# On Characterizing the Capacity of Neural Networks using Algebraic Topology

**William H. Guss** [1]   **Ruslan Salakhutdinov** [1]
{wguss, rsalakhu}@cs.cmu.edu

## Abstract

The learnability of different neural architectures can be characterized directly by computable measures of data complexity. In this paper, we reframe the problem of architecture selection as understanding how data determines the most expressive and generalizable architectures suited to that data, beyond inductive bias. After suggesting algebraic topology as a measure for data complexity, we show that the power of a network to express the topological complexity of a dataset in its decision region is a strictly limiting factor in its ability to generalize. We then provide the first empirical characterization of the topological capacity of neural networks. Our empirical analysis shows that at every level of dataset complexity, neural networks exhibit topological phase transitions. This observation allowed us to connect existing theory to empirically driven conjectures on the choice of architectures for fully-connected neural networks.

## 1. Introduction

Deep learning has rapidly become one of the most pervasively applied techniques in machine learning. From computer vision (Krizhevsky et al., 2012) and reinforcement learning (Mnih et al., 2013) to natural language processing (Wu et al., 2016) and speech recognition (Hinton et al., 2012), the core principles of hierarchical representation and optimization central to deep learning have revolutionized the state of the art; see Goodfellow et al. (2016). In each domain, a major difficulty lies in selecting the architectures of models that most optimally take advantage of structure in the data. In computer vision, for example, a large body of work ((Simonyan & Zisserman, 2014), (Szegedy et al., 2014), (He et al., 2015), etc.) focuses on improving the initial architectural choices of Krizhevsky et al. (2012) by developing novel network topologies and optimization schemes specific

to vision tasks. Despite the success of this approach, there are still not general principles for choosing architectures in arbitrary settings, and in order for deep learning to scale efficiently to new problems and domains without expert architecture designers, the problem of architecture selection must be better understood.

Theoretically, substantial analysis has explored how various properties of neural networks, (eg. the depth, width, and connectivity) relate to their expressivity and generalization capability ((Raghu et al., 2016), (Daniely et al., 2016), (Guss, 2016)). However, the foregoing theory can only be used to determine an architecture in practice if it is understood how expressive a model need be in order to solve a problem. On the other hand, neural architecture search (NAS) views architecture selection as a compositional hyperparameter search ((Saxena & Verbeek, 2016), (Fernando et al., 2017), (Zoph & Le, 2017)). As a result NAS ideally yields expressive and powerful architectures, but it is often difficult to interpret the resulting architectures beyond justifying their use from their empirical optimality.

We propose a third alternative to the foregoing: data-first architecture selection. In practice, experts design architectures with some inductive bias about the data, and more generally, like any hyperparameter selection problem, the most expressive neural architectures for learning on a particular dataset are solely determined by the nature of the true data distribution. Therefore, architecture selection can be rephrased as follows: *given a learning problem (some dataset), which architectures are suitably regularized and expressive enough to learn and generalize on that problem?*

A natural approach to this question is to develop some objective measure of data complexity, and then characterize neural architectures by their ability to learn subject to that complexity. Then given some new dataset, the problem of architecture selection is distilled to computing the data complexity and choosing the appropriate architecture.

For example, take the two datasets $\mathcal{D}_1$ and $\mathcal{D}_2$ given in Figure 1(a,b) and Figure 1(c,d) respectively. The first dataset, $\mathcal{D}_1$, consists of positive examples sampled from two disks and negative examples from their compliment. On the right, dataset $\mathcal{D}_2$ consists of positive points sampled from two disks and two rings with hollow centers. Under some ge-

ometric measure of complexity $\mathcal{D}_2$ appears more 'complicated' than $\mathcal{D}_1$ because it contains more holes and clusters. As one trains single layer neural networks of increasing hidden dimension on both datasets, *the minimum number of hidden units required to achieve zero testing error is ordered according to this geometric complexity.* Visually in Figure 1, regardless of initialization no single hidden layer neural network with $\leq 12$ units, denoted $h_{\leq 12}$, can express the two holes and clusters in $\mathcal{D}_2$. Whereas on the simpler $\mathcal{D}_1$, both $h_{12}$ and $h_{26}$ can express the decision boundary perfectly. Returning to architecture selection, one wonders if this characterization can be extrapolated; that is, is it true that for datasets with 'similar' geometric complexity to $\mathcal{D}_1$, any architecture with $\geq 12$ hidden learns perfectly, and likewise for those datasets similar in complexity to $\mathcal{D}_2$, architectures with $\leq 12$ hidden units can never learn to completion?

