

1 **STATISTICAL INFERENCE FOR PERSISTENT HOMOLOGY**
2 **APPLIED TO SIMULATED FMRI TIME SERIES DATA**

HASSAN ABDALLAH*

Department of Mathematics,
Wayne State University, MI 48202, USA

ADAM REGALSKI¹, MOHAMMAD BEHZAD KANG¹, MARIA BERISHAJ¹, NKECHI NNADI¹, ASADUR CHOWDURY²,
VAIBHAV A. DIWADKAR², ANDREW SALCH¹

¹ Department of Mathematics,

² Dept. of Psychiatry & Behavioral Neuroscience,
Wayne State University, MI 48202, USA

(Communicated by the associate editor name)

ABSTRACT. Time-series data are amongst the most widely-used in biomedical sciences, including domains such as functional Magnetic Resonance Imaging (fMRI). Structure within time series data can be captured by the tools of topological data analysis (TDA). Persistent homology is the mostly commonly used data-analytic tool in TDA, and can effectively summarize complex high-dimensional data into an interpretable 2-dimensional representation called a *persistence diagram*. Existing methods for statistical inference for persistent homology of data depend on an independence assumption being satisfied. While persistent homology can be computed for each time index in a time-series, time-series data often fail to satisfy the independence assumption. This paper develops a statistical test that obviates the independence assumption by implementing a multi-level block sampled Monte Carlo test with sets of persistence diagrams. Its efficacy for detecting task-dependent topological organization is then demonstrated on simulated fMRI data. This new statistical test is therefore suitable for analyzing persistent homology of fMRI data, and of non-independent data in general.

3 **1. Introduction.** Functional magnetic resonance imaging (fMRI) is a tool that
4 provides a rich avenue for studying brain activity via the hemodynamic responses.
5 Making sense of the complex spatio-temporal relationships in fMRI data can provide
6 insight into the functional and structural organization of the brain. A common goal
7 in fMRI experiments is to establish associations between changes in the fMRI signal
8 induced by the given task used specifically to evoke changes in the signal. Statistical
9 and data-analytic methods play a pivotal role in identifying and evaluating the
10 validity of such associations.

2020 *Mathematics Subject Classification.* 62R40, 55N31,

Key words and phrases. Topological Data Analysis, Statistics, Persistent Homology, Neuroscience, fMRI.

The seventh author is supported by Lycaki-Young Funds (State of Michigan), Mark Cohen Neuroscience Endowment, and 627 National Institutes of Mental Health (MH111177, MH059299).

* Corresponding author: Hassan Abdallah.

1 Statistical methods for fMRI analyses differentiate themselves by the type of
 2 effects they evaluate, and the manner by which they incorporate spatial informa-
 3 tion in the analysis. The general linear model (GLM), a predominant statistical
 4 workhorse, mostly evaluates effects under the assumption that the relationship be-
 5 tween the signal and any underlying variables or co-variates is linear. However, such
 6 effects or associations are frequently non-linear, and not straightforwardly discover-
 7 able. The fMRI signal is high dimensional, and characterized by hidden properties,
 8 the nature of which are not always known *a priori*. In this vein, topological data
 9 analysis (TDA) is in fact a viable option for exploring associations related to the
 10 topological or geometric characteristics of fMRI [21]. Within TDA, persistent ho-
 11 mology is one of the best known tools for characterizing topological features of a set
 12 of points in a relatively high-dimensional space, such as the four-dimensional space
 13 (three spatial dimensions, and one signal amplitude dimension) in which fMRI data
 14 naturally sit.

15 Persistent homology is a technique for discovery, but many scientific applications
 16 of any technique demand hypothesis testing to validate discoveries. Thus, if tools
 17 like persistent homology are to be widely adopted for fMRI research (and, more
 18 broadly, by scientists in biomedical and other fields), it is necessary to incorporate
 19 conventional statistical ideas for *hypothesis testing* into their application. In partic-
 20 ular, it is necessary to have a statistical test which can be applied to the results of
 21 using persistent homology to summarize fMRI data (or, more generally, to any class
 22 of time series data). The statistical test should ideally yield a numerical measure,
 23 such as a p-value, that is informative of the statistical significance of any kind of
 24 topological phenomena encoded in persistence diagrams. Our explicit goal herein is
 25 to provide a unique statistical measure that, when applied to persistent homology,
 26 permits conventional hypothesis testing to derive the significance of differences in
 27 topological properties summarized across experimental conditions or groups. Be-
 28 cause we focus on time series data, we describe a method for determining whether
 29 topological characteristics within a set of time series intervals are *significantly dif-*
 30 *ferent* from those in another set of time series intervals, where these intervals are
 31 related to different experimental conditions. This investigation builds on our previ-
 32 ous work [21] where we demonstrated the use of persistent homology to characterize
 33 structure in fMRI data, though without a framework for statistical inference.

34 We begin with the original motivating example for this paper: suppose we are
 35 given the data of an fMRI (functional magnetic resonance imaging) scan for a single
 36 participant in a study. This data set consists of, for each time index t and each
 37 spatial coordinate (x, y, z) in some representative set of spatial locations within the
 38 physical space of the brain, a number $f(x, y, z, t)$, the **fMRI signal amplitude**,
 39 which varies with the ratio of oxygenated hemoglobin to deoxygenated hemoglobin
 40 within the brain tissues near spatial location (x, y, z) at time t . The fMRI signal
 41 amplitude $f(x, y, z, t)$ is understood to vary, in an indirect and highly nonlinear
 42 way, with neuronal activity in the brain near (x, y, z) shortly preceding time t . In
 43 a task-based fMRI acquisition, the participant is engaged in a controlled cognitive
 44 experiment while fMRI data are being contemporaneously acquired. In an epoch
 45 structure for an associative memory experiment for instance:

46 **Epoch 1:** the person is asked to memorize associations between different classes
 47 of memoranda,

48 **Epoch 2:** the person’s memory for those associations is tested using cued recall.

49 A typical data-analytic approach might involve:

- 1 • the use of a spatial “mask” to the data, focusing analyses on a specific region
- 2 (e.g. the hippocampus),
- 3 • before asking whether the masked fMRI data collected in the two epochs is
- 4 statistically significantly different from each other.

5 If statistically significant differences in activity in the brain region (e.g., the
6 hippocampus) differs across each of the epoch types or from some baseline, then by
7 inference the task exerts significant effects on the region. Of course, this deductive
8 method is not restricted to fMRI data, but generalizes to any time series data
9 originating from *any* empirical study¹.

10 With the rise of the use of topological methods in data analysis in the past ten
11 years (see [23] for an introduction and brief survey) and in fMRI in particular (see
12 [21] for an introduction and brief survey), here, we motivate a combination of sta-
13 tistical inference with topological methods. The idea is to calculate the *persistence*
14 *diagram* (see [16] for an introduction) of the time-series data at each time index
15 separately, and then to ask whether the temporal organization of the data into
16 epochs can be recovered from the persistence diagrams in some statistically signif-
17 icant way. A Monte Carlo test for statistically significant clustering of persistence
18 diagrams was given in [20] and generalized in [4], but in both of those references,
19 an independence hypothesis on the persistence diagrams makes the resulting test
20 unsuited to time-series data. In particular, fMRI time-series data usually fails to
21 satisfy an independence assumption, since ongoing state-based processes in a given
22 brain region can cause the collected fMRI signal in that region at a given time index
23 to be dependent on the collected signal at the previous time index. More generally,
24 it is well-accepted that the fMRI signal is an index of dynamic continuing processes,
25 the state of the signal at any time t is dependent on the signal at time $t - 1$, and
26 will be predictive to some degree of the signal at time $t + 1$.

27 In the current paper we lift the independence assumption by describing a multi-
28 level block-sampled version of that Monte Carlo test. We demonstrate the utility of
29 our version on simulated fMRI time series data but reiterate its suitability for hy-
30 pothesis testing relating to any time-series data. We provide the R software package,
31 that our group developed for this test, at [https://github.com/hassan-abdallah/](https://github.com/hassan-abdallah/TimeSeriesTDA)
32 [TimeSeriesTDA](https://github.com/hassan-abdallah/TimeSeriesTDA). Furthermore, while time-series data is the main area of applica-
33 tion for this test, it is also useful on any other set of observations in which the
34 independence hypothesis fails.

35 As input, our analytic methods takes a) a set of points of a time-series (each
36 of which is a point cloud²) b) a labelling of the points (i.e., which epoch do they

¹As an example which is far removed from fMRI, we might consider average property value $pv(x, y, t)$ in some city, as a function of time t and of longitude-latitude coordinate pairs (x, y) . At each individual time index t , the persistent H_1 of the point cloud of triples $(x, y, pv(x, y, t))$ in \mathbb{R}^3 is sensitive to pockets of significantly higher or significantly lower property value than their surroundings. This persistent H_1 changes over time as the property values change, and over long periods of time, one might imagine that certain economic policies might have a statistically significant impact on the presence and distinctness of these pockets of higher or lower property value. This yields a labelling scheme, in the sense of our Definition 4.2, by labelling each time index with the economic policies in effect during that time. Our statistical test yields a way to determine whether the economic policies indeed have a statistically significant effect on the presence and distinctness of these pockets of higher or lower property value, insofar as these pockets are visible in persistent homology.

²A *point cloud* is a finite subset of the Euclidean space \mathbb{R}^n for some n . Consequently our analytic method requires that the observations have some kind of spatial organization to begin with.

1 belong to), c) a grouping of the points into exchangeability blocks, d) a labelling
 2 scheme for those blocks. A careful definition of this kind of structure is given in
 3 Definitions 4.1 and 4.2. The output of the analytic method is a p-value which
 4 reports on whether the persistent homology of the point clouds of each given label
 5 are statistically significantly distinct from the point clouds with other labels.

6 A brief introduction to persistent homology is given, though we direct interested
 7 readers to a thorough introduction and overview given by [6]. For an extended
 8 discussion of the relevance of topological summaries to fMRI, and a detailed com-
 9 parison of the kinds of insights about fMRI data obtainable via persistent homology,
 10 but not by classical statistical methods (e.g. regression analysis), we refer the reader
 11 to the paper [21], which is devoted to that topic. There are a variety of papers on
 12 TDA applied to time-series data, including [22], [14], and [18]; see [17] and [19] for
 13 nice surveys of some current ideas. The questions about hypothesis testing which
 14 motivate the present paper are not taken up in those works, however.

