 0\%$ , which is typical in other classification/segmentation problems [66, 52]. The mean accuracy scores are robust up to  $b_p = 20\%$ , corresponding to a 13-fold compression. At  $b_p > 40\%$ , a high mean accuracy (84%) is still observed, which is in agreement with the well-known observation [67] that ML algorithms are reasonably robust to feature noise.

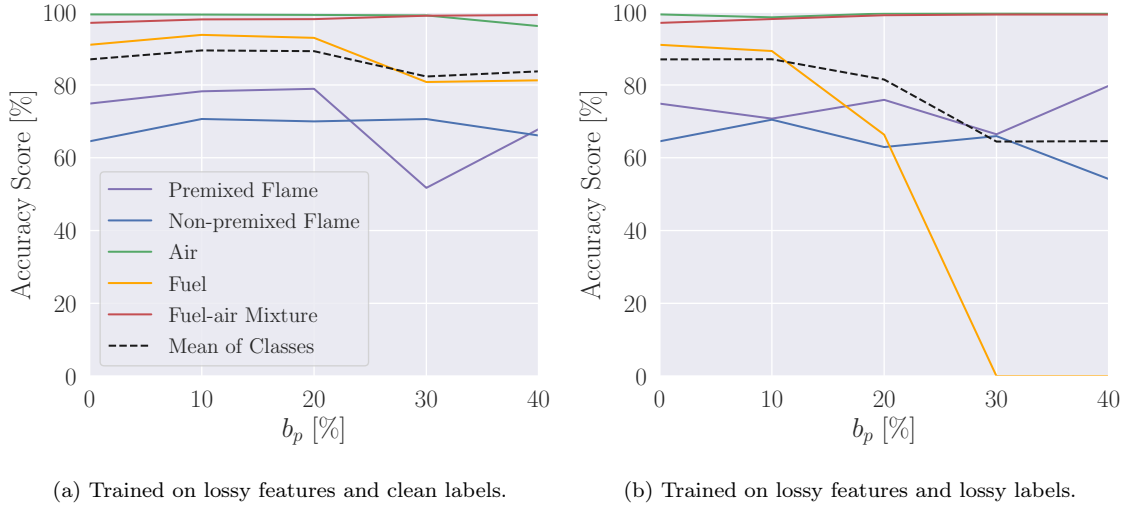


Figure 9: Class accuracy score at different levels of maximum point-wise error specified during compression.

Figure 9b compares class accuracy scores for different  $b_p$ , when training with both lossy features and (post-processed) labels generated from lossy data. The mean accuracy scores are robust to errors up to only  $b_p = 10\%$ , which still corresponds to a 11-fold compression in the original data. At  $b_p \geq 20\%$ , class accuracy for fuel begins to decrease towards 0. This is caused by the distorted fuel labels shown in Figure 7c. Nevertheless, the deep learning model demonstrates reasonably robust behavior in the other classes, especially in the flame regions, up until  $b_p = 40\%$ .

Figure 10 visualizes predictions from the deep learning model trained on lossy features and post-processed labels. Figure 10b shows that the model predictions at  $b_p = 0\%$  are in

agreement with the ground truth labels in Figure 10a. A slightly higher class accuracy of 89% is seen when training with clean labels and lossy features at  $b_p = 10\%$ , as shown in Figure 10c. Increase in accuracy is commonly observed in ML models with the introduction of small amounts of noise during data-augmentation [26], which is well-known to improve neural network models [25]. However, Figure 10d shows that non-premixed flame samples are misclassified as premixed flame near the flame boundary with air when the ML model is trained with lossy labels and lossy features at  $b_p = 10\%$ . This is likely caused by the excessive label noise between the premixed and non-premixed flame regions, as seen in Figure 7b. Similarly, misclassification is seen in the air and premixed flame labels in Figure 10e, where feature noise exceeds  $b_p = 40$ . The aforementioned failure in classifying fuel is clearly observed when the ML model is trained with lossy labels and lossy features at  $b_p = 40\%$  (Figure 10f). Nevertheless, coherent classification is still observed in the flame regions at  $b_p = 40\%$ , despite the high label noise seen in Figure 7c.

