



Anomaly detection in scientific data using joint statistical moments

Konduri Aditya^a, Hemanth Kolla^{a,*}, W. Philip Kegelmeyer^a, Timothy M. Shead^b, Julia Ling^c, Warren L. Davis IV^b

^a Sandia National Laboratories, Livermore, CA 94550, United States

^b Sandia National Laboratories, Albuquerque, NM 87123, United States

^c Citrine Informatics, Redwood City, CA 94063, United States

ARTICLE INFO

Article history:

Received 22 September 2018

Received in revised form 3 March 2019

Accepted 4 March 2019

Available online 13 March 2019

Keywords:

Anomaly detection

Scientific computing

Co-kurtosis

Tensor decomposition

Hellinger distance

Auto-ignition

ABSTRACT

We propose an anomaly detection method for multi-variate scientific data based on analysis of high-order joint moments. Using kurtosis as a reliable measure of outliers, we suggest that principal kurtosis vectors, by analogy to principal component analysis (PCA) vectors, signify the principal directions along which outliers appear. The inception of an anomaly, then, manifests as a change in the principal values and vectors of kurtosis. Obtaining the principal kurtosis vectors requires decomposing a fourth order joint cumulant tensor for which we use a simple, computationally less expensive approach that involves performing a singular value decomposition (SVD) over the matricized tensor. We demonstrate the efficacy of this approach on synthetic data, and develop an algorithm to identify the occurrence of a spatial and/or temporal anomalous event in scientific phenomena. The algorithm decomposes the data into several spatial sub-domains and time steps to identify regions with such events. Feature moment metrics, based on the alignments of the principal kurtosis vectors, are computed at each sub-domain and time step for all features to quantify their relative importance towards the overall kurtosis in the data. Accordingly, spatial and temporal anomaly metrics for each sub-domain are proposed using the Hellinger distance of the feature moment metric distribution from a suitable nominal distribution. We apply the algorithm to two turbulent auto-ignition combustion cases and demonstrate that the anomaly metrics reliably capture the occurrence of auto-ignition in relevant spatial sub-domains at the right time steps.

© 2019 Published by Elsevier Inc.

1. Introduction

Anomaly detection is such a widely studied topic, and has found numerous applications in various contexts, that it defies easy generalization. Nonetheless, the vast majority of applications that have embraced anomaly detection methods have characteristics that may not be representative of scientific data. Chandola et al. [1] emphasize that the key aspects of anomaly detection include the nature of input data, type(s) of anomaly and output of anomaly detection. In all these aspects, scientific data have distinctly different attributes compared to all other domains. As the scale of scientific investigations

* Corresponding author.

E-mail addresses: akondur@sandia.gov (K. Aditya), hnkolla@sandia.gov (H. Kolla), wkp@sandia.gov (W.P. Kegelmeyer), tshead@sandia.gov (T.M. Shead), jling@citrine.io (J. Ling), wldavis@sandia.gov (W.L. Davis).

<https://doi.org/10.1016/j.jcp.2019.03.003>

0021-9991/© 2019 Published by Elsevier Inc.

keeps ever increasing, robust anomaly detection is becoming increasingly critical. One of the key findings of a Department of Energy Workshop on mathematics of data [2] is “*near real-time identification of anomalies in streaming and evolving data is needed in order to detect and respond to phenomena that are either short-lived or urgent*”.

Some of the challenges of anomaly detection in scientific data stem from the following attributes:

- Multi-variate, multi-physics phenomena: the observations are of numerous variables (tens to hundreds) that represent coupled non-linear physics and hence elude easy assumptions about statistical (in)dependence.
- Multi-scale dynamics: the different observed variables span vastly different orders of magnitude since they are active at different scales in space and time.
- The observed data are not discrete but rather continuous and smoothly varying over multiple decades. With computational power ever increasing, the numerical resolution of investigations, e.g. scientific computing, is increasingly finer over a broader range of scales.
- In most cases the field variables are not Gaussian. Examples of non-Gaussian distributions include beta or bi-modal shape distributions for reacting scalars in turbulent combustion and log-normal shape for dissipation in turbulence.
- We are focusing on data rich, not data-sparse, scenarios like extreme-scale computational simulations or extremely well resolved measurements/observations.

There are many scenarios where detecting scientific anomalies may be critically important. Many scientific investigations involve large quantities of rapidly streaming data: massively parallel computational fluid dynamics (CFD) simulations, large-scale particle accelerator data, real-time climate and meteorological simulations, etc. Identifying anomalies as they occur can help judicious steering of these investigations (e.g. trigger analyses, data check pointing, mesh/time-step refinement, model parameter refinements, etc.). Across these varied scientific domains a general definition of an “anomaly” may be elusive and even within one domain it could be problem or regime dependent. Yet physics-driven anomalies, as opposed to investigative/measurement anomalies, occur in these settings and are important to detect.

We clarify that our focus is on identifying anomalous events in scientific investigations, if and when they occur in streaming scientific data, and not whether a specific observation is anomalous relative to the rest. To this end, we want to develop an unsupervised methodology for detecting *statistically identifiable but anomalous scientific phenomena e.g. ignition events in turbulent combustion and cyclones in climate simulations*. We had previously proposed a semi-supervised method based on random forests [3]. Our intuition is that these phenomena have a statistical signature measurable in the higher statistical joint moments. Accordingly, we propose a methodology that is centered on analyzing high-order joint moments in multi-variate scientific datasets. The anomalous event may occur at any time on some subset(s) of the spatial domain.

The central hypothesis of the approach we propose is as follows:

- In data rich settings like scientific computing, anomalies have a discernible statistical signature since the data samples are usually large in number.
- Higher statistical moments, in particular kurtosis, are good indicators of outliers.
- For multi-variate data, by analogy to principal component analysis (PCA), principal vectors of the fourth order joint moment tensor – principal kurtosis vectors – signify directions along which outliers lie.
- The occurrence of an anomalous event manifests as a *sudden detectable change in the principal kurtosis vectors*.

We also note that, in choosing to analyze statistical joint moments, we assume that the large data sizes ensure that higher moments are reasonably converged. Moreover, joint moments can be computed in a computationally efficient manner and fast, single-pass algorithms for computing arbitrary order joint moments over distributed datasets are readily available in literature [4]. From a scientific computing perspective these algorithms have a high compute intensity and have regular, contiguous memory access patterns. They are scalable and likely to be efficient compared to algorithms that may be based on, say, decision trees.

An outline for the remainder of the paper is as follows. In section 2, we present the background necessary to develop the anomaly detection algorithm, which includes a brief summary of previous work, mathematical idea behind the algorithm, a survey of tensor decomposition methods, and test of the mathematical idea on synthetic datasets. The algorithm is described in section 3. Section 4 describes the combustion test cases on which the algorithm is employed to detect auto-ignition (anomalous) events and the results from their analyses. Conclusions and further work are presented in section 5.

2. Background

An anomaly can be loosely defined as an occurrence of something that is “abnormal”, “atypical” or “unexpected” [5]. This is predicated on the assumption that what may be normal/typical/expected is well known and well defined, which may not always be the case. Accordingly, a wide variety of methods across various domains have come to be considered as anomaly detection methods. We first present a brief summary of anomaly detection methods from recent review papers that survey the literature on this topic. We will then discuss the applicability of existing methods, or lack thereof, for our domain of interest, anomalous events in scientific investigations.

2.1. Previous work

Recent review papers by Chandola et al. [1], Campos et al. [6] and Goldstein and Uchida [7] provide a useful survey and overview of various anomaly detection methods. They all conclude that it may be difficult to generalize, and compare, the various methods which stems from the inherent difficulty in defining, across domains, reasonable measures of what may be deemed as “normal” and “anomalous”. Nonetheless, a broad categorization may be discerned. The types of anomalies usually sought are broadly three:

- *point anomalies*: whether a given individual sample or observation is anomalous,
- *collective anomalies*: whether a collection of samples/observations when considered together is anomalous, even if individual samples are not,
- *contextual anomalies*: whether a point or collection of samples is otherwise normal, but anomalous given a specific context.

The vast majority of existing methods deal with the first kind, *point anomalies*. Accordingly the methods result in the identification of an anomaly at the individual sample level. Far fewer methods are devoted specifically to *collective* and *contextual anomalies*, and sometimes these problems are transformed to or posed as a *point* anomaly detection problem. The methods can also be distinguished based on whether the result is an anomaly score or a binary label/classification. Furthermore, the methods span the full spectrum of *supervised*, *semi-supervised* and *unsupervised* paradigm. While we will not attempt an overview here, suffice it to say that the most commonly used *unsupervised* methods (see [6,7]) are clustering or nearest-neighbor based and involve measures of distance (with respect to nearest neighbors, global/local clusters), or density, the intuition being outlier samples have a large distance from normal samples or occur in regions of low density. The more comprehensive survey of Chandola et al. establishes other prominent categories such as classification-based, statistical (parametric and non-parametric), information theoretic and spectral anomaly detection methods.

The specific setting of interest to us, as detailed in the next subsection, is *unsupervised* detection of *anomalous events* in scientific data, which may be considered closest to the *collective* anomaly detection class of methods. Such methods have been developed in the context of sequential anomaly detection, spatial anomaly detection and graph anomaly detection [1]. However, none of these methods appear to be readily applicable for detecting anomalous scientific events, as described below, motivating the need for a new method.

2.2. Anomalous or extreme events in scientific data

Our interest is specifically in identifying anomalous events in scientific data which are sometimes also referred to as extreme events. These events represent genuine scientific phenomena that may be rare and/or extreme but are not necessarily spurious, as the interpretation may be for *point* anomalies. Some examples are ignition/extinction events in combustion data, tornadoes in climate data, crack propagation in fracture mechanics, etc. In a scientific investigation such extreme events are, by definition, not indicated by any individual sample observation, but by a group of observations, rendering *point* anomaly detection methods ineffective. Of the *collective* anomaly detection methods described by Chandola et al. [1], graph based techniques are not obviously applicable since few scientific data sets are represented using graphs. Continuum mechanics (e.g. solid/fluid mechanics) scientific data, which are representative of a broad class of scientific applications, consider joint spatio-temporal domains (discretized spatial mesh and discrete time instances) and hence a purely sequence-based (time domain) or spatial-anomaly based techniques will not capture a joint spatio-temporal anomaly. This is the very setting for our work.