15 **2. Background on persistent homology.** Topological data analysis involves
 16 computations of *homology* and *persistent homology*. In this section, we offer a primer
 17 on those ideas, but for reasons of space, restrict its scope. For a more complete in-
 18 troductory treatment of persistent homology, see [16]. For a more comprehensive
 19 introduction to topological data analysis in general (rather than specifically persis-
 20 tent homology), we refer the reader to [6]. Even more generally, a more completely
 21 introductory treatment of homology can be found in any textbook on algebraic
 22 topology, such as the widely used book [11].

23 Before we begin, we note that persistent homology is defined on a choice of
 24 “point cloud” (see footnote for the definition of this term) together with a choice of
 25 coefficient ring. In most practical applications of persistent homology, the coefficient
 26 ring has been chosen to be the field with two elements, $\mathbb{F}_2 = \{0, 1\}$; see for example
 27 Table 3.1 in [15] for a 2015 list of commonly-used persistent homology software
 28 libraries which use \mathbb{F}_2 as either the default coefficient ring or as the only supported
 29 coefficient ring, e.g. Perseus, Dionysus, and GUDHI. We adhere to that convention
 30 in this paper: throughout, all homology is taken with coefficients in \mathbb{F}_2 .

31 Now we sketch the definition of a simplicial complex and its homology. We begin
 32 with a set of points v_0, v_1, \dots, v_k in \mathbb{R}^n such that the vectors $v_1 - v_0, v_2 - v_0, \dots, v_k - v_0$
 33 are linearly independent. Taking the convex hull $[v_0, v_1, \dots, v_k]$ of this set, we form
 34 its *k-simplex*. A *face* of that *k-simplex* is then the convex hull of a proper subset
 35 of $\{v_0, v_1, \dots, v_k\}$. So, for example, a 1-simplex is a line segment, and its faces are
 36 the endpoints of that line segment. Similarly, a 2-simplex is a solid triangle, and its
 37 faces are the edges of the triangle. A 3-simplex is a solid tetrahedron, and its faces
 38 are the triangles comprising the surface of the tetrahedron.

39 Next, consider a countable set K of simplices in \mathbb{R}^n such that:

- 40 • for each simplex in K , each of its faces are also contained in K , and
- 41 • the intersection of two simplices in K is either a face of both simplices, or is
 42 empty.

43 Such a set K is known as a *simplicial complex*. The intuition here is that a simplicial
 44 complex K is a geometric object which is “built” by taking a union of simplices,
 45 allowing any two to intersect only along a common face. If a simplicial complex K
 46 has only finitely many simplices, then K is a *finite simplicial complex*.

Let K be a simplicial complex, and for each integer k , consider the vector space
 $V_k(K)$ of formal \mathbb{F}_2 -linear combinations of *k*-simplices in K . That is, $V_k(K)$ is the

vector space of *simplicial k -chains*. Then the *boundary map*, extending to k -chains by linearity, is given by

$$\begin{aligned} \delta_k(K) : V_k(K) &\longrightarrow V_{k-1}(K) \\ [v_0, v_1, \dots, v_k] &\longmapsto \sum_{j=0}^k [v_0, v_1, \dots, \hat{v}_j, \dots, v_k], \end{aligned}$$

1 where \hat{v}_j indicates that v_j is omitted from the simplex. The *simplicial chain complex*
 2 of the finite simplicial complex K is the sequence of \mathbb{F}_2 -vector spaces and \mathbb{F}_2 -linear
 3 functions

$$\dots \xrightarrow{\delta_{k+1}} V_k(K) \xrightarrow{\delta_k} V_{k-1}(K) \xrightarrow{\delta_{k-1}} \dots \xrightarrow{\delta_2} V_1(K) \xrightarrow{\delta_1} V_0(K) \xrightarrow{\delta_0} 0.$$

4 The image (that is, range) of $\delta_{k+1} : V_{k+1}(K) \rightarrow V_k(K)$ is called the vector space of
 5 *k -boundaries of K* , while the kernel (that is, nullspace) of $\delta_k : V_k(K) \rightarrow V_{k-1}(K)$
 6 is called the vector space of *k -cycles of K* .

7 The boundary maps in the simplicial complex satisfy $\delta_k \circ \delta_{k+1} = 0$ for each integer
 8 k , that is, every k -boundary is also a k -cycle. Consequently we have a well-defined
 9 quotient vector space $\ker \delta_k / \text{im } \delta_{k+1}$ which is trivial if and only if every k -cycle is
 10 a k -boundary. The vector space $\ker \delta_k / \text{im } \delta_{k+1}$ is called the *k th homology of K* ,
 11 written $H_k(K)$. When it is important to remember that the coefficient ring has
 12 been taken to be the field \mathbb{F}_2 , we write $H_k(K; \mathbb{F}_2)$ instead of $H_k(K)$.

13 Now, given a simplicial complex K , consider a family $\{K_a : a \in \mathbb{R}\}$ of simplicial
 14 sub-complexes of K such that $K_m \subseteq K_n$ whenever $m \leq n$. That is, for each real
 15 number a , K_a is a simplicial sub-complex of K , and if a, b are real numbers with
 16 $a < b$, then every simplex in K_b is also in K_a . (So, as the subscript a gets smaller,
 17 the simplicial complex K_a also gets smaller.) The simplicial complex K together
 18 with the family $\{K_a : a \in \mathbb{R}\}$ is known as a *filtered simplicial complex*. For $a \leq b$,
 19 denoting the boundary maps on $V_k(K_a)$ and $V_k(K_b)$ by δ_k^a and δ_k^b , respectively, we
 20 naturally have inclusion maps $\iota : K_a \rightarrow K_b$, which, in turn, gives inclusion maps
 21 $\iota : \text{Im}(\delta_{k+1}^a) \rightarrow \text{Im}(\delta_{k+1}^b)$ and $\iota : \text{Ker}(\delta_k^a) \rightarrow \text{Ker}(\delta_k^b)$.

22 If the simplicial complex K is finite, then for most pairs of real numbers $a < b$
 23 with a sufficiently close to b , the subcomplex K_a of K_b is simply the entirety of
 24 K_b . There is only a *finite* list of real numbers b such that K_a differs from K_b for all
 25 $a < b$, no matter how close a is to b . Writing b_1, b_2, \dots, b_m for that finite sequence
 26 of real numbers, we have a sequence of \mathbb{F}_2 -linear functions

$$0 \rightarrow H_k(K_{b_1}) \rightarrow H_k(K_{b_2}) \rightarrow \dots \rightarrow H_k(K_{b_m})$$

27 called the *persistent homology groups of K* .

28 An element z of $H_k(K_{b_i})$ has a *birth radius*, that is, the least real number b_h
 29 such that z is in the image of the function $H_k(K_{b_h}) \rightarrow H_k(K_{b_i})$. Similarly, z has a
 30 *death radius*, that is, the least real number b_j such that z maps to zero under the
 31 function $H_k(K_{b_i}) \rightarrow H_k(K_{b_j})$. If the image of z is nonzero in $H_k(K_{b_j})$ for all b_j ,
 32 then the death radius of z is defined to be ∞ . (The birth radius of z , however, is
 33 always finite.)

34 We are now prepared to define the *persistence diagram*. The *k th persistence*
 35 *diagram* of the k th persistence module is a multiset³ of points in $\mathbb{R} \times (\mathbb{R} \cup \{\infty\})$.
 36 Each point in the diagram represents a homology class; the x -coordinate of the point

³Recall that a “multiset” is a set with (unordered) multiplicities, that is, an element of a multiset can be contained in that multiset “multiple times.” A typical way to make this intuitive idea rigorous is to simply think of a multiset as an ordinary set S equipped with an equivalence

1 representing a homology class z is the birth radius of z , while the y -coordinate
 2 of that point is the death radius of z . By convention, we include (with infinite
 3 multiplicity) all points such that $x = y$ (that is, the points lying along the diagonal).
 4 The further a point is from the diagonal of a persistence diagram, the longer the
 5 homology class *persists* (i.e., is nonzero) as the filtration parameter ranges over the
 6 real numbers. The intuition here, then is that the closer a point in the persistence
 7 diagram is to the diagonal, the more we think of the topological feature represented
 8 by that homology class as a kind of “topological noise,” rather than a meaningful
 9 topological pattern involving and organizing a large part of the data set.

10 The typical intended use of persistent homology for the sake of data analysis
 11 is that one begins with a point cloud, one builds a finite filtered simplicial com-
 12 plex whose structure reflects the geometry of the point cloud in some desired way,
 13 and then one calculates the persistent homology groups of that filtered simplicial
 14 complex. We have explained the last step, but we have not yet explained how to
 15 build a finite filtered simplicial complex from a point cloud. There are several ways
 16 to do this: a point cloud has an associated Čech complex, Vietoris-Rips complex,
 17 Delaunay complex, witness complexes, and others, each of which is a finite filtered
 18 simplicial complex whose structure “encodes” the geometry of the point cloud in
 19 some particular way. See [8] for discussion and comparison of the Čech and Vietoris-
 20 Rips complexes, for example. For brevity, here we do not attempt a survey of these
 21 various filtered simplicial complexes, but we at least give a definition of the Čech
 22 complex, since it is the most geometrically straightforward: given a point cloud
 23 $X \subseteq \mathbb{R}^n$ and a subset U of X , the *diameter of X* is the least real number ϵ such
 24 that every element of U is contained in a closed ball of radius ϵ in \mathbb{R}^n . The *Čech*
 25 *complex of X* is the filtered simplicial complex $\{K_a : a \in \mathbb{R}\}$ such that K_a is the
 26 union of the convex hulls of each of the subsets of X of diameter $< a$.

