

Topology-enhanced machine learning for consonant recognition

Yifei Zhu

zhu yf@sustech.edu.cn

Southern University of Science and Technology <https://orcid.org/0000-0001-8918-1896>

Pingyao Feng

Southern University of Science and Technology

Siheng Yi

Southern University of Science and Technology

Qingrui Qu

Southern University of Science and Technology

Zhiwang Yu

Southern University of Science and Technology

Article

Keywords:

Posted Date: March 12th, 2024

DOI: <https://doi.org/10.21203/rs.3.rs-3978261/v1>

License:   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Additional Declarations: There is **NO** Competing Interest.

Topology-enhanced machine learning for consonant recognition

Pingyao Feng , Siheng Yi, Qingrui Qu, Zhiwang Yu, Yifei Zhu 



Abstract—In artificial-intelligence-aided signal processing, existing deep learning models often exhibit a black-box structure. The integration of topological methods serves a dual purpose of making models more interpretable as well as extracting structural information from time-dependent data for smarter learning. Here, we provide a transparent and broadly applicable methodology, TopCap, to capture topological features inherent in time series for machine learning. Rooted in high-dimensional ambient spaces, TopCap is capable of capturing features rarely detected in datasets with low intrinsic dimensionality. Compared to prior approaches, we obtain descriptors which probe finer information such as the vibration of a time series. This information is then vectorised and fed to multiple machine learning algorithms. Notably, in classifying voiced and voiceless consonants, TopCap achieves an accuracy exceeding 96%, significantly outperforming traditional convolutional neural networks in both accuracy and efficiency, and is geared towards designing topologically enhanced convolutional layers for deep learning speech and audio signals.

1 INTRODUCTION

IN 1966, Mark Kac asked the famous question: “Can you hear the shape of a drum?” To hear the shape of a drum is to infer information about the shape of the drumhead from the sound it makes, using mathematical theory. In this article, we venture to flip and mirror the question across senses and address instead: “Can we see the sound of a human speech?”

The artificial intelligence (AI) advancements have led to a widespread adoption of voice recognition technologies, encompassing applications such as speech-to-text conversion and music generation. The rise of topological data analysis (TDA) [1] has integrated topological methods into many areas including AI [2, 3], which makes neural networks more interpretable and efficient, with a focus on structural information. In the field of voice recognition [4, 5], more specifically consonant recognition [6, 7, 8, 9, 10], prevalent methodologies frequently revolve around the analysis of energy and spectral information. While topological approaches are still rare in this area, we combine TDA and machine learning to obtain a classification for speech data, based on geometric patterns hidden within phonetic segments. The method we propose, TopCap (referring to capturing topological structures of data), is not only applicable to audio data but also to general-purpose time series data that require extraction of structural information for machine learning algorithms. Initially, we endow

phonetic time series with point-cloud structure in a high-dimensional Euclidean space via time-delay embedding (TDE, see Fig. 1a) with appropriate choices of parameters. Subsequently, 1-dimensional persistence diagrams are computed using persistent homology (see Sec. S.2.2 for an explanation of the terminologies). We then conduct evaluations with nine machine learning algorithms, in comparison with a convolutional neural network (CNN) without topological inputs, to demonstrate the significant capabilities of TopCap in the desired classification.

Conceptually, TDA is an approach that examines data structure through the lens of topology. This discipline was originally formulated to investigate the *shape* of data, particularly point-cloud data in high-dimensional spaces [11]. Characterised by a unique insensitivity to metrics, robustness against noise, invariance under continuous deformation, and coordinate-free computation [1], TDA has been combined with machine learning algorithms to uncover intricate and concealed information within datasets [12, 3, 13, 14, 15, 16]. In these contexts, topological methods have been employed to extract structural information from the dataset, thereby enhancing the efficiency of the original algorithms. Notably, TDA excels in identifying patterns such as clusters, loops, and voids in data, establishing it as a burgeoning tool in the realm of data analysis [17]. Despite being a nascent field of study, with its distinctive emphasis on the shape of data, TDA has led to novel applications in various far-reaching fields, as evidenced in the literature. These include image recognition [18, 19, 20], time series forecasting [21] and classification [22], brain activity monitoring [23, 24], protein structural analysis [25, 26], speech recognition [27], signal processing [28, 29], neural networks [30, 31, 32, 2], among others. It is anticipated that further development of TDA will pave a new direction to enhance numerous aspects of daily life.

The task of extracting features that pertain to structural information is both intriguing and formidable. This process is integral to a multitude of practical applications [33, 34, 35, 36], as scholars strive to identify the most effective representatives and descriptors of shape within a given dataset. Despite the fact that TDA is specifically designed for shape capture, there are several hurdles that persist in this newly developed field of study. These include (1) the nature and sensitivity of descriptors obtained by methods in TDA, (2) the dimensionality of the data and other parameter

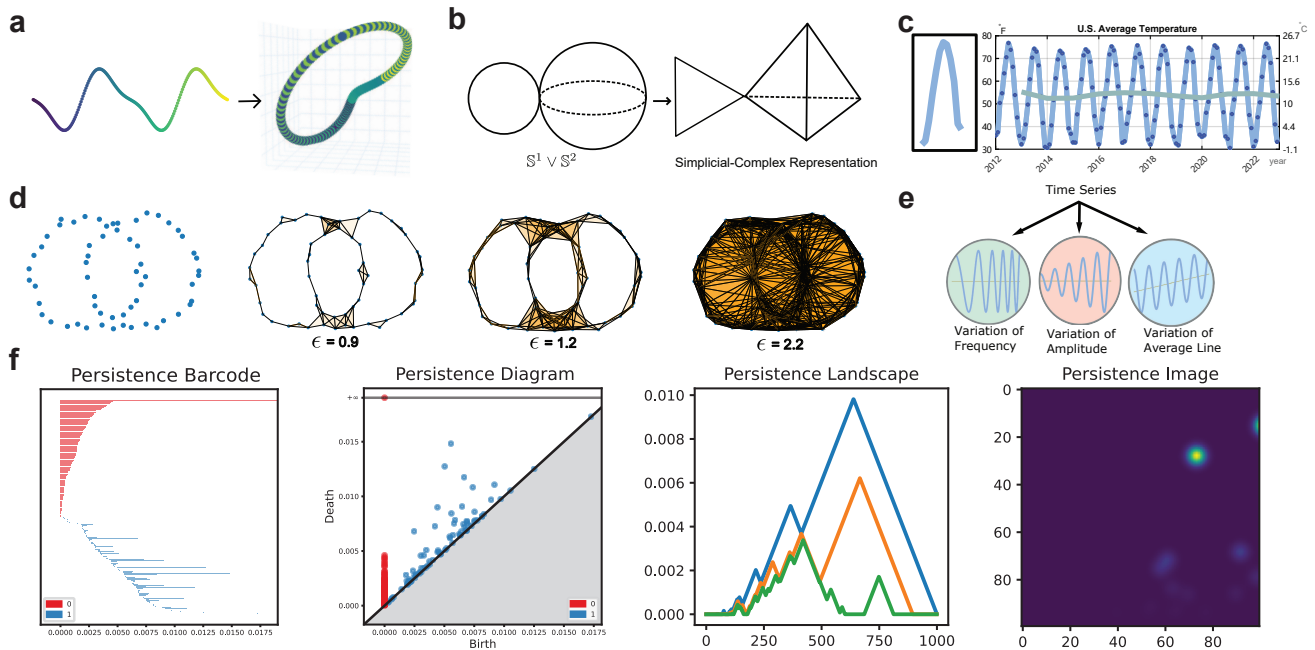


Fig. 1: Illustrations of methodology. **a**, Time-delay embedding (dimension=3, delay=10, skip=1) of $f(t_n) = \sin(2t_n) - 3\sin(t_n)$, with $t_n = \frac{\pi}{50}n$ ($0 \leq n \leq 200$). Resulting point clouds lay on a closed curve in 3-dimensional Euclidean space. The colour indicates their original locations in the time series. **b**, A topological space and its triangulation. On the left is a topological space consisting of a 1-dimensional sphere (i.e., a circle) and a 2-dimensional sphere with a single point of contact, denoted as $S^1 \vee S^2$. The right depicts a triangulation of this topological space. **c**, Average temperature in the U.S. with monthly values (dark blue dots) and yearly values (green curve). The left panel shows a single-year section of average temperature. **d**, Computing PH. The four plots consecutively show how a diagram or a barcode is computed: Connect each pair of points with a distance less than ϵ by a line segment, fill in each triple of points with mutual distances less than ϵ with a triangular region, etc., and compute the corresponding homology groups. In this way, as “time” ϵ increases, points in the diagram or intervals in the barcode record the “birth” and “death” of each generator of a homology group, i.e., the occurrence and disappearance of a loop (or a higher-dimensional hole), thereby revealing the essential topological features of the point cloud that persist. **e**, Characterising the vibration of a time series in terms of its variability of frequency, amplitude, and average line. **f**, Commonly used representations for PH, with an example of 100 points uniformly distributed over a bounded region in 2D Euclidean space.

90 choices, (3) the vectorisation of topological features, and (4)
 91 computational cost. These challenges will be elaborated in
 92 the following paragraphs within this section. Subsequently,
 93 we will demonstrate how our proposed methodology, Top-
 94 Cap, addresses these challenges through an application to
 95 consonant classification.

96 When applying TDA, the most imminent question is to
 97 comprehend the characteristics and nature of descriptors
 98 extracted via topological methods. TDA is grounded in the
 99 pure-mathematical field of algebraic topology (AT) [37, 38],
 100 with persistent homology (PH) being its primary tool [39,
 101 40]. While AT can quantify topological information to a
 102 certain extent [38, 1, 17], it is vitally important to understand
 103 both the capabilities and limitations of TDA. Generally
 104 speaking, TDA methods distinguish objects based on con-
 105 tinuous deformation. For example, PH cannot differentiate a
 106 disk from a filled rectangle, given that one can continuously
 107 deform the rectangle into a disk by pulling out its four
 108 edges. In contrast, PH can distinguish between a filled rect-
 109 angle and an unfilled one due to the presence of a “hole” in
 110 the latter, preventing a continuous deformation between the
 111 two. In certain circumstances, these methods are considered

excessively ambiguous to capture the structural information
 in data, thereby necessitating a more precise descriptor of
 shapes. To draw an analogy, TDA can be conceptualised
 as a scanner with diverse inputs encompassing time series,
 graphs, pictures, videos, etc. The output of this scanner is a
 multiset of intervals in the extended real line, referred to as a
 persistence diagram (PD)¹ or a persistence barcode (PB) [11,
 41, 42] (cf. Fig. 1f). In particular, by *maximal persistence* (MP)
 we mean the maximal length of the intervals. The precision
 of the topological descriptor depends on two factors: (1)
 the association of a topological space, i.e., the process of
 transforming the input data into a topological space (see
 Fig. 1b for a simplicial-complex representation of spaces;
 typically, the original datasets are less structured, and one
 should find a suitable representation of the data), and (2)
 the vectorisation of PD or PB, i.e., how to perform statistical
 inference with PD/PB. Despite there are many theoretical
 results which provide a solid foundation for TDA, few can
 elucidate the practical implications of PD and PB. For exam-
 112
 113
 114
 115
 116
 117
 118
 119
 120
 121
 122
 123
 124
 125
 126
 127
 128
 129
 130

¹In this article, we shall freely use the usual birth-by-death PDs and their birth-by-lifetime variants, whichever better serve our purposes. See Sec. S.2.2 for details.

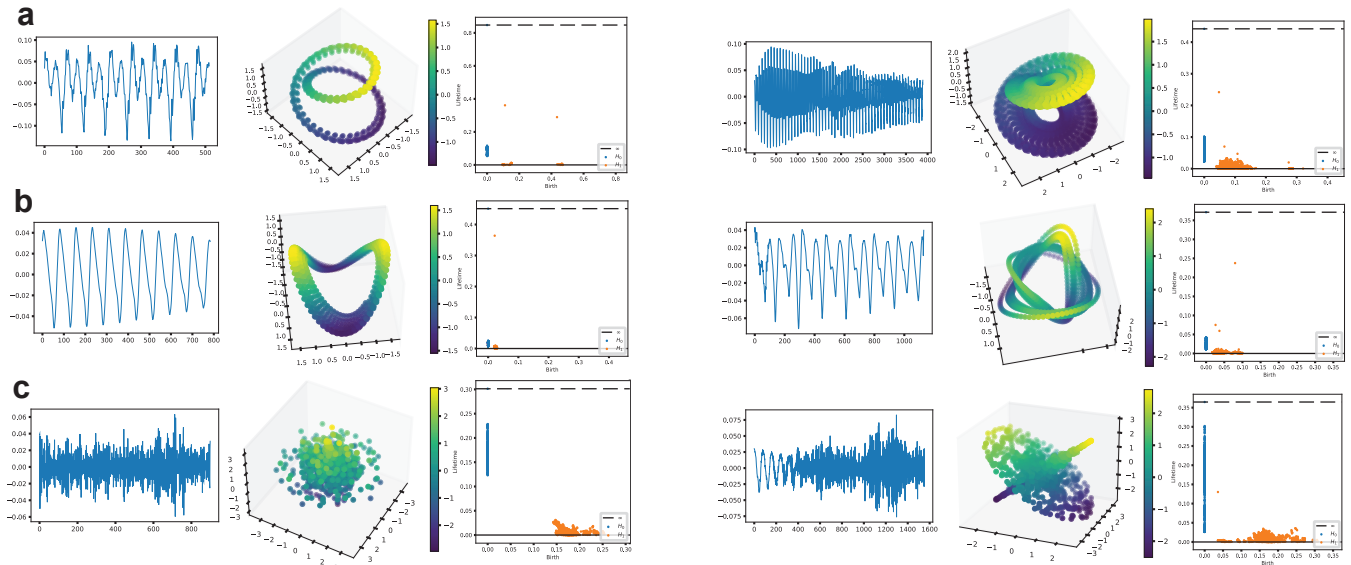


Fig. 2: The varied shapes of vowels, voiced consonants, and voiceless consonants. **a**, the left 3 panels and the right 3 panels depict 2 vowels, respectively. For each, the first picture is the time series of the vowel, the second picture corresponds to the 3-dimensional principal component analysis of the point cloud resulting from performing TDE (dimension=100, delay=1, skip=1) on this time series, and the third picture is the PD of this point cloud. **b**, The analogous features for 2 voiced consonants. **c**, Those for 2 voiceless consonants.