## 1.1. Our Contribution

In this paper, we formalize the above notion of geometric complexity in the language of algebraic topology. We show that questions of architecture selection can be answered by understanding the 'topological capacity' of different neural networks. In particular, a geometric complexity measure, called persistent homology, characterizes the capacity of neural architectures in direct relation to their ability to generalize on data. Using persistent homology, we develop a method which gives the first empirical insight into the learnability of different architectures as data complexity increases. In addition, our method allows us to generate conjectures which tighten known theoretical bounds on the expressivity of neural networks. Finally, we show that topological characterizations of architectures areuseful in practice by presenting a new method, topological architecture selection, and applying it to several OpenML datasets.

## 2. Background

### 2.1. General Topology

In order to more formally describe notions of geometric complexity in datasets, we will turn to the language of topology. Broadly speaking, topology is a branch of mathematics that deals with characterizing shapes, spaces, and sets by their *connectivity*. In the context of characterizing neural networks, we will work towards defining the topological complexity of a dataset in terms of how that dataset is 'connected', and then group neural networks by their capacity to produce decision regions of the same connectivity.

In topology, one understands the relationships between two different spaces of points by the *continuous maps* between them. Informally, we say that two topological spaces $A$ and $B$ are *equivalent* ($A \cong B$) if there is a continuous function $f : A \to B$ that has an inverse $f^{-1}$ that is also continuous.
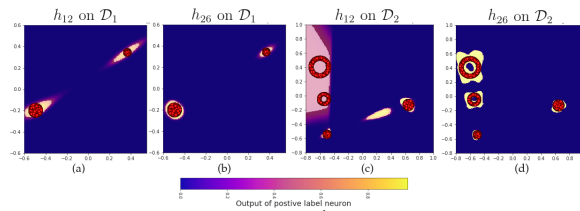


*Figure 1.* The positive label outptus of single hidden layer neural networks, $h_{12}$ and $h_{26}$, of 2 inputs with 12 and 26 hidden units respectively after training on datasets $\mathcal{D}_1$ and $\mathcal{D}_2$ with positive examples in red. Highlighted regions of the output constitute the positive decision region.

When $f$ exists, we say that $A$ and $B$ are *homeomorphic* and $f$ is their *homeomorphism*; for a more detailed treatment of general topology see Bredon (2013). In an informal way, $\mathcal{D}_1 \ncong \mathcal{D}_2$ in Figure 1 since if there were a homeomorphism $f : \mathcal{D}_1 \to \mathcal{D}_2$ at least one of the clusters in $\mathcal{D}_1$ would need to be split discontinuously in order to produce the four different regions in $\mathcal{D}_2$.

The power of topology lies in its capacity to differentiate sets (topological spaces) in a meaningful geometric way that discards certain irrelevant properties such as rotation, translation, curvature, etc. For the purposes of defining geometric complexity, non-topological properties[1] like curvature would further fine-tune architecture selection–say if $\mathcal{D}_2$ had the same regions but with squigly (differentially complex) boundaries, certain architectures might not converge–but as we will show, grouping neural networks by 'topological capacity' provides a powerful minimality condition. That is, we will show that if a certain architecture is incapable of expressing a decision region that is equivalent in topology to training data, then there is no hope of it ever generalizing to the true data.

### 2.2. Algebraic Topology

Algebraic topology provides the tools necessary to not only build the foregoing notion of topological equivalence into a measure of geometric complexity, but also to compute that measure on real data ((Betti, 1872), (Dey et al., 1998), (Bredon, 2013)). At its core, algebraic topology takes topological spaces (shapes and sets with certain properties) and assigns them algebraic objects such as *groups*, *chains*, and other more exotic constructs. In doing so, two spaces can be shown to be topologically equivalent (or distinct) if the algebraic objects to which they are assigned are isomorphic (or not). Thus algebraic topology will allow us to compare the complexity of decision boundaries and datasets by the objects to which they are assigned.