For a vast majority of statistical learning applications the assumption of normally distributed data is reasonable and hence all the statistical information is encapsulated in the covariance matrix. However, for scientific data the distributions are not often Gaussian and there is relevant statistical information in moments higher than the second (variance). Furthermore, scientific data often represent tightly coupled physical processes and hence the observed variables are, more often than not, statistically dependent. Both aspects suggest that joint distributions, and joint moments higher than second order (covariance), are relevant, as opposed to marginal distributions and moments.

For the specific case of anomaly detection, then, it is a matter of analyzing the appropriate higher joint moment. By definition, the outlier samples have an increasingly greater contribution the higher the moment. However, from a practical perspective, we suggest that the fourth moment – kurtosis – is appropriate. While kurtosis has been often held as a measure of “peakiness” or “flatness”, we refer to the paper by Westfall [8] that establishes kurtosis as an unambiguous measure of “either existing outliers (for the sample kurtosis) or propensity to produce outliers (for the kurtosis of a probability distribution).” Westfall illustrates [8] that as kurtosis increases, the contribution to it from the portion of data centered around the mean becomes vanishingly small, no matter how the center is defined (i.e. no matter how many multiples of standard deviation). In the limit of infinite kurtosis the contribution from any finite portion centered around the mean is zero no matter how large this portion. These results suggest that kurtosis is a reliable enough measure of outliers. Accordingly, our approach is centered around analyzing the joint fourth moments – the co-kurtosis – which is a fourth order tensor (co-variance is a second order tensor i.e. a matrix). For a vector \mathbf{V} of size N_v , the co-kurtosis tensor can be constructed as:

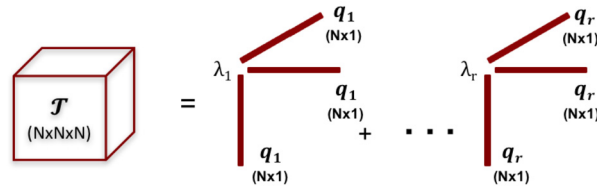


Fig. 1. Symmetric CP decomposition of a third order symmetric tensor \mathcal{T} .

$$\mathcal{T}_{ijkl} = v_i v_j v_k v_l, 1 \leq i, j, k, l \leq N_v, \quad (1)$$

where $v_i \in \mathbf{V}$. To reiterate, we seek, by analogy to PCA, principal vectors of the co-kurtosis tensor which can be interpreted as principal directions along which outliers lie. This effectively becomes a symmetric tensor decomposition problem.

2.3. Joint moment tensor decomposition

Joint moment tensors have been analyzed in different settings, but are usually considered expensive due to the curse of dimensionality, e.g., the size of the fourth moment tensor is N_f^4 , where N_f is the number of random variables or features. By definition, the tensor is symmetric which means that the number of unique entries is smaller than N_f^4 , but to leading order the scaling is still a fourth power of N_f . Jondeau et al. [9] recognizing that non-normality is important in analyzing financial market data, consider decompositions of the co-skewness and co-kurtosis tensors for analysis of stock market data, and conclude that the first few factors of the decomposition of these tensors contain useful information about market returns. The fourth moment tensor is also at the heart of mathematical underpinnings of Independent Component Analysis (ICA), which we will review shortly.

Since most of the properties of matrix factorizations do not extend in general to higher order tensors [10], the appropriate decomposition technique, and its associated interpretation, needs to be chosen carefully. We briefly present an overview of major classes of tensor decompositions that may be applicable to the present problem, before explaining the particular method chosen in this study. We use third order tensors, only for illustration purposes, in the following discussion.

2.3.1. Symmetric CP decomposition

Canonical polyadic (CP) decomposition seeks a factorization of a tensor as a sum of outer products of real-valued vectors. Mathematically, for a third order tensor \mathcal{T} , this can be expressed as

$$\mathcal{T} = \sum_{i=1}^r \lambda_i x_i \otimes y_i \otimes z_i, \quad (2)$$

where \otimes denotes the outer product between the sets of vectors x_i, y_i, z_i and the number of such sets sought, r , is the rank of the decomposition. Even when the tensor \mathcal{T} is symmetric, it is possible to seek a decomposition which is asymmetric, that is, x_i, y_i and z_i are not equal. However, a symmetric decomposition always exists for a symmetric tensor [11] such that

$$\mathcal{T} = \sum_{i=1}^r \lambda_i x_i \otimes x_i \otimes x_i, \quad (3)$$

as illustrated in Fig. 1, and r in this case is referred to as the *symmetric rank*. While seemingly analogous to eigenvalue decomposition of a symmetric matrix (PCA can be interpreted as the eigenvalue decomposition of the co-variance matrix), a few key differences remain. The rank r in the above decomposition is not a unique number, but lies within certain bounds [11]. Moreover, it is not necessary that the decomposition be orthogonal, which is the case for eigenvalue decomposition of a matrix. Orthogonal decompositions of symmetric tensor do not exist in general [12], but if it is known to exist for a specific tensor the decomposition problem can be reduced to a matrix factorization problem [12,13].

2.3.2. Higher order Singular Value Decomposition (HOSVD)

An alternate decomposition, which extends the matrix singular value decomposition (SVD) concept to higher order tensors is Higher-Order Singular Value Decomposition (HOSVD) [14]. For a general third order tensor, this is written mathematically as

$$\mathcal{T} = \mathcal{S} \times_1 \mathbf{U}^{(1)} \times_2 \mathbf{U}^{(2)} \times_3 \mathbf{U}^{(3)} \quad (4)$$

where the core tensor, \mathcal{S} , has the same dimensions as \mathcal{T} . The factor matrices $\mathbf{U}^{(1)}, \mathbf{U}^{(2)}$ and $\mathbf{U}^{(3)}$ are orthogonal matrices that result from “unfolding” (matricizing) \mathcal{T} along modes 1, 2 and 3, respectively, and performing SVD over the resulting matrix. The symbol \times_k denotes a “mode- k ” tensor-matrix product (see [14] for details). For a symmetric tensor the result of “unfolding” is the same along any mode and hence the factor matrices become identical. This special case is illustrated

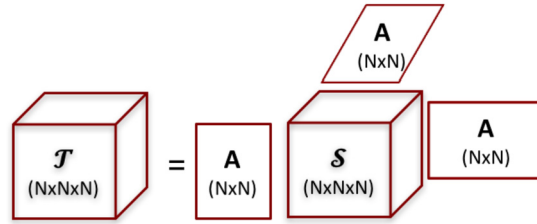


Fig. 2. HOSVD of a third order symmetric tensor \mathcal{T} .

in Fig. 2. For more details on the mathematical properties and conventions of HOSVD, the analogies to and differences with matrix SVD, the reader is referred to De Lathauwer et al. [14]. The column vectors of the (identical) factor matrices may be interpreted by analogy to principal vectors in PCA. However, one key departure between HOSVD and matrix SVD is that the core tensor \mathcal{S} is not purely diagonal but dense. This makes the interpretation of the column vectors of factor matrices, by analogy, not straightforward.

2.4. Independent Component Analysis

In the present study the most useful approach for the fourth moment tensor decomposition appeared to be by way of analogy to Independent Component Analysis (ICA). ICA specifically deals with non-Gaussian random variables and considers a scenario where a set of statistically independent non-Gaussian random variables $s = [s_1 \ s_2 \ \dots \ s_q]^T$ are linearly mixed, superimposed with independent white noise, and observed as the set of random variables $y = [y_1 \ y_2 \ \dots \ y_p]^T$. Mathematically,

$$y = As + n \quad (5)$$

where $A \in \mathbb{R}^{p \times q}$ is the mixing matrix and $n \in \mathbb{R}^p$ is independent white noise, and ICA aims to identify the source set s given the observed set y .

The connection between ICA and higher moment tensor decomposition is explained in many papers (see De Lathauwer and Moore [15] or Anandkumar et al. [13]). Formally, given the statistical model in Eq. (5), the fourth order cumulant tensor of y (observed variables) is related to the outer product of column vectors, a_i , of the mixing matrix A and the excess kurtosis, κ_i , of the individual sources s_i [16]:

$$C_4^y = \sum_{i=1}^q \kappa_i a_i \otimes a_i \otimes a_i \otimes a_i. \quad (6)$$

The cumulant tensor is related to the fourth and second moment tensors as

$$\begin{aligned} [C_4^y]_{i_1 i_2 i_3 i_4} &= \mathbb{E}[y \otimes y \otimes y \otimes y] - \mathbb{E}[y_{i_1} y_{i_2}] \mathbb{E}[y_{i_3} y_{i_4}] \\ &\quad - \mathbb{E}[y_{i_1} y_{i_3}] \mathbb{E}[y_{i_2} y_{i_4}] - \mathbb{E}[y_{i_1} y_{i_4}] \mathbb{E}[y_{i_2} y_{i_3}], \quad 1 \leq i_1 \dots i_4 \leq q, \end{aligned} \quad (7)$$

where \mathbb{E} is the expectation operator. These authors have also proposed algorithms for performing the tensor decomposition in the form of Eq. (6). De Lathauwer and Moore [15] propose a method based on “simultaneous third-order tensor diagonalization”, Anandkumar et al. [13] propose a robust variant of a tensor power method, while Kolda [12] proposes a method involving eigenvalue decomposition of a matrix that is linear combination of slices of the tensor.