131 ple, what does it mean if many points are distributed near
 132 the birth–death diagonal line in a PD? In most cases, these
 133 points are regarded as descriptors of noise and are often
 134 disregarded if possible. Consequently, the TDA scanner can
 135 be seen as an imprecise observer, overlooking much of the
 136 information contained in less significant regions. In this
 137 article, we present an example of simulated time series to
 138 demonstrate that points distributed in such regions indeed
 139 encode vibration patterns of the time series, and a different
 140 distribution in these regions leads to a different pattern
 141 of vibration. This serves as a motivation for proposing
 142 TopCap and is further discussed in Sec. 2.1. It turns out that
 143 topological descriptors can be sharpened by noting patterns
 144 in these regions.

145 In view of the capability of topological methods to discern
 146 vibration patterns in time series, we apply them to classify
 147 consonant signals into voiced and voiceless categories. As a
 148 first demonstration of our findings, to *visualise* vowels,
 149 voiced consonants, and voiceless consonants in TDE and
 150 PD, see Fig. 2 (cf. Sec. S.1 for details of phonetic categories).

151 The first challenge, as many researchers may encounter
 152 when applying topological methods, is to determine the
 153 dimension of point clouds derived from input data [43, 44,
 154 45]. This essentially involves transforming the input into a
 155 topological space. In situations where the dimensionality
 156 of the data is large, researchers often project the data into
 157 a lower-dimensional topological space to facilitate visualisation
 158 and reduce computational cost [23, 24, 46]. On the other
 159 hand, as in this study and other applications with time
 160 series analysis [47, 48, 49, 50, 22, 51, 27], low-dimensional
 161 data are embedded into a higher-dimensional space. In
 162 both scenarios, deciding on the data dimensionality is both
 163 critical and challenging. Often, tuning the dimension is a
 164 tremendous task. In Sec. 3 of Discussion below, we delve

165 into the issue of data dimensionality. In our case, as it might
 166 seem counterintuitive compared to most algorithms, when
 167 the data are embedded into a higher-dimensional space, the
 168 computation will be a little faster, the point cloud appears
 169 smoother and more regular, and most importantly, more
 170 salient topological features can be spotted, which seldom
 171 happen in lower-dimensional spaces. When encountering
 172 the dimensionality of data, researchers would think of the
 173 well-known curse of dimensionality [52]: As a typical algo-
 174 rithm grapple, with the increase of dimension, more data
 175 are needed to be involved, often growing exponentially
 176 and thereby escalating computational cost. Even worse, the
 177 computational cost of the algorithm itself normally rises as
 178 the dimension goes higher. However, topological methods
 179 do not necessarily prefer data of lower dimension. For com-
 180 puting PH (see Fig. 1d for the process of computing PD/PB
 181 from point clouds), a commonly used algorithm [53, 54]
 182 sees complexity grow with an increase in the number n of
 183 simplices during the process, with a worst-case polynomial
 184 time-complexity of $O(n^3)$. As such, the computational cost
 185 is directly related to the number of simplices formed during
 186 filtration. Our observation shows that computation time
 187 may not increase much given an increase of dimension of
 188 data, because the latter may have little effect on the size
 189 (i.e., number of points) of the point cloud and thus neither
 190 on the number of simplices formed during filtration.

191 Having obtained a suitable topological space from input
 192 data, one can derive a PD/PB from the topological space,
 193 which constitutes a multiset of intervals. The subsequent
 194 challenge lies in the vectorisation of the PD/PB for its
 195 integration into a machine-learning algorithm. The vector-
 196 isation process is essentially linked to the construction
 197 of the topological space, as the combination of different
 198 methods for constructing the topological space and vectori-

sation together determine the descriptor utilised in machine learning. A plethora of vectorisation methods exist, such as persistence landscape (PL) [55] and persistence image (PI) [56], among others, as documented in various studies [40, 57] (cf. Fig. 1f). The selection of these methods requires careful consideration. In Sec. 4 of Methods, we employ MP and its corresponding birth time as two features. These have been integrated into nine traditional machine learning algorithms to classify voiced and voiceless consonants, yielding an accuracy that exceeds 96% with each algorithm. This vectorisation method is quite simple, primarily due to our construction of topological spaces from phonetic time series, as detailed in the Method section. This construction enables PH to capture significant topological features within the time series. In Sec. 2.1, we also observe a pattern of vibration which could potentially be vectorised by PI into a matrix. As one of its strengths, PI emphasises regions where the weighting function scores are high, which makes it a computationally flexible method. Future work may involve a more precise recognition of such patterns using PI.

An outline for the remainder of this article goes as follows. Sec. 1.1 gives an overview of closely related works in the field, with an extended commentary relegated to Sec. S.4. Sec. 2 of Results provides in more detail the motivations for TopCap, presents final results of classifying voiced and voiceless consonants, including a comparison with traditional deep learning neural networks, and explains our purposes in practical use. Sec. 3 of Discussion highlights important parameter setups and indicates potential directions for future work, with further discussion in Sec. S.3. Sec. 4 of Methods contains a detailed template of TopCap. Sec. 5 gives the data and code sources for our experiments.

1.1 Related works

Time series analysis [58] is a prevalent tool for various applied sciences. The recent surge in TDA has opened new avenues for the integration of topological methods into time series analysis [21, 59, 60]. Much literature has contributed to the theoretical foundation in this area. For example, theoretical frameworks for processing periodic time series have been proposed by Perea and Harer [61], followed by their and their collaborators' implementation in discovering periodicity in gene expressions [62]. Their article [61] studied the geometric structure of truncated Fourier series of a periodic function and its dependence on parameters in time-delay embedding (TDE), providing a solid background for TopCap. In addition to periodic time series, towards more general and complex scenarios, quasi-periodic time series have also been the subject of scholarly attention. Research in this direction has primarily concentrated on the selection of parameters for geometric space reconstruction [63] and extended to vector-valued time series [64].

In this article, a topological space is constructed from data using TDE, a technique that has been widely employed in the reconstruction of time series (see Fig. 1a and cf. Sec. S.2.1 for more background). Thanks to the topological invariance of TDE, the general construction of simplicial-complex representation (see Fig. 1b) and computation of PH from point clouds (see Fig. 1d) apply to time series data,

although this transformation involves subtle technical issues in practice. For instance, Emrani et al. utilised TDE and PH to identify the periodic structure of dynamical systems, with applications to wheeze detection in pulmonology [47]. They selected the embedded dimension d as 2, and their delay parameter τ was determined by an autocorrelation-like (ACL) function, which provided a range for the delay between the first and second critical points of the ACL function. Pereira and de Mello proposed a data clustering approach based on PD [48]. The data were initially reconstructed by TDE, with $d = 2$ and $\tau = 3$, so as to obtain the corresponding PD, which was then subjected to k -means clustering. The delay τ was determined using the first minimum of an auto mutual information, and the embedded dimension d was set to be 2 as using 3 dimensions did not significantly improve the results. Khasawneh and Munch introduced a topological approach for examining the stability of a class of nonlinear stochastic delay equations [49]. They used false nearest neighbours to determine the embedded dimension $d = 3$ and chose the delay to equal the first zeros of the ACL function. Subsequently, the longest persistence lifetime in PD was used as a vectorisation to quantify periodicity. Umeda focused on a classification problem for volatile time series by extracting the structure of attractors, using TDA to represent transition rules of the time series [22]. He assigned $d = 3$, $\tau = 1$ in his study and introduced a novel vectorisation method, which was then applied to a convolutional neural network (CNN) to achieve classification. Gidea and Katz employed TDA to detect early signs prior to financial crashes [51]. They studied multi-dimensional time series with $\tau = 1$ and used persistence landscape as a vectorisation method. For speech recognition, Brown and Knudson examined the structure of point clouds obtained via TDE of human speech signals [27]. The TDE parameters were set as $d = 3$, $\tau = 20$, after which they examined the structure of point clouds and their corresponding PB.

Upon reviewing the relevant literature, we see that currently there is no general framework for systematically choosing d and τ , and researchers often have to make choices in an ad hoc fashion for practical needs. While the TDE-PH topological approach to handling time series data is not new, TopCap extracts features from high-dimensional spaces. For example, in our experiment $d = 100$. It happens in some cases that in a low-dimensional space, regardless of how optimal the choice of τ is, the structure of the time series cannot be adequately captured. In contrast, given a high-dimensional space, feature extraction from data becomes simpler. Of course, operating in a high-dimensional space comes with its own cost. For example, the adjustment of τ then requires careful consideration. Nonetheless, it also offers advantages, which we will elucidate step by step in the subsequent sections.

2 RESULTS

This research drew inspiration from Carlsson and his collaborators' discovery of the Klein-bottle distribution of high-contrast, local patches of natural images [20], as well as their subsequent recent work on topological CNNs for learning image and even video data [2]. By analogy, we aim to understand a distribution space for speech data, even a

directed graph structure on it modeling the complex network of speech-signal sequences for practical purposes such as speaker diarisation, and how these topological inputs may enable smarter learning (cf. Sec. S.1). Here are some of our first findings in this direction, set in the context of topological time series analysis.

2.1 Detection of vibration patterns

The impetus behind TopCap lies in an observation of how PD can capture vibration patterns within time series. To begin with, our aim is to determine which sorts of information can be extracted using topological methods. As the name indicates, topological methods quantify features based on topology, which distinguishes spaces that cannot continuously deform to each other. In the context of time series, we conduct a series of experiments to scrutinise the performance of topological methods, their limitations as well as their potential.

Given a periodic time series, its TDE target is situated on a closed curve (i.e., a loop) in a sufficiently high-dimensional Euclidean space (see Fig. 1a). Despite the satisfactory point-cloud representation of a periodic time series, it remains rare in practical measurement and observation to capture a truly periodic series. Often, we find ourselves dealing with time series that are not periodic yet exhibit certain patterns within some time segments. For instance, Fig. 1c portrays the average temperature of the United States from the year 2012 to 2022, as documented in [65]. Although the temperature does not adhere strictly to a periodic pattern, it does display a noticeable cyclical trend on an annual basis. Typically, the temperature tends to rise from January to July and fall from August to December, with each year approximately comprising one cycle of the variation pattern. One strength of topological methods is their ability to capture “cycles”. A question then arises naturally: Can these methods also capture the cycle of temperature as well as subtle variations within and among these cycles? To be more precise, we first observe that variations occur in several ways. For instance, the amplitude (or range) of the annual temperature variation may fluctuate slightly, with the maximum and minimum annual temperatures varying from year to year. Additionally, the trend line for the annual average temperature also shows fluctuations, such as the average temperature in 2012 surpassing that of 2013. Despite each year’s temperature pattern bearing resemblance to that depicted in the left panel in Fig. 1c (representing a single cycle of temperature within a year), it may be more beneficial for prediction and response strategies to focus on the evolution of this pattern rather than its specific form. In other words, attention should be directed towards how this cycle varies over the years. This leads to several questions. How can we consistently capture these subtle changes in the pattern’s evolution, such as variations in the frequency, amplitude, and trend line of cycles? How can we describe the similarities and differences between time series that possess distinct evolutionary trajectories? In applications, these are crucial inquiries that warrant further exploration.

To address these questions, we propose three kinds of “fundamental variations” which are utilised for depicting the evolutionary trace of a time series. Consider a series of a periodic function $f(t_n) = f(t_n + T)$, where T is a period.

- (1) *Variation of frequency.* Denote the frequency by $F = T^{-1}$. Note that the series is not necessarily periodic in the mathematical sense. Rather, it exhibits a recurring pattern after the period T . For instance, the average temperature from Fig. 1c is not a periodic series, but we consider its period to be one year since it follows a specific pattern, i.e., the one displayed in the left panel of Fig. 1c. This 1-year pattern always lasts for a year as time progresses. Hence, there is no frequency variation in this example. This type of variations can be represented as $g_1(t_n) = f(F(t_n) \cdot t_n)$, where $F(t_n)$ is a series representing the changing frequency. This type of variation occurs, for example, when one switches their vocal tone or when one’s heartbeats experience a transition from walking mode to running mode.
- (2) *Variation of amplitude.* The amplitudes of temperature in the years 2014 and 2015 are 42.73°F and 40.93°F , respectively. So the variation of amplitude from 2014 to 2015 is -1.80°F . This can be represented by $g_2(t_n) = A(t_n) \cdot f(t_n)$, where $A(t_n)$ is a series of the changing amplitude. This type of variation is observed when a particle vibrates with resistance or when there is a change in the volume of a sound.
- (3) *Variation of average line.* The average temperatures through the years 2012 and 2013 are 55.28°F and 52.43°F , respectively. The variation of average line from 2012 to 2013 is -2.85°F . Let $g_3(t_n) = f(t_n) + L(t_n)$, where $L(t_n)$ is a series representing the variation of average line. This type of variation is observed when a stock experiences a downturn over several days or when global warming causes a year-by-year increase in temperature.

To summarise, Fig. 1e provides a visual representation of the three fundamental variations. It is important to note that these variations are not utilised to depict the pattern itself but rather to illustrate the variation within the pattern or how the time series oscillates over time. This approach offers a dynamic perspective on the evolution of the time series, capturing changes in patterns that static analyses may overlook.