---

[1]A *topological property* or *invariant* is one that is preserved by a homeomorphism. For example, the number of holes and regions which are disjoint from one another are topological properties, whereas curvature is not.

Although there are many flavors of algebraic topology, a powerful and computationally realizable tool is homology.

**Definition 2.1** (Informal, ([Bredon](), [2013])). *If $X$ is a topological space, then $H_n(X) = \mathbb{Z}^{\beta_n}$ is called **the** $n^{th}$ **homology group** of $X$ if the power $\beta_n$ is the number of 'holes' of dimension $n$ in $X$. Note that $\beta_0$ is the number of separate connected components. We call $\beta_n(X)$ the nth Betti number of $X$. Finally, the homology[2] of $X$ is defined as $H(X) = \{H_n(X)\}_{n=0}^\infty$.*

Immediately homology brings us closer to defining the complexity of $\mathcal{D}_1$ and $\mathcal{D}_2$. If we assume that $\mathcal{D}_1$ is not actually a collection of $N$ datapoints, but really the union of 2 solid balls, and likewise that $\mathcal{D}_2$ is the union of 2 solid balls and 2 rings, then we can compute the homology directly. In this case $H_0(\mathcal{D}_1) = \mathbb{Z}^2$ since there are two connected components[3]; $H_1(\mathcal{D}_1) = \{0\}$ since there are no circles (one-dimensional holes); and clearly, $H_n(\mathcal{D}_1) = \{0\}$ for $n \geq 2$. Performing the same computation in the second case, we get $H_0(\mathcal{D}_2) = \mathbb{Z}^4$ and $H_1(\mathcal{D}_2) = \mathbb{Z}^2$ as there are 4 seperate clusters and 2 rings/holes. With respect to any reasonable ordering on homology, $\mathcal{D}_2$ is more complex than $\mathcal{D}_1$. The measure yields non-trivial differentiation of spaces in higher dimension. For example, the homology of a hollow donut is $\{\mathbb{Z}^1, \mathbb{Z}^2, \mathbb{Z}^1, 0, \dots\}$.

Surprisingly, the homology of a space contains a great deal of information about its topological complexity[1]. The following theorem suggests the absolute power of homology to group topologically similar spaces, and therefore neural networks with topologically similar decision regions.

**Theorem 2.2** (Informal). *Let $X$ and $Y$ be topological spaces. If $X \cong Y$ then $H(X) = H(Y)$.[4]*

Intuitively, Theorem 2.2 states that number of 'holes' (and in the case of $H_0(X)$, connected components) are topologically invariant, and can be used to show that two shapes (or decision regions) are different.

### 2.3. Computational Methods for Homological Complexity

In order to compute the homology of both $\mathcal{D}_1$ and $\mathcal{D}_2$ we needed to assume that they were actually the geometric shapes from which they were sampled. Without such assumptions, *for any dataset $\mathcal{D}$ a $H(\mathcal{D}) = \{\mathbb{Z}^N, 0, \dots\}$* where $N$ is the number of data points. This is because, at small enough scales each data point can be isolated as

---

[2]This definition of homology makes many assumptions on $X$ and the base field of computation, but for introductory purposes, this informality is edifying.

[3]Informally, a *connected component* is a set which is not contained in another connected set except for itself.

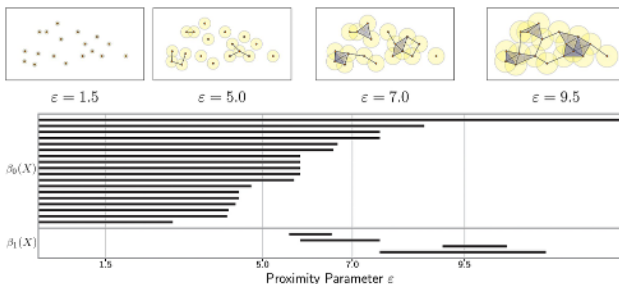[4]Equality of $H(X)$ and $H(Y)$ should be interpreted as isomorphism between each individual $H_i(X)$ and $H_i(Y)$.