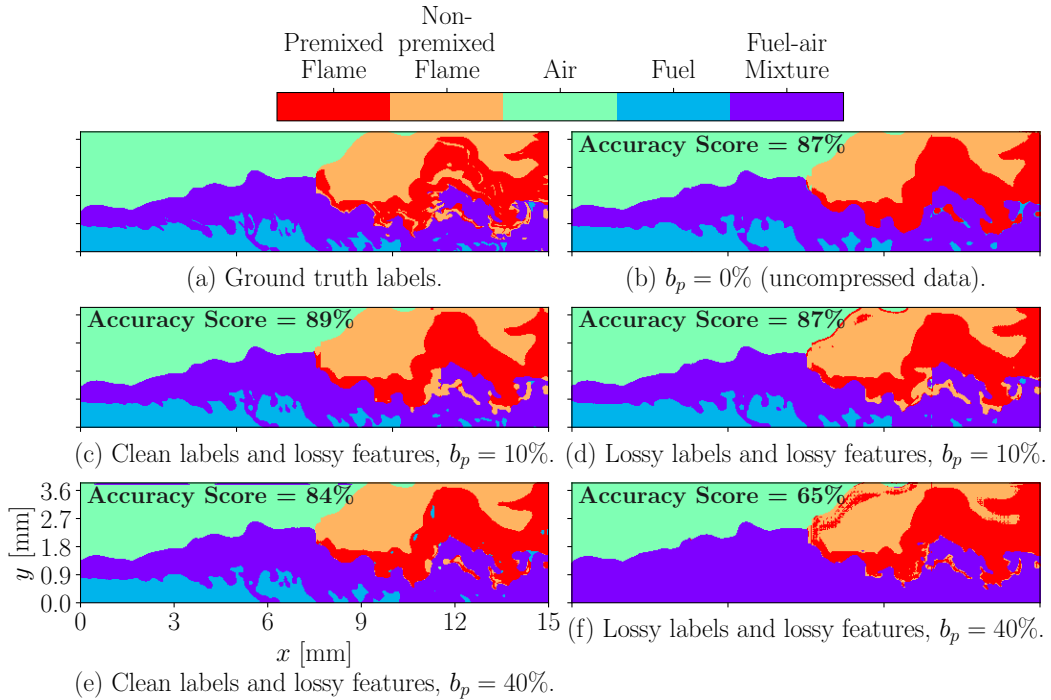


Figure 10: Visualization of ground truth and predictions from ML models in the classification problem.

#### 4.2.2. Regression

Figure 11a compares regression accuracy and error metrics, namely SSIM and the normalized MSE, respectively, for different levels of point-wise error bounds  $b_p$ , for ML models trained on lossy features and lossy pre-processed labels  $\tilde{\omega}_C^l$ . For  $b_p = 0\%$ , a normalized MSE of 22% is similar to results from another study [60]. These values and the high SSIM  $\approx 0.92$  are reasonably consistent up to  $b_p = 40\%$ , corresponding to a 20-fold compression.

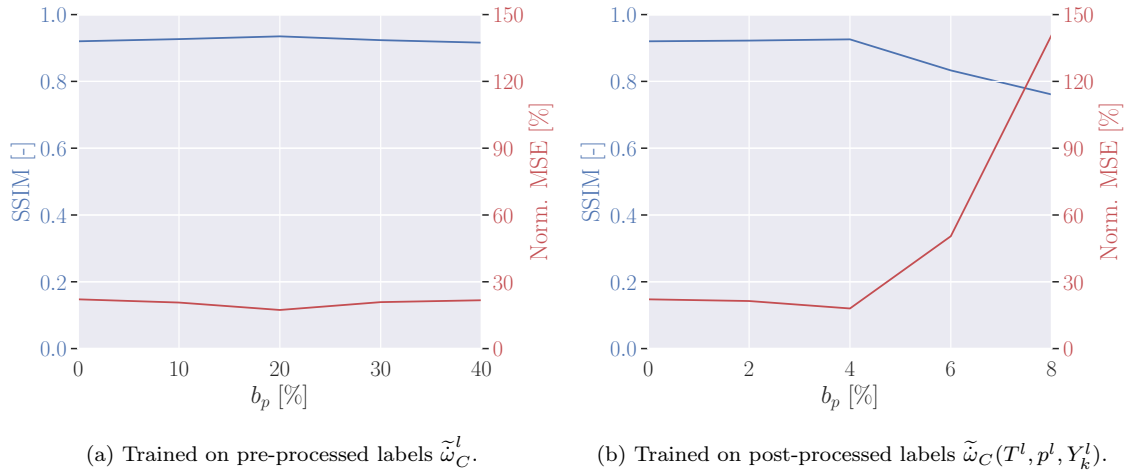


Figure 11: Structural similarity index measure (SSIM) and normalized mean-squared-error (MSE) at different values of point-wise error bounds  $b_p$ . Tested on clean features and clean labels.

Figure 11b compares SSIM and normalized MSE when training with both lossy features and post-processed labels  $\tilde{\omega}_C(T^l, p^l, Y_k^l)$  generated from lossy data, as a function of  $b_p$ . SSIM and normalized MSE are robust to errors up to only  $b_p = 4\%$ , which still corresponds to a 10-fold compression. At  $b_p \geq 4\%$ , SSIM begins to increase while normalized MSE increases significantly due to the magnification of errors during label generation, as shown in Figure 8.

Figure 12 visualizes the predictions from the ML model trained on lossy features and post-processed labels  $\tilde{\omega}_C(T^l, p^l, Y_k^l)$ , along with the filtered DNS. Figure 12b,c shows that the model predictions at  $b_p = 0\%$  and  $b_p = 4\%$  are in reasonable agreement with the ground truth labels in Figure 12a. Over-prediction and under-prediction of  $\tilde{\omega}_C$  is observed in Figure 12d, where  $b_p = 8\%$ .