Using three simulated time series corresponding to the above three fundamental types of variation (see Sec. 4.1 for detailed construction), we demonstrate that PD can distinguish these variations and detect how significant they are. See Fig. 3, where a smaller value of c indicates a more rapid fundamental variation. Here, regardless of which value c takes, each individual diagram features a prominent single point at the top and a cluster of points with relatively short duration, except when $F(t_n) = 1$ (i.e., $c = 4$). In this case, the series represents a cosine function, and thus the diagram consists of a single point. Normally, one tends to overlook the points in a PD that exhibit a short duration as they are sometimes inferred as noise. However, in this example, the distribution of those points holds valuable information regarding the three fundamental variations. As shown in Fig. 3, each fundamental variation has its distinct pattern of distribution in the lower region of a diagram, which leads to refined inferences: If the points spiral along the vertical axis of lifetime, it is probably due to a variation of amplitude; if every two or four points stay close to form a “shuttle”, it probably indicates a variation of average line;

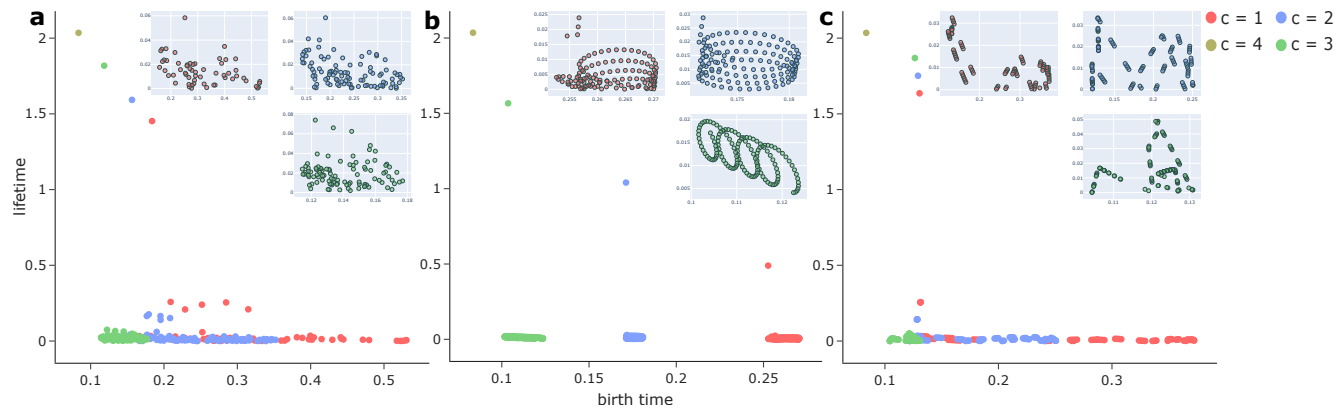


Fig. 3: 1-dimensional PH reveals three fundamental variations. **a**, Detecting variation of frequency. Upper-right panels zoom in to show the barcode distribution in the lower dense region, where the position and colour of each value of c in the main legend corresponds to those of its panel. Note that when $c = 4$, there is a single point, and so the panel for this value is omitted. **b**, Detecting variation of amplitude. **c**, Detecting variation of average line.

otherwise the points just seem to randomly spread over, which more likely results from a variation of frequency. It is also straightforward to distinguish the values of c for a specific fundamental variation, by their most significant point in the diagram: Longer lifetime for the barcode of the solitary point indicates slower variation. The lower region of a diagram also gives some hints in this respect.

In this simulated example, we demonstrated how PD could be utilised as a uniform means to distinguish three fundamental variations of the cosine series and their respective rates of change. However, it is important to note that in general scenarios, identifying the fundamental variations in a time series using topological methods may encounter significant challenges. Although topological methods are indeed capable of capturing this information, vectorising this information for subsequent utilisation remains a complex task at this stage. Having recognised the potential of topological methods, we resort to an alternative algorithm for handling time series. Specifically, despite the difficulty in vectorising PD to measure each fundamental variation, we have developed a simplified algorithm to measure the vibration of time series as a whole. This approach provides a comprehensive understanding of the overall behaviour of a time series, bypassing the need for complex vectorisation.

2.2 Traditional machine learning methods with novel topological features

Using datasets comprising human speech, we initially employ the Montreal Forced Aligner to align natural speech into phonetic segments. Following preprocessing of these phonetic segments, TDE is conducted with dimension parameter $d = 100$ and delay parameter τ set to equal $6T/d$, where T approximates the (minimal) period of the time series. Following additional refinement procedures, PDs are computed for these segments and are then vectorised based on MP and its corresponding birth time. The comprehensive procedural framework is expounded in Secs. 4.2 and 4.3, while the corresponding workflow is shown in Fig. 4e. In the applications of TDE, the dimension parameter d is usually determined through some algorithms designed to

identify the minimal appropriate dimension [45, 66]. The delay parameter τ is determined by an ACL function with no specific rule, but in many cases $\tau = mT/d$ for some positive integer m . In our pursuit of enhanced extraction of topological features, a relatively high dimension is chosen (see Sec. 3 for more discussion on dimension in TDE). Given this higher dimension, the usual case of $\tau = T/d$ with $m = 1$ may prove excessively diminutive, particularly in light of the time series only taking values in discrete time steps. Consequently, in TopCap we adopt an adjusted parametrisation for $\tau = mT/d$ with a relatively large value $m = 6$.

We input the pair of MP and birth time from 1-dimensional PD for each sound record to multiple traditional classification algorithms: Tree, Discriminant, Logistic Regression, Naive Bayes, Support Vector Machine, k -Nearest Neighbours, Kernel, Ensemble, and Neural Network. We use the application of the MATLAB (R2022b) Classification Learner, with 5-fold cross-validation, and set aside 30% records as test data. This application performs machine learning algorithms in an automatic way. There are a total of 1016 records, with 712 training samples and 304 test samples. Among them, 694 records are voiced consonants and the remaining are voiceless consonants. The models we choose in this application are Optimizable Tree, Optimizable Discriminant, Efficient Logistic Regression, Optimizable Naive Bayes, Optimizable SVM, Optimizable KNN, Kernel, Optimizable Ensemble, and Optimizable Neural Network. Our results are compared with those obtained from a CNN, for which we compute the short-time Fourier transform of phones (implemented in Python with `signal.stft` or `scipy.signal.spectrogram`) and directly classify the resulting spectrograms using CNN, without extracting any topological features.

The results are shown in Fig. 4a–d. The receiver operating characteristic curve (ROC), area under the curve (AUC), and accuracy metrics collectively demonstrate the efficacy of these topological features as inputs for a variety of machine learning algorithms. Each of the algorithms incorporating topological inputs attains AUC and accuracy surpassing 96%, whereas CNN without topological inputs

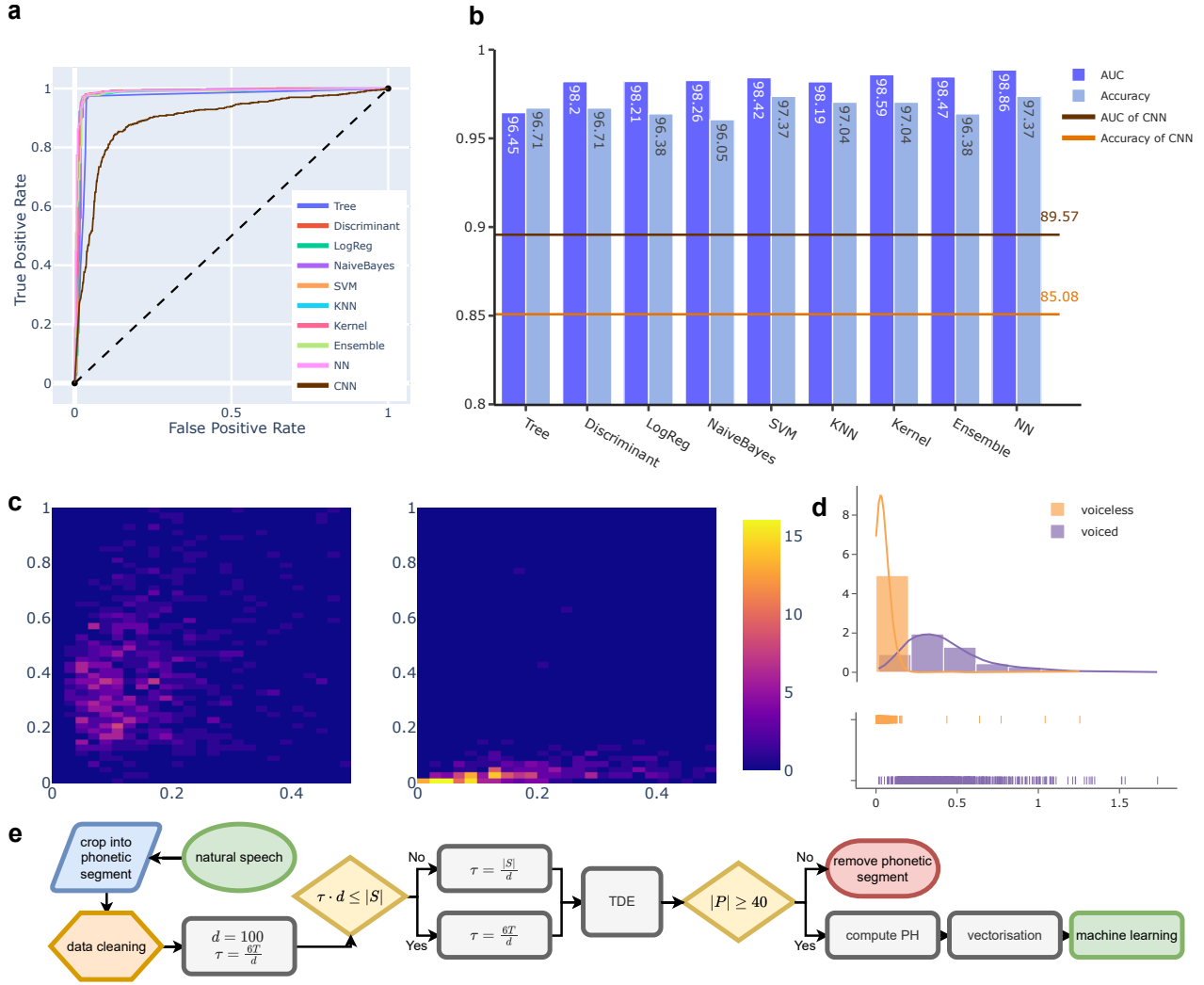


Fig. 4: Machine learning results with topological features. **a**, ROCs of TopCap’s traditional machine learning algorithms with topological inputs and of CNN without topological inputs. **b**, Accuracy and AUC of TopCap versus CNN. **c**, Diagrams of records represented as (birth time, lifetime) for voiced consonants (left) and voiceless consonants (right), where voiced consonants exhibit relatively higher birth time and lifetime. The colour represents the density of points in each unit grid box. **d**, Histograms of records represented by their lifetime for voiced and voiceless consonants, together with kernel density estimation and rug plot. The distributions of MP can distinguish voiced and voiceless consonants. **e**, Flow chart of experiment. Here $|S|$ denotes the number of samples in a time series, $|P|$ denotes the number of points in the point cloud, and T denotes the (minimal) period of the time series computed by the ACL function.

516 merely yields an AUC of 90% and an accuracy of 85%. The
 517 ROC and AUC together depict the high performance of our
 518 classification model across all classification thresholds. The
 519 2D histograms depicted in Fig. 4c–d collectively illustrate
 520 the distinct distributions of voiced and voiceless consonants.
 521 Voiced consonants tend to exhibit a relatively higher birth
 522 time and lifetime, which provides an explanation for the
 523 high performance of these algorithms. Despite the intricate
 524 structure that a PD may present, appropriately extracted
 525 topological features enable traditional machine learning al-
 526 gorithms to separate complex data effectively. This high-
 527 lights the potential of TDA in enhancing the performance
 528 of machine learning models.

529 It is noteworthy that the CNN we use as a compar-
 530 ative, which comprises 5 layers with more than 43 million

parameters, is considerably more intricate than traditional
 machine learning algorithms with TopCap. Nonetheless, in
 effect, this CNN requires 2 hours for sufficient training (1602
 spectrograms in total). In contrast, learning with topological
 inputs achieves both higher accuracy as in Fig. 4a–b and
 higher efficiency, under 5 minutes including topological
 feature extraction on the same device (mere seconds for
 machine learning alone).

539 In summary, from our topological detection results, the
 540 most significant distinction between voiced and voiceless
 541 consonants is that the former exhibit higher MP. This can
 542 scarcely be detected in lower dimensions regardless of how
 543 we tune the delay parameter τ . Besides the figure above, see
 544 also Fig. 2 for a sample of the recognition of vowels as well
 545 as consonants in terms of their *shapes*.

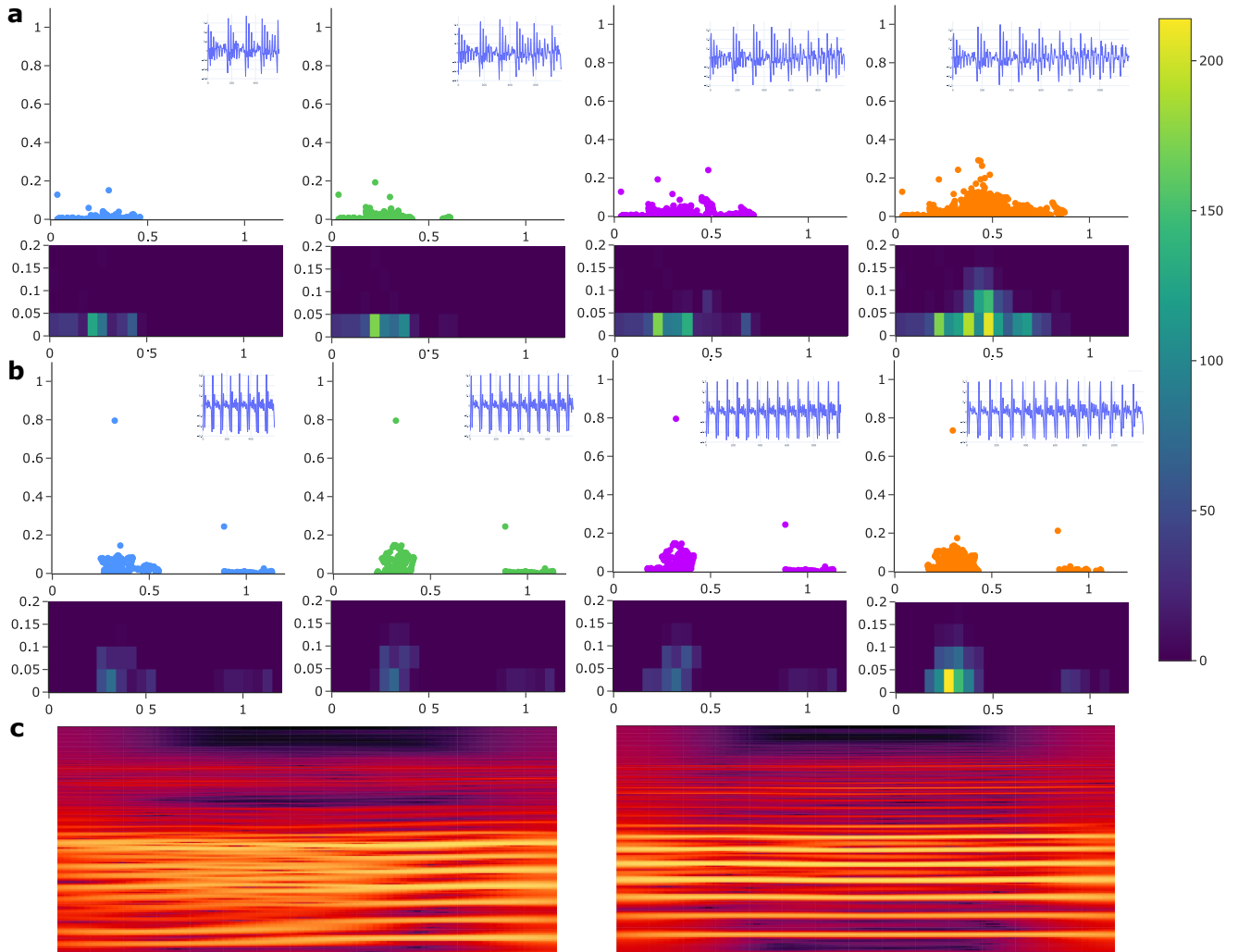


Fig. 5: Variation of 1-dimensional PDs due to the fundamental variations of time series. **a**, PDs of drastic fundamental variations. The small panel on top right of each diagram shows the original time series, with 4 segments extracted from the same record of [a], each starting from time 0 and ending at time 600, 800, 1000, 1200, respectively. It can directly be seen from the time series that the variation of amplitude in (a) is bigger than (b); for frequency, see c; normally, we do not discuss the average line of phonetic data as it is assumed to be constant. Below, each diagram shows the clustering density of points in the lower region of the PD. **b**, PDs of mild fundamental variations for 4 time-series segments extracted from the other record of [a], with the same ending and starting times as in (a). The lower density diagrams demonstrate that unstable time series are characterised by a higher density of points in the lower region of PD. Moreover, stable series tend to attain high MP. **c**, Spectral frequency plots of the time series with rapid variations (left) and with mild variations (right).