27 The persistent homology groups have intuitive geometric significance, of which
 28 we now give a very brief account. The dimension of the vector space $H_0(K)$ counts
 29 the *connected components* in the geometric realization of the simplicial complex K .
 30 Similarly, the dimension of the vector space $H_1(K)$ counts *noncontractible loops* (up
 31 to “homology”, a certain equivalence relation) in the geometric realization of K .
 32 The dimension of the vector space $H_2(K)$ counts *noncontractible spheres* (again,
 33 up to “homology”) in the geometric realization of K ; one often thinks of such
 34 noncontractible spheres as being wrapped around three-dimensional *voids* in the
 35 geometric realization. As applied to a point cloud arising from real-world data,
 36 persistent H_0 measures clustering at different scales, while persistent H_1 measures
 37 loop-shaped “gaps” at different scales in the point cloud, and persistent H_2 measures
 38 open “voids” at different scales in the point cloud. The persistent H_n for $n > 2$
 39 measures higher-dimensional analogues of loop-shaped gaps, voids, etc

40 As an example, consider a point cloud formed by sampling points from an an-
 41 nulus. Figure 1 shows a side-by-side comparison of balls of radius $\frac{1}{3}$ around each
 42 point and a visualization of the associated simplicial complex. Figure 2 shows the
 43 persistence diagram computed from that point cloud. The triangles in the persis-
 44 tence diagram represent 1-dimensional homological features which occur as a result
 45 of non-contractible loops in the filtered simplicial complex. The triangles closer to
 46 the diagonal (i.e. lower persistence features) are a result of smaller loops in the
 47 point cloud that exist because of our noisy sampling of the annulus. The single

relation. Given an element s of S , the multiplicity of s in S is understood to be the number of elements in the equivalence class of s .

- 1 triangle far from the diagonal (i.e. a high persistence feature) is a result of the large
- 2 hole in the center of the point cloud.
- 3 For more details, the reader can consult the references cited at the start of this
- 4 section.

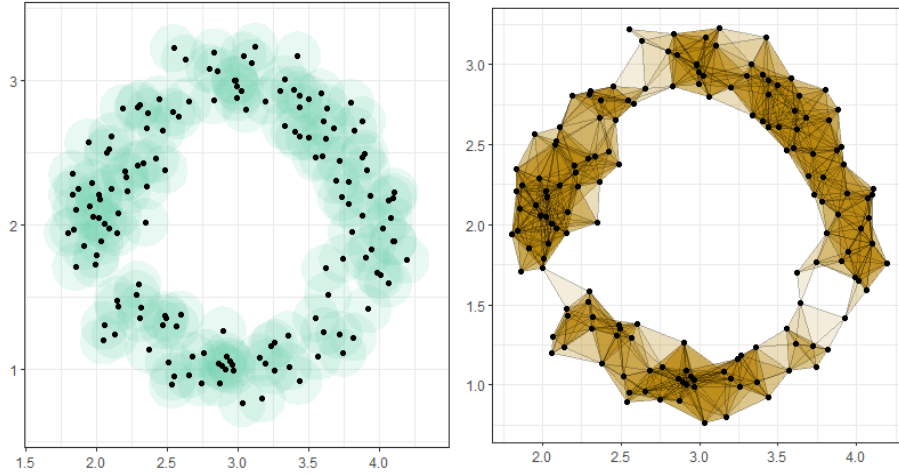


FIGURE 1. On the left is a plot of the point cloud with balls of radius $\frac{1}{3}$ around each point. On the right is a visualization of the simplicial complex for filtration= $\frac{1}{3}$.

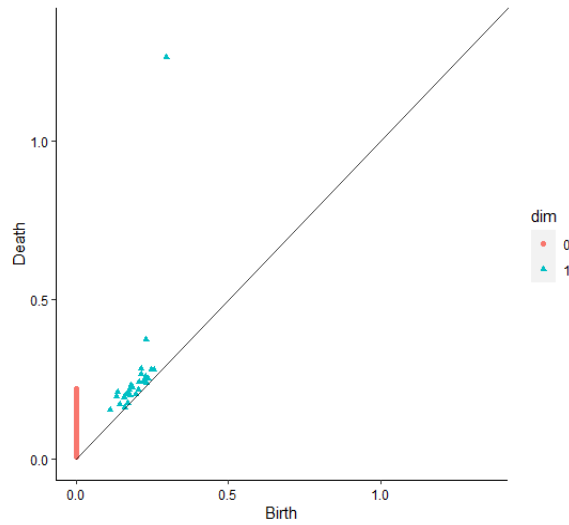


FIGURE 2. The persistence diagram computed from the point cloud in Figure 1.

1 **3. Summary of Existing Hypothesis Testing Methods for Topological**
 2 **Data Analysis.** The purpose of this section is to discuss (and extend via [4]) the
 3 methods used in [20] which we will apply to our fMRI data. To clarify, everything
 4 we discuss in this section is due to [20] and [4].

5 We begin by involving some ideas from statistics alongside the basic notions in
 6 persistent homology. The idea here is as follows: Imagine that we collect fMRI data
 7 using a task that oscillates between blocks of multiple task active conditions and rest
 8 (i.e., what is known as a typical “block design”). We want to calculate the persistent
 9 homology of all the data acquired in each block and develop a statistical test (i.e.,
 10 a hypothesis test) to determine whether the persistence diagrams generated from
 11 one condition are distinguishable from persistence diagrams generated from another
 12 condition within the same acquisition. Thus, in order to assess the strength of evi-
 13 dence against the claim that the two conditions elicit indistinguishable topological
 14 organization, we can study the distributions of persistence diagrams associated with
 15 each condition. The goal of Robinson and Turner’s work in [20] is to use hypothe-
 16 sis testing to compare two groups of persistence diagrams. The methods discussed
 17 in [20] are extended in [4] in order to use hypothesis testing to compare multiple
 18 groups of persistence diagrams. The need for us to extend the comparisons between
 19 persistence diagrams to 3 or more groups of persistence diagrams comes from the
 20 multi-level block sampling framework that we apply to our time-series data in the
 21 next section of this paper, where we freely permute multiple blocks (and, hence,
 22 multiple groups of persistence diagrams) to carry out our hypothesis test.

23 Our hypothesis test begins with a set of n persistence diagrams divided into
 24 s groups $\beta_1 = \{X_{1,1}, X_{1,2}, \dots, X_{1,n_1}\}$, $\beta_2 = \{X_{2,1}, X_{2,2}, \dots, X_{2,n_2}\}$, ... , $\beta_s =$
 25 $\{X_{s,1}, X_{s,2}, \dots, X_{s,n_s}\}$ containing n_1, n_2, \dots, n_s diagrams, respectively, with this di-
 26 vision into multiple groups done according to some initially-chosen labeling scheme.
 27 The hypothesis test corresponding to the case $s = 2$ is the subject of [20], while
 28 the generalization to arbitrary finite s was the focus of [4]. The null hypothesis is
 29 that the underlying distribution of β_1 is the same as the underlying distribution
 30 of β_2 . The alternative hypothesis is that the underlying distributions are different.
 31 An observed test statistic is computed using the initial labeling scheme, and com-
 32 puted further for each permutation of labels in the permutation test. The key to
 33 computing the final p-value, which assesses the strength of evidence against the null
 34 hypothesis, then, is to compute the ratio of permutations that yield a test statis-
 35 tic more extreme than the observed statistic to the total number of permutations.
 36 We note that a necessary assumption for the test is that observations (respectively,
 37 persistence diagrams) are independent.⁴ Also, the permutation test we carry out is
 38 a randomization test. As mentioned in Section 2.6 of [20], using a randomization
 39 test avoids any need to hypothesize a distribution model from which persistence
 40 diagrams are drawn under the null hypothesis .

41 **3.1. Metric on Persistence Diagrams.** In order to carry out our hypothesis
 42 test, we first need to introduce a metric, i.e., a distance function, on the space of
 43 persistence diagrams. This metric allows us to compare two persistence diagrams
 44 and is a key piece of the test statistic that we’ll utilize in this hypothesis test.

⁴Since our goal is to have a statistical test that can be applied to the persistence diagrams of *non-independent* time series data, in the next section, we apply a multi-level block sampling framework to satisfy the exchangeability criteria for our permutation test, thereby removing the requirement of our observations being independent.

The appropriate distance metric between persistence diagrams X and Y that we consider is the bottleneck distance

$$d_\infty(X, Y) = \inf_{\text{bij. } \phi: X \rightarrow Y} \sup_{x \in X} \|x - \phi(x)\|_\infty$$

which occurs as the limit of the metric

$$d_p(X, Y) = \left(\inf_{\text{bij. } \phi: X \rightarrow Y} \sum_{x \in X} \|x - \phi(x)\|_p^p \right)^{1/p}$$

1 as p goes to infinity, where $\|x - \phi(x)\|_p^p$ is the L^p -norm between x and $\phi(x)$ raised to
 2 the p -th power, and the infimum is taken over all bijections ϕ between the points of
 3 X and the points of Y . Note that, as a metric on the space of persistence diagrams,
 4 the bottleneck distance $d_\infty(X, Y)$ between X and Y is indeed symmetric. This
 5 follows since a bijection $\phi : X \rightarrow Y$ also defines a bijection $\phi^{-1} : Y \rightarrow X$. We
 6 now unpack the construction of these metrics.

7 The metrics take into account an optimal bijection $\phi : X \rightarrow Y$ between the
 8 points of X and the points of Y . A bijection $\phi : X \rightarrow Y$ is said to be “optimal”
 9 if it minimizes the total cost $\sum_{x \in X} \|x - \phi(x)\|^2$. Optimal bijections are found by
 10 using the Hungarian algorithm. Given two sets of elements $S = \{s_1, \dots, s_n\}$ and
 11 $T = \{t_1, \dots, t_n\}$, and a square matrix A , where the i th row of A is represented by
 12 the element s_i and the j th column is represented by the element t_j , one can apply
 13 the *Hungarian algorithm* to A to find the optimal bijection between elements of S
 14 and elements of T . (The original reference for the Hungarian algorithm is the 1955
 15 paper [12], but today the Hungarian algorithm is a standard topic covered in many
 16 discrete mathematics textbooks, so many modern expositions are available.)