*Figure 2.* An illustration of computing persistent homology on a collection of points (([Topaz et al.](), [2015]))

its own connected component; that is, as sets each pair of different positive points $d_1, d_2 \in \mathcal{D}$ are disjoint. To properly utilize homological complexity in better understanding architecture selection, we need to be able to compute the homology of the data directly and still capture meaningful topological information.

Persistent homology, introduced in [Zomorodian & Carlsson]() [(2005)](), avoids the trivialization of computation of dataset homology by providing an algorithm to calculate the homology of a *filtration* of a space. Specifically, a filtration is a topological space $X$ equipped with a sequence of subspaces $X_0 \subset X_1 \subset \cdots \subset X$. In Figure 2 one such particular filtration is given by growing balls of size $\epsilon$ centered at each point, and then letting $X_\epsilon$ be the resulting subspace in the filtration. Define $\beta_n(X)$ to be the $n$th Betti number of the homology $H(X_\epsilon)$ of $X_\epsilon$. Then for example at $\epsilon = 1.5$, $\beta_0(X_\epsilon) = 19$ and $\beta_1(X_\epsilon) = 0$ as every ball is disjoint. At $\epsilon = 5.0$ some connected components merge and $\beta_0(X_\epsilon) = 12$ and $\beta_1(X_\epsilon) = 0$. Finally at $\epsilon = 7$, the union of the balls forms a hole towards the center of the dataset and $\beta_1(X_\epsilon) > 0$ with $\beta_0(X_\epsilon) = 4$.

All together the change in homology and therefore Betti numbers for $X_\epsilon$ as $\epsilon$ changes can be summarized succinctly in the *persistence barcode diagram* given in Figure 2. Each bar in the section $\beta_n(X)$ denotes a 'hole' of dimension $n$. The left endpoint of the bar is the point at which homology detects that particular component, and the right endpoint is when that component becomes indistinguishable in the filtration. When calculating the persistent homology of datasets we will frequently use these diagrams.

With the foregoing algorithms established, we are now equipped with the tools to study the capacity of neural networks in the language of algebraic topology.

## 3. Homological Characterization of Neural Architectures

In the forthcoming section, we will apply persistent homology to empirically characterize the power of certain neural

architectures. To understand why homological complexity is a powerful measure for differentiating architectures, we present the following principle.

Suppose that $\mathcal{D}$ is some dataset drawn from a joint distribution $F$ with continuous CDF on some topological space $X \times \{0, 1\}$. Let $X^+$ denote the support of the distribution of points with positive labels, and $X^-$ denote that of the points with negative labels. Then let $H_S(f) := H[f^{-1}((0, \infty))]$ denote the *support homology* of some function $f : X \to \{0, 1\}$. Essentially $H_S(f)$ is homology of the set of $x$ such that $f(x) > 0$. For a binary classifier, $f$, $H_S(f)$ is roughly a characterization of how many 'holes' are in the positive decision region of $f$. We will sometimes use $\beta_n(f)$ to denote the $n$th Betti number of this support homology. Finally let $\mathcal{F} = \{f : X \to \{0, 1\}\}$ be some family of binary classifiers on $X$.

**Theorem 3.1** (Homological Generalization)**.** *If $X = X^- \sqcup X^+$ and for all $f \in \mathcal{F}$ with $H_S(f) \neq H(X^+)$, then for all $f \in \mathcal{F}$ there exists $A \subset X^+$ so $f$ misclassifies every $x \in A$.*

Essentially, Theorem 3.1 says that if an architecture (a family of models $\mathcal{F}$) is incapable of producing a certain homological complexity, then for any model using that architecture there will always be a set $A$ of true datapoints on which the model will fail. Note that the above principle holds regardless of how $f \in \mathcal{F}$ is attained, learned or otherwise. The principle implies that no matter how well some $f$ learns to correctly classify $\mathcal{D}$ there will always be counter examples in the true data.

In the context of architecture selection, the foregoing minimality condition significantly reduces the size of the search space by eliminating smaller architectures which cannot even express the 'holes' (persistent homology) of the data $H(\mathcal{D})$. This allows us to return to our original question of finding suitably expressive and generalizable architectures but in the very computable language of homological complexity: Let $\mathcal{F}_A$ the set of all neural networks with 'architecture' $A$, then

*Given a dataset $\mathcal{D}$, for which architectures $A$ does there exist a neural network $f \in \mathcal{F}_A$ such that $H_S(f) = H(\mathcal{D})$?*

We will resurface a contemporary theoretical view on this question, and thereafter make the first steps towards an empirical characterization of the capacity of neural architectures in the view of topology.