546 2.3 The three fundamental variations gleaned from a 547 persistence diagram

548 A PD for 1-dimensional PH encodes much more information
549 beyond the birth time and lifetime of the point of MP.
550 The three fundamental variations examined in Sec. 2.1 also
551 manifest themselves in certain regions of the PD, which can
552 in turn be vectorised.

553 To capture these variations, we perform an experiment
554 with two records of the vowel [a]. Specifically, we demon-
555 strate the fundamental variations by comparing the PDs
556 of (a) the record of [a] relatively unstable with respect to
557 the fundamental variations and (b) the other record of the
558 same vowel that is relatively stable. To better illustrate the

559 results, we crop each record into 4 overlapping intervals,
560 each starting from time 0 and ending at 600, 800, 1000, 1200,
561 respectively. When adding a new segment of 200 units into
562 the original sample each time, the amplitude and frequency
563 of the series altered more drastically in case (a). A more
564 rapid changing rate may lead to more points distributed
565 in the lower region of the diagram. The outcomes are
566 presented in Fig. 5. The plots in Fig. 5c show that the spectral
567 frequency of (a) indeed varies faster than that of (b).

568 We should also mention that the 1-dimensional PD here
569 serves as a profile for the collective effect of the fundamen-
570 tal variations. Currently, it is unclear how the points in the
571 lower region change in response to a specific variation.

3 DISCUSSION

In the realm of applying topological methods to analyse time series [47, 48, 49, 50, 22, 51, 27], the determination of parameters for TDE emerges as a pivotal aspect. This stems from the significant impact that the selection of parameters has on the resulting topological spaces and their corresponding PDs. There exist several convenient algorithms for parameter selection. For example, the False Nearest Neighbours algorithm (FNN), a widely utilised tool, provides a method for deciding the minimal embedded dimension [66]. However, in the context of PH, usually the objective is not to achieve a *minimal* dimension. Contrarily, a dimension of substantial magnitude may be desirable due to certain advantages it offers.

In this section, as a main novel feature of TopCap, we reveal and leverage the relationship between embedded dimension and maximal persistence. We relegate further aspects of parameter selection to Sec. S.3.

In the TDE-PH approach, the determination of dimension in a TDE can be complex. However, it plays a pivotal role in the extraction of topological descriptors such as MP. It is observed that a larger dimension can significantly enhance the theoretically optimal MP of a time series. In TopCap, the dimension of TDE is set to be 100, a relatively large dimension for the experiment. On the other hand, several factors also constrain this choice. These include the length of the sampled time series, since the dimension cannot exceed the length (otherwise it would render the resulting point cloud literally pointless). The constraints also include the periodicity of the time series, as the time-delay window size should be compatible with the approximate period of the time series, which is to be elaborated below.

According to Perea and Harer [61, Proposition 5.1], there is no information loss for trigonometric polynomials if and only if the dimension of TDE exceeds twice the maximal frequency. Here, no information loss implies that the original time series can be fully reconstructed from the embedded point cloud. In general, for a periodic function, a higher dimension of TDE can yield a more precise approximation by trigonometric polynomials. Although there are no absolutely periodic functions in real data, each time series exhibits its own pattern of vibration, as discussed in Sec. 2.1, and a higher dimension of embedding may be employed to capture a more accurate vibration pattern in the time series. Furthermore, an increased embedded dimension may result in reduced computation time for PD. For instance, computation times for a voiced consonant [ŋ] are 0.2671, 0.2473, and 0.2375 seconds, corresponding to embedded dimensions 10, 100, and 1000 (see Fig. 6a). This is attributed to the reduction due to a higher dimension on the number of points in the embedded point cloud. While this reduction in computation time may not be considered substantial compared to the impact of changing skip (see Fig. 6d), it may become significant when handling large datasets. More importantly, an increased embedded dimension can yield benefits such as enhanced MP, which serves as a major motivation for higher dimensions, as well as a smoother shape of resulting point clouds obtained through TDE, which makes the embedding visibly reasonable. Typically, for most algorithms, a lower dimension is preferred due to factors

such as those associated with curse of dimensionality and computation cost. By contrast, in TopCap, we opt instead for a higher dimension.

However, the embedded dimension cannot be arbitrarily large. As illustrated in Fig. 6c, when the embedded dimension escalates to 1280, it becomes unfeasible to capture a significant MP in the phonetic time series. This results from a break of the point cloud. When the embedded dimension further reaches 1290, an empty 1-dimensional barcode is obtained due to the lack of points necessary to form even a single cycle. In this way, the dimension of TDE is related to the length of the time series.

Using a sound record of the voiced consonant [ŋ] as an exemplar, we delineate the correlation between MP and embedded dimension in Fig. 6a–c. As depicted in Fig. 6b, MP tends to escalate rapidly and nonlinearly with the increase in dimension, signifying that a more substantial MP is captured in higher-dimensional TDE. Notably, two precipitous drops in MP are observed, corresponding to embedded dimensions 600 and 1190. When $d = 600$, this time series can theoretically attain its optimal MP when $\tau = 2$ (see Sec. S.2.1). However, given the length of the series is 1337 and the window size is $d \cdot \tau = 1200$, with the skip set as 5, only 28 points are in the resulting point cloud for PD computation. The sparse point cloud fails to represent the original series adequately, leading to a decrease in MP. A similar phenomenon occurs when the dimension reaches 1190. The principal component analysis for dimension 1280 is shown in Fig. 6c. In this scenario, as observed above, the hypothetical cycle fails to form as there is a break in the point cloud, resulting in a free-fall in MP. In contrast, when $d = 630$, this series has a significant MP when $\tau = 1$, resulting in a window size of $d \cdot \tau = 630$. There are 142 points in the point cloud for the persistence diagram if skip equals 5, ensuring that the MP rises again without any breakdown. The embedded dimension also contributes significantly to the geometric property of time-delay embedding, as the shape becomes smoother in higher dimensions and the point cloud more structural.

As mentioned above, there are three crucial parameters in TDE, namely, d , τ , and skip. However, it is worth noting that the TDE-PH approach encompasses many other significant variables and choices. These include the construction of underlying topological space of the point clouds (i.e., the distance function for pairwise points), and the type of complexes utilised in filtering PH, among others. Some of these choices, despite their importance, were seldom addressed in the literature. Here, we propose a method for determining delay in order to capture the theoretically optimal MP of a time series in high-dimensional TDE. In future research, we aim at more systematic approaches for determining other parameters, particularly dimension of the TDE.

4 METHODS

4.1 Constructing vibrating time series

There are three kinds of fundamental variations mentioned in Sec. 2.1. In order to substantiate our argument, let $t_n =$

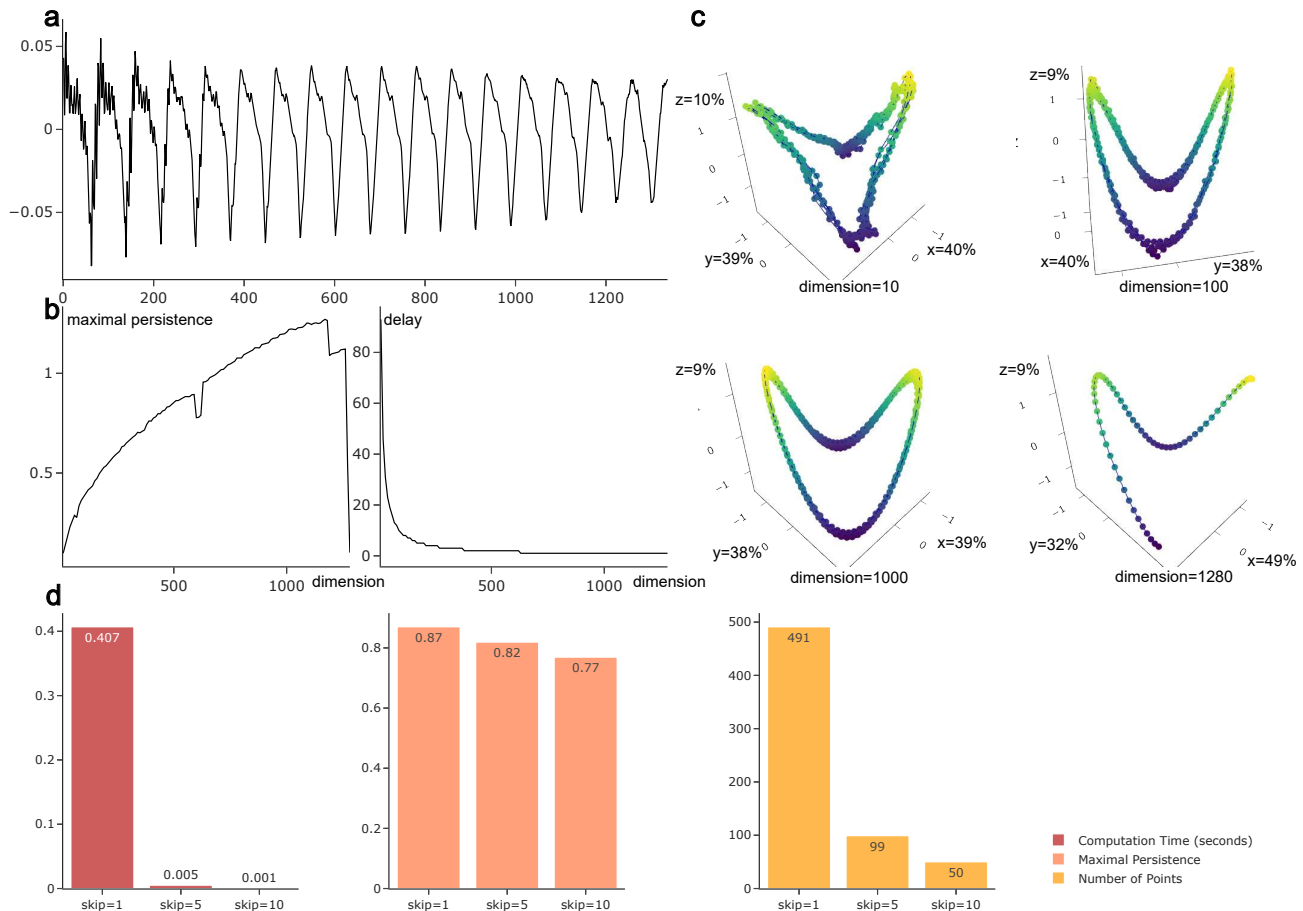


Fig. 6: Point-cloud behaviour with increasing embedded dimension. **a**, Original .wav file of a record of [ɹ] (voiced consonant). **b**, MP of the series after TDE as dimension increases (left) and the corresponding delay that ensures the time series to reach theoretically optimal MP (right). Skip equals 5 when computing PD. **c**, Visualisation of the embedded point clouds, which shows principal component analysis (PCA) of the embedded point clouds in 3D as projected from various dimensions. Skip equals 1 when performing PCA. The percentage along each axis indicates the PCA explained variance ratio. **d**, Given a sound record of the voiced consonant [m], computation time, MP, and the size of point clouds as skip increases (see Sec. S.3.1 for details). An increase in skip can lead to a significant reduction in computation time, owing to the reduced size of the point cloud. However, MP remains resilient to an increase in the skip parameter.

688 $0.01n$ with $0 \leq t_n \leq 7\pi$ and for each $c \in \{1, 2, 3, 4\}$ define

$$\begin{aligned} f(t_n) &= \cos(t_n) \\ F(t_n) &= \frac{c}{4} + \frac{1-c}{7\pi} \cdot t_n \\ g_1(t_n) &= f(F(t_n) \cdot t_n) \end{aligned}$$

689 Note that $F(t_n) = c/4$ when $t_n = 0$ and $F(t_n) = 1$ when
690 $t_n = 7\pi$. In fact, $F(t_n)$ is a sequence of line segments con-
691 necting $(0, c/4)$ and $(7\pi, 1)$. Correspondingly, the frequency
692 of $g_1(t_n)$ changes more slowly as c increases. In the extreme
693 case when $c = 4$, we have $F(t_n) = 1$, so

$$g_1(t_n) = f(F(t_n) \cdot t_n) = f(t_n) = \cos(t_n)$$

694 which is a periodic function. For each value of c , we applied
695 TDE to the series $g_1(t_n)$ with dimension 3, delay 100, skip
696 10 and computed the 1-dimensional PD of the embedded
697 point cloud. See Fig. 3a for the results. Replacing $F(t_n)$ by
698 $A(t_n)$ and $L(t_n)$, we obtained the diagrams in Figs. 3b and
699 3c, respectively.