17 Here is a bit more detail about what the bottleneck distance between two persis-
 18 tence diagrams is. If X has points x_1, \dots, x_n and Y has points y_1, \dots, y_m , one takes
 19 copies x_{n+1}, \dots, x_{n+m} and y_{m+1}, \dots, y_{m+n} of the diagonal in a persistence diagram,
 20 where this diagonal is the line of slope 1 in the birth-death plane, and constructs
 21 the $(n+m) \times (n+m)$ matrix in which the (i, j) entry is the cost $\|x_i - y_j\|_2^2$. When
 22 one of x_i or y_j is a copy of the diagonal, this is the perpendicular distance between
 23 x_i and y_j . When both x_i and y_j are copies of the diagonal, the cost is simply 0.

3.2. Test Statistic and p-value for Comparing Groupings of Persistence Diagrams. Now that we have established an appropriate metric on two persistence diagrams, we can formulate a test statistic for our hypothesis test. The test statistic is the joint loss function given by

$$F'_{p,q}(\{X_{1,i}\}, \{X_{2,i}\}, \dots, \{X_{s,i}\}) := \sum_{m=1}^s \frac{1}{2n_m(n_m - 1)} \sum_{i=1}^{n_m} \sum_{j=1}^{n_m} d_p(X_{m,i}, X_{m,j})^q,$$

24 where $p \in [1, \infty)$, $q \in [1, \infty)$. This joint loss function, as a test statistic, was
 25 introduced in [4] as a generalization of the $s = 2$ case considered in [20]. (In [4],
 26 only the case of $p = 2$ and $q = 2$ is considered, but the extension to other values of
 27 p and q is straightforward. In our application of these ideas, we take p to be infinity
 28 and q to be 1.) Since the groups $\beta_1, \beta_2, \dots, \beta_s$ are determined by a choice of
 29 labeling L , we will use the notation $F'(L)$ to mean the joint loss function computed
 30 on the s groups of persistence diagrams determined by L , and $F'(L_{\text{observed}})$ to mean
 31 the joint loss function computed on the s groups of diagrams determined by the
 32 initial choice of labeling. When implemented in software, the pairwise distances

1 between persistence diagrams are only computed once and stored in a table. Note
 2 that the test statistic given by the joint loss function takes into account distances
 3 between observations (respectively, persistence diagrams), rather than distances
 4 between observations and the mean. This is because the latter consideration is very
 5 computationally expensive⁵.

6 Taking α to be the proportion of all labelings L such that $F'(L) \leq F'(L_{observed})$,
 7 where all of the possible labelings L are determined by the permutation test carried
 8 out to permute the labels on persistence diagrams, we can now generalize (via
 9 [4]) the algorithm developed in [20] to compute the proportion α to be taken as
 10 the p-value (after a standard modification to α in order to avoid a p-value of 0).
 11 The difference is that, rather than having $n_1 + n_2$ persistence diagrams with labels
 12 $L_{observed}$ in disjoint sets of size n_1 and n_2 and randomly shuffling the group labels
 13 into disjoint sets of size n_1 and n_2 to give the labeling L , we now have $n_1 + n_2 + \dots + n_s$
 14 persistence diagrams with labels $L_{observed}$ in disjoint sets of sizes n_1, n_2, \dots, n_s and
 15 we randomly shuffle the group labels into disjoint sets of sizes n_1, n_2, \dots, n_s to give
 16 the labeling L . It is shown in [20] that the modified α is a true p-value, and by
 17 Lemma 1 of [20], α is an unbiased estimator of the permutation p-value under the
 18 assumption that the persistence diagrams are i.i.d. As mentioned before, our goal in
 19 this paper is to adapt the Robinson-Turner test to the common real-world situation
 20 of time series data which is not independent, and consequently our persistence
 21 diagrams, regarded as observations, are not independent observations. However,
 22 again, we're able to correct for this using our methods in the next section.

23 **4. Hypothesis Testing for Topological Data Analysis extended to Non-**
 24 **Independent Data.** In this section, we describe a single and multi-level block vari-
 25 ation of the original Monte Carlo test that allows for the analysis of non-independent
 26 data sets. The primary idea involves accommodating the unique *exchangeability*
 27 structure of a particular set of data.

28 **4.1. Exchangeability.** A sequence of random variables X_1, X_2, \dots, X_n is *exchange-*
 29 *able* under a set of permutations Π of $\{1, 2, \dots, n\}$ if it has the same joint distribution
 30 as the sequence $X_{\pi(1)}, X_{\pi(2)}, \dots, X_{\pi(n)}$ for every $\pi \in \Pi$. Determining the set Π for
 31 which exchangeability holds is critical to perform sound statistical inference via a
 32 permutation test. If the joint distribution of a set of data changes under particular
 33 permutations of labels, then the distribution of a test statistic under those per-
 34 mutations is not suitable to be compared to the observed test statistic and could
 35 elicit spurious results. In the case of independent and identically-distributed ran-
 36 dom variables, the set Π contains all permutations of $\{1, 2, \dots, n\}$, meaning labels
 37 may be freely exchanged during a permutation test. As a result, the hypothesis
 38 testing framework described in Section 2 did not require any considerations of ex-
 39 changeability. In many cases, however, the set of permutations that satisfy the
 40 above exchangeability criterion is far more restrictive. Fortunately, by restricting
 41 permutations to the set Π while generating the distribution of a test statistic, a
 42 permutation test may proceed without the iid requirement.

43 In practice, implementing a restrictive set of permutations is done via a multi-
 44 level block shuffling scheme, as in [26]. Instead of exchanging the label of one ob-
 45 servation with another, shuffling takes places across blocks of data called *exchange-*
 46 *ability blocks*. Exchangeability blocks can either be shuffled as a whole (defined as

⁵Also, there is not a clear notion of the mean of a set of persistence diagrams.

1 whole-block exchangeability) or labels may be shuffled within a block(defined as
 2 within-block exchangeability). The block sizes and attributes are chosen in accor-
 3 dance with the permitted set of permutations.

4 For example, consider an fMRI experimental design that consists of 120 total
 5 scans. Suppose a stimulus or task is administered every 10 scans and lasts for 10
 6 scans. Our whole-block exchangeable level in this scenario would be defined as
 7 each contiguous set of 10 scans starting at 1, accounting for 12 blocks in total.
 8 These first-level exchangeability blocks ensure that any label shuffling would result
 9 in a set of labels exhibiting a similar contiguity to the initial set of labels (those
 10 assigned stimulus/no stimulus during the experiment) and retain structure related
 11 to the experiment design. They do not, however, account for a temporal dependence
 12 structure across the whole of the experiment. Without additional restrictions, labels
 13 may be shuffled into two groups where one group is the data associated with scans
 14 1-60 and the other is the data associated with 61-120. A distinguishing pattern
 15 across early versus late stages of fMRI experiments have been noted so the test
 16 statistic computed for this set of labels could display an extremeness resulting from
 17 this temporal phenomena [13]. As such, a within-block exchangeability level is
 18 necessary.

19 A within-block exchangeability level would consist of two blocks, one covering
 20 scans 1-60 and another covering scans 61-120. In this design, the whole-block ex-
 21 changeable blocks present in labels 1-60 could only be exchanged amongst them-
 22 selves and not with their corresponding blocks in labels 61-120. This ensures that
 23 the first half of the experiment and the second half of the experiment would have
 24 equal representation in any set of labels shuffled under this scheme, accounting for
 25 early versus late confounding.

26 In the context of topological hypothesis testing, we define a *two-level point cloud*
 27 *grouping* and associated *labelling scheme* to encode the multi-level block shuffling
 28 technique described above:

29

30 **Definition 4.1.** A (*two-level*) *point cloud grouping* is the following data:

- 31 1. A set T of indexes for each point cloud.
- 32 2. A function pc from T to the set of all (observed) point clouds.
- 33 3. A partition T into subsets T_1, \dots, T_n and,
- 34 4. A partition T into subsets T'_1, \dots, T'_m which is finer than the partition T_1, \dots, T_n
 35 of T .

36 **Definition 4.2.** Given a two-level point-cloud grouping X , a 2-group *labelling*
 37 *scheme* on X is a partition of T into subsets S_1 and S_2 by the following:

- 38 1. For each $i \in \{1, \dots, n\}$, $\exists j_1, \dots, j_{k_i} \in \{1, \dots, m\}$ such that $\bigcup_{a=1}^{k_i} T'_{j_a} = T_i$.
- 39 2. For each i , choose $k_i/2$ elements without replacement, i.e. without repetition,
 40 from the set $\{j_1, \dots, j_{k_i}\}$, denoted $\{s^i_1, \dots, s^i_{k_i/2}\}$.

- 41 3. Define $S_1 = \bigcup_{i=1}^n \bigcup_{k=1}^{k_i/2} T'_{s^i_k}$ and $S_2 = T - S_1$.

42 The partition T_1, \dots, T_n represents the whole-block exchangeable level and the
 43 partition T'_1, \dots, T'_m represents the within-block exchangeable level. Permutations of
 44 labels are then obtained by generating distinct labelling schemes as defined above.

1 In the fMRI example described above, T would be the set of a time indices from 1
2 to 120 and the partitions would be defined as follows:

$$\begin{aligned} T_1 &= \{1 : 10\} & T_2 &= \{11 : 20\} & T_3 &= \{21 : 30\} & T_4 &= \{31 : 40\} \\ T_5 &= \{41 : 50\} & T_6 &= \{51 : 60\} & T_7 &= \{61 : 70\} & T_8 &= \{71 : 80\} \\ T_9 &= \{81 : 90\} & T_{10} &= \{91 : 100\} & T_{11} &= \{101 : 110\} & T_{12} &= \{111 : 120\} \\ T'_1 &= \{1 : 60\} & T'_2 &= \{61 : 120\} \end{aligned}$$

3 The function pc would map an element of T to its corresponding point cloud.

4 **4.2. Overview of Analysis Pipeline.** In this section, a broad overview of the
5 steps to go from a data set to a p-value for a hypothesis is given. Several of the
6 computational tasks involved can be accomplished using our R package “TimeSeriesTDA”.
7 Beginning with a data set of interest, carry out the following:

- 8 1. Compute persistent homology on all observations of your data set to produce a
9 collection of persistence diagrams. Each observation should be a point cloud.
- 10 2. Generate a hypothesis that conjectures a significant difference between two
11 sub-collections of your collection of persistent diagrams, called *groupings*. The
12 null hypothesis is that the two groupings are not significantly different from
13 each other. Choose an α level for rejecting the null hypothesis. For example,
14 for $\alpha = 0.05$, the null hypothesis will be rejected if the resulting p-value of
15 this hypothesis test is less than 0.05.
- 16 3. Compute the value of the appropriate joint loss function given labels for the
17 above groupings, called $F'(L_{observed})$ as defined in Section 3.
- 18 4. Determine the exchangeability structure of your observations and encode it in
19 a two-level point cloud grouping. In particular, define the set T and partitions
20 of T corresponding to whole-block and within-block exchangeability levels, as
21 in Definitions 4.1 and 4.2.
- 22 5. Generate distinct labelling schemes and recompute the joint loss function value
23 for each new set of labels. Compute a p-value by taking the proportion of
24 permuted labels L such that $F'(L) \leq F'(L_{observed})$. Compare the p-value to
25 the pre-determined α -threshold to evaluate whether the null hypothesis will
26 be rejected or not.