Figure 13 compares mean conditional filtered progress variable  $\langle \tilde{\omega}_C | \tilde{Z} \rangle$  from the ML

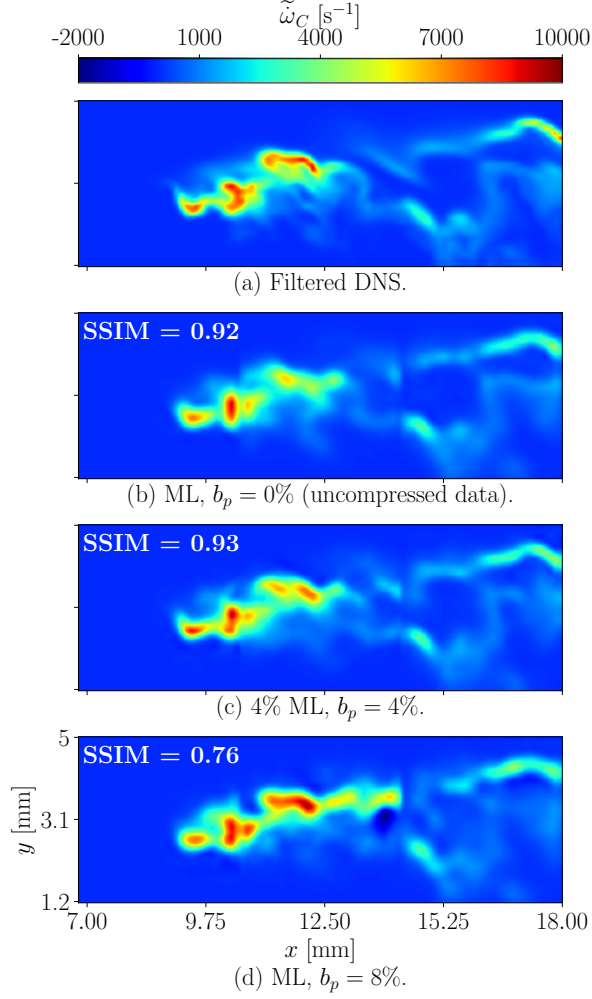


Figure 12: Visualization of filtered DNS and predictions from model trained on lossy features and post-processed labels  $\tilde{\omega}_C(T^l, p^l, Y_k^l)$ , and tested on clean features and clean labels.

model (trained with post-processed labels) with ground truth labels from the filtered DNS. The misprediction observed at  $b_p = 8\%$  in Figure 8d can also be observed here, where a 2-fold over-prediction in  $\langle \tilde{\omega}_C | \tilde{Z} \rangle$  occurs at  $\tilde{Z} = 0.24$ .  $\langle \tilde{\omega}_C | \tilde{Z} \rangle$  at  $b_p = 4\%$  is seen to be in better agreement with the filtered DNS than at  $b_p = 2\%$  and  $b_p = 0\%$ . Introducing small amounts of noise is also seen to improve the classification model, as discussed with Figure 10c.

We note that while the addition of noise can improve training (as commonly done via data augmentation [26]), an excessive amount of noise can lead to bad predictions as shown

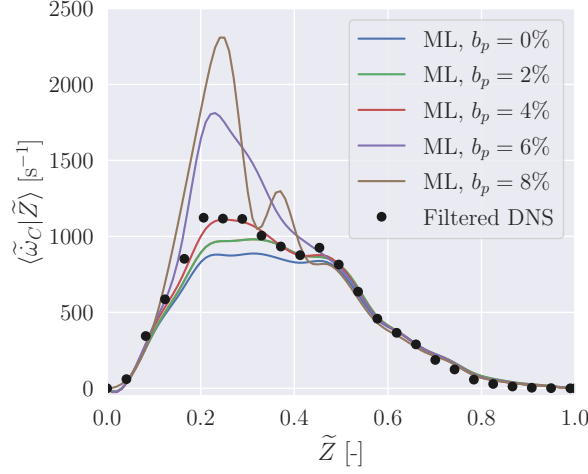


Figure 13: Comparison of mean conditional filtered progress variable  $\langle \tilde{\omega}_C | \tilde{Z} \rangle$  from the filtered DNS and predictions from the ML model, trained with post-processed labels  $\tilde{\omega}_C(T^l, p^l, Y_k^l)$ .

by  $b_p = 8\%$ , and thus caution should be exercised when dealing with noisy data. Hence, for the purposes of BLASTNet, we recommend a  $b_p = 1\%$  (7-fold compression), unless necessary to achieve higher compression ratio in very large datasets. Any further augmentation with noise should be performed after downloading and during training at a user’s discretion. Nevertheless, these results demonstrate that controlled amounts of noise does not affect deep learning models, and in some cases can even be beneficial.

## 5. Conclusions

In this paper, we propose BLASTNet, a realistic framework that combines (i) community involvement, (ii) public data repositories, and (iii) lossy compression algorithms for accessing the wealth of combustion data that already exists in the form of high-fidelity simulations and detailed measurements. Alongside this, we introduce a web-platform, at <https://blastnet.github.io/>, for consolidating and curating the proposed network-of-datasets.

Given the limitations in public storage capacity, a key component of this framework involves the use of lossy compression algorithms for enabling access to petascale simulation data and large experimental measurements. Thus, we evaluate effects of lossy compression algorithms on data quality and deep learning performance on a  $\text{H}_2$ -air lifted flame DNS. To

this end, we train CNN models with labels and features, extracted from lossy DNS data in two completely different regression and classification problems.

In scientific supervised learning, two broad categories of datasets can be encountered:

(i) a dataset targeted at a specific problem, and (ii) raw simulation data or measurements.