4.2 Obtaining phonetic data from natural speech

700

We used speech files sourced from SpeechBox [67],
701 ALLSTAR Corpus, task HT1 language English L1 file,
702 retrieved on 28th January 2023. SpeechBox is a web-based
703 system providing access to an extensive collection of digital
704 speech corpora developed by the Speech Communication
705 Research Group in the Department of Linguistics at North-
706 western University. This section contains a total of 25 indi-
707 vidual files, comprising 14 files from women and 11 files
708 from men. The age range of these speakers spans from 18 to
709 26 years, with an average of 19.92. Each file is presented in
710 the WAV format and is accompanied by its corresponding
711 aligned file in Textgrid format, which features three tiers of
712 sentences, words, and phones. Collectively, these 25 speech
713 files amount to a total duration of 41.21 minutes. The speech
714 file contains each individual reading the same sentences
715 consecutively for a duration ranging from 80 to 120 seconds,
716 contingent upon each person's pace. The original .wav file
717 has a sampling frequency of 22050 and comprises only
718 one channel. Since the Montreal Forced Aligner (MFA) [68]
719

is trained in a sampling frequency of 16000, we opted to adjust the sampling frequency of the .wav files accordingly. We then extracted the “words” tier from Textgrid and aligned words into phones using English_MFA dictionary and acoustic model (MFA version 2.0.6). Thus we obtained corresponding phonetic data from these speech files.

Subsequently, we used voiced and voiceless consonants in those segments as our dataset. Voiced consonants are consonants for which vocal cords vibrate in the throat during articulation, while voiceless consonants are pronounced otherwise (see also Sec. S.1). Specifically, using Praat [69], we extracted voiced consonants [ɹ], [m], [n], [j], [l], [v], and [ʒ]; for voiceless consonants, we selected [f], [k], [θ], [t], [s], and [tʃ]. These phones were then read as time series. Our selection was limited to these voiced and voiceless consonants, as we aimed to balance the ratio of voiced and voiceless consonant records in these speech files. Additionally, some consonants, such as [d] and [h], appeared difficult to classify by our methods.

4.3 Deriving topological features from phonetic data

Prior to the extraction of topological features from a time series, we first imbued this 1-dimensional time series with a (Euclidean) topological structure through TDE. It is noteworthy that this technique also applies to multi-dimensional time series. The ambient space throughout this article is always a Euclidean space. By establishing the topological structure there, or more precisely, the distance matrices, we subsequently calculated PH. We elaborate on the following main steps. See Fig. 4e for the flow chart of this section.

4.3.1 Data cleaning

This involved eliminating the initial and final segments of a time series until the first point with an amplitude exceeding 0.03 occurred. This approach was aimed at mitigating the impact of environmental noise at the beginning and end of a phone. Any resulting series with fewer than 500 points will be disregarded, as such series were considered insufficiently long or to contain excessive environmental noise.

4.3.2 Parameter selection for time-delay embedding

We selected suitable parameters for TDE to capture the theoretically optimal MP of a given time series. The dimension of the embedding was fixed to be 100. Our principle for determining an appropriate dimension is that we want to choose the embedded dimension to be large for a time series of limited length. As discussed in Sec. 3 and cf. Sec. S.2.1, a higher dimension results in a more accurate approximation. This approach also aimed to enhance computational efficiency and the occurrence of more prominent MP. Nonetheless, it is imperative to exercise caution when selecting the dimension, as excessively large dimensions may lead to empty point clouds and other uncontrollable factors.

With a proper dimension, we then computed the delay for the embedding. According to Perea and Harer [61], in the case of a periodic function, the optimal delays τ can be expressed as

$$\tau = m \cdot \frac{T}{d}$$

where T denotes the (minimal) period, d represents the dimension of the embedding, and m is a positive integer.

Under these conditions, we could obtain the theoretically optimal MP. The time series under consideration in our case was far from periodic, however, so we used the first peak of the ACL function to represent the period T and set $m = 6$, thus obtaining a relatively proper delay τ . The common choice of τ is to let window size equal the (minimal) period. However, in the case of a discrete time series, one often obtains $\tau = 0$ or $\tau = 1$ in this way, since the dimension of TDE is too large in comparison. Therefore, one strategy is to increase m to get a relatively reasonable τ . The performance of delay obtained in this way is presented in Sec. 3.

Then τ was rounded to the nearest integer (if it equals 0, take 1 instead). It was common that $\tau \cdot d$ exceeded the number of points in the series, resulting in an empty embedding. In this case, we adopted $\tau = \lfloor |S|/d \rfloor$, where $|S|$ denotes the number of points (i.e., the point capacity of the time series), and then rounded it downwards. This enabled us to obtain the appropriate delay for each time series, thereby facilitating the attainment of significant MP for the specified dimension.

Lastly, we let skip equal to 5. We chose this skip mainly to reach a satisfactory computation time. The impact of the skip parameter in TDE on MP and computation time is expounded upon in Sec. S.3.1.

Once the parameters were set, the time series were transformed into point clouds. If the number $|P|$ of points in a point cloud was less than 40, we excluded this time series from further analysis, considering that there were too few points to represent the original structure of the time series. The problem of lacking points is also discussed in Sec. 3.

4.3.3 Computing persistent homology

Using Ripser [70, 71], we could compute the PDs of the point clouds in a fast and efficient way. We then extracted MP from each 1-dimensional PD, using persistence birth time and lifetime as two features of a time series. The process of vectorising a PD presents a challenge due to the indeterminate (and potentially large) number of intervals in the barcode, coupled with the ambiguous information they contain. This ambiguity arises from our lack of knowledge about the types of information that can be derived from different parts of the PD. Here we only extracted the MP and corresponding birth time. This decision was informed by our prior selection of an appropriate set of parameters, which ensured that the MP reached its optimal.

5 DATA AND CODE AVAILABILITY

The data that support the findings of this study are openly available in SpeechBox [67], ALLSTAR Corpus, L1-ENG division at <https://speechbox.linguistics.northwestern.edu>.

The source code and supplementary materials for Top-Cap can be accessed on the GitHub page at https://github.com/AnnFeng233/TDA_Consonant_Recognition.

REFERENCES

- [1] Gunnar Carlsson. "Topology and data". In: *Bulletin of The American Mathematical Society* 46 (Apr. 2009), pp. 255–308. DOI: 10.1090/S0273-0979-09-01249-X.
- [2] Ephy R. Love et al. "Topological convolutional layers for deep learning". In: *Journal of Machine Learning Research* 24.59 (2023), pp. 1–35.
- [3] Gunnar Carlsson and Rickard Brül Gabrielsson. "Topological approaches to deep learning". In: *Topological Data Analysis: The Abel Symposium 2018*. Springer. 2020, pp. 119–146.
- [4] Ken W Grant, Brian E Walden, and Philip F Seitz. "Auditory-visual speech recognition by hearing-impaired subjects: Consonant recognition, sentence recognition, and auditory-visual integration". In: *The Journal of the Acoustical Society of America* 103.5 (1998), pp. 2677–2690. ISSN: 0001-4966.
- [5] Yichen Shen et al. "Deep learning with coherent nanophotonic circuits". In: *Nature Photonics* 11.7 (2017), pp. 441–446. ISSN: 1749-4893.
- [6] Eric W Healy et al. "Speech-cue transmission by an algorithm to increase consonant recognition in noise for hearing-impaired listeners". In: *The Journal of the Acoustical Society of America* 136.6 (2014), pp. 3325–3336. ISSN: 0001-4966.
- [7] Thierry Nazzi and Anne Cutler. "How consonants and vowels shape spoken-language recognition". In: *Annual Review of Linguistics* 5 (2019), pp. 25–47. ISSN: 2333-9683.
- [8] Dianne J Van Tasell et al. "Speech waveform envelope cues for consonant recognition". In: *The Journal of the Acoustical Society of America* 82.4 (1987), pp. 1152–1161. ISSN: 0001-4966.
- [9] DeLiang Wang and Guoning Hu. "Unvoiced speech segregation". In: *2006 IEEE International Conference on Acoustics Speech and Signal Processing Proceedings*. Vol. 5. 2006, pp. V–V. DOI: 10.1109/ICASSP.2006.1661435.
- [10] Philip Weber et al. "Consonant recognition with continuous-state hidden Markov models and perceptually-motivated features". In: *Sixteenth Annual Conference of the International Speech Communication Association*. 2015.
- [11] Gunnar Carlsson et al. "Persistence barcodes for shapes". In: *International Journal of Shape Modeling* 11.02 (2005), pp. 149–187. DOI: 10.1142/S0218654305000761.
- [12] Oleksandr Balabanov and Mats Granath. "Unsupervised learning using topological data augmentation". In: *Physical Review Research* 2.1 (2020), p. 013354.
- [13] Azadeh Hadadi et al. "Prediction of cybersickness in virtual environments using topological data analysis and machine learning". In: *Frontiers in Virtual Reality* 3 (2022), p. 973236. ISSN: 2673-4192.
- [14] Firas A. Khasawneh, Elizabeth Munch, and Jose A. Perea. "Chatter classification in turning using machine learning and topological data analysis". In: *IFAC-PapersOnLine* 51.14 (2018), pp. 195–200. ISSN: 2405-8963. DOI: 10.1016/j.ifacol.2018.07.222.
- [15] Daniel Leykam and Dimitris G. Angelakis. "Topological data analysis and machine learning". In: *Advances in Physics: X* 8.1 (2023). ISSN: 2374-6149. DOI: 10.1080/23746149.2023.2202331.
- [16] Grzegorz Muszynski et al. "Topological data analysis and machine learning for recognizing atmospheric river patterns in large climate datasets". In: *Geoscientific Model Development* 12.2 (2019), pp. 613–628. ISSN: 1991-9603. DOI: 10.5194/gmd-12-613-2019.
- [17] Frédéric Chazal and Bertrand Michel. "An introduction to topological data analysis: Fundamental and practical aspects for data scientists". In: *Frontiers in Artificial Intelligence* 4 (2021), p. 108. ISSN: 2624-8212.
- [18] Aras Asaad and Sabah Jassim. "Topological data analysis for image tampering detection". In: *Digital Forensics and Watermarking: 16th International Workshop, IWDW 2017, Magdeburg, Germany, August 23-25, 2017, Proceedings 16*. Springer. 2017, pp. 136–146.
- [19] Lander Ver Hoef et al. "A primer on topological data analysis to support image analysis tasks in environmental science". In: *Artificial Intelligence for the Earth Systems* 2.1 (2023), e220039.
- [20] Gunnar Carlsson et al. "On the local behavior of spaces of natural images". In: *International Journal of Computer Vision* 76 (Jan. 2008), pp. 1–12. DOI: 10.1007/s11263-007-0056-x.
- [21] Sebastian Zeng et al. "Topological attention for time series forecasting". In: *Advances in Neural Information Processing Systems* 34 (2021), pp. 24871–24882.
- [22] Yuhei Umeda. "Time series classification via topological data analysis". In: *Information and Media Technologies* 12 (2017), pp. 228–239. ISSN: 1881-0896.
- [23] Manish Saggarr et al. "Towards a new approach to reveal dynamical organization of the brain using topological data analysis". In: *Nature Communications* 9.1 (2018), p. 1399. ISSN: 2041-1723. DOI: 10.1038/s41467-018-03664-4.
- [24] Jessica L Nielson et al. "Uncovering precision phenotype-biomarker associations in traumatic brain injury using topological data analysis". In: *PloS One* 12.3 (2017), e0169490. ISSN: 1932-6203.
- [25] Tamal K Dey and Sayan Mandal. "Protein classification with improved topological data analysis". In: *18th International Workshop on Algorithms in Bioinformatics (WABI 2018)*. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik. 2018.
- [26] Alessio Martino, Antonello Rizzi, and Fabio Massimo Frattale Mascioli. "Supervised approaches for protein function prediction by topological data analysis". In: *2018 International Joint Conference on Neural Networks (IJCNN)*. IEEE. 2018, pp. 1–8.
- [27] Kenneth Brown and Kevin Knudson. "Nonlinear statistics of human speech data". In: *International Journal of Bifurcation and Chaos* 19 (July 2009), pp. 2307–2319. DOI: 10.1142/S0218127409024086.
- [28] Sergio Barbarossa and Stefania Sardellitti. "Topological signal processing over simplicial complexes". In: *IEEE Transactions on Signal Processing* 68 (2020), pp. 2992–3007.