However, De Lathauwer and Moore [15] and Anandkumar et al. [13] also observe the connections to a much simpler matrix SVD problem. Specifically, in a discussion on computational complexity, Anandkumar et al. [13] note that if the cumulant tensor C_4^y is unfolded into a matrix M^y (the unfolding is invariant to choice of mode since C_4^y is symmetric), then Eq. (6) can be transformed to

$$\text{mat}(C_4^y) := M^y = \sum_{i=1}^q \kappa_i a_i \otimes \text{vec}(a_i \otimes a_i \otimes a_i) \quad (8)$$

and hence the vectors a_i can be determined from an SVD of M^y (mat and vec denote operations that convert a tensor to a matrix and a vector, respectively). For its simplicity, ease of interpretation and computational cost considerations, we adopt this approach to decomposing the fourth moment tensor in this study.

2.5. Tests with synthetic data

The method chosen to identify principal directions of joint moment tensor has three simple steps: (a) construct the joint fourth cumulant tensor (C_4^y in Eq. (7)), (b) matricize the tensor (M^y) along any of the modes, and (c) perform SVD of the

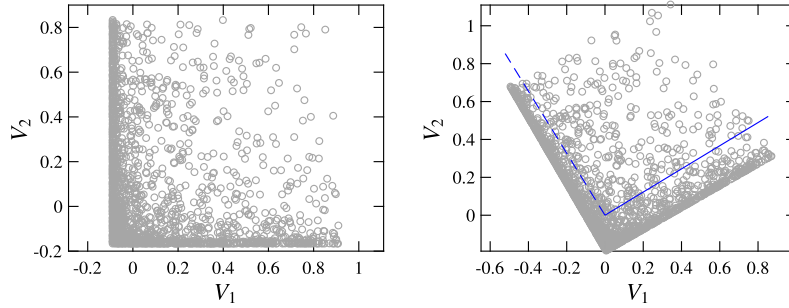


Fig. 3. (Left) Synthetic bi-variate statistically independent dataset with beta marginal distributions. (Right) The independent bi-variate is “mixed” (rotated by 30°) to transform to a dataset with non-zero higher joint moments. The first (solid blue line) and second (dashed blue line) principal kurtosis vectors of this dataset are also shown. (For interpretation of the colors in the figure(s), the reader is referred to the web version of this article.)

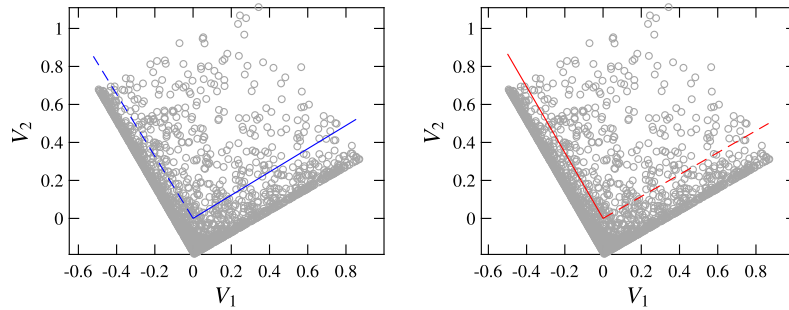


Fig. 4. Comparison between the principal kurtosis vectors shown in blue (left) and the PCA vectors shown in red (right) for the synthetic dataset from Fig. 3.

matrix M^y . The resulting singular vectors of the SVD (whether left or right singular vectors depends on how the tensor is matricized) are the desired principal vectors. We present tests of this method on synthetic datasets. All tests were conducted in Matlab and the tensor operations were performed using tensor toolbox [17]. Two bi-variate datasets are chosen and the main purpose is to illustrate that:

- the method recovers the principal vectors for datasets when they are known, by construction,
- compare the principal kurtosis vectors with PCA vectors to show that they need not be equal even for simple datasets.

For the first test a bi-variate dataset with known statistical joint moments is constructed in two steps. In the first step a bi-variate Gaussian copula with zero cross-correlation is sampled to ensure a statistically independent dataset. The inverse beta cumulative distribution function is applied over each copula generated vector such that the resulting dataset has beta marginal distributions, with known moments, and all the joint moments are approximately zero within sampling error. More importantly, the dataset is non-Gaussian and has non-zero higher order moments. The resulting dataset is shown in the left panel of Fig. 3. The $[\alpha, \beta]$ parameters of the beta distributions are chosen to be $[1.0, 0.1]$ for the first variable V_1 (x-axis) and $[1.0, 0.2]$ for the second variable V_2 (y-axis). These parameters were chosen carefully such that the excess kurtosis of V_1 is greater than that of V_2 , whereas the variance of V_2 is greater than that of V_1 . In the second step this independent dataset is “mixed” by performing an Euler rotation, and the resulting dataset now has non-zero joint moments. The first and second principal kurtosis vectors for this dataset are extracted using the described methodology and are verified to be equal to the vectors of the rotation matrix. These are shown in blue in the right panel of Fig. 3. Crucially, the parameters of this dataset were so chosen such that the principal directions of variance (PCA vectors) are different from that of the kurtosis i.e. the first PCA vector and the first kurtosis vector are orthogonal to each other, as are the second. This is illustrated in Fig. 4.

For the second test a dataset with a more visually obvious set of outliers is constructed. An uncorrelated bi-variate Gaussian dataset, with zero means and variances $[1.0, 0.25]$ is generated, and then mixed by rotation. Fig. 5 (a) shows the dataset, along with the PCA vectors shown in red. Since the data has Gaussian structure, the higher moments and their principal directions are ill-defined. Then, a small subset of samples in this dataset are chosen at random and placed far away from the rest at different locations, to mimic outliers, as shown in Fig. 5 (b) – (d). The PCA vectors (red) and principal kurtosis vectors (blue) for the data with outliers are also shown. In comparison to the part (a), it is evident that the PCA vectors remain nearly unchanged. However, the principal kurtosis vectors are different from the PCA vectors, and clearly align in the direction of outliers. The results from both test datasets illustrate that principal kurtosis vectors need not be the same as PCA vectors, and are more sensitive to the change in distribution due to the presence of outliers.

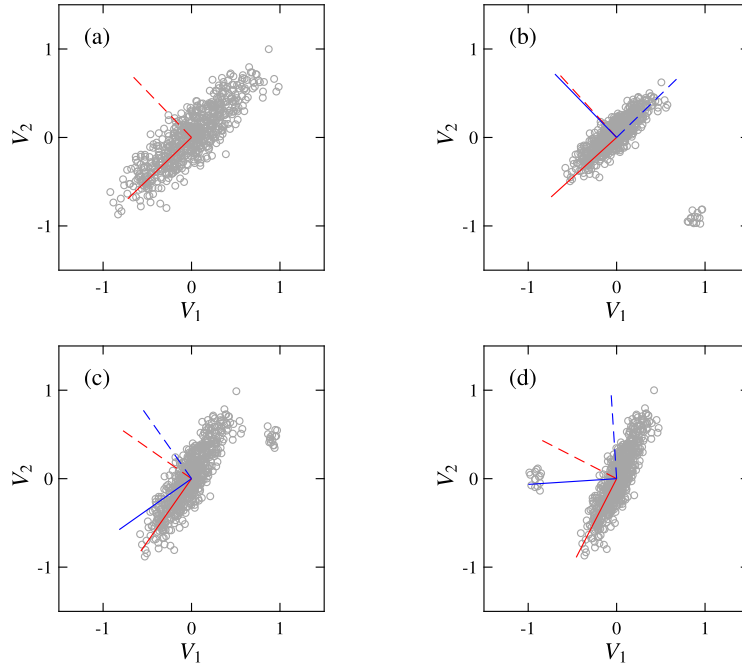


Fig. 5. (a) An uncorrelated bi-variate Gaussian dataset with zero means and variances $[1.0, 0.25]$ is mixed by rotation. (b) – (d) Gaussian dataset with a few samples converted to outliers in different locations. The first (solid line) and second (dashed line) PCA vectors are shown in red color. Similarly, the corresponding principal kurtosis vectors are shown in blue color.

3. Anomaly detection algorithm

In the previous section, we have shown that the variance or the kurtosis in data can be characterized in terms of the principal values and vectors of the joint moment tensors. Anomalies, as they occur, result in a change in the distribution of the data reflecting a change in the magnitude of the principal values and the orientation of the principal vectors. We base our anomaly detection algorithm for smoothly varying scientific data on this concept, as explained in this section.

As mentioned earlier, anomalous events in scientific phenomena can appear locally in space and/or time. We choose to decompose the data into several spatial sub-domains and time steps to explicitly flag the locality of these events. For distributed streaming data, e.g. data from massively parallel simulations, such a decomposition is inherently present by way of domain decomposition. Let N_d and N_t be the number of spatial sub-domains and the number of time steps, respectively. In regard to the feature space, depending on the nature of anomaly and guidance from the scientific domain, not all the features may be relevant to identify its occurrence. Let N_f be the number of relevant features to be used in the anomaly detection algorithm. The first step in the algorithm involves a data preprocessing stage, where we scale the data from each feature by subtracting its mean and dividing by the absolute spatial maximum. It is quite common in multi-scale simulations that value ranges of different features are decades apart. Hence, scaling the data ensures an equitable contribution to the joint moments from all the features.