27 **5. Application to fMRI data.** fMRI imaging is a rich source for obtaining non-
28 independent time-series data. The data obtained from an fMRI scan is time-series
29 data consisting of, at each time index t , a real number $f(x, y, z, t)$ at each point
30 (x, y, z) in a certain set of lattice points in \mathbb{R}^3 . The number $f(x, y, z, t)$ is the
31 *fMRI signal amplitude*, which is understood to vary (non-linearly) with the ratio
32 of oxygenated hemoglobin to deoxygenated hemoglobin in the blood in the tissues
33 near physical location (x, y, z) at time t . That is, fMRI data is time series data,
34 such that at each time index, we have a point cloud in \mathbb{R}^4 : three spatial dimensions,
35 and one signal amplitude dimension. The fMRI signal amplitude has a relationship
36 to unfolding biological processes in the brain. These processes are, at each moment
37 in time, potentially dependent on their states at prior moments in time. [2]

38 Before applying persistent homology, a suitable normalization technique for the
39 fMRI signal needs to be identified such that the 4-dimensional point clouds ob-
40 tained from fMRI data are organized in such a way that persistent homology is
41 adequately sensitive to evolving topological structure. Additionally, parameters

1 related to persistence, such as maximum birth and death radiuses to compute per-
 2 sistent homology to, and cutoffs for persistence diagram feature selection, need to
 3 be explored in the context of fMRI data. In this section, we discuss each of these
 4 decisions (normalization method, and parameter choices) in turn.

5 **5.1. Normalization.** At each individual time index, the structure of fMRI data
 6 consists of three spatial coordinates and a signal amplitude coordinate. When
 7 discovering topological features in the 4-dimensional point cloud, one would hope
 8 that the same features would be obtained regardless of the choice of units. This
 9 poses an issue as the three spatial coordinates are measured in millimeters, whereas
 10 the signal amplitude is unitless. A change in units of distance would rescale the three
 11 spatial dimensions but not the fourth (signal amplitude), changing the topological
 12 features and persistence diagrams acquired. Consequently, in order to yield results
 13 that are invariant under changing units, fMRI data must be normalized before
 14 calculating persistence diagrams. Different choices of how to normalize fMRI data
 15 may yield different persistence diagrams, so topological structure in fMRI data is
 16 impacted by *how* we normalize the data. Below, two methods of normalizing fMRI
 17 data are discussed. In section 5.4, we report on which of these two normalization
 18 methods, when applied to our simulated fMRI data, allow our statistical test to
 19 achieve greater statistical power.

Definition 5.1 (Normalization Scheme 1). Define the following notation:

$$\begin{aligned}
 S_{\min} &= \min \left\{ \min\{x - \text{coordinates}\}, \min\{y - \text{coordinates}\}, \min\{z - \text{coordinates}\} \right\} \\
 S_{\max} &= \max \left\{ \max\{x - \text{coordinates}\}, \max\{y - \text{coordinates}\}, \max\{z - \text{coordinates}\} \right\} \\
 A_{\min} &= \min \left\{ \text{signal amplitude} \right\} \\
 A_{\max} &= \max \left\{ \text{signal amplitude} \right\}
 \end{aligned}$$

where the minimums and maximums are taken for each time slice and each subject individually. For any given coordinate (x, y, z, ϵ) , replace ϵ (the signal amplitude) with

$$\left[\frac{\epsilon - A_{\min}}{A_{\max} - A_{\min}} \cdot (S_{\max} - S_{\min}) \right] + S_{\min}.$$

Definition 5.2 (Normalization Scheme 2). Define the following notation:

$$\begin{aligned}
 S_{\min} &= \frac{\min \{x - \text{coordinates}\} + \min \{y - \text{coordinates}\} + \min \{z - \text{coordinates}\}}{3} \\
 S_{\max} &= \frac{\max \{x - \text{coordinates}\} + \max \{y - \text{coordinates}\} + \max \{z - \text{coordinates}\}}{3},
 \end{aligned}$$

where the minimums and maximums are taken for each time slice and each subject individually. We let A_{\min} and A_{\max} be as in Definition 5.1. For any given coordinate (x, y, z, ϵ) , replace ϵ (here, the fMRI signal amplitude) with

$$\left[\frac{\epsilon - A_{\min}}{A_{\max} - A_{\min}} \cdot (S_{\max} - S_{\min}) \right] + S_{\min}.$$

20 Normalization Scheme 2 is preferred due to the fMRI signal amplitude simi-
 21 larity to the spatial coordinates range, maximums, minimums, and magnitudes.

1 This occurs because Scheme 2 utilizes the averages of the spatial coordinates. We
 2 demonstrate this preference in the following example:

3 **Example 5.1.** Suppose we consider fMRI data whose spatial coordinates have
 4 x, y , and z coordinates in the ranges $x \in (35, 57), y \in (65, 90), z \in (32, 63)$. Using
 5 Normalization Scheme 1 linearly rescales the fMRI signal amplitudes, at each time
 6 slice, so that the normalized signal amplitudes lie in the range $(32, 90)$. On the other
 7 hand, Normalization Scheme 2 linearly rescales the fMRI signal amplitudes so that
 8 the normalized signal amplitudes lie in the range $(44, 70)$. The reason an average
 9 is preferred is due to the sizes of the ranges: the normalized amplitudes under the
 10 first method lie in an interval of length 58, while the normalized amplitudes under
 11 the second method lie an interval of length 26, which is closer to the ranges of the
 12 spatial coordinates, since the lengths of the ranges for the x, y , and z coordinates are
 13 22, 25, and 31, respectively. In this example, we see that the second normalization
 14 scheme yields a normalized fMRI amplitude whose properties more closely mirror
 15 the properties of the spatial coordinates. See section 5.4 for empirical calculations
 16 of the power of our statistical test when applied to simulated fMRI data with each
 17 of the two normalization schemes.

18 **5.2. Parameter considerations.** To compute persistent homology, a choice of
 19 maximum filtration parameter and maximum dimension of homology is made. Vary-
 20 ing these choices does not introduce external artifacts, but instead varies how com-
 21 prehensive of a view of the data is obtained. We emphasize that computational
 22 constraints typically play the biggest role in selecting these parameters.

23 The first consideration is the maximum filtration parameter for which to compute
 24 the persistent homology. For the Čech filtration (defined in section 2), the ideal
 25 choice for this parameter is half the distance of the two farthest points in a data
 26 set, since there are no non-trivial changes in the topology of the space beyond that
 27 radius (the topology is that of a single convex body). This is an example of a
 28 “canonical choice” of the filtration parameter.

29 However, in most fMRI data sets, computing persistent homology up to that
 30 distance is not computationally feasible. Instead, a threshold value is chosen such
 31 that it has the potential to capture nontrivial topology, and the process completes
 32 in a reasonable amount of time. For example, using maximum filtration parameter
 33 1 or 2 with the two normalization techniques previously discussed yields virtually
 34 no one-dimensional homological features in fMRI data. This is not because those
 35 features are not present, but rather because the birth radius or death radius of those
 36 features is greater than 1 or 2. Using maximum radius 3 or 4, on the other hand, is
 37 large enough to capture interesting topological information. In practice, one should
 38 choose a value as close to the “canonical choice” as your time and computational
 39 resources allow. It is important to determine whether tweaking this choice of pa-
 40 rameter alters results (and we present some conclusions to this effect in section 5.4),
 41 since as this parameter changes, so does the hypothesis and statistical conclusion
 42 of our test. For example, rejecting the null hypothesis would show that there is
 43 enough evidence to support the claim that the two groups of persistence diagrams,
 44 *up to persistence* (= maximum radius), are statistically significantly different from
 45 each other. This not only indicates differing topological structure, it also indicates
 46 the maximum size and scale of the topological structures that influence the result.

1 The second consideration is the maximum dimension of homology for which
 2 to compute persistent homology. Recall that 0-dimensional homology (H_0) is re-
 3 lated to connected components, 1-dimensional homology (H_1) is related to non-
 4 contractible loops, and higher dimensional homology is related to voids and their
 5 higher-dimensional analogues. If a point cloud is n -dimensional, a “canonical” up-
 6 per bound for dimension of homology to calculate is $n - 1$. This is because it is not
 7 possible for there to be nontrivial homology in dimensions greater than $n - 1$; see
 8 Corollary 2.2 of [7] and surrounding discussion for a nice exposition of why this is
 9 true (technically this discussion handles only Čech homology; for the close relation-
 10 ship between Čech homology and Vietoris-Rips homology, see (6.5) in section 6.1
 11 of [7]). Consider the point cloud sampled from an annulus in Figure 1. Each point
 12 is a point of \mathbb{R}^2 and the one-dimensional hole in the point cloud is captured by H_1 .
 13 It is not possible for there to be nontrivial H_2 because a two-dimensional void is
 14 not possible in \mathbb{R}^2 . Therefore, computing persistence homology up to dimension 1
 15 is satisfactory. Unfortunately, the canonical upper bound is not always achievable.
 16 For example, with fMRI data, our point cloud consists of points in 4-dimensional
 17 space (i.e. in \mathbb{R}^4). In this case, it would be ideal to compute persistent homology
 18 up to dimension 3. In reality, at present, computing up to dimension 1 is all that
 19 is possible for the computation to finish in a reasonable amount of time and with
 20 modest computational resources. As such, persistent homology is only computed
 21 up to H_1 in our simulation.