Thus, in the 5-class classification problem, two corresponding scenarios are investigated: (i) where lossy features are shared into the repository with clean labels, and (ii) where lossy labels are generated from raw lossy data obtained from the repository, respectively. In the clean label scenario, the classification model is robust to lossy errors in the features up to a point-wise error bound of  $b_p = 20\%$ , which corresponds to a 13-fold compression. In the case of lossy labels and features, the CNN is robust up to  $b_p = 10\%$ , corresponding to a 11-fold compression ratio.

For the regression problem where the filtered progress variable source term is modeled, the two corresponding scenarios are: (i) where lossy features are trained with lossy pre-processed labels  $\tilde{\omega}_c^l$ , and (ii) where lossy features are trained with post-processed labels  $\tilde{\omega}_c(T^l, p^l, Y_k^l)$  generated from lossy simulation data. In the pre-processed scenario, the performance of the regression model is unhindered even at  $b_p = 40\%$ . Due to the magnification of the lossy errors by Arrhenius term calculations, large lossy errors are seen in the post-processed training labels even at  $b_p = 4\%$ . Nevertheless, the regression model still predicts  $\tilde{\omega}_c$  accurately, and the presence of small amounts of noise is even seen to improve the performance of the deep learning model. However, model predictive accuracy drops sharply at  $b_p > 4\%$ . In both regression and classification problems, our results demonstrate that deep learning models applied to combustion can be robust to small amounts of noise.

We now summarize the findings from all sections of this paper towards recommendations for standards in BLASTNet. Based on the requirements a good training datasets listed in Section 1.2, we envision DNS and LES data, covering  $\sim 100$  different configurations with a total of  $\sim 1000$  different snapshots for the first iteration of BLASTNet, with later versions considering experimental data. Since this work demonstrates that deep learning models can train on labels that are post-processed from lossy data, the flowfield in these datasets should contain at least  $[\rho, \mathbf{u}, T, p, Y_k]^T$ , with additional files required for evaluating thermodynamic



1  
2  
3 and transport properties provided to BLASTNet as metadata, so users can recreate any  
4 labels required for training in the wide range of supervised learning problems that are of  
5 interest to CombML. For the choice of the public repository for BLASTNet, we recommend  
6 the use of Kaggle [22], due to the platform’s command-line interface that can enable data  
7  
8  
9  
10  
11 495 access from computing clusters, and ability to provide each data contribution a unique  
12 digital-object-identifier (DOI),  
13

14  
15 We recommend the use of a consistent lossy compressor (SZ2 [43]) to allow for a consistent  
16 data format that would expedite the construction of a data pipeline during training. Since  
17 caution should be exercised as the performance of deep learning algorithms are seen here to  
18  
19  
20  
21 500 degrade rapidly in the presence of excessive noise, lossy compression should only be applied  
22 when necessary, which is largely relevant for bigger DNS cases that exceed 100 GB per  
23 snapshot. In these cases, we recommend a soft-constraint of  $b_p = 1\%$  with SZ2 so that  
24 a few snapshots from a petascale simulation can be stored onto Kaggle. Since the total  
25 compression ratio observed in this study (7-fold compression) is limited by compression of  
26  
27  
28  
29  
30 505 the chemical species, we expect that a 5 to 10 compression ratio would also be observed in  
31 other flame configurations since the volumetric ratio of reacting to non-reacting gases should  
32 be relatively similar across different simulation configurations.  
33

34  
35 To help facilitate these standards and guidelines, tutorials on Kaggle, SZ2, and read-  
36 ing/writing with the recommended data format are provided in BLASTNet. BLASTNet  
37  
38  
39  
40 510 also curates information (boundary conditions, initial conditions, fuel composition, chemical  
41 mechanism, DOI) regarding individual simulation configurations, and provides a centralized  
42 search interface that enables users to download individual cases, along with scripts that  
43 enable batch access to all shared data. In this web-platform, a BLASTNet discussion forum  
44 is also hosted in order to receive community feedback and to provide user support.  
45  
46  
47  
48

49 515 We remind the readers that each BLASTNet contributor will be included to the list-  
50 of-authors in order to cultivate a truly community-involved big training database for com-  
51 bustion. Thus, we call on the combustion community to contribute to this bearable large  
52 accessible scientific training network-of-datasets.  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

## Supplementary Material

The web-platform for consolidating BLASTNet [46] is found in <https://blastnet.github.io/>, which also provides standards for contributing data and tutorials on reading and accessing shared data. The code and models used for this study can found in [https://github.com/IhmeGroup/lossy\\_ml](https://github.com/IhmeGroup/lossy_ml).

## Acknowledgments

The authors acknowledge financial support and computing resources from the Department of Energy (DoE), under award DE-NA0003968. We are also thankful for funding support from the DoE Office of Basic Energy Sciences under award DE-SC0022222. The work at Sandia National Laboratories was supported by the DoE, Office of Basic Energy Sciences, Division of Chemical Sciences, Geosciences, and Biosciences. Sandia National Laboratories is a multimission laboratory managed and operated by National Technology and Engineering Solutions of Sandia, LLC., a wholly owned subsidiary of Honeywell International, Inc., for DoE National Nuclear Security Administration under contract DE-NA0003525.