Once the data is preprocessed, for each sub-domain j at a given time step n , the joint moment tensor $\mathcal{T}^{j,n}$ is constructed. This symmetric tensor is decomposed using the two-step (matricize and perform SVD) method described in section 2.3, to obtain the principal values λ_j and the principal vectors \hat{v}_j . For a given phenomena, in the absence of anomalous events, the principal values and vectors would remain nearly the same in all the sub-domains. As mentioned above, the occurrence of an anomalous event will result in a significant change in the magnitude of the principal values and/or the orientation of the principal vectors. We, now, define a feature moment metric $F_i^{j,n}$ for each feature i in a given sub-domain j and time step n , which can be used to quantify the changes in the principal values and vectors.

$$F_i^{j,n} = \frac{\sum_{k=1}^{N_f} \lambda_k (\hat{e}_i \cdot \hat{v}_k)^2}{\sum_{k=1}^{N_f} \lambda_k} \quad (9)$$

It should be noted that $\hat{e}_i \cdot \hat{v}_k$ is effectively the i -th entry in the k -th vector \hat{v}_k . Since the set of vectors \hat{v}_k are all unit vectors, by construction, the set of feature moment metrics in every spatial (j) and temporal (n) sub-domain sum to unity, i.e. $\sum_{i=1}^{N_f} F_i^{j,n} = 1, \forall j, n$. Accordingly, the moment metric for a given i can be interpreted as a measure of the fraction of

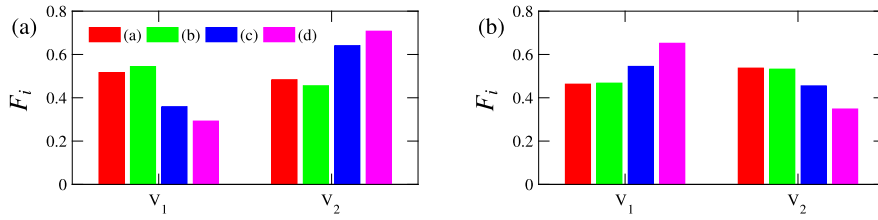


Fig. 6. Distribution of feature moment metrics (FMMs) for data in Fig. 5 obtained from the decomposition of (a) second moment co-variance matrix and (b) fourth moment co-kurtosis tensor. Different colors in the graphs are labeled according to the sub-figure labels in the Fig. 5.

the overall moment (variance or excess-kurtosis as case may be) contained in feature i , in other words a distribution of the moment in the feature space. To illustrate this further, the moment metrics for the data sets in Fig. 5 are shown in Fig. 6. The principal values and vectors from both the second moment (PCA) as well as the fourth moment (cumulant tensor) are used to compute the feature moment metrics. The data in 5 (a) appears to be aligned along the diagonal of the graph, which implies that the spread in the data has nearly equal contribution from both the features (V_1 and V_2). Accordingly, metrics (red bars in the figure) from both the moments are nearly equal. Note that this data set by construction does not possess an anomaly. For the other data sets (moving from (b) to (d) in Fig. 5), the position of data points that correspond to an anomalous event are more aligned, relative to the normal data, with feature V_1 . While the metrics from the second moment incorrectly show a greater contribution of V_2 to the moment, the metrics from the fourth moment accurately capture the greater contribution from V_1 towards the anomaly. This distribution interpretation allows comparing the feature moment metrics between different sub-domains in space and different steps in time, and flag the occurrence of anomalous events in space and/or time. For the first data set

If, as hypothesized, the statistical signature of anomalies is such that the distribution of feature moment metrics measurably changes, then distribution divergence metrics, such as f -divergence, can be used to quantify the change. We use the Hellinger distance, a symmetric measure of difference between two discrete distributions P and Q :

$$D_{PQ} = \frac{1}{\sqrt{2}} \sqrt{\sum_i (\sqrt{p_i} - \sqrt{q_i})^2}. \quad (10)$$

The Hellinger distance lies between 0 and 1, and for a discrete distribution the distance is 1 when the two distributions being compared are exact complements of each other i.e. if $\forall i$ when $p_i \neq 0$ and $q_i = 0$, and vice-versa. Intuitively, for the anomaly containing spatial/temporal sub-domain, the Hellinger distance of the $p_i \equiv F_i^{j,n}$ from a nominal distribution q_i would be large, and a suitable threshold of this distance can be used as an anomaly metric. The nominal set, q_i , can be chosen to be the spatial average such that the distance quantifies a spatial anomaly (at every time instance), whereas, setting q_i to be the previous time distribution quantifies a temporal anomaly (in each spatial sub-domain). Accordingly, we define a spatial anomaly metric

$$M_1^n(j) = \frac{1}{\sqrt{2}} \sqrt{\sum_{i=1}^{N_f} \left(\sqrt{F_i^{j,n}} - \sqrt{\bar{F}_i^n} \right)^2}, \quad (11)$$

where \bar{F}_i^n denotes the spatial (over j) average of $F_i^{j,n}$. A corresponding temporal anomaly metric can be defined as

$$M_2^j(n) = \frac{1}{\sqrt{2}} \sqrt{\sum_{i=1}^{N_f} \left(\sqrt{F_i^{j,n}} - \sqrt{F_i^{j,n-1}} \right)^2}. \quad (12)$$

We, now, proceed to describe the implementation of the anomaly detection algorithm, which is outline in Algorithm 1. The line 1 decomposes the data into different sub-domains and time steps. The relevant features for the algorithm are selected in line 2. Once these initialization steps are complete, the computations enter the time step loop in line 3. For each time step, the data is accessed with a sub-domain loop which begins in line 4. For a given time step and sub-domain, the data is scales in line 5 before computing the joint moment tensor in line 6. The tensor is then matricized in line 7 to perform SVD in line 8 to obtain the principal values λ_j and principal vectors \hat{v}_j . The feature moment metrics and the anomaly metrics are computed in line 9 and 10, respectively. After these computations are performed over all the sub-domains, the anomaly metrics are compared with a threshold value in line 12 to flag the occurrence of any anomalous event.

Algorithm 1: Anomaly detection algorithm.

```

// initialization
1  $N_t, N_d \leftarrow$  decompose data;
2  $N_f \leftarrow$  select features;
// time step loop
3 for  $n \leftarrow 1$  to  $N_t$  do
    // sub-domain loop
4     for  $j \leftarrow 1$  to  $N_d$  do
5         scale data;
6          $\mathcal{T}^{j,n} \leftarrow$  compute joint moment tensor;
7         matricize tensor  $\mathcal{T}^{j,n}$ ;
8          $\lambda_j, v_j \leftarrow$  perform SVD;
9          $F^{j,n} \leftarrow$  compute feature importance;
10         $M_1^n(j), M_2^n(j) \leftarrow$  compute anomaly metrics;
11    end
12    flag anomalous sub-domains;
13 end

```

4. Auto-ignition detection in combustion simulations

Auto-ignition is a spontaneous combustion process where a small temperature rise due to exothermic reactions results in the self-ignition of a fuel in the presence of an oxidizer. This phenomenon plays a critical role in the operation of devices such as the compression ignition engines and gas turbines [18–20]. The characteristics of auto-ignition strongly depend on the fuel-oxidizer composition, thermodynamic and flow conditions, among others. Hence, several experimental and simulation studies are devoted to understand the auto-ignition phenomena at conditions relevant to practical devices.

As the chemical phenomena resulting in auto-ignition possess an exponential behavior, the inception of ignition kernels (parcels of burning mixtures) is very sensitive to the initial conditions. Also, the ignition events are short-lived in time. In practical devices, where spatial inhomogeneities persist, e.g. due to turbulence and mixing, these events are often highly localized in space. This makes it challenging to detect and analyze the phenomena [21]. In the combustion community, researchers often use simple threshold based metrics which are ad hoc in nature, or mathematically rigorous formulations such as chemical mode explosive analysis (CEMA) [22] which are computationally expensive, to detect the auto-ignition events in complex turbulent reacting flows. We will use the anomaly detection algorithm described in section 3 to identify the occurrence of auto-ignition events in two representative combustion simulations: (a) a canonical spatially one-dimensional (1D) simulation of a turbulent auto-ignitive mixture, and (b) a spatially two-dimensional (2D) simulation of turbulent auto-ignition representing compression ignition conditions. Both simulations are performed using a direct numerical simulation solver, S3D [23], which solves the reacting compressible flow governing equations. The governing partial differential equations are approximated using explicit eight-order central difference schemes for spatial derivatives and six-stage fourth-order Runge-Kutta method for time integration. Explicit tenth-order filters are used to remove spurious numerical oscillations.

4.1. Canonical 1D simulation

A 1D time-varying simulation of an auto-ignitive mixture of syngas is performed using a 12-species 29-reactions mechanism [24]. The nuances of the anomaly detection method are first illustrated in this canonical simulation before being deployed in the more realistic two-dimensional simulation in the following section. The spatial domain is 1 cm long and is discretized using 1024 grid points. The domain is initialized with premixed fuel-oxidizer mixture comprised of species concentrations $0.6\text{CO} + 0.4\text{H}_2 + 0.5(\text{O}_2 + 3.76\text{N}_2)$. The initial pressure is set to 4 atm and the flow is quiescent. A spatially inhomogeneous temperature field with a mean of 1200 K is imposed using a linear combination of sinusoidal waves with random phases. In addition, a spike in the temperature is superimposed in the first quarter of the domain, which ensures that this portion ignites first, resulting in a spatially anomalous ignition event. The initial temperature profile is shown in Fig. 7. A periodic boundary condition is imposed for solving the equations. The simulation is time advanced up to 20 μs using fixed time-steps of 0.001 μs . The data from the simulation is stored at equal intervals of 1 μs , which is one-tenth of the ignition delay time based on a homogeneous mixture at a temperature of 1200 K. The time evolution of the temperature profiles are shown in Fig. 8. The contour plots of temperature and mass fractions of HO_2 , H and O are also shown in Fig. 9 to illustrate production and consumption of key radicals in the auto-ignition process. Initially, the temperature gradients reduce as a consequence of the diffusion process, resulting in a decrease in the peak value. The diffusion process dominates until a thermal runaway is set off due to a minor consumption of the fuel and production of HO_2 in exothermic reactions. This is first observed in the higher temperature regions which have greater propensity to trigger reactions. The thermal runaway then leads to an inception of an ignition kernel where oxidation of the fuel occurs with a rapid temperature rise and heat release. Several intermediate species such as H and O are produced and consumed during the oxidation process. Eventually, regions with lower initial temperature also ignite and the combustion process completes in the entire domain.

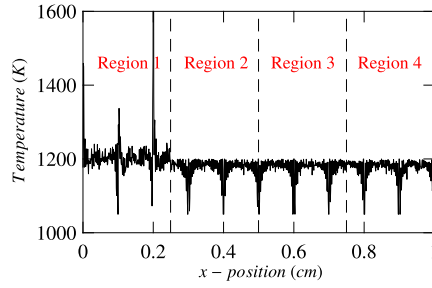


Fig. 7. Profile of the temperature initial condition. Dashed lines represent the sub-domain boundaries of the decomposed spatial domain.