- [29] Eduard Tulchinskii et al. "Topological data analysis for speech processing". In: *Proc. Interspeech 2023*, pp. 311–315. DOI: 10.21437/Interspeech.2023-1861.
- [30] Deli Chen et al. "Measuring and relieving the over-smoothing problem for graph neural networks from the topological view". In: *Proceedings of the AAAI Conference on Artificial Intelligence*. Vol. 34. 2020, pp. 3438–3445.
- [31] Woong Bae, Jaejun Yoo, and Jong Chul Ye. "Beyond deep residual learning for image restoration: Persistent homology-guided manifold simplification". In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*. 2017, pp. 145–153.
- [32] Christoph Hofer et al. "Deep learning with topological signatures". In: *Advances in Neural Information Processing Systems* 30 (2017).
- [33] Emilie Gerardin et al. "Multidimensional classification of hippocampal shape features discriminates Alzheimer's disease and mild cognitive impairment from normal aging". In: *Neuroimage* 47.4 (2009), pp. 1476–1486. ISSN: 1053-8119.
- [34] Mingqiang Yang, Kidiyo Kpalma, and Joseph Rossin. "A survey of shape feature extraction techniques". In: *Pattern Recognition* 15.7 (2008), pp. 43–90.
- [35] Zenggang Xiong et al. "Research on image retrieval algorithm based on combination of color and shape features". In: *Journal of Signal Processing Systems* 93 (2021), pp. 139–146. ISSN: 1939-8018.
- [36] Dengsheng Zhang and Guojun Lu. "Review of shape representation and description techniques". In: *Pattern Recognition* 37.1 (2004), pp. 1–19. ISSN: 0031-3203. DOI: 10.1016/j.patco.2003.07.008.
- [37] John Lee. *Introduction to topological manifolds*. Vol. 202. Springer Science & Business Media, 2010. ISBN: 1441979409.
- [38] Allen Hatcher. *Algebraic topology*. Cambridge: Cambridge University Press, 2002.
- [39] Afra Zomorodian and Gunnar Carlsson. "Computing persistent homology". In: *Discrete & Computational Geometry* 33.2 (2005), pp. 249–274. ISSN: 1432-0444. DOI: 10.1007/s00454-004-1146-y.
- [40] Herbert Edelsbrunner and John Harer. "Persistent homology – a survey". In: *Contemporary Mathematics* 453.26 (2008), pp. 257–282.
- [41] Robert Ghrist. "Barcodes: The persistent topology of data". In: *Bulletin of the American Mathematical Society* 45.1 (2008), pp. 61–75. ISSN: 0273-0979.
- [42] David Cohen-Steiner, Herbert Edelsbrunner, and John Harer. "Stability of persistence diagrams". In: *Proceedings of the Twenty-First Annual Symposium on Computational geometry*. 2005, pp. 263–271.
- [43] Charu C Aggarwal, Alexander Hinneburg, and Daniel A Keim. "On the surprising behavior of distance metrics in high dimensional space". In: *Database Theory—ICDT 2001: 8th International Conference London, UK, January 4–6, 2001 Proceedings 8*. Springer. 2001, pp. 420–434.
- [44] Vladislav Polianskii and Florian T. Pokorný. "Voronoi graph traversal in high dimensions with applications to topological data analysis and piecewise linear interpolation". In: *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. KDD '20. Virtual Event, CA, USA: Association for Computing Machinery, 2020, pp. 2154–2164. ISBN: 9781450379984. DOI: 10.1145/3394486.3403266.
- [45] Matthew B Kennel, Reggie Brown, and Henry DI Abarbanel. "Determining embedding dimension for phase-space reconstruction using a geometrical construction". In: *Physical Review A* 45.6 (1992), p. 3403.
- [46] Baihan Lin. "Topological data analysis in time series: Temporal filtration and application to single-cell genomics". In: *Algorithms* 15.10 (2022), p. 371. ISSN: 1999-4893.
- [47] Saba Emrani, Thanos Gentimis, and Hamid Krim. "Persistent homology of delay embeddings and its application to wheeze detection". In: *IEEE Signal Processing Letters* 21.4 (2014), pp. 459–463. ISSN: 1070-9908.
- [48] Cássio MM Pereira and Rodrigo F de Mello. "Persistent homology for time series and spatial data clustering". In: *Expert Systems with Applications* 42.15-16 (2015), pp. 6026–6038. ISSN: 0957-4174.
- [49] Firas A Khasawneh and Elizabeth Munch. "Chatter detection in turning using persistent homology". In: *Mechanical Systems and Signal Processing* 70 (2016), pp. 527–541. ISSN: 0888-3270.
- [50] Lee M Seversky, Shelby Davis, and Matthew Berger. "On time-series topological data analysis: New data and opportunities". In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*. 2016, pp. 59–67.
- [51] Marian Gidea and Yuri Katz. "Topological data analysis of financial time series: Landscapes of crashes". In: *Physica A: Statistical Mechanics and its Applications* 491 (2018), pp. 820–834. ISSN: 0378-4371.
- [52] Richard Bellman. "Dynamic programming". In: *Science* 153.3731 (1966), pp. 34–37.
- [53] Afra Zomorodian and Gunnar Carlsson. "Computing persistent homology". In: *Discrete & Computational Geometry* 33.2 (Feb. 2005), pp. 249–274. ISSN: 1432-0444. DOI: 10.1007/s00454-004-1146-y.
- [54] Nikola Milosavljević, Dmitriy Morozov, and Primoz Skraba. "Zigzag persistent homology in matrix multiplication time". In: *Proceedings of the Twenty-Seventh Annual Symposium on Computational Geometry*. 2011, pp. 216–225.
- [55] Peter Bubenik. "Statistical topological data analysis using persistence landscapes". In: *Journal of Machine Learning Research* 16.1 (2015), pp. 77–102.
- [56] Henry Adams et al. "Persistence images: A stable vector representation of persistent homology". In: *Journal of Machine Learning Research* 18 (2017).
- [57] D. Ali et al. "A survey of vectorization methods in topological data analysis". In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2023), pp. 1–14. ISSN: 1939-3539. DOI: 10.1109/TPAMI.2023.3308391.
- [58] James D Hamilton. *Time series analysis*. Princeton University Press, 2020.
- [59] Yu-Min Chung et al. "A persistent homology approach to heart rate variability analysis with an

- 1069 application to sleep-wake classification". In: *Frontiers*
 1070 *in Physiology* 12 (2021), p. 637684. ISSN: 1664-042X. 1128
- 1071 [60] Jose A Perea. "Topological time series analysis". 1129
 1072 In: *Notices of the American Mathematical Society* 66.5
 1073 (2019), pp. 686–694.
- 1074 [61] Jose A. Perea and John Harer. "Sliding windows and
 1075 persistence: An application of topological methods
 1076 to signal analysis". In: *Foundations of Computational*
 1077 *Mathematics* 15.3 (June 2015), pp. 799–838. ISSN: 1615-
 1078 3383. DOI: 10.1007/s10208-014-9206-z.
- 1079 [62] Jose Perea et al. "SW1PerS: Sliding windows and 1-
 1080 persistence scoring; discovering periodicity in gene
 1081 expression time series data". In: *BMC bioinformatics*
 1082 16 (Aug. 2015), p. 257. DOI: 10.1186 / s12859-015-
 1083 0645-6.
- 1084 [63] Christopher Tralie and Jose Perea. "(Quasi) periodic-
 1085 ity quantification in video data, using topology". In:
 1086 *SIAM Journal on Imaging Sciences* 11 (Apr. 2017). DOI:
 1087 10.1137/17M1150736.
- 1088 [64] Hitesh Gakhar and Jose A. Perea. "Sliding window
 1089 persistence of quasiperiodic functions". In: *Journal of*
 1090 *Applied and Computational Topology* (2023). DOI: 10.
 1091 1007/s41468-023-00136-7.
- 1092 [65] NOAA. *Climate at a glance: National time series*. Na-
 1093 tional Centers for Environmental Information, re-
 1094 trieved 28th January 2023 from <https://www.citedrive.com/overleaf>.
- 1095 [66] Holger Kantz and Thomas Schreiber. *Nonlinear time*
 1096 *series analysis*. Vol. 7. Cambridge University Press,
 1097 2004.
- 1098 [67] A. R. Bradlow. SpeechBox, retrieved from [https://](https://speechbox.linguistics.northwestern.edu)
 1100 speechbox.linguistics.northwestern.edu.
- 1101 [68] Michael McAuliffe et al. "Montreal Forced Aligner:
 1102 Trainable text-speech alignment using Kaldi". In:
 1103 *Proc. Interspeech 2017*, pp. 498–502. DOI: 10.21437/
 1104 Interspeech.2017-1386.
- 1105 [69] Paul Boersma and David Weenink. *Praat: Doing*
 1106 *phonetics by computer*. Version 6.3.09, retrieved 2nd
 1107 March 2023 from <http://www.praat.org/>. 2023.
- 1108 [70] Christopher Tralie, Nathaniel Saul, and Rann Bar-On.
 1109 "Ripser.py: A Lean persistent homology library for
 1110 Python". In: *The Journal of Open Source Software* 3.29
 1111 (Sept. 2018), p. 925. DOI: 10.21105/joss.00925.
- 1112 [71] Ulrich Bauer. "Ripser: Efficient computation of
 1113 Vietoris-Rips persistence barcodes". In: *Journal of Ap-*
 1114 *plied and Computational Topology* 5.3 (2021), pp. 391–
 1115 423. DOI: 10.1007/s41468-021-00071-5.

1116 ACKNOWLEDGEMENTS

1117 The authors would like to thank Meng Yu for his invaluable
 1118 mentorship on audio and speech signal processing, on neu-
 1119 ral networks and deep learning, as well as for his comments
 1120 and suggestions on an earlier draft of this manuscript.
 1121 The authors would also like to thank Andrew Blumberg,
 1122 Fangyi Chen, Haibao Duan, Houhong Fan, Fuquan Fang,
 1123 Wenfei Jin, Fengchun Lei, Yanlin Li, Andy Luchuan Liu,
 1124 Tao Luo, Zhi Lü, Jianxin Pan, Jie Wu, Kelin Xia, Jiang
 1125 Yang, Jin Zhang and Zhen Zhang for helpful discussions
 1126 and encouragement. This work was partly supported by
 1127 the National Natural Science Foundation of China grant

12371069 and the Guangdong Provincial Key Laboratory of
 Interdisciplinary Research and Application for Data Science.

AUTHOR INFORMATION

Authors and Affiliations

Department of Mathematics, Southern University of Sci-
 ence and Technology, Shenzhen, China

Pingyao Feng, Siheng Yi, Qingrui Qu, Zhiwang Yu, Yifei
 Zhu

Contributions

Y.Z. planned the project. P.F. and S.Y. constructed the the-
 oretical framework. P.F. designed the sample, built the
 algorithms, and analysed the data. S.Y. assisted with the
 algorithms. P.F., S.Y., Q.Q., Z.Y., and Y.Z. wrote the paper
 and contributed to the discussion.

Corresponding authors

Correspondence to Pingyao Feng or Yifei Zhu.

SUPPLEMENTARY INFORMATION

S.1 Generalities on phonetic data

As a research field of linguistics, phonetics studies the production as well as the classification of human speech sounds from the world’s languages. In phonetics, a *phone* is the smallest basic unit of human speech sounds. It is a short speech segment possessing distinct physical or perceptual properties. Phones are generally classified into two principal categories: vowels and consonants. A *vowel* is defined as a speech sound pronounced by an open vocal tract with no significant build-up of air pressure at any point above the glottis, and at least making some airflow escape through the mouth. In contrast, a *consonant* is a speech sound that is articulated with a complete or partial closure of the vocal tract and usually forces air through a narrow channel in one’s mouth or nose.

Unlike vowels which must be pronounced by vibrated vocal cords, consonants can be further categorised into two classes according to whether the vocal cords vibrate or not during articulation. If the vocal cords vibrate, the consonant is known as a *voiced* consonant. Otherwise, the consonant is *voiceless*. Since vocal cord vibration can produce a stable periodic signal of air pressure, voiced consonants tend to have more periodic components than voiceless consonants, which can in turn be detected by PH as topological characteristics from phonetic time series data.

Indeed, one of the more heuristic motivations for our research project is to reexamine (and even revise) the linguistic classifications of phones through the mathematical lens of topological patterns and shape of speech data, analogous to Carlsson and his collaborators’ seminal work [S1] on the distribution of image data (cf. Fig. S1).

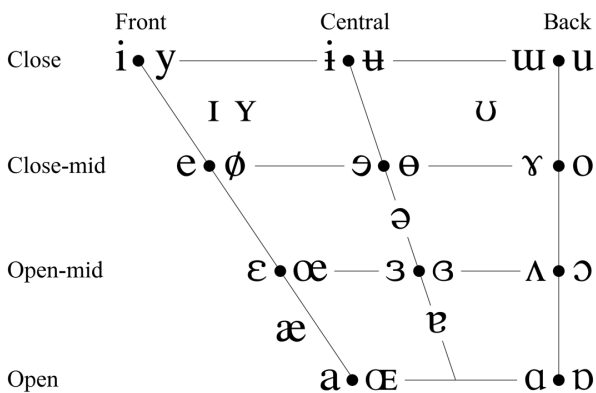


Fig. S1: A charted “distribution space” of vowels created by linguists [S2]. The vertical axis of the chart denotes vowel height. Vowels pronounced with the tongue lowered are located at the bottom and those raised are at the top. The horizontal axis of this chart denotes vowel backness. Vowels with the tongue moved towards the front of the mouth are in the left of the chart, while those with to the back are placed in the right. The last parameter is whether the lips are rounded. At each given spot, vowels on the right and left are rounded and unrounded, respectively.

S.2 Mathematical generalities of the TDE–PH approach to time series data

S.2.1 Time-delay embedding

Time-delay embedding (TDE) is also known as sliding window embedding, delay embedding, and delay coordinate embedding. For simplicity, we focus on 1-dimensional time series. TDE of a real-valued function $f: \mathbb{R} \rightarrow \mathbb{R}$, with parameters positive integer d and positive real number τ , is defined to be the vector-valued function

$$SW_{d,\tau}f: \mathbb{R} \rightarrow \mathbb{R}^d \\ t \mapsto \left(f(t), f(t+\tau), \dots, f(t+(d-1)\tau) \right)$$

Here, d is the *dimension* of the target space for the embedding, τ is the *delay*, and their product $d \cdot \tau$ is called the *window size*. According to the Manifold Hypothesis, a time series lies on a manifold. The method then reconstructs this topological space from the input time series, when d is at least twice the dimension of the latent manifold M . Given a trajectory $\gamma: \mathbb{R} \rightarrow M$ whose image is dense in M , the embedding property holds for the time series $f(t_n)$ (generically, in a technical sense we omit here) via an “observation” function $G: M \rightarrow \mathbb{R}$, i.e., $f(t_n) = G(\gamma(t_n))$.

In [S3, Sec. 5], Perea and Harer established that the N -truncated Fourier series expansion

$$S_N f(t) = \sum_{n=0}^N a_n \cos(kt) + b_n \sin(kt)$$

of a periodic time series f can be reconstructed into a circle when $d \geq 2N$, i.e.,

$$SW_{d,\tau}f(\mathbb{R}) \cong \mathbb{S}^1$$

Moreover, let L be a constant such that

$$f\left(t + \frac{2\pi}{L}\right) = f(t)$$

Then the 1-dimensional MP of the resulting point cloud is the largest when the window size $d \cdot \tau$ is integrally proportional to $2\pi/L$, i.e.,

$$d \cdot \tau = m \frac{2\pi}{L}$$

for a positive integer m . Intuitively, an increase in the dimension of TDE results in a better approximation when truncating the Fourier series, and the MP of the point cloud becomes the most significant when the window size equals a period.

This methodology also proves particularly advantageous in scenarios where the system under investigation exhibits nonlinear dynamics, precluding straightforward analysis of the time series data. Via a suitable embedding, the inherent geometric configuration of the system emerges, enabling deeper comprehension and refined analysis.

S.2.2 Persistent homology

Topology is a subject area that studies the properties of geometric objects that remain unchanged under continuous transformations or smooth perturbations. It focuses on the intrinsic features of a space that regardless of its rigid shape or size. Algebraic topology (AT) provides a quantitative description of these topological properties.