## References

- [1] M. Ihme, W. T. Chung, A. A. Mishra, Combustion machine learning: Principles, progress and prospects, Prog. Energy Combust. Sci. 91 (2022) 101010.
- [2] C. Sun, A. Shrivastava, S. Singh, A. K. Gupta, Revisiting unreasonable effectiveness of data in deep learning era, Proc. IEEE Int. Conf. Comput. Vis. (2017) 843–852.
- [3] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, L. Fei-Fei, ImageNet: A large-scale hierarchical image database, Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit. (2009) 248–255.
- [4] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al., ImageNet large scale visual recognition challenge, Int. J. Comput. Vis. 115 (2015) 211–252.
- [5] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (2016) 770–778.
- [6] Y. Li, E. Perlman, M. Wan, Y. Yang, C. Meneveau, R. Burns, S. Chen, A. Szalay, G. Eyink, A public turbulence database cluster and applications to study Lagrangian evolution of velocity increments in turbulence, J. Turbul. 9 (2008) No. 31.

- [7] E. Agustsson, R. Timofte, NTIRE 2017 challenge on single image super-resolution: Dataset and study, IEEE Conf. Comput. Vis. Pattern Recognit. Workshop (2017).
- [8] G. E. Karniadakis, I. G. Kevrekidis, L. Lu, P. Perdikaris, S. Wang, L. Yang, Physics-informed machine learning, Nat. Rev. Phys. 3 (2021) 422–440.
- [9] S. Ioffe, C. Szegedy, Batch normalization: Accelerating deep network training by reducing internal covariate shift, Proc. Int. Conf. Mach. Learn. 37 (2015) 448–456.
- [10] V. Nair, G. Hinton, Rectified linear units improve restricted Boltzmann machines, Proc. Int. Conf. Mach. Learn. 27 (2010) 807–814.
- [11] S. Yuan, H. Zhao, S. Zhao, J. Leng, Y. Liang, X. Wang, J. Yu, X. Lv, Z. Shao, J. He, et al., A roadmap for big model, arXiv pre-print 2203.14101 (2022).
- [12] R. Bommasani, D. A. Hudson, E. Adeli, R. Altman, S. Arora, S. von Arx, M. S. Bernstein, J. Bohg, A. Bosselut, E. Brunskill, et al., On the opportunities and risks of foundation models, arXiv pre-print 2108.07258 (2021).
- [13] S. Thrun, Lifelong learning algorithms, in: S. Thrun, L. Pratt (Eds.), Learning to Learn, Springer US, Boston, MA, 1998, pp. 181–209.
- [14] J. H. Chen, Petascale direct numerical simulation of turbulent combustion—fundamental insights towards predictive models, Proc. Combust. Inst. 33 (2011) 99–123.
- [15] S. Treichler, M. Bauer, A. Bhagatwala, G. Borghesi, R. Sankaran, H. Kolla, P. S. McCormick, E. Slaughter, W. Lee, A. Aiken, et al., S3D-Legion: An exascale software for direct numerical simulation of turbulent combustion with complex multicomponent chemistry, in: T. P. Straatsma, K. B. Antypas, T. J. Williams (Eds.), Exascale Scientific Applications, Chapman and Hall/CRC, New York, NY, 2017, pp. 257–278.
- [16] J. H. Frank, Advances in imaging of chemically reacting flows, J. Chem. Phys. 154 (2021) 040901.
- [17] J. Li, Z. Zhao, A. Kazakov, F. L. Dryer, An updated comprehensive kinetic model of hydrogen combustion, Int. J. Chem. Kinet. 36 (2004) 566–575.
- [18] X. Wang, K. Yu, S. Wu, J. Gu, Y. Liu, C. Dong, Y. Qiao, C. C. Loy, ESRGAN: Enhanced super-resolution generative adversarial networks, Proc. IEEE Euro. Conf. Comput. Vis. (2018) 63–79.
- [19] M. Bode, M. Gauding, Z. Lian, D. Denker, M. Davidovic, K. Kleinheinz, J. Jitsev, H. Pitsch, Using physics-informed enhanced super-resolution generative adversarial networks for subfilter modeling in turbulent reactive flows, Proc. Combust. Inst. 38 (2021) 2617–2625.
- [20] I. Foster, Globus online: Accelerating and democratizing science through cloud-based services, IEEE Internet Comput. 15 (2011) 70–73.
- [21] M. R. Blanton, M. A. Bershad, B. Abolfathi, F. D. Albareti, C. A. Prieto, A. Almeida, J. Alonso-García, F. Anders, S. F. Anderson, B. Andrews, et al., Sloan digital sky survey IV: Mapping the Milky