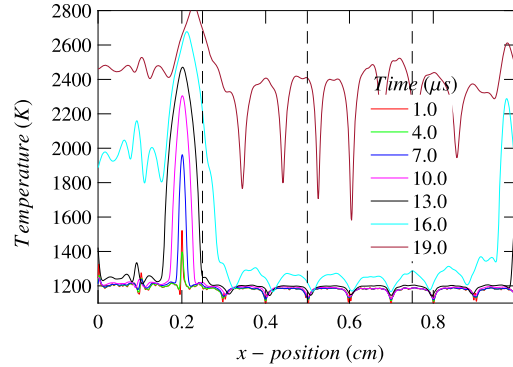


Fig. 8. Time evolution of the temperature profiles. Dashed lines represent the sub-domain boundaries of the decomposed spatial domain.

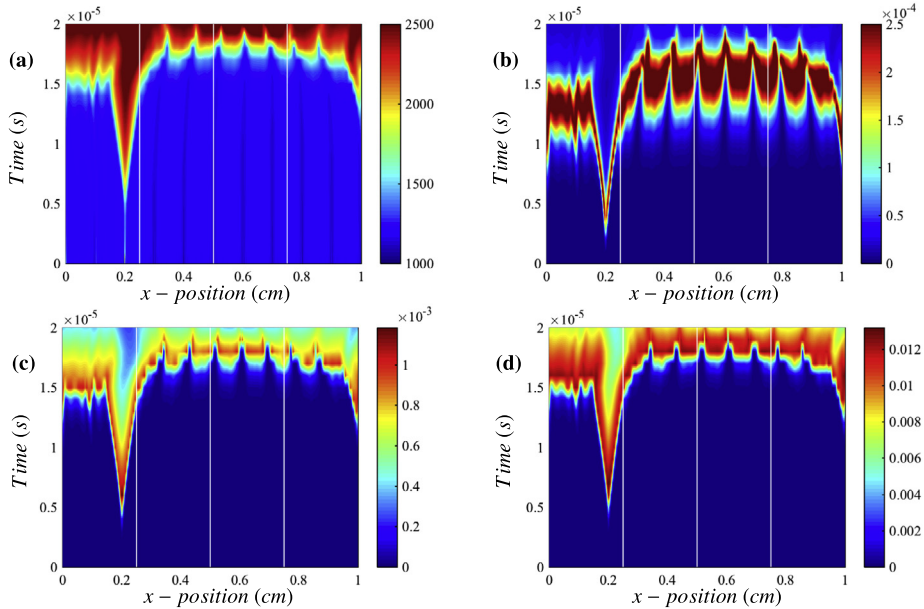


Fig. 9. Contour plot of (a) temperature (K), (b) mass fraction of HO_2 , (c) H and (d) O, illustrating time evolution of the solution in 1D domain. White vertical lines represent sub-domain boundaries.

To apply the anomaly detection algorithm, outlined in Algorithm 1, we decompose the spatial domain into four sub-domains ($N_d = 4$), which are labeled Regions 1–4, as shown in Fig. 7. In time, the dataset consists of 21 save files, which makes the number of time steps $N_t = 21$. A total of 17 features are stored in the datasets, of which 13 are the relevant features for the algorithm (12 species mass fractions and temperature). The data are pre-processed according to the scaling described in section 3. Fig. 10 shows scatter plots of scaled mass fraction of H_2 and temperature for Region 1 at four different times. The solid and dashed lines in the plots represent the first and second principal kurtosis vectors, respectively. Note that for this illustration, the joint moment tensor is constructed only for these two features, which results only in

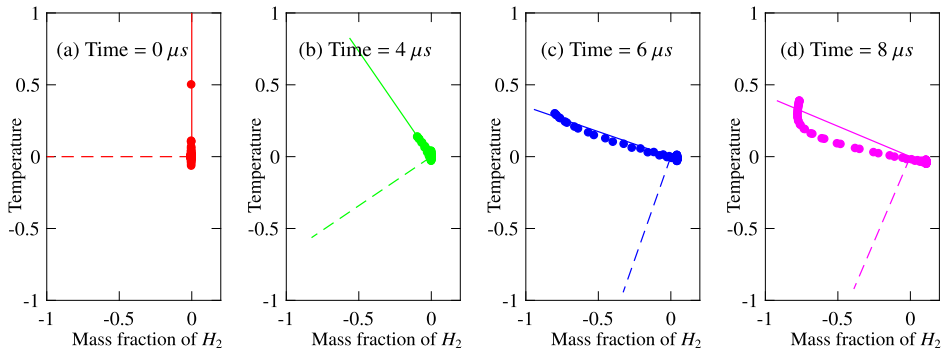


Fig. 10. Scatter plots of scaled mass fraction of H_2 and temperature in Region 1 at different times. Solid line: first principal vector, dashed line: second principal vector.

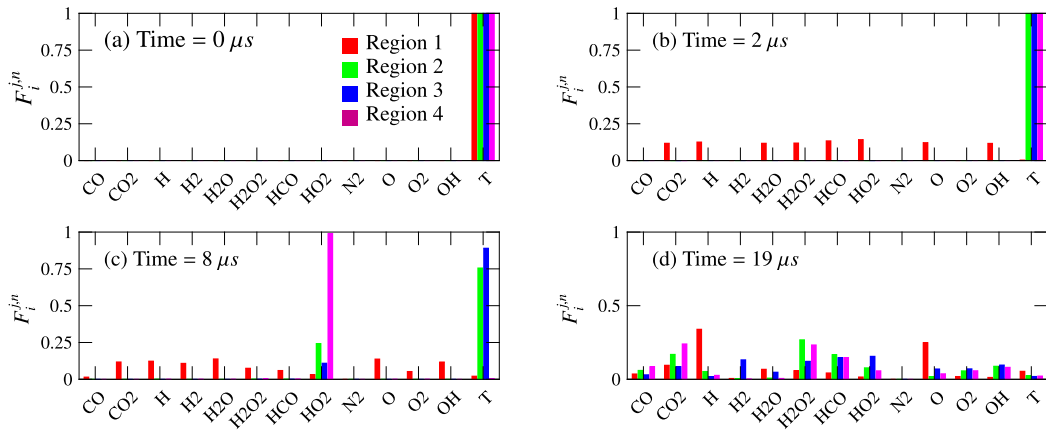


Fig. 11. Distribution of feature moment metrics (FMMs). Different colors correspond to different sub-domains.

two principal vectors. At initial condition the inhomogeneity in the data is present in temperature alone. Hence, the first principal vector is aligned along the temperature axis. As time advances, an ignition kernel appears in the peak temperature region, which leads to the consumption of the H_2 -fuel and an increase in the temperature. This is depicted by some points migrating along the negative x-axis and positive y-axis. Clearly, the orientation of the principal vectors has also changed with time and the first principal vector aligns along the extreme points which represent the ignition event.

The feature moment metrics (FMMs), $F_i^{j,n}$, which quantify the relative importance of each feature towards the kurtosis in the data, are plotted at different times in Fig. 11. The histograms are colored separately for the different spatial sub-domains. The x-axis consists the labels of the 13 relevant features, which include the 12 chemical species and temperature (T). At the initial time, as temperature is the only feature having any spread in all sub-domains, the temperature-FMM is equal to 1 for all of them. As time advances to $2 \mu s$ (part (b)), the FMM in temperature disappears in Region 1, and spreads to the other features, signifying a change in the FMM distribution. This is attributed to an early auto-ignition in the Region 1 due to the significantly greater peak temperature, which leads to locally active chemical reactions and an inhomogeneity in the mass fractions of species in the sub-domain. At the later time ($8 \mu s$), ignition also begins in the other three sub-domains accompanied by a similar decrease in temperature-FMM. A careful look at the temperature profiles in Fig. 8 will indicate that, relatively, Region 4 has the second highest temperature, followed by Region 2 and Region 3. Hence, ignition events as well as the corresponding decrease in the temperature-FMM, in these three sub-domains, will also appear in the same order in time. Eventually, ignition occurs in all the sub-domains, leading to the non-zero FMMs appearing across different features, as seen at time $19 \mu s$. The FMM distributions in Fig. 11 confirm the hypothesis that at the inception of auto-ignition event the FMMs, which were initially significant only for temperature, begin to spread among other features, which is the change in the statistical signature anticipated.

An anomalous auto-ignition event, which results in a significant change in the FMM distribution, can be detected based on the Hellinger distance, as described in section 3. The spatial and temporal anomalies are measured in terms of the distance metrics $M_1^n(j)$ and $M_2^j(n)$, respectively. Since the metrics lie between 0 and 1 we choose a threshold of 0.5 for these metrics to identify an anomalous event. A sub-domain j at a given time step n can be flagged to contain an anomaly if maximum of $M_1^n(j)$ and $M_2^j(n)$ is greater than 0.5. Fig. 12 shows the evolution of the $M_1^n(j)$ and $M_2^j(n)$ anomaly metrics, as well the greater of the two, with time. At earlier times ($< 5 \mu s$), an ignition kernel appears in the Region 1, as mentioned

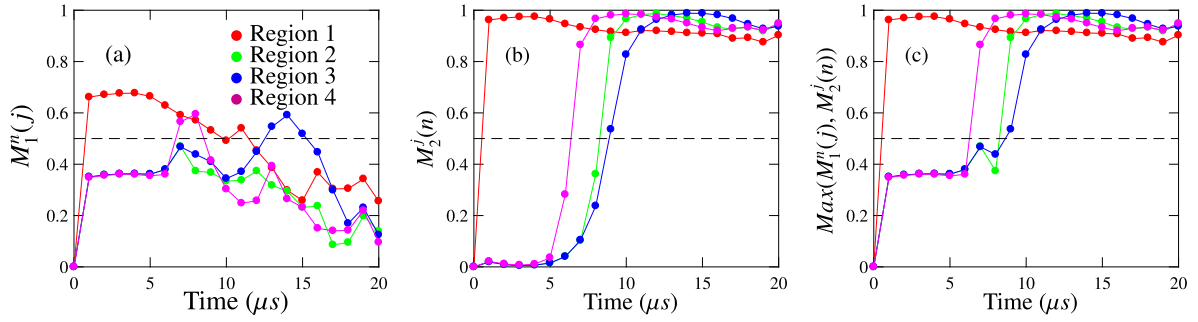


Fig. 12. Time evolution of (a) spatial anomaly metric, (b) temporal anomaly metric, and (c) their maximum anomaly metrics. Different colors correspond to different sub-domains. Dashed line represents a constant threshold value of 0.5.