22 After sets of persistence diagrams are in hand, the next parameter to consider is
 23 the number of features in persistence diagrams that are retained for our analysis.
 24 A distance matrix of the persistence diagrams is necessary to compute the test
 25 statistic in our Monte Carlo test. Ideally, one would not remove any features from
 26 the persistence diagrams when computing this distance matrix, however that is not
 27 always possible. For example, it has been found that reasonably sized sets of fMRI
 28 data (>1000 4-dimensional points) contain potentially thousands of 1-dimensional
 29 homological features. It is not tractable to compute a distance matrix between
 30 more than a few dozen persistence diagrams when each has that many features.
 31 Fortunately, there are methodological considerations for filtering out a large subset
 32 of features. Masked fMRI data is composed on a lattice with distance 1 between
 33 adjacent points. Setting signal to zero for all points, computing persistent homology
 34 on such a space would result in an abundance of features with persistence $\sqrt{(2)}/2$
 35 (≈ 0.707). It is then reasonable to infer that features at that persistence and below
 36 are likely more related to small-scale “topological noise” rather than large-scale,
 37 meaningful topological organization within the data. Thus, our initial cutoff for
 38 minimum persistence threshold is 0.8. In our results section, we determine whether
 39 increasing that cutoff gives a more powerful test or not. Although we have just
 40 given a logical explanation for why choosing this cutoff is reasonable, we again
 41 emphasize that this choice of parameter should be made to be as close to zero (i.e.
 42 not removing any features) as is computationally feasible.

43 **5.3. fMRI Data Simulation.** Here, we discuss how we generated simulated fMRI
 44 data in order to test the power, accuracy, and reliability of the proposed method.
 45 We used the R package neuRosim [24] to simulate the data.

46 **5.3.1. Experimental Design.** We generated simulated fMRI data with a repetition
 47 time (TR) of two seconds in a spatial region (i.e., a region of stereotactic space) in
 48 the shape of a standard fMRI mask of the hippocampus. Throughout each simulated

1 run, the signal amplitudes in this spatial region are first given by a standard simu-
 2 lation of physiological noise. Physiological noise is intended to mimic noise caused
 3 by heart beat and respiratory rate. It is modelled by sine and cosine functions with
 4 the addition of Gaussian noise to increase variability across voxels.

5 The simulated data is structured so that, in each simulated run, there are six
 6 “epochs,” consisting of 20 seconds each. At the onset of each epoch, the signal
 7 amplitudes are increased in a sphere-shaped region within the hippocampus-shaped
 8 region, depicted in the figures below. Activation is greatest at the beginning of each
 9 epoch and fades throughout.

10 5.3.2. *Simulation Characteristics.* With this experimental design we varied the char-
 11 acteristics of both noise and signal in the interest of deciding what, if any, topological
 12 structure this method might detect. By varying noise and signal characteristics, the
 13 topology of the data will vary with it.

14 Each simulated data set was generated by the following process:

- 15 1. For each time index t from $t = 1$ to $t = 120$, set the signal amplitude in
 16 each voxel in a standard hippocampal mask to the values given by simulated
 17 physiological noise, depicted in Figure 3.
- 18 2. Choose a radius r (we considered the values $r = 1, 3, 5, 7$, and 15 , in separate
 19 runs) and a point p in the mask. In a spherical region of radius r (measured in
 20 voxel edge lengths) with center p , replace the physiological noise signal with an
 21 “activated signal” of high amplitude at the start of each epoch, and decaying
 22 in amplitude throughout the epoch. We used a standard amplitude curve for
 23 simulated fMRI provided by neuRosim, depicted in Figure 4. It is necessary
 24 to choose the initial effect size (which can be thought of as a measure of the
 25 magnitude of activation) of the activation in the sphere. We generated data
 26 with initial effect sizes $2, 5, 10$, and 20 , to compare the results. Figure 5
 27 contains images of the resulting mask, with the spheres indicated in yellow.

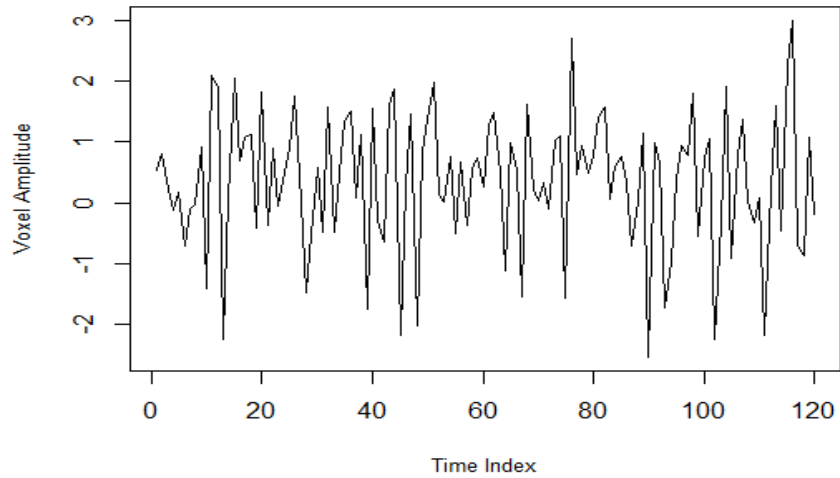


FIGURE 3. This shows the amplitude of a voxel outside of the embedded sphere that does not respond to the experimental task and has been simulated with physiological noise. Simulated with effect size = 5.

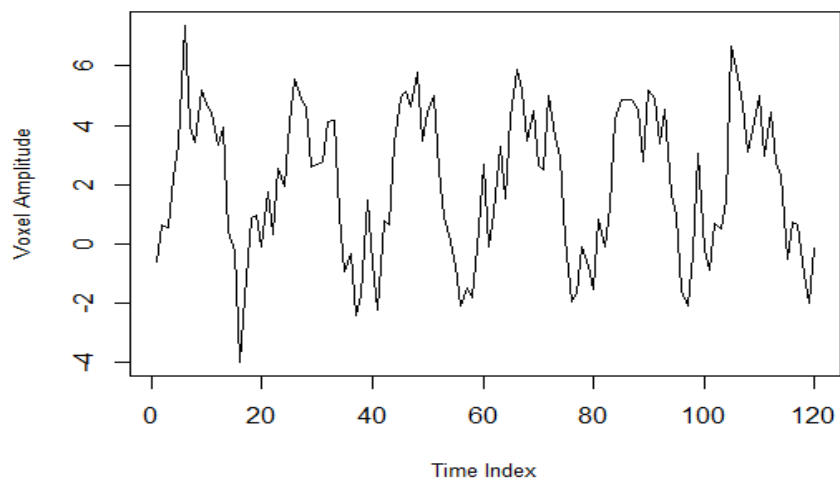
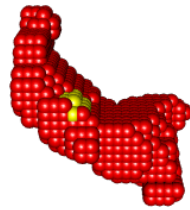
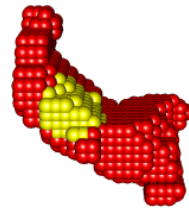


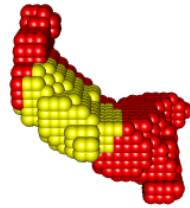
FIGURE 4. This shows the amplitude of a voxel within the embedded sphere that does respond to the periodic experimental task. Simulated with effect size = 5.



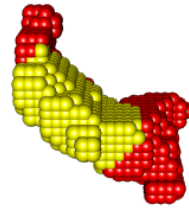
(A) Radius
= 1



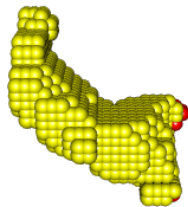
(B) Radius
= 3



(C) Radius
= 5



(D) Radius
= 7



(E) Radius
= 15

FIGURE 5. Simulation Volumes: Spheres of various sizes embedded in a mask of the right hippocampus (lateral view).