1221 A simplicial complex (and its numerous variants and
 1222 analogues) is a powerful tool in AT which enables us to
 1223 represent a topological space using discrete data. Unlike
 1224 the original space, which can be challenging to compute
 1225 and analyse, a simplicial complex provides a combinatorial
 1226 description that is much more amenable to computation.
 1227 We can use algebraic techniques to study the properties of a
 1228 simplicial complex, such as its homology and cohomology
 1229 groups, which encode and reveal information about the
 1230 topology of the underlying space.

1231 Formally, a *simplicial complex* with *vertices* in a set V is
 1232 a collection K of nonempty finite subsets $\sigma \subset V$ such that
 1233 any nonempty subset τ of σ always implies $\tau \in K$ (called a
 1234 *face* of σ) and that σ intersecting σ' implies their intersection
 1235 $\sigma \cap \sigma' \in K$. A set $\sigma \in K$ with $(i + 1)$ elements is called an
 1236 i -*simplex* of the simplicial complex K . For instance, consider
 1237 $\mathbb{S}^1 \vee \mathbb{S}^2$, a circle kissing a sphere at a single point, as a
 1238 topology space. It can be approximated by the simplicial
 1239 complex K with 6 vertices a, b, c, d, e, f . This simplicial
 1240 complex can be enumerated as

$$K = \{\{a\}, \{b\}, \{c\}, \{d\}, \{e\}, \{f\}, \\ \{a, b\}, \{a, c\}, \{b, c\}, \{c, d\}, \{c, f\}, \{d, f\}, \{c, e\}, \\ \{d, e\}, \{f, e\}, \\ \{c, d, f\}, \{c, e, f\}, \{c, d, e\}, \{d, e, f\}\}$$

1241 which is a combinatorial avatar for $\mathbb{S}^1 \vee \mathbb{S}^2$ via a “triangulation”
 1242 operation on the latter. See Fig. S2.

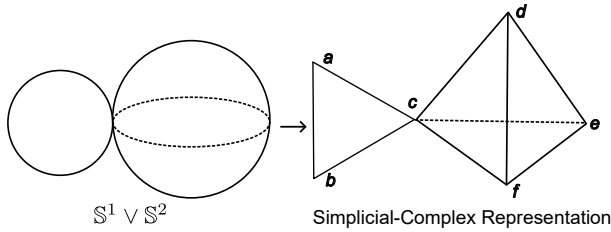


Fig. S2: From a topological space to its triangulation.

1243 Given a simplicial complex K , let p be a prime number
 1244 and \mathbb{F}_p be the finite field with p elements. Define $C_i(K; \mathbb{F}_p)$
 1245 to be the \mathbb{F}_p -vector space with basis the set of i -simplices in
 1246 K . To keep track of the order of vertices within a simplex,
 1247 we use the alternative notation with square brackets in the
 1248 following. If $\sigma = [v_0, v_1, \dots, v_i]$ is an i -simplex, define the
 1249 *boundary* of σ , denoted by $\partial\sigma$, to be the alternating sum of
 1250 the $(i - 1)$ -dimensional faces of σ given by

$$\partial\sigma := \sum_{k=0}^i (-1)^k [v_0, \dots, \hat{v}_k, \dots, v_i]$$

1251 where $[v_0, \dots, \hat{v}_k, \dots, v_i]$ is the k -th $(i - 1)$ -dimensional
 1252 face of σ missing the vertex v_k . We can extend ∂ to $C_i(K; \mathbb{F}_p)$
 1253 as an \mathbb{F}_p -linear operator so that $\partial: C_i(K; \mathbb{F}_p) \rightarrow C_{i-1}(K; \mathbb{F}_p)$.
 1254 The composition of boundary operators satisfies $\partial \circ \partial = 0$.
 1255 The elements in $C_i(K; \mathbb{F}_p)$ with boundary 0 are called i -
 1256 *cycles*. They form a subspace of $C_i(K; \mathbb{F}_p)$, denoted by
 1257 $Z_i(K; \mathbb{F}_p)$. The elements in $C_i(K; \mathbb{F}_p)$ that are the images
 1258 of elements of $C_{i+1}(K; \mathbb{F}_p)$ under ∂ are called i -*boundaries*.

They form a subspace too, denoted by $B_i(K; \mathbb{F}_p)$. It follows
 from $\partial \circ \partial = 0$ that

$$B_i(K; \mathbb{F}_p) \subset Z_i(K; \mathbb{F}_p)$$

Then define the quotient space

$$H_i(K; \mathbb{F}_p) := Z_i(K; \mathbb{F}_p) / B_i(K; \mathbb{F}_p)$$

to be the i -th homology group of K with \mathbb{F}_p -coefficients. We call
 $\dim(H_i(K; \mathbb{F}_p))$ the i -th Betti number, denoted by $\beta_i(K)$,
 which counts the number of i -dimensional holes in the
 corresponding topological space. As such, these homology
 groups are also called the homology groups of the space (it
 can be shown that they are independent of the particular
 ways in which the space is triangulated). For example, the
 Betti numbers of $\mathbb{S}^1 \vee \mathbb{S}^2$ from above are $\beta_1 = 1, \beta_2 = 1$, and
 $\beta_i = 0$ when $i \geq 3$.

The usefulness of these invariants, besides their com-
 putability (essentially Gaussian elimination in linear algebra),
 lies in their tractability along deformations. Given two
 simplicial complexes K and L , a simplicial map $f: K \rightarrow L$
 (that preserves the simplicial structure) induces an \mathbb{F}_p -linear
 map $H_i(f; \mathbb{F}_p): H_i(K; \mathbb{F}_p) \rightarrow H_i(L; \mathbb{F}_p)$. Thus, if two spaces
 are topologically equivalent (in fact, “homotopy equivalent”
 suffices), their homology groups must be isomorphic and
 the Betti numbers match up.

Let (X, d) be a finite point cloud with metric d . Define a
 family of simplicial complexes, called *Rips complexes*, by

$$R_\epsilon(X) := \{\sigma \subset X \mid d(x, x') \leq \epsilon \text{ for all } x, x' \in \sigma\}$$

The family

$$\mathcal{R}(X) := \{R_\epsilon(X)\}_{\epsilon \geq 0}$$

is known as the Rips filtration of X . Clearly, if $\epsilon_1 \leq \epsilon_2$, then
 $R_{\epsilon_1}(X) \subset R_{\epsilon_2}(X)$. Thus, for each i we obtain a sequence

$$H_i(R_{\epsilon_0}(X); \mathbb{F}_p) \rightarrow H_i(R_{\epsilon_1}(X); \mathbb{F}_p) \rightarrow \dots \\ \rightarrow H_i(R_{\epsilon_m}(X); \mathbb{F}_p)$$

where $0 = \epsilon_0 < \epsilon_1 < \dots < \epsilon_m < \infty$. As ϵ varies, the
 topological features in the simplicial complexes $R_\epsilon(X)$ vary,
 resulting in the emergence and disappearance of holes.

Given the values of ϵ , record the instances of emergence
 and disappearance of holes, which correspond to cycle
 classes in the homology groups along the above sequence.
 Each class has a descriptor $(b, d) \in \mathbb{R}^2$, where b represents
 the *birth time*, d represents the *death time*, and $b - d$ represents
 the *lifetime* of the holes. In this way, we obtain a multiset

$$\{(b_j, d_j)\}_{j \in J} =: \text{dgm}_i(\mathcal{R}(X))$$

which encodes the “persistence” of topological features of
 X . This multiset can be represented as a multiset of points
 in the 2-dimensional coordinate system called a *persistence*
diagram for the i -th PH or as an array of interval segments
 called a *persistence barcode*. In particular, we use *maximal*
persistence to refer to the maximal lifetime among all the
 points in a persistence diagram.

dimension = 10 desired delay = 40			dimension = 50 desired delay = 8			dimension = 100 desired delay = 4		
delay	skip	MP	delay	skip	MP	delay	skip	MP
1	1	0.0610	1	1	0.2834	1	1	0.4270
10	1	0.1299	3	1	0.3021	2	1	0.4337
20	1	0.1312	4	1	0.3054	2	5	0.4146
30	1	0.1281	5	1	0.3058	3	1	0.4357
39	1	0.1229	6	1	0.3042	3	5	0.4120
39	5	0.1134	7	1	0.3052	4	1	0.4381
40	1	0.1290	7	5	0.2886	4	5	0.4139
40	5	0.1195	8	1	0.3093	5	1	0.4375
41	1	0.1200	8	5	0.2928	5	5	0.4105
41	5	0.1153	9	1	0.3091	6	1	0.4347
45	1	0.0940	9	5	0.2913	6	5	0.4114
50	1	0.1226	10	1	0.3069	7	1	0.4380
60	1	0.1315	15	1	0.3070	8	1	0.4378
94	1	empty	18	1	empty	9	1	empty

Tab. S1: MP for choices of dimension, delay, and skip in TDE. The desired delay is computed by the algorithm in Sec. 4 of Methods. Empty in MP means the delay is too large to obtain point-cloud data.

S.3 More specifics on parameter selection with TopCap

S.3.1 Skip, maximal persistence, and persistence execution time

Computation time assumes a critical role when processing a substantial volume of data. In this context, the parameter skip in TDE is considered, as it significantly influences the number of points within the point clouds, thereby directly impacting the number of simplices during persistent filtration and thus the computation time for PD. In this subsection, we demonstrate that an appropriate increment in the skip parameter can markedly reduce computation time. However, it is noteworthy that MP exhibits resilience to an increase in skip to a certain extent. Consequently, in this case, it is feasible to augment skip in TDE to expedite the computation of PD. For details on the complexity of computing persistent homology, the interested reader may refer to Zomorodian and Carlsson [S4, Sec. 4.3] as well as Edelsbrunner et al. [S5, Sec. 4].

Using an example of a sound record of the voiced consonant [m], we elucidate the relationship between skip, computation duration, and size of the resulting point clouds obtained via TDE in Fig. 6d. Computation duration is measured each time after restarting the Jupyter notebook, on Dell Precision 3581, with CPU Intel® Core™ i7-13800H of basic frequency 2.50 GHz and 14 cores. Computation time means the time for executing the code `ripser(Points, maxdim=1)`. As depicted in Fig. 6d, a substantial reduction in computation time is observed with an increase in the skip parameter. In contrast, our computation's output MP appears stable.

S.3.2 Multiple dependency of maximal persistence

As mentioned in the main text, there are three crucial parameters in TDE, namely, d , τ , and skip. In this subsection, we present a table that delineates the topological descriptor MP in relation to these from TopCap.

The experiment is executed on a record of the voiced consonant [ŋ], which comprises 887 sampled points as the length of this time series. Theoretically, given a periodic function, one obtains the optimal MP of the function in a fixed dimension under the condition that the TDE window size (i.e., the product of dimension and delay) equals a period (cf. Sec. S.2.1). However, the phonetic time series

that we typically handle deviate far from being periodic. Despite our approach to calculating the period of time series by ACL functions, we cannot assure that the (theoretically derived) desired delay will indeed yield the optimal MP of a time series in general. Nevertheless, this desired delay usually gives relatively good MP. For instance, as illustrated in Tab. S1, when the dimension is 10, the desired delay is 40. This corresponds to an MP of 0.1290, which is marginally lower than the MP of 0.1315 achieved at a delay of 60. However, as the dimension rises, the point clouds from TDE become more regular. It becomes increasingly probable that at the desired delay, one can indeed obtain the optimal MP of the time series. For example, when the dimension is either 50 or 100, the MP of the time series is achieved at the desired delay. This provides additional justification for preferring higher dimensions: The table reveals that an augmentation in dimension may lead to a more substantial enhancement in the MP of a time series than simply tuning delay.

S.4 Review and outlook on topology-enhanced machine learning

Here we present a general review of literature on the topics (1) TDA and its applications, which encompasses genesis of the subject, recommended resources, and practical applications; (2) vectorisation of PH, wherein we summarize topological methods geared towards machine learning.

S.4.1 Topological data analysis and its applications

The evolution of TDA is relatively nascent when juxtaposed with other enduring fields, and its applications are still somewhat delimited. The genesis of the concept of invariants of filtered complexes can be traced back to Baranikov in 1994, which are nowadays referred to as PD/PB (persistence diagram/barcode) [S6]. These invariants were conceived with the objective of quantifying some specific critical point within some ambit of an extension of function. In 1999, Robins pioneered the concept of *persistent Betti numbers* of inverse systems and underscored their stability in Hausdorff distance [S7].

The modern incarnation of persistent homology was established in the first decade of the 21st century. Zomorodian, under the tutelage of Edelsbrunner, completed his doctoral

thesis in 2001, wherein he employed persistence to distinguish between topological noise and inherent features of a space [S8]. After that, the term *persistent homology group* first appeared in the work by Edelsbrunner et al. in 2002 [S9]. This seminal work formalised topological methodologies to chronicle the evolution of an expanding complex originating from a point set in Euclidean 3-space, a process they termed as topological simplification. The expansion process is recognised as filtration. They classified topological modifications based on the lifetime of topological features during filtration and proposed an algorithm to compute this simplification process. Subsequently, in 2005, Carlsson et al. applied persistent homology to generate a barcode as a shape descriptor [S10]. Their methodology was able to distinguish between shapes with varying degrees of “sharp” features, such as corners. In the same year, Zomorodian and Carlsson presented an algebraic interpretation of persistent homology and developed a natural algorithm for computing persistent homology of spaces in any dimension over any field [S11]. Cohen-Steiner et al. considered the stability property of persistence algorithm [S12]. Robustness is measured by the bottleneck distance between persistence diagrams.

In 2008, Carlsson, Singh, and Sexton founded Ayasdi, a company that combines mathematics and finance to truly put theory into practice. The inception of TDA may be complex, as it originates from some pure mathematical fields such as Morse theory and PH. However, the underlying principle remains steadfast: to identify topological features that can quantify the shape of the data to certain degrees, which is robust against noise and perturbations.

An abundance of materials is available that offer a thorough understanding of TDA for both specialists and general audience. In 2009, Carlsson wrote an extensive survey on the applications of geometry and topology to the analysis of various types of data [S13]. This work introduced topics such as the characteristics of topological methods, persistence, and clusters. A recent publication by Carlsson and Vejdemo-Johansson discussed practical case studies of topological methods, such as their applications to image data and time series [S14]. For nonspecialists seeking to delve into TDA, the introductory article [S15] by Chazal and Michel may be more accessible. It provides explicit explanations and hands-on guidance on both the theoretical and practical aspects of the subject.