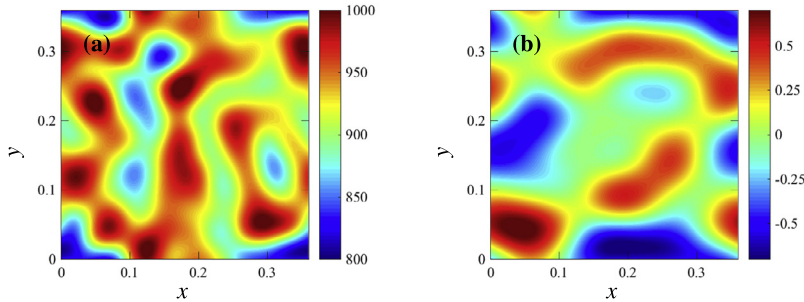


Fig. 13. Initial conditions of the HCCI simulation: contour plot of temperature in K (left) and x-component velocity in m/s (right).

earlier, due to a significantly larger value of temperature. This is an anomaly both in space as well as time, because the distribution of FMMs is distinctively different from (compare parts (a) and (b) in Fig. 11) other sub-domains and also varies quickly in time. Hence, both $M_1^n(j)$ and $M_2^j(n)$ shoot up above the threshold value. In this time period, the metrics are nearly the same and below the threshold in other sub-domains. As time progresses, beyond 5 μs , the temporal anomaly metric goes above the threshold in other sub-domains. For these sub-domains, the occurrence of ignition event is an anomaly in time, but not in space, and hence $M_1^n(j)$ is observed to be consistently lower than $M_2^j(n)$ in Fig. 12. The order in which the events appear i.e. earliest in Region 1, then in Region 4 and then in Regions 2 and 3, is also in agreement with the observed change in the FMM distribution in Fig. 11. It appears that the metrics, and the suggested threshold, are consistent with the expected evolution of the auto-ignitive system, and are able to flag the occurrence of auto-ignition in the right sub-domains and at the right times.

4.2. 2D simulation of compression ignition

We next demonstrate the anomaly detection method on a 2D simulation representing conditions of homogeneous charge compression ignition (HCCI) of ethanol. The simulation is part of a parametric study investigating the influence of initial temperature and composition stratification on ignition timings under realistic internal combustion (IC) engine conditions [25]. All the simulations in the parametric study have the occurrence of ignition events, but the precise time and location of these events is difficult to ascertain, which can be accomplished by our anomaly detection algorithm. We demonstrate this for the chosen simulation, described next. A premixture of ethanol and air mixed with combustion products of the same mixture, representing “exhaust gas recirculation” (EGR), is initialized in a 2D spatial domain with periodic boundaries. The nominal initial conditions are pressure of 45 atm, equivalence ratio of 0.4 and mean temperature of 924 K. Temperature inhomogeneity, with an r.m.s of 40 K, is superimposed over the mean value to mimic uneven mixing by EGR. A divergence-free turbulent velocity field, with an r.m.s of 0.61 m/s, is specified from an auxiliary simulation of homogeneous isotropic decaying turbulence [25]. Compression heating, pertaining to piston motion in an IC engine cylinder, is also imposed on the mixture and the simulation is performed for a total duration spanning 45 crank angle degrees of the piston movement. Fig. 13 shows the initial conditions: contour plots of the initial temperature illustrating the inhomogeneity, and initial x-component velocity showing the turbulent flow field. All other quantities, except y-component velocity, are spatially uniform at initialization. The simulation was performed in a domain of $0.36 \text{ cm} \times 0.36 \text{ cm}$ with a uniform grid of 672×672 grid points, decomposed into a set of 12×12 sub-domains for MPI parallelism i.e. each sub-domain was a square subset of the domain comprising 56×56 grid points. Accordingly, the algorithm was employed in each of these 144 sub-domains/regions. The mechanism used for ethanol combustion comprised of 28 chemical species, which are the key features for constructing the principal kurtosis vectors and moment metrics. For the spatio-temporal analysis, 201 time instants of the solution

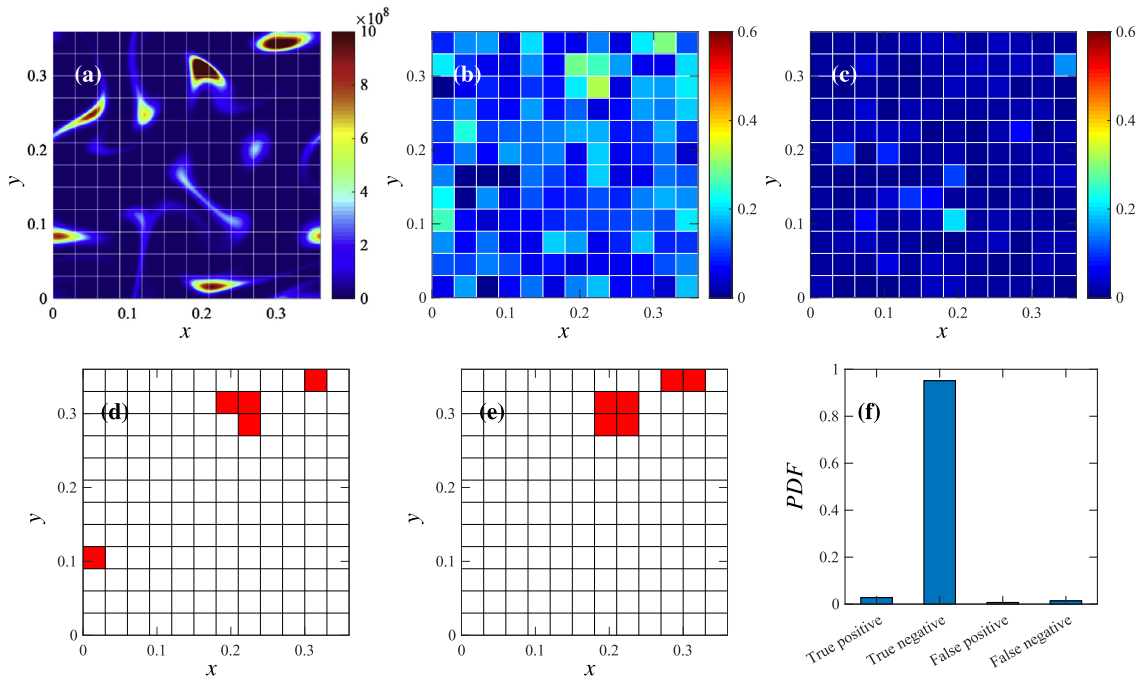


Fig. 14. Instantaneous plots of (a) heat release rate, (b) spatial anomaly metric ($M_1^n(j)$), (c) temporal anomaly metric ($M_2^j(n)$), (d) anomalous event map from proposed method, (e) anomalous event map from heat release rate threshold, and (f) anomaly detection performance from 2D simulation of compression ignition at $t = 0.845$ ms. Tiles in parts (a)–(e) represent the sub-domains/regions.

equi-spaced at intervals of $1.5 \mu\text{s}$ are considered. From a combustion perspective, the total heat release rate, and not any particular species concentration, maybe a better indicator of exothermic chemical reaction events such as ignition kernels or flame fronts. Accordingly, we use heat release rate as the “domain guided” true marker of ignition, and use it to evaluate the prediction accuracy the proposed anomaly detection method.

Contour plots from the earliest time instant when ignition kernels appear are shown in Fig. 14. The contours of heat release rate, Fig. 14 (a), show the presence of a few ignition kernels that occupy a fraction of the domain. For the conditions of this simulation, a heat release rate of $1 \times 10^9 \text{ J/m}^3/\text{s}$ is an appropriate threshold to denote ignition (see Fig. 5b of [25]). Among the kernels in Fig. 14 (a), two of them appear with dark red contours and are above this threshold. The values the M_1 and M_2 anomaly metrics in different subdomains are shown in Figs. 14 (b) and (c). Clearly, the spatial metric (M_1) has high values in the regions with the ignition kernels, whereas the temporal metric (M_2) appears low. This can be attributed to the finely spaced time steps in the evaluation of $M_2^j(n)$, where the changes in FMM distribution are gradual, and this is precisely why a joint criterion based on both spatial and temporal anomaly metrics is suggested. For this 2D compression ignition simulation, we identify a threshold value of 0.25 as appropriate to detect anomalous events, i.e., a region is anomalous if $M_1^n(j)$ or $M_2^j(n)$ is greater than 0.25. Using this threshold the identified anomalous sub-domains containing the ignition kernels are shown in the event map in Fig. 14 (d) (red marks the regions containing the anomaly). The map in Fig. 14 (d) compares reasonably well to the *true* ignition event map computed using the heat release rate threshold in Fig. 14 (e). Comparing the two event maps, each region (a tile in the grid) can be assigned an attribute: a combination of “true”/“false” and “positive”/“negative”. A histogram of the four resulting attributes, over all regions considered at this instant, is shown in Fig. 14 (f). As evident from the histogram, the regions with “false” predictions are insignificant compared to those with “true” predictions.

We next show, in Fig. 15, contours of heat release rate, predicted anomalous event map, true event map and the prediction histograms at three different time instants representative of various combustion stages in the HCCI simulation: first (top row) at the inception of auto-ignition in the domain (same as in Fig. 14), second (middle row) at an instant when auto-ignition has occurred in roughly half of the domain, and third (bottom row) at a time when combustion has completed in most of the domain and only three flame kernels remain. The kernels are physically anomalous at the first and last time instants, as they occupy only a fraction of the domain. Accordingly, a better performance is observed at these instants. At the second instant, where ignition kernels are present in half of the domain, ignition occurrence is not anomalous and hence negative predictions are dominant, many of them being false negatives.