- Way, nearby galaxies, and the distant universe, *Astron. J.* 154 (2017) 28.
- [22] A. Goldbloom, B. Hamner, Kaggle: Your machine learning and data science community, 2010. <https://www.kaggle.com>.
- [23] C. G. Northcutt, A. Athalye, J. Mueller, Pervasive label errors in test sets destabilize machine learning benchmarks, *Proc. Neural Inf. Process. Syst. Track Datasets Benchmarks 1* (2021).
- [24] D. Mahajan, R. Girshick, V. Ramanathan, K. He, M. Paluri, Y. Li, A. Bharambe, L. van der Maaten, Exploring the limits of weakly supervised pretraining, *Proc. Euro. Conf. Comput. Vis.* (2018) 185–201.
- [25] C. M. Bishop, Training with noise is equivalent to Tikhonov regularization, *Neural Comput.* 7 (1995) 108–116.
- [26] C. Shorten, T. M. Khoshgoftaar, A survey on image data augmentation for deep learning, *J. Big Data* 6 (2019) 60.
- [27] I. T. Jolliffe, J. Cadima, Principal component analysis: A review and recent developments, *Phil. Trans. R. Soc. A* 374 (2016) 20150202.
- [28] M. R. Malik, P. Obando Vega, A. Coussement, A. Parente, Combustion modeling using principal component analysis: A posteriori validation on Sandia flames D, E and F, *Proc. Combust. Inst.* 38 (2021) 2635–2643.
- [29] K. M. Gitushi, R. Ranade, T. Echekki, Investigation of deep learning methods for efficient high-fidelity simulations in turbulent combustion, *Combust. Flame* 236 (2022) 111814.
- [30] S. P. Burke, T. E. W. Schumann, Diffusion flames, *Ind. Eng. Chem.* 20 (1928) 998–1004.
- [31] O. Gicquel, N. Darabiha, D. Thévenin, Laminar premixed hydrogen/air counterflow flame simulations using flame prolongation of ILDM with differential diffusion, *Proc. Combust. Inst.* 28 (2000) 1901–1908.
- [32] J. van Oijen, L. de Goey, Modelling of premixed laminar flames using flamelet-generated manifolds, *Combust. Sci. Technol.* 161 (2000) 113–137.
- [33] C. D. Pierce, P. Moin, Progress-variable approach for large-eddy simulation of non-premixed turbulent combustion, *J. Fluid Mech.* 504 (2004) 73–97.
- [34] M. Ihme, C. M. Cha, H. Pitsch, Prediction of local extinction and re-ignition effects in non-premixed turbulent combustion using a flamelet/progress variable approach, *Proc. Combust. Inst.* 30 (2005) 793–800.
- [35] G. E. Hinton, R. R. Salakhutdinov, Reducing the dimensionality of data with neural networks, *Science* 313 (2006) 504–507.
- [36] A. Glaws, R. King, M. Sprague, Deep learning for in situ data compression of large turbulent flow simulations, *Phys. Rev. Fluids* 5 (2020) 114602.
- [37] Y. Lu, K. Jiang, J. A. Levine, M. Berger, Compressive neural representations of volumetric scalar fields, *Comput. Graph. Forum* 40 (2021) 135–146.

- [38] T. Liu, J. Wang, Q. Liu, S. Alibhai, T. Lu, X. He, High-ratio lossy compression: Exploring the autoencoder to compress scientific data, *IEEE Trans. Big Data* (2021). In press.
- [39] D. Meister, J. Kaiser, A. Brinkmann, T. Cortes, M. Kuhn, J. Kunkel, A study on data deduplication in HPC storage systems, *Proc. Int. Conf. High Perform. Comput. Netw. Storage Anal.* 7 (2012).
- [40] J.-L. Gailly, M. Adler, GNU gzip, 1992. <https://www.gnu.org/software/gzip/>.
- [41] M. Burtscher, P. Ratanaworabhan, FPC: A high-speed compressor for double-precision floating-point data, *IEEE Trans. Comput.* 58 (2009) 18–31.
- [42] S. Lakshminarasimhan, N. Shah, S. Ethier, S.-H. Ku, C. S. Chang, S. Klasky, R. Latham, R. Ross, N. F. Samatova, ISABELA for effective in situ compression of scientific data, *Concurr. Comput.* 25 (2013) 524–540.
- [43] X. Liang, S. Di, D. Tao, S. Li, S. Li, H. Guo, Z. Chen, F. Cappello, Error-controlled lossy compression optimized for high compression ratios of scientific datasets, *Proc. IEEE Int. Conf. Big Data* (2018) 438–447.
- [44] P. Lindstrom, Fixed-rate compressed floating-point arrays, *IEEE Trans. Vis. Comput. Graph.* 20 (2014) 2674–2683.
- [45] R. Ballester-Ripoll, P. Lindstrom, R. Pajarola, TTHRESH: Tensor compression for multidimensional visual data, *IEEE Trans. Vis. Comput. Graph.* 26 (2019) 2891–2903.
- [46] W. T. Chung, K. S. Jung, J. H. Chen, M. Ihme, J. Guo, D. Brouzet, M. Talei, BLASTNet simulation dataset, 2022. <https://blastnet.github.io/>.
- [47] D. G. Goodwin, H. K. Moffat, I. Schoegl, R. L. Speth, B. W. Weber, Cantera: An object-oriented software toolkit for chemical kinetics, thermodynamics, and transport processes, 2022. <https://www.cantera.org>.
- [48] K. S. Jung, S. O. Kim, T. Lu, J. H. Chen, C. S. Yoo, On the flame stabilization of turbulent lifted hydrogen jet flames in heated coflows near the autoignition limit: A comparative DNS study, *Combust. Flame* 233 (2021) 111584.
- [49] C. S. Yoo, Y. Wang, A. Trouvé, H. G. Im, Characteristic boundary conditions for direct simulations of turbulent counterflow flames, *Combust. Theor. Model.* 9 (2005) 617–646.
- [50] C. S. Yoo, H. G. Im, Characteristic boundary conditions for simulations of compressible reacting flows with multi-dimensional, viscous and reaction effects, *Combust. Theor. Model.* 11 (2007) 259–286.
- [51] J. H. Chen, A. Choudhary, B. de Supinski, M. DeVries, E. R. Hawkes, S. Klasky, W. K. Liao, K. L. Ma, J. Mellor-Crummey, N. Podhorszki, R. Sankaran, S. Shende, C. S. Yoo, Terascale direct numerical simulations of turbulent combustion using S3D, *Comput. Sci. Discov.* 2 (2009) 015001.
- [52] W. T. Chung, A. A. Mishra, N. Perakis, M. Ihme, Data-assisted combustion simulations with dynamic submodel assignment using random forests, *Combust. Flame* 227 (2021) 172–185.