### 3.1. Theoretical Basis for Neural Homology

Theoretically, the homological complexity of neural network can be framed in terms of *the sum of the number of holes* expressible by certain architectures. In particular, Bianchini et al. (2014) gives an analysis of how the maximum sum of
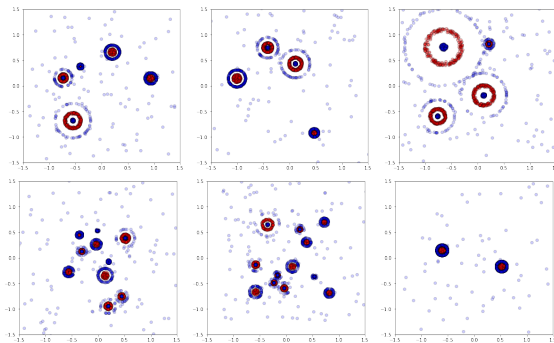


*Figure 3.* Scatter plots of 6 different synthetic datasets of varying homological complexity.

Betti numbers grows as $\mathcal{F}_A$ changes. The results show that the width and activation of a fully connected architecture effect its topological expressivity to varying polynomial and exponential degrees.

What is unclear from this analysis is how these bounds describe expressivity or learnability in terms of *individual* Betti numbers. From a theoretical perspective Bianchini et al. (2014) stipulated that a characterization of individual homology groups require the solution of deeper unsolved problems in algebraic topology. However, for topological complexity to be effective in architecture selection, understanding *each* Betti number is essential in that it grants direct inference of architectural properties from the persistent homology of the data. Therefore we turn to an empirical characterization.

### 3.2. Empirical Characterization

To understand how the homology of data determines expressive architectures we characterize the capacities of architectures with an increasing number of layers and hidden units to *learn* and *express* homology on datasets of varying homological complexity.

Restricting[5] our analysis to the case of $n = 2$ inputs, we generate binary datasets of increasing homological complexity by sampling $N = 5000$ points from mixtures of $\text{Unif}(\mathbb{S}^1)$ and $\text{Unif}(B^2)$, uniform random distributions on solid and empty circles with known support homologies. The homologies chosen range contiguously from $H(\mathcal{D}) = \{\mathbb{Z}^1, 0\}$ to $H(\mathcal{D}) = \{\mathbb{Z}^{30}, \mathbb{Z}^{30}\}$ and each sampled distribution is geometrically balanced *i.e.* each topological feature occupies the same order of magnitude. Additionally, margins were induced between the classes to aid learning. Examples are shown in Figure 3.

---

[5] Although we chose to study the low dimensional setting because it allows us to compute the persistent homology of the decision region directly, the convergence analysis extends to any number of dimensions.
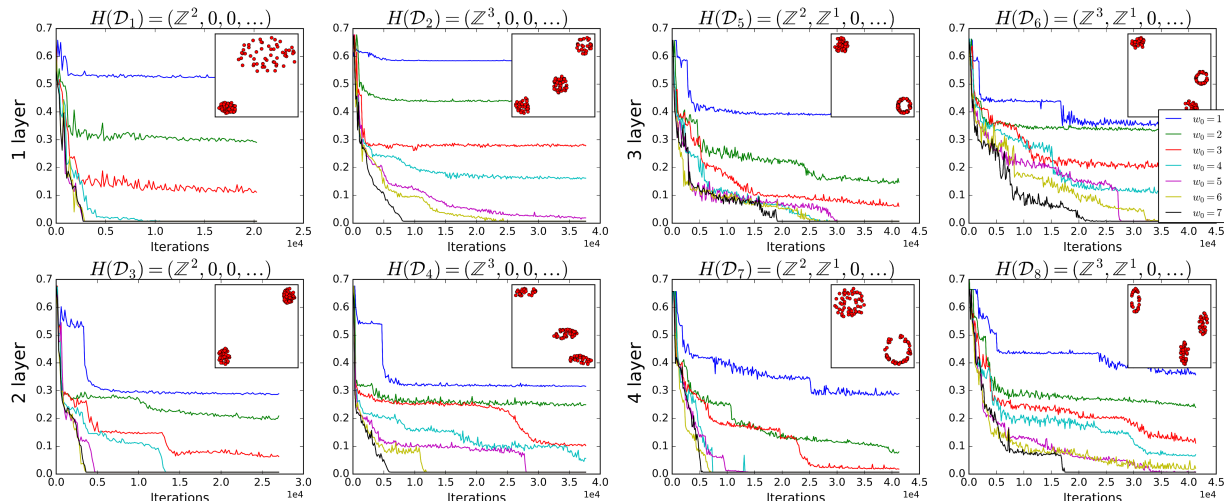
*Figure 4.* Topological phase transitions in low dimensional neural networks as the homological complexity of the data increases. The upper right corner of each plot is a dataset on which the neural networks of increasing first layer hidden dimension are trained. Each plot gives the minimum error for each architecture versus the number of minibatches seen.