To compare the performance of the proposed anomaly detection method against an existing algorithm, we choose the well known local outlier factor (LOF) method [26]. The LOF method compares the “reachability density” of each data sample with the average reachability density of its neighbors and assigns a score. Outliers are characterized by a LOF score greater than 1. For the distributed scientific data setting, we reduce the data in each sub-domain using a maximum or a minimum

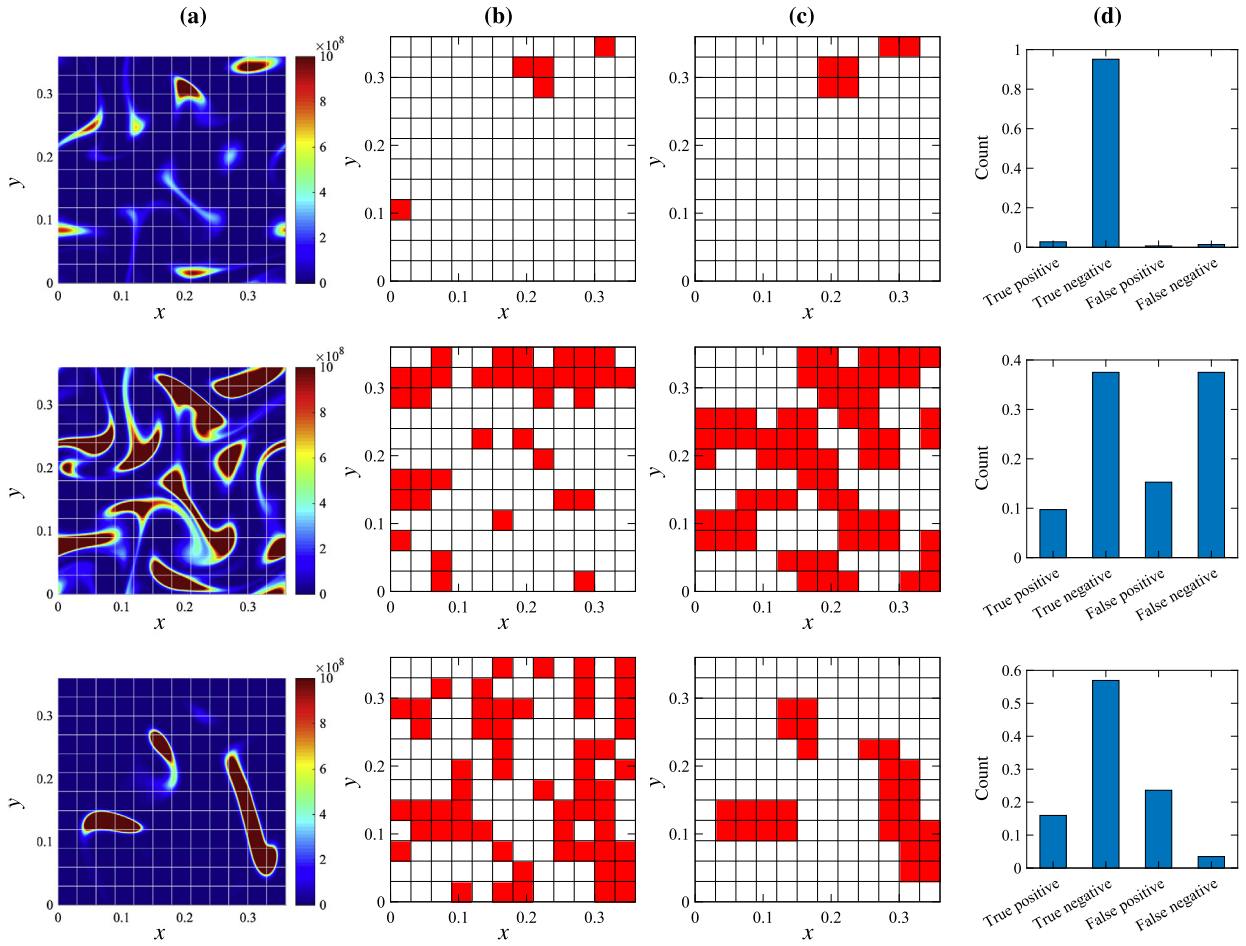


Fig. 15. Plots of (a) heat release rate, (b) predicted anomalous event map based on proposed method, (c) true anomalous event map based on heat release rate threshold, and (d) anomaly detection performance from 2D simulation of compression ignition at three different time instants. Top row: $t = 0.845$ ms, middle row: $t = 1.2$ ms, and bottom row $t = 3$ ms. Tiles in parts (a)–(c) represent the sub-domains/regions.

or an average for each feature, representing the feature signature similar to the feature moment metric. We then employ the LOF method on the feature signatures from all the sub-domains to identify the anomalous sub-domains. The algorithm is implemented using the widely used Scikit-learn Python package [27]. We first consider the subset labeled as anomalous by Scikit-learn's `fit_predict` function. We then take the LOF scores of this subset (returned by the `negative_outlier_factor` attribute) and normalize these scores to a 0-1 range, to apply a threshold for decisions. Results from this evaluation for the same time instants used in Fig. 15 at two different threshold values (0.01 and 0.5) are shown in Fig. 16. Similar to the procedure followed earlier, we use the true ignition event map to compute the distribution of the combination of “true”/“false” and “positive”/“negative” attributes. Although the LOF method is able to capture the true positives well, its true negative prediction is poor, resulting in a large number of false positives. The results, as evident from the figure, are insensitive to the threshold used for the LOF method. In other words, the method falsely predicts most sub-domains to be anomalous or consisting of an ignition event. Clearly, for scientific data the proposed anomaly detection method outperforms the popular LOF method.

The heat release rate threshold, used to denote true ignition events and compare predicted events against, is arguably arbitrary although it is informed by the conditions expected in the simulation. Likewise, the M_1/M_2 metric threshold of 0.25 used for the prediction is a hyper parameter. In order to evaluate the effect of these threshold values, we plot the receiver operating characteristic (ROC) curves in Fig. 17, which show the diagnostic ability of a binary classifier system as its thresholds are varied. The graph consists of false positive rate or (1-specificity) on x-axis and true positive rate or sensitivity on y-axis. The specificity and sensitivity are the ratios of true positives over total positives and true negatives over total negatives, respectively. For reference we show a dashed line corresponding to the performance of a random guess for a binary classification. The graph shows the ROC curves for three different threshold values of the heat release rate. Each curve is obtained by varying the anomaly metric threshold between 0 and 1. All the curves appear above the dashed line, indicating that the performance of the proposed algorithm is considerably better than a random guess, as is desirable. As the heat release rate threshold is increased, the curve become flatter, possessing a greater true positive rate and a lower

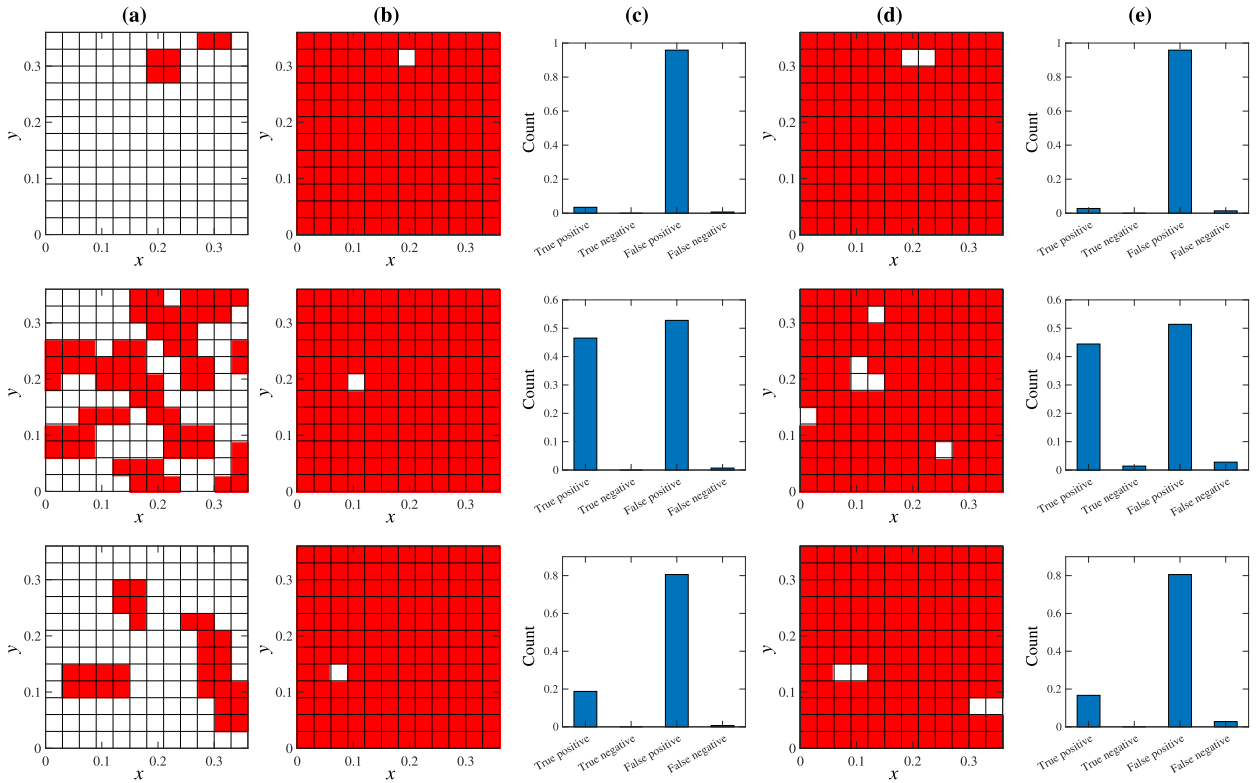


Fig. 16. Plots of (a) true anomalous event map based on heat release rate threshold, (b) predicted anomalous event map based on normalized local outlier factor (LOF) scores with $threshold = 0.01$, (c) LOF performance with $threshold = 0.01$, (d) predicted anomalous event map based on LOF method with $threshold = 0.5$, and (e) LOF performance with $threshold = 0.5$, obtained from 2D simulation of compression ignition at three different time instants. Top row: $t = 0.845$ ms, middle row: $t = 1.2$ ms, and bottom row $t = 3$ ms. Tiles in parts (a), (b) and (d) represent the sub-domains/regions.