- [53] A. Cellier, C. Lapeyre, G. Öztarlik, T. Poinso, T. Schuller, L. Selle, Detection of precursors of combustion instability using convolutional recurrent neural networks, *Combust. Flame* 233 (2021) 111558.
- [54] K. Wan, S. Hartl, L. Vervisch, P. Domingo, R. S. Barlow, C. Hasse, Combustion regime identification from machine learning trained by Raman/Rayleigh line measurements, *Combust. Flame* 219 (2020) 268–274.
- [55] H. Yamashita, M. Shimada, T. Takeno, A numerical study on flame stability at the transition point of jet diffusion flames, *Proc. Combust. Inst.* 26 (1996) 27–34.
- [56] R. Bilger, Turbulent jet diffusion flames, *Prog. Energy Combust. Sci.* 1 (1976) 87–109.
- [57] W. T. Chung, A. A. Mishra, M. Ihme, Interpretable data-driven methods for subgrid-scale closure in LES for transcritical LOX/GCH<sub>4</sub> combustion, *Combust. Flame* 239 (2022) 111758.
- [58] R. Nakazawa, Y. Minamoto, N. Inoue, M. Tanahashi, Species reaction rate modelling based on physics-guided machine learning, *Combust. Flame* 235 (2022) 111696.
- [59] Z. M. Nikolaou, C. Chrysostomou, L. Vervisch, S. Cant, Progress variable variance and filtered rate modelling using convolutional neural networks and flamelet methods, *Flow Turbul. Combust.* 103 (2019) 485–501.
- [60] S. Yellapantula, M. T. H. de Frahan, R. King, M. Day, R. Grout, Machine learning of combustion les models from reacting direct numerical simulation, in: H. Pitsch, A. Attili (Eds.), *Data Analysis for Direct Numerical Simulations of Turbulent Combustion: From Equation-Based Analysis to Machine Learning*, Springer International Publishing, Cham, Switzerland, 2020, pp. 273–292.
- [61] A. Glaws, R. King, M. Sprague, Deep learning for in situ data compression of large turbulent flow simulations, *Phys. Rev. Fluid* 5 (2020) 114602.
- [62] X. Glorot, Y. Bengio, Understanding the difficulty of training deep feedforward neural networks, *Proc. Int. Conf. Artif. Intell. Stat.* (2010) 249–256.
- [63] X. Liang, S. Di, D. Tao, Z. Chen, F. Cappello, An efficient transformation scheme for lossy data compression with point-wise relative error bound, *Proc. IEEE Int. Conf. Clust. Comput.* (2018) 179–189.
- [64] Z. Wang, A. Bovik, H. Sheikh, E. Simoncelli, Image quality assessment: From error visibility to structural similarity, *IEEE Trans. Image Process.* 13 (2004) 600–612.
- [65] A. Horé, D. Ziou, Image quality metrics: PSNR vs. SSIM, *Proc. IEEE Int. Conf. Pattern Recognit.* (2010) 2366–2369.
- [66] O. Ronneberger, P. Fischer, T. Brox, U-Net: Convolutional networks for biomedical image segmentation, *Proc. Med. Image Comput. Comput.-Assist. Interv.* (2015) 234–241.
- [67] X. Zhu, X. Wu, Class noise vs. attribute noise: A quantitative study, *Artif. Intell. Rev.* 22 (2004)

177–210.

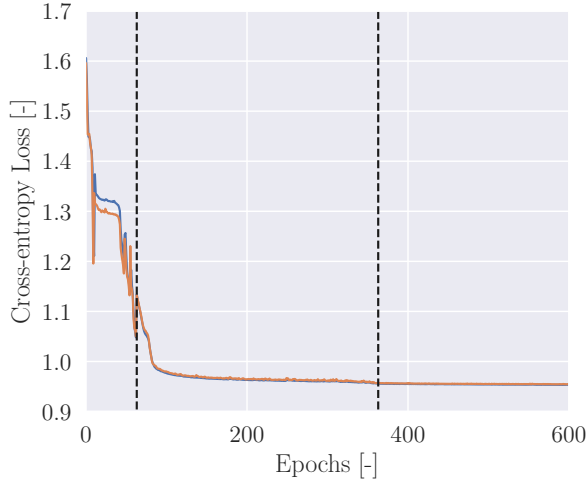
685 [68] D. P. Kingma, J. Ba, Adam: A method for stochastic optimization, arXiv pre-print 1412.6980 (2014).