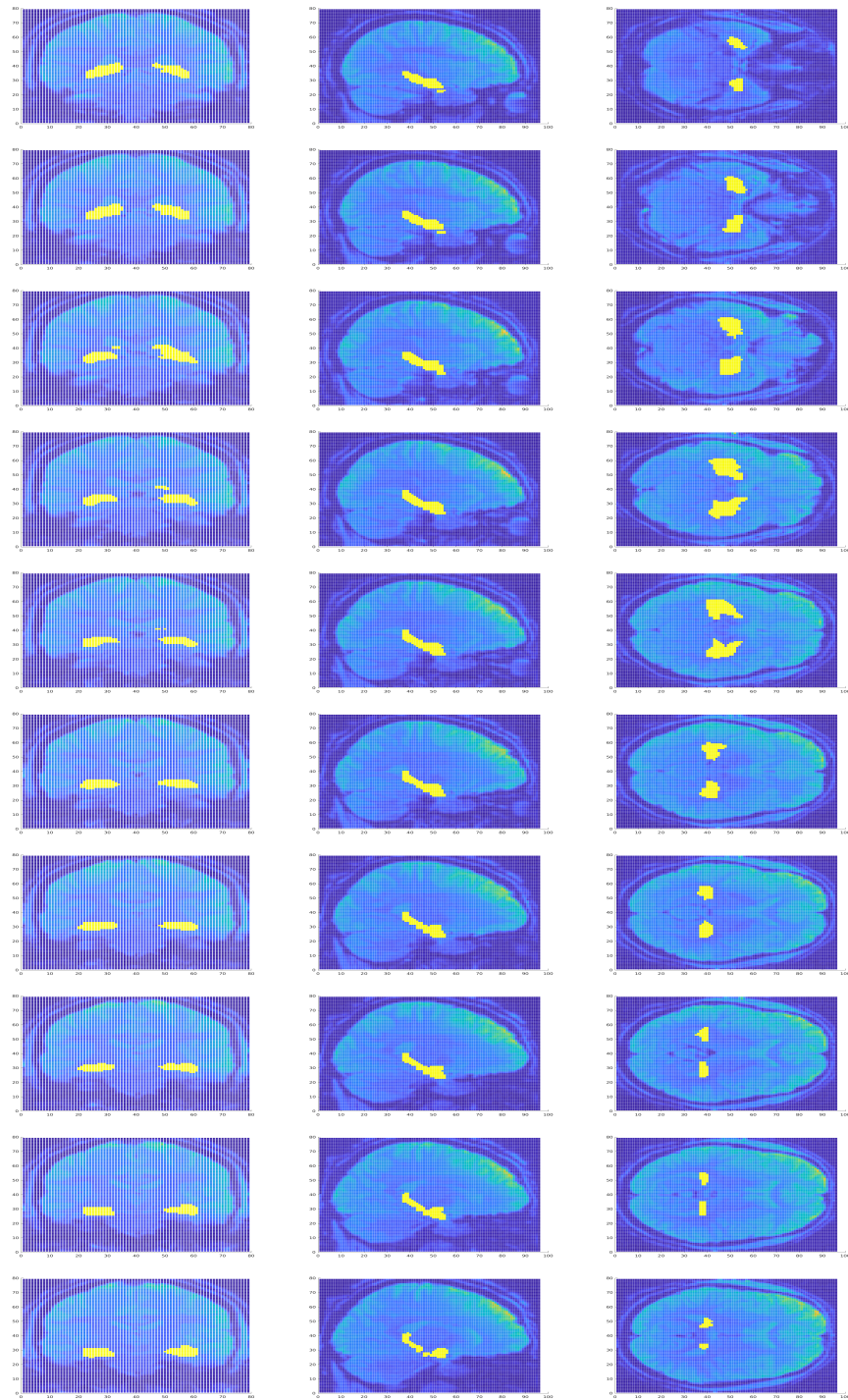


FIGURE 6. Hippocampus mask overlaid onto a brain image.

1 In order to clearly indicate how the simulation volumes in Figure 5 correspond
 2 to the shape of one half of a standard hippocampus mask, in Figure 6 we show
 3 sagittal, coronal, and transverse cross-sections of a hippocampus mask overlaid
 4 onto a brain image. The yellow-highlighted voxels are those in the hippocampus
 5 mask. Our simulated volumes, as pictured in Figure 5, are precisely the voxels in
 6 the right hippocampus.

7 We see from Figure 5 that, in the simulation volumes, the activated regions
 8 (pictured in yellow) do not form tunnel or ring-like shapes. In particular, if we
 9 regard the red regions of the simulation volumes pictured in Figure 5 as subsets of
 10 \mathbb{R}^3 , the classical singular homology group H_1 is trivial. Consequently one expects
 11 to find that the persistent H_1 of these data sets consists of relatively low-persistence
 12 features. This expectation about the simulated data pictured in Figure 5 is borne
 13 out: see below, in Figure 9. (See [3] for an influential study of the sensitivity of
 14 low-persistence features in persistent homology to geometric structure in a data
 15 set.)

16 Signal to noise ratio (SNR) is the magnitude of the signal over the magnitude
 17 of the noise. The SNR establishes the rough amplitude of noise only after the
 18 amplitude of the non-noise signal has already been established. SNR is defined a
 19 variety of ways in the literature. NeuRosim defines average SNR as the following:

$$\text{SNR} = \frac{\bar{S}}{\sigma_N}$$

20 where \bar{S} is the average signal magnitude and σ_N is the standard deviation of
 21 the noise. For this particular definition of SNR, an overview of fMRI studies found
 22 its value to range from 1 to 1000 in the literature [25]. As such, our simulations
 23 included SNR values of 2, 5, 10, and 20.

24 Minimum persistence was also investigated at values of .8, 1, and 1.2. Recall
 25 from section 5.2 that, for fMRI data, we see .8 as a canonical choice to remove noise
 26 from the persistence diagrams.

27 The two normalization functions discussed earlier were also compared, with re-
 28 sults explained in the Results section, below.

29 **5.4. Results.** Our method was evaluated on its ability to identify the task-based
 30 activation of embedded spheres of various radii. This was accomplished by calcu-
 31 lating statistical power. Statistical power is the probability that a method rejects
 32 the null hypothesis when the alternative hypothesis is correct. In this case, the null
 33 hypothesis is that the persistence diagrams of observations during the “resting”
 34 phases of our simulated experiment are no different than the persistence diagrams
 35 of those during the “task” phases. Power was empirically estimated by first simu-
 36 lating each set of parameters 500 times and conducting the permutation test with
 37 2000 permutations for each simulation. The proportion of tests that rejected the
 38 null hypothesis is then our empirical estimate for power. The figures at the end of
 39 this section summarize the empirical power estimates across sphere radius, mini-
 40 mum persistence threshold, and effect size. In addition, Figure 7 gives an example
 41 of a persistence diagram from a “rest” epoch and a persistence diagram from an
 42 “activation” epoch. Though similar, we point out the band of higher-persistence 1-
 43 dimensional features (triangles far from the diagonal) that is present for birth radius
 44 greater than 2 in the “activation” epoch persistence diagram that is not present in
 45 the “rest” epoch persistence diagram. Furthermore, there appears to be a denser

1 cluster of 1-dimensional features in the “rest” epoch persistence diagram compared
2 to the “activation” epoch persistence diagram.

3 The *lower* the minimum persistence threshold for features considered, the more
4 powerful the method became (see Figure 4). This indicates that the incorporation
5 of lower persistence features provides information useful for identifying the activity
6 of the embedded sphere. Performance also improved when increasing maximum
7 radius of homology computed from 3 to 4, informing us that more information im-
8 proved results rather than overwhelming the method. Additionally, the canonical
9 normalization scheme in Definition 4.1 outperformed the Definition 4.2 normaliza-
10 tion scheme for constants 10 and 100, and performed comparably with constant 50.
11 Thus, normalizing the signal to have a similar spread to the spatial coordinates per-
12 forms better than either having a smaller variation in the signal or larger variation
13 in the signal relative to the spatial coordinates.

14 For effect size 5 and above, our method displayed power > 0.85 for all radii
15 except $r=15$ (see Figure 5). The sensitivity of our method to task-activated spheres
16 as small as radius 1 without sub-setting the data is evidence that, even for more
17 subtle patterns of activity, persistence diagrams record differentiating topological
18 structure. The drop-off in power for radius=15 is likely because, as the embedded
19 sphere at that radius made up most of the hippocampus-shaped data, it likely was
20 not as detectable via one-dimensional homological features. Perhaps including zero-
21 dimensional homological features(which represent connected components) would
22 improve sensitivity to larger clusters.

23 Our simulations demonstrate the efficacy of statistical inference using persistent
24 homology to capture associations in task-based fMRI experiments.

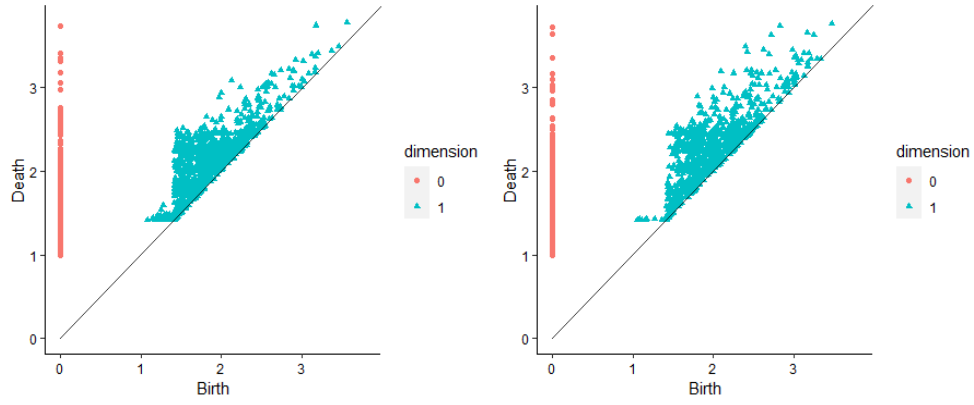


FIGURE 7. On the left is a persistence diagram from a “rest” epoch of our simulation and on the right is a persistence diagram from an “activation” epoch of our simulation. This is for effect size=5, sphere radius=5, and SNR=2.

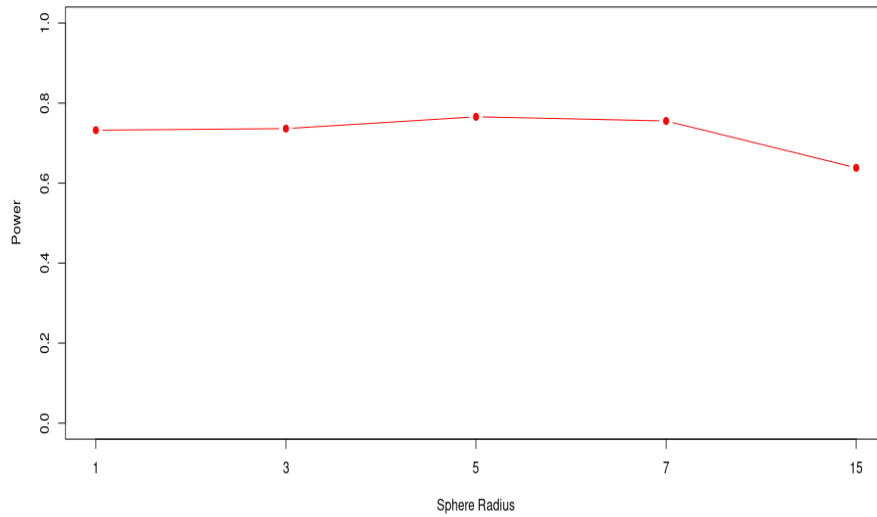


FIGURE 8. Empirical power estimates by radius of embedded sphere.