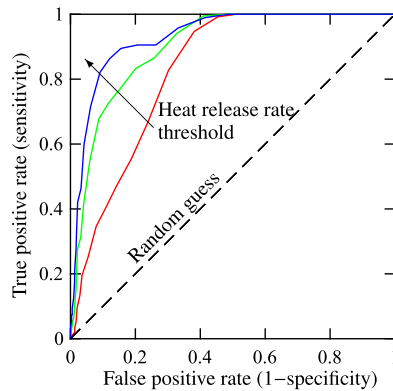


Fig. 17. Receiver operating characteristic (ROC) curves indicating the performance of the anomaly detection algorithm under varying heat release rate and anomaly metric thresholds. The values of heat release rate thresholds are 1×10^9 J/m³/s (red curve), 5×10^9 J/m³/s (green curve) and 1×10^{10} J/m³/s (blue curve).

false positive rate. In terms of the physics, a greater heat release rate occurs over a relatively narrower region in space, making the ignition kernels more localized and anomalous. Hence, a better prediction is observed for greater heat release rate thresholds.

5. Conclusions

Detection of anomalous events in scientific phenomena remains a key challenge, particularly, due to the availability of increasingly high resolution data which is enabled by the advances in measurement techniques and computing power. Scientific datasets are characterized by smoothly varying multi-variate data which represent complex non-linear multi-physics

processes. Hence, commonly used anomaly detection algorithms developed for the use in other domains may not be readily applicable for scientific data. It is well known that occurrence of anomalies manifest as extreme values in the distribution of at least some of the features, and significantly affect their higher order joint statistical moments. In this paper, we have used this idea to develop a robust anomaly detection algorithm for scientific data.

For normally distributed data, the statistical information is fully encapsulated in the second order co-variance matrix. However, scientific data are often non-Gaussian and need further higher order joint moment tensors to characterize the anomalous events. By analogy to principal component analysis, a general joint moment tensor can be analyzed in terms of its principal values and vectors which can be computed using symmetric tensor decomposition methods. We have reviewed commonly used methods and identified that a simple singular value decomposition of the matricized tensor works best for our purpose. Using synthetic data, we have shown that vectors computed from the decomposition of cumulant fourth moment tensor capture the outlier data perfectly.

In general, anomalous events in scientific phenomena appear in space and/or time. To identify the locality of an anomaly, the proposed algorithm decomposes the data into spatial sub-domains and time steps. Feature moment metrics are computed for each sub-domain and time step to quantify the magnitude and alignment of principal values and vectors, respectively. With an inception of an anomaly the distribution of the feature moment metrics significantly changes. The algorithm, then, uses Hellinger distance to compare the feature moment metrics between different sub-domains and successive time steps to flag anomalies in space and time, respectively. The algorithm has been tested using a turbulent combustion test cases to detect auto-ignition events.

In comparison to the widely used local outlier factor method, the unsupervised anomaly detection algorithm presented in this paper has been shown to robustly identify the spatial and temporal anomalies. The statistical approach in the algorithm uses a decomposed layout of the data which is inherent to large streaming distributed scientific datasets. Our future work includes two aspects. First, an in-situ implementation of the algorithm into the massively parallel direct numerical simulation solver (S3D) and evaluate its scalability. Second, to apply the algorithm to detect anomalies in other scientific phenomena.

Acknowledgements

We would like to acknowledge helpful discussions with, and suggestions from, Dr. Jacqueline H. Chen and Dr. Tamara G. Kolda at Sandia National Laboratories. This work was funded through U.S. Department of Energy Advanced Scientific Computing Research (ASCR) grant FWP16-019471. Sandia National Laboratories is a multi-mission laboratory managed and operated by National Technology and Engineering Solutions of Sandia, LLC., a wholly owned subsidiary of Honeywell International, Inc., for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-NA-0003525. The views expressed in the article do not necessarily represent the views of the U.S. Department of Energy or the United States Government.

References

- [1] V. Chandola, A. Banerjee, V. Kumar, Anomaly detection: a survey, *ACM Comput. Surv.* 41 (2009) 15:1–15:58.
- [2] W.P. Kegelmeyer, R. Calderbank, T. Critchlow, L. Jameson, C. Kamath, J. Meza, N. Samatova, A. Wilson, Mathematics for Analysis of Petascale Data, Report on a Department of Energy Workshop, Technical Report 2007-2008-4349P, Sandia National Laboratories, 2008.
- [3] J. Ling, W.P. Kegelmeyer, K. Aditya, H. Kolla, K.A. Reed, T.M. Shead, W.L. Davis, Using feature importance metrics to detect events of interest in scientific computing applications, in: 2017 IEEE 7th Symposium on Large Data Analysis and Visualization (LDAV), IEEE, pp. 55–63.
- [4] P. Pébay, T.B. Terriberry, H. Kolla, J. Bennett, Numerically stable, scalable formulas for parallel and online computation of higher-order multivariate central moments with arbitrary weights, *Comput. Stat.* 31 (2016) 1305–1325.
- [5] F.E. Grubbs, Procedures for detecting outlying observations in samples, *Technometrics* 11 (1969) 1–21.
- [6] G.O. Campos, A. Zimek, J. Sander, R.J.G.B. Campello, B. Micenkova, E. Schubert, I. Assent, M.E. Houle, On the evaluation of unsupervised outlier detection: measures, datasets, and an empirical study, *Data Min. Knowl. Discov.* 30 (2016) 891–927.
- [7] M. Goldstein, S. Uchida, A comparative evaluation of unsupervised anomaly detection algorithms for multivariate data, *PLoS ONE* 11 (2016) 1–31.
- [8] P.H. Westfall, Kurtosis as Peakedness, 1905–2014, *R.I.P.*, *Am. Stat.* 68 (2014) 191–195.
- [9] E. Jondeau, E. Jurczenko, M. Rockinger, Moment Component Analysis: An Illustration with International Stock Markets, Research Paper 10-43, Swiss Finance Institute, 2015.
- [10] T.G. Kolda, B.W. Bader, Tensor decompositions and applications, *SIAM Rev.* 51 (2009) 455–500.
- [11] P. Comon, G. Golub, L.-H. Lim, B. Mourrain, Symmetric tensors and symmetric tensor rank, *SIAM J. Matrix Anal. Appl.* 30 (2008) 1254–1279.
- [12] T.G. Kolda, Symmetric Orthogonal Tensor Decomposition is Trivial, 2015.
- [13] A. Anandkumar, R. Ge, D. Hsu, S.M. Kakade, M. Telgarsky, Tensor decompositions for learning latent variable models, *J. Mach. Learn. Res.* 15 (2014) 2773–2832.
- [14] L.D. Lathauwer, B.D. Moor, J. Vandewalle, A multilinear singular value decomposition, *SIAM J. Matrix Anal. Appl.* 21 (2000) 1253–1278.
- [15] L.D. Lathauwer, B.D. Moor, Independent component analysis and (simultaneous) third-order tensor diagonalization, *IEEE Trans. Signal Process.* 49 (2001) 2262–2271.
- [16] P. Comon, Tensor decompositions, in: J.G. McWhirter, I.K. Proudler (Eds.), *Mathematics and Signal Processing V*, Clarendon Press, Oxford, UK, 2002, pp. 1–24.
- [17] B.W. Bader, T.G. Kolda, et al., Matlab tensor toolbox version 2.6, 2015, Available online.
- [18] J.E. Dec, Advanced compression-ignition engines—understanding the in-cylinder processes, *Proc. Combust. Inst.* 32 (2009) 2727–2742.
- [19] F. Güthe, J. Hellat, P. Flohr, The reheat concept: the proven pathway to ultra-low emissions and high efficiency and flexibility, *J. Eng. Gas Turbines Power* 131 (2009) 021503.
- [20] K. Aditya, A. Gruber, C. Xu, T. Lu, A. Krisman, M.R. Bothien, J.H. Chen, Direct numerical simulation of flame stabilization assisted by autoignition in a reheat gas turbine combustor, *Proc. Combust. Inst.* 37 (2019) 2635–2642.

- [21] J. Bennett, A. Bhagatwala, J. Chen, A. Pinar, M. Slalom, C. Seshadhri, Trigger detection for adaptive scientific workflows using percentile sampling, *SIAM J. Sci. Comput.* 38 (2016) 240–260.
- [22] T.F. Lu, C.S. Yoo, J.H. Chen, C.K. Law, Three-dimensional direct numerical simulation of a turbulent lifted hydrogen jet flame in heated coflow: a chemical explosive mode analysis, *J. Fluid Mech.* 652 (2010) 45–64.
- [23] 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, et al., Terascale direct numerical simulations of turbulent combustion using s3d, *Comput. Sci. Discov.* 2 (2009) 015001.
- [24] E.R. Hawkes, R. Sankaran, J.C. Sutherland, J.H. Chen, Scalar mixing in direct numerical simulations of temporally evolving plane jet flames with skeletal co/h 2 kinetics, *Proc. Combust. Inst.* 31 (2007) 1633–1640.
- [25] A. Bhagatwala, J.H. Chen, T. Lu, Direct numerical simulations of hcci/saci with ethanol, *Combust. Flame* 161 (2014) 1826–1841.
- [26] M.M. Breunig, H.-P. Kriegel, R.T. Ng, J. Sander, Lof: identifying density-based local outliers, in: *ACM Sigmod Record*, vol. 29, ACM, 2000, pp. 93–104.
- [27] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, et al., Scikit-learn: machine learning in python, *J. Mach. Learn. Res.* 12 (2011) 2825–2830.