## Appendix A. Training and Validation

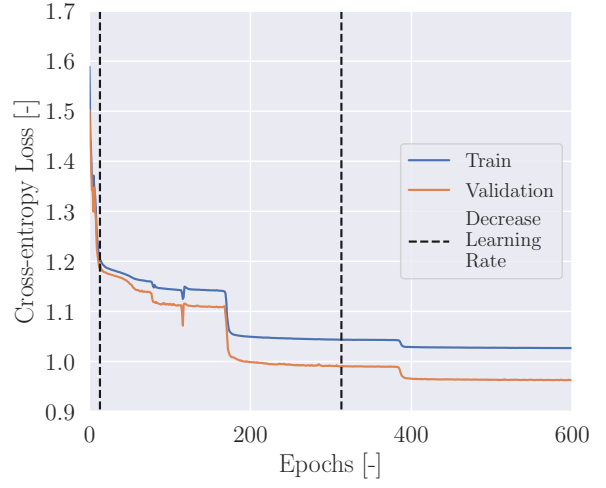
In both classification and regression problems, training is performed with the Adam [68] optimizer. In the classification problem, we employ raw learning rates of 1E-4, 1E-5, and 1E-6 for 100, 300, and 300 epochs, respectively, and early-stopping is employed when necessary. Prior to training, the raw learning rates are multiplied by the square root of the batch size. Here, the batch size is 24.

In the regression problem, we employ raw learning rates of 1E-4, 5E-5, and 1E-5, for 300 epochs each, with batch size of 36. Training both regression and classification models on four Tesla V100 GPUs requires a total of approximately 4 hours of wall-clock-time for each case.

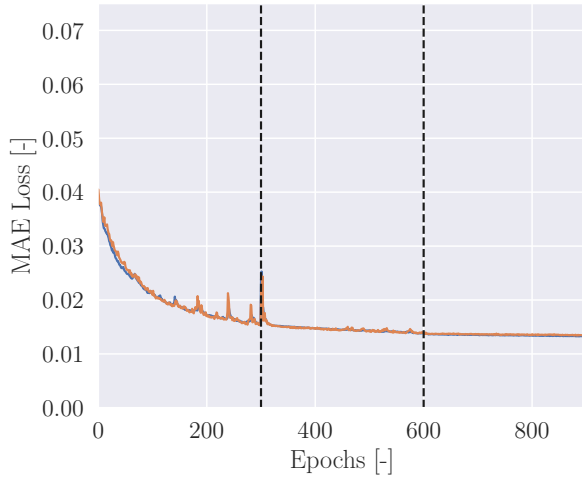
Training and validation losses for selected cases are shown in Figure A.14. In Figures A.14b and A.14d, the converged validation loss can be lower than the training loss, leading to higher validation accuracy than training accuracy. This is caused by the absence of lossy errors in the validation set, as described in Section 4. Otherwise, training shows no sign of overfitting in Figures A.14a and A.14c



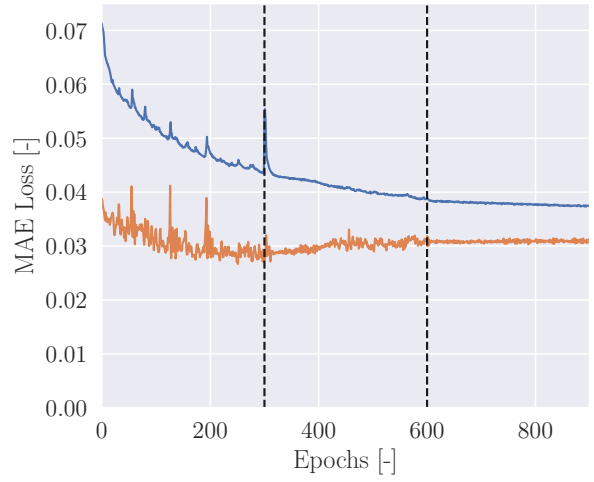
(a) Classification: Lossy and clean labels,  $b_p = 20\%$ .



(b) Classification: Lossy features and lossy labels,  $b_p = 10\%$ .



(c) Regression: Lossy features and lossy pre-processed labels,  $b_p = 10\%$ .



(d) Regression: Lossy features and post-processed labels,  $b_p = 8\%$ .

Figure A.14: Loss during training.



– Highlights –

## BLASTNet: A Call for Community-Involvement Big Data in Combustion Machine Learning

Wai Tong Chung, Ki Sung Jung, Jacqueline H. Chen, Matthias Ihme

- Present weakly centralized framework for enabling access to diverse scientific data for combustion machine learning
- BLASTNet: Bearable Large Accessible Scientific Training Network-of-Datasets combines community involvement, public data repository, lossy compression, and consolidation through community-hosted webpage: <https://blastnet.github.io>
- Demonstrate framework and data compression in application of DNS data to two diverse CombML-problems: regression and classification
- Provide recommendation for community-contribution to shared database for deep learning algorithms