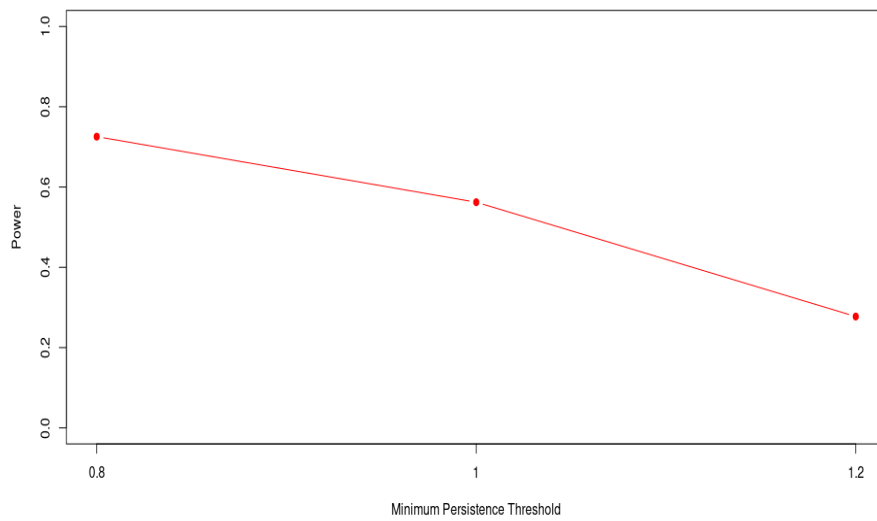


FIGURE 9. Empirical power estimates by minimum persistence threshold of homological features.

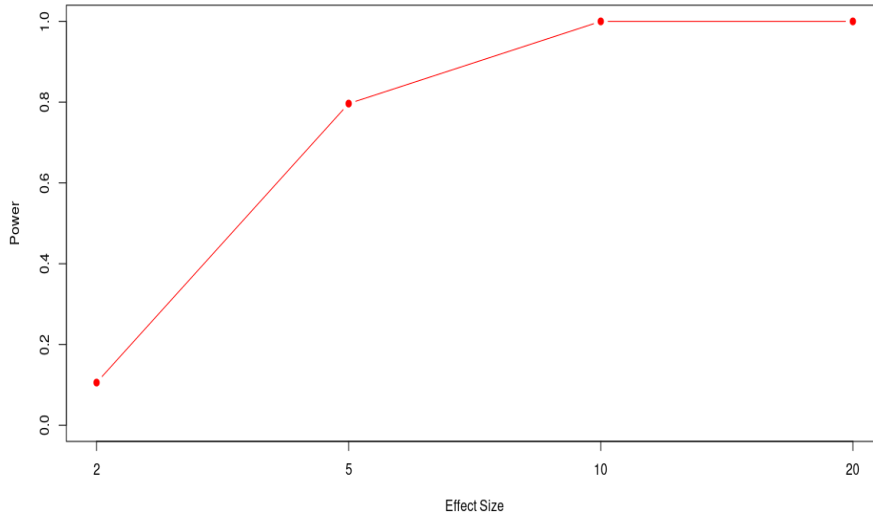


FIGURE 10. Empirical power estimates by effect size for embedded sphere's response to task.

1 **6. Discussion.** We have described a statistical test intended for use on time series
 2 data which comes equipped with a point cloud grouping and a labelling scheme,
 3 as defined in 4.1 and 4.2. Such time series data are typical of task-based fMRI
 4 studies. More generally, any time-series data collected from any source with tempo-
 5 rally distinct “epochs” serves as an appropriate target application. Our statistical
 6 test yields a p-value which provides information on the statistically reliability of
 7 the difference between the veridical labelling scheme of the persistent homology of
 8 the observations against a randomly-assigned labelling scheme. By incorporating
 9 a multi-level block sampling protocol, our test does not have the independence re-
 10 quirement that prevented earlier analogues (in [20] and [4]) from being applied to
 11 time-series data.

12 Our simulated fMRI data was of a simple and typical (in real-world fMRI data)
 13 pattern, consisting of an activated region inside a convex ROI mask (Figure 5),
 14 providing a typical representation of real-world fMRI data. Figures 8 through
 15 10 indicate that our test was able to distinguish between the topological “sig-
 16 nature” (i.e., the persistence diagram) of the simulated fMRI signal during ac-
 17 tive against during resting epochs. Our approach will be valuable to other sci-
 18 entists working with labelled time-series data who a) chose to apply persistent
 19 homology to capture the topological properties of distinct parts of the time se-
 20 ries before b) exploring whether topological properties of the epochs are statis-
 21 tically significant from each other. We invite researchers to apply these meth-
 22 ods (<https://github.com/hassan-abdallah/TimeSeriesTDA>) to their time series
 23 data of choice.

1

REFERENCES

- 2 [1] Sebastian Benzekry, Jack A. Tuszynski, Edward A. Rietman, and Giannoula Lakka Klement.
3 Design principles for cancer therapy guided by changes in complexity of protein-protein in-
4 teraction networks. *Biology Direct*, 10(1), 2015.
- 5 [2] Steven L. Bressler and Anil K. Seth. Wiener–granger causality: A well established method-
6 ology. *NeuroImage*, 58(2):323–329, 2011.
- 7 [3] Peter Bubenik, Michael Hull, Dhruv Patel, and Benjamin Whittle. Persistent homology de-
8 tects curvature. *Inverse Problems*, 36(2):025008, jan 2020.
- 9 [4] Christopher Cericola, Inga Jo Johnson, Joshua Kiers, Mitchell Krock, Jordan Purdy, and
10 Johanna Torrence. Extending hypothesis testing with persistent homology to three or more
11 groups. *Involve*, 11(1):27–51, 2018.
- 12 [5] Frédéric Chazal, Marc Glisse, Catherine Labruère, and Bertrand Michel. Convergence rates
13 for persistence diagram estimation in topological data analysis. *Journal of Machine Learning
14 Research*, 16(110):3603–3635, 2015.
- 15 [6] Frédéric Chazal and Bertrand Michel. An introduction to topological data analysis: funda-
16 mental and practical aspects for data scientists. 10 2017.
- 17 [7] Tamal Krishna Dey and Yusu Wang. *Computational topology for data analysis*. Cambridge
18 University Press, Cambridge, 2022.
- 19 [8] Herbert Edelsbrunner and John L. Harer. *Computational topology*. American Mathematical
20 Society, Providence, RI, 2010. An introduction.
- 21 [9] Anders Eklund, Thomas Nichols, and Hans Knutsson. Can parametric statistical methods be
22 trusted for fMRI based group studies? *arXiv preprint arXiv:1511.01863*, 11 2015.
- 23 [10] Lohmann G, Stelzer J, Lacosse E, Kumar VJ, Mueller K, Kuehn E, Grodd W, and Scheffler
24 K. LISA improves statistical analysis for fMRI. *Nature Communications*, 10 2018.
- 25 [11] Allen Hatcher. *Algebraic topology*. Cambridge University Press, Cambridge, 2002.
- 26 [12] H. W. Kuhn. The Hungarian method for the assignment problem. *Naval Res. Logist. Quart.*,
27 2:83–97, 1955.
- 28 [13] Bradley Macintosh, Richard Mraz, William McIlroy, and Simon Graham. Brain activity during
29 a motor learning task: An fMRI and skin conductance study. *Human brain mapping*, 28:1359–
30 67, 12 2007.
- 31 [14] Christopher Oballe, Alan Cherne, Dave Boothe, Scott Kerick, Piotr J. Franaszczuk, and
32 Vasileios Maroulas. Bayesian topological signal processing. *Discrete Contin. Dyn. Syst. Ser.
33 S*, 15(4):797–, 2022.
- 34 [15] Nina Otter, Mason A Porter, Ulrike Tillmann, Peter Grindrod, and Heather A Harrington.
35 A roadmap for the computation of persistent homology. *EPJ Data Science*, 6:1–38, 2017.
- 36 [16] Jose A Perea. A brief history of persistence. *Morfismos*, 23(1):1–16, 2019.
- 37 [17] Jose A. Perea. Topological times series analysis. *Notices Amer. Math. Soc.*, 66(5):686–694,
38 2019.
- 39 [18] Jose A. Perea and John Harer. Sliding windows and persistence: an application of topological
40 methods to signal analysis. *Found. Comput. Math.*, 15(3):799–838, 2015.
- 41 [19] Nalini Ravishanker and Renjie Chen. An introduction to persistent homology for time series.
42 *Wiley Interdiscip. Rev. Comput. Stat.*, 13(3):Paper No. e1548, 25, 2021.
- 43 [20] Andrew Robinson and Katharine Turner. Hypothesis testing for topological data analysis.
44 *Journal of Applied and Computational Topology*, 10 2013.
- 45 [21] Andrew Salch, Adam Regalski, Hassan Abdallah, Raviteja Suryadevara, Michael J Catanzaro,
46 and Vaibhav A Diwadkar. From mathematics to medicine: A practical primer on topological
47 data analysis (TDA) and the development of related analytic tools for the functional discovery
48 of latent structure in fMRI data. *PLoS ONE*, 16(8):e0255859, 2021.
- 49 [22] Lee M Seversky, Shelby Davis, and Matthew Berger. On time-series topological data analysis:
50 New data and opportunities. In *Proceedings of the IEEE conference on computer vision and
51 pattern recognition workshops*, pages 59–67, 2016.
- 52 [23] Larry Wasserman. Topological data analysis. *Annu. Rev. Stat. Appl.*, 5:501–535, 2018.
- 53 [24] Marijke Welvaert, Joke Durnez, Beatrijs Moerkerke, Geert Verdoolaege, and Yves Rosseel.
54 neuRosim: An R package for generating fMRI data. *Journal of Statistical Software*, 44(10):1–
55 18, 2011.
- 56 [25] Marijke Welvaert and Yves Rosseel. On the definition of signal-to-noise ratio and contrast-
57 to-noise ratio for fmri data. *PloS one*, 8:e77089, 11 2013.

- 1 [26] Anderson Winkler, Matthew Webster, Diego Vidaurre, Thomas Nichols, and Stephen Smith.
2 Multi-level block permutation. *NeuroImage*, 62, 06 2015.

3 Received xxxx 20xx; revised xxxx 20xx.

- 4 *E-mail address:* hassan@wayne.edu
5 *E-mail address:* adam.regalski@wayne.edu
6 *E-mail address:* mohammad.behzad.kang@wayne.edu
7 *E-mail address:* maria.berishaj@wayne.edu
8 *E-mail address:* nkechinnadi@wayne.edu
9 *E-mail address:* er4974@wayne.edu
10 *E-mail address:* vdiwadka@med.wayne.edu
11 *E-mail address:* asalch@wayne.edu