Several software tools assist researchers in building case studies on data. The GUDHI library [S16], an open-source C++ library with a Python interface, includes a comprehensive set of tools involving different complexes and vectorisation tools. Ripser [S17], also a C++ library with a Python binding, surpasses GUDHI in computing Vietoris–Rips PD/PB, especially when high-dimensional cases or large quantities of PD/PB are present. TTK [S18] is both a library and software designed for topological analysis with a focus on scientific visualisation. Other standard libraries include Dionysus, PHAT, DIPHA, and Giotto². Additionally,

²In order, they are available at

<https://mrzv.org/software/dionysus2>
<https://bitbucket.org/phant-code/phant>
<https://github.com/DIPHA/dipha>
<https://giotto-ai.github.io/gtda-docs/0.4.0>

an R interface named TDA [S19] is available for the libraries GUDHI, Dionysus, and PHAT.

The recent proliferation of TDA has established it as an effective instrument in numerous studies. Owing to the characteristics of topological methods [S13], a multitude of applications have been discovered, particularly in the realm of recognition. In the field of biomedicine, Nicolau et al. utilised the topological method Mapper [S20] to analyse transcriptional data related to breast cancer [S21]. This method is used due to its high performance in shape recognition in high dimensions. The book [S22] authored by Rabadán and Blumberg provides an introduction to TDA techniques and their specific applications in biology, encompassing topics such as evolutionary processes and cancer genomics.

In signal processing, Emrani et al. introduced a topological approach for the analysis of breathing sound signals for the detection of wheezing, which can distinguish abnormal wheeze signals from normal breathing signals due to the periodic patterns within wheezing [S23]. Robinson’s monograph [S24] offers a systematic exploration of the intersection between topology and signal processing.

In the context of deep learning, Bae et al. proposed a PH-based deep residual learning algorithm for image restoration tasks [S25]. Hofer et al. incorporated topological signatures into deep neural networks to learn unusual structures that are typically challenging for most machine learning techniques [S26]. More recently, having extracted statistical features of images and videos through topological means, Love et al. input these features to the kernel of convolutional layers [S27, S1]. In their case, manifolds in relation to the natural-image space are used to parametrise image filters, which also parametrise slices in layers of neural networks. These signify a new phase of development for the subject.

For complex networks, an early application of PH on sensor networks is presented in the work [S28] by de Silva and Ghrist. They applied topological methods to graphs representing the distance estimation between nodes and a proximity sensor. Subsequently, Horak et al. discussed PH in different networks, observing that persistent topological attributes are related to the robustness of networks and reflect deficiencies in certain connectivity properties [S29]. Additionally, Jonsson’s book [S30] provides insights on how to construct a simplicial complex from a graph. Recently, Wu et al. applied a persistent variant of the GLMY homology for directed graphs of Grigor’yan, Lin, Muranov, and Yau to the study of networks of complex diseases [S31, S32].

S.4.2 Vectorising persistent homology for machine learning

When executing PH on point-cloud data, one typically obtains PD/PB, which is a set of intervals on the (extended real) line. Indeed, PD/PB can be considered a form of vectorisation of the original data. However, they may not be sufficiently accessible for further applications, such as integration into machine learning algorithms for future model development. Since the intervals exist on the extended line, some may involve $+\infty$ as their terminal point, which can pose challenges for certain algorithms. This issue can be mitigated by setting a threshold for the maximal lifetime, which is a relatively straightforward solution. However, there are more intrinsic challenges embedded in the vectorisation of

1498 PD/PB that are not easily resolved and may pose difficulties
 1499 for researchers attempting to leverage this powerful tool.
 1500 For example, the number of intervals in PD/PB is not fixed;
 1501 sometimes, there may be 10, and other times there may be
 1502 100. Moreover, PD is too sparse to put into machine learning
 1503 algorithms. Researchers may extract the top five longest
 1504 intervals from the set as a method of vectorisation, or
 1505 remove intervals with a length less than a certain threshold
 1506 from the set, or implement the distance functions and kernel
 1507 methods of PD/PB to achieve vectorisation. In this article,
 1508 vectorisation in TopCap is relatively simple, as we extract
 1509 the MP and its corresponding birth time as two topological
 1510 features to feed into machine learning algorithms.

1511 There is no definitive rule to determine that one method
 1512 of vectorisation is superior to another, as the performance
 1513 of vectorisation methods largely depends on the data and
 1514 how they are transformed into a topological space. Indeed,
 1515 there are a great many creative methods for vectorising PH.
 1516 Persistence Landscapes (PL) [S33], developed by Bubenik,
 1517 is one popular method. Bubenik’s work introduces both
 1518 theoretical and experimental aspects of PL in a statistical
 1519 manner. Generally speaking, PL maps PD into a function
 1520 space that is stable and invertible [S34]. A toolbox [S35] is
 1521 also available for implementing PL. Persistence Image [S36],
 1522 another vectorisation method developed by Adams et al.,
 1523 stably maps PD to a finite-dimensional vector representation
 1524 depending on resolution, weight function, and distribution
 1525 of points in PD. For additional vectorisation methods, one
 1526 may consider the article [S37] by Ali et al., which presents
 1527 13 ways to vectorise PD.

1528 References for supplementary information

- 1529
- 1530 [S1] Gunnar Carlsson et al. “On the local behavior of
 1531 spaces of natural images”. In: *International Journal*
 1532 *of Computer Vision* 76 (Jan. 2008), pp. 1–12. DOI: 10.
 1533 1007/s11263-007-0056-x.
- 1534 [S2] IPA Chart. *The international phonetic alphabet (re-*
 1535 *vised to 2020)*. International Phonetic Association,
 1536 retrieved 16th January 2024 from [https://www.](https://www.internationalphoneticassociation.org/content/ipa-chart)
 1537 [internationalphoneticassociation.org/content/ipa-](https://www.internationalphoneticassociation.org/content/ipa-chart)
 1538 [chart](https://www.internationalphoneticassociation.org/content/ipa-chart).
- 1539 [S3] Jose A. Perea and John Harer. “Sliding windows and
 1540 persistence: An application of topological methods
 1541 to signal analysis”. In: *Foundations of Computational*
 1542 *Mathematics* 15.3 (June 2015), pp. 799–838. ISSN: 1615-
 1543 3383. DOI: 10.1007/s10208-014-9206-z.
- 1544 [S4] Afra Zomorodian and Gunnar Carlsson. “Comput-
 1545 ing persistent homology”. In: *Discrete & Computa-*
 1546 *tional Geometry* 33.2 (Feb. 2005), pp. 249–274. ISSN:
 1547 1432-0444. DOI: 10.1007/s00454-004-1146-y.
- 1548 [S5] Edelsbrunner, Letscher, and Zomorodian. “Topolog-
 1549 ical persistence and simplification”. In: *Discrete &*
 1550 *Computational Geometry* 28.4 (Nov. 2002), pp. 511–533.
 1551 ISSN: 1432-0444. DOI: 10.1007/s00454-002-2885-2.
- 1552 [S6] Serguei Barannikov. “The framed Morse complex
 1553 and its invariants”. In: *Advances in Soviet Mathematics*.

- Singularities and Bifurcations 21 (Apr. 1994), pp. 93–
 116. DOI: 10.1090/advsov/021/03.
- [S7] Vanessa Robins. “Towards computing homology
 from finite approximations”. In: *Topology Proceedings*.
 Vol. 24. 1999, pp. 503–532.
- [S8] Afra Joze Zomorodian. *Computing and comprehending*
topology: Persistence and hierarchical Morse complexes.
 University of Illinois at Urbana-Champaign, 2001.
- [S9] Edelsbrunner, Letscher, and Zomorodian. “Topo-
 logical persistence and simplification”. In: *Discrete*
& Computational Geometry 28.4 (2002), pp. 511–533.
 ISSN: 1432-0444. DOI: 10.1007/s00454-002-2885-2.
- [S10] Gunnar Carlsson et al. “Persistence barcodes for
 shapes”. In: *International Journal of Shape Model-*
ing 11.02 (2005), pp. 149–187. DOI: 10 . 1142 /
 s0218654305000761.
- [S11] Afra Zomorodian and Gunnar Carlsson. “Comput-
 ing persistent homology”. In: *Discrete & Computa-*
tional Geometry 33.2 (2005), pp. 249–274. ISSN: 1432-
 0444. DOI: 10.1007/s00454-004-1146-y.
- [S12] David Cohen-Steiner, Herbert Edelsbrunner, and
 John Harer. “Stability of persistence diagrams”. In:
Proceedings of the twenty-first annual symposium on
Computational geometry. 2005, pp. 263–271.
- [S13] Gunnar Carlsson. “Topology and data”. In: *Bulletin*
of The American Mathematical Society 46 (Apr. 2009),
 pp. 255–308. DOI: 10.1090/S0273-0979-09-01249-X.
- [S14] Gunnar Carlsson and Mikael Vejdemo-Johansson.
Topological data analysis with applications. Cambridge
 University Press, 2021. ISBN: 1108983944.
- [S15] Frédéric Chazal and Bertrand Michel. “An introduc-
 tion to topological data analysis: Fundamental and
 practical aspects for data scientists”. In: *Frontiers in*
Artificial Intelligence 4 (2021), p. 108. ISSN: 2624-8212.
- [S16] Clément Maria et al. “The gudhi library: Simplicial
 complexes and persistent homology”. In: *Mathemat-*
ical Software–ICMS 2014: 4th International Congress,
Seoul, South Korea, August 5-9, 2014. Proceedings 4.
 Springer. 2014, pp. 167–174.
- [S17] Ulrich Bauer. “Ripser: Efficient computation of
 Vietoris-Rips persistence barcodes”. In: *Journal of Ap-*
plied and Computational Topology 5.3 (2021), pp. 391–
 423. DOI: 10.1007/s41468-021-00071-5.
- [S18] Julien Tierny et al. “The topology toolkit”. In: *IEEE*
Transactions on Visualization and Computer Graphics
 24.1 (2017), pp. 832–842. ISSN: 1077-2626.
- [S19] Brittany Terese Fasy et al. *Introduction to the R package*
TDA. 2015. arXiv: 1411.1830 [cs.LG].
- [S20] Gurjeet Singh, Facundo Mémoli, and Gunnar E
 Carlsson. “Topological methods for the analysis of
 high dimensional data sets and 3D object recog-
 nition”. In: *Eurographics Symposium on Point-Based*
Graphics 2 (2007), pp. 91–100.
- [S21] Monica Nicolau, Arnold J Levine, and Gunnar Carls-
 son. “Topology based data analysis identifies a sub-
 group of breast cancers with a unique mutational
 profile and excellent survival”. In: *Proceedings of the*
National Academy of Sciences 108.17 (2011), pp. 7265–
 7270. ISSN: 0027-8424.
- [S22] Raúl Rabadán and Andrew J Blumberg. *Topologi-*
cal data analysis for genomics and evolution: Topology

- 1615 *in biology*. Cambridge University Press, 2019. ISBN:
1616 1108753396.
- 1617 [S23] Saba Emrani, Thanos Gentimis, and Hamid Krim.
1618 “Persistent homology of delay embeddings and its
1619 application to wheeze detection”. In: *IEEE Signal*
1620 *Processing Letters* 21.4 (2014), pp. 459–463. ISSN: 1070-
1621 9908.
- 1622 [S24] Michael Robinson. *Topological signal processing*.
1623 Vol. 81. Springer, 2014.
- 1624 [S25] Woong Bae, Jaejun Yoo, and Jong Chul Ye. “Beyond
1625 deep residual learning for image restoration: Persis-
1626 tent homology-guided manifold simplification”. In:
1627 *Proceedings of the IEEE Conference on Computer Vision*
1628 *and Pattern Recognition Workshops*. 2017, pp. 145–153.
- 1629 [S26] Christoph Hofer et al. “Deep learning with topolog-
1630 ical signatures”. In: *Advances in Neural Information*
1631 *Processing Systems* 30 (2017).
- 1632 [S27] Ephy R. Love et al. “Topological convolutional layers
1633 for deep learning”. In: *Journal of Machine Learning*
1634 *Research* 24.59 (2023), pp. 1–35.
- 1635 [S28] Vin De Silva and Robert Ghrist. “Coverage in sensor
1636 networks via persistent homology”. In: *Algebraic &*
1637 *Geometric Topology* 7.1 (2007), pp. 339–358. ISSN: 1472-
1638 2739.
- 1639 [S29] Danijela Horak, Slobodan Maletić, and Milan Ra-
1640 jković. “Persistent homology of complex networks”.
1641 In: *Journal of Statistical Mechanics: Theory and Experi-*
1642 *ment* 2009.03 (2009), P03034. ISSN: 1742-5468.
- 1643 [S30] Jakob Jonsson. *Simplicial complexes of graphs*. Vol. 3.
1644 Springer, 2008.
- 1645 [S31] Shuang Wu et al. “The metabolomic physics of com-
1646 plex diseases”. In: *Proceedings of the National Academy*
1647 *of Sciences* 120.42 (2023), e2308496120. DOI: 10.1073/
1648 pnas.2308496120.
- 1649 [S32] Alexander Grigor’yan. “Advances in path homology
1650 theory of digraphs”. In: *Notices of the International*
1651 *Congress of Chinese Mathematicians* 10.2 (2022), pp. 61–
1652 124. ISSN: 2326-4810,2326-4845. DOI: 10.4310/ICCM.
1653 2022.v10.n2.a7.
- 1654 [S33] Peter Bubenik. “Statistical topological data analysis
1655 using persistence landscapes”. In: *Journal of Machine*
1656 *Learning Research* 16.1 (2015), pp. 77–102.
- 1657 [S34] Peter Bubenik. “The persistence landscape and some
1658 of its properties”. In: *Topological Data Analysis: The*
1659 *Abel Symposium 2018*. Springer. 2020, pp. 97–117.
- 1660 [S35] Peter Bubenik and Paweł Dłotko. “A persistence
1661 landscapes toolbox for topological statistics”. In:
1662 *Journal of Symbolic Computation* 78 (2017), pp. 91–114.
1663 ISSN: 0747-7171.
- 1664 [S36] Henry Adams et al. “Persistence images: A stable
1665 vector representation of persistent homology”. In:
1666 *Journal of Machine Learning Research* 18 (2017).
- 1667 [S37] D. Ali et al. “A survey of vectorization methods
1668 in topological data analysis”. In: *IEEE Transactions*
1669 *on Pattern Analysis and Machine Intelligence* (2023),
1670 pp. 1–14. ISSN: 1939-3539. DOI: 10.1109/TPAMI.2023.
1671 3308391.