To characterize the difficulty of learning homologically complex data in terms of both depth and width, we consider fully connected architectures with ReLu activation functions (Nair & Hinton, 2010) of depth $\ell = \{1, \ldots, 6\}$ and width $h_l = \beta_0(\mathcal{D})$ when $1 \leq l \leq \ell$ and $h_l \in \{1, \ldots, 500\}$ when $l = 0$. We will denote individual architectures by the pair $(\ell, h_0)$. We vary the number of hidden units in the first layer, as they form a half-space basis for the decision region. The weights of each architecture are initialized to samples from a normal distribution $\mathcal{N}(0, \frac{1}{\beta_0})$ with variance respecting the scale of each synthetic dataset. For each homology we take several datasets sampled from the foregoing procedure and optimize 100 initializations of each architecture against the standard cross-entropy loss. To minimize the objective we use the Adam optimizer (Kingma & Ba, 2014) with a fixed learning rate of 0.01 and an increasing batch size schedule (Smith et al., 2017).

We compare each architectures average and best performance by measuring misclassification error over the course of training and homological expressivity at the end of training. The latter quantity, given by

$$E_H^p(f, \mathcal{D}) = \min\left\{ \frac{\beta_p(f)}{\beta_p(\mathcal{D})}, 1 \right\},$$

measures the capacity of a model to exhibit the true homology of the data. We compute the homology of individual decision regions by constructing a filtration on Heaviside step function of the difference of the outputs, yielding a persistence diagram with the exact homological components of the decision regions. The results are summarized in Figures 4-6

The resulting convergence analysis indicates that neural networks exhibit a statistically significant *topological phase transition* during learning which depends directly on the homological complexity of the data. For any dataset and any random homeomorphism applied thereto, the best error of architectures with $\ell$ layers and $h$ hidden units (on the first layer) is *strictly* limited in magnitude and convergence time by $h_{phase}$. For example in Figure 4, $\ell = 3$ layer neural networks fail to converge for $h < h_{phase} = 4$ on datasets with homology $H(\mathcal{D}_5) = (\mathbb{Z}^2, \mathbb{Z}^1, \ldots)$.

More generally, homological complexity directly effects the efficacy of optimization in neural networks. As shown in Figure 5, taking any increasing progression of homologies against the average convergence time of a class of architectures yields an approximately monotonic relationship; in this case, convergence time at $h_{phase}$ increases with increasing $\beta_p(\mathcal{D})$, and convergence time at $h > h_{phase}$ decreases with fixed $\beta_p(\mathcal{D})$. A broader analysis for a varying number of layers is given in the appendix.

Returning to the initial question of architecture selection, the analysis of empirical estimation of $E_H^p(f, \mathcal{D})$ provides the first complete probabilistic picture of the homological expressivity of neural architectures. For architectures with $\ell \in \{1, 2, 3, 4\}$ and $h_0 \in \{0, \ldots, 30\}$ Figure 6 displays the estimated probability that $(\ell, h_0)$ expresses the homology of the decision region after training. Specifically, Figure 6(top) indicates that, for $\ell = 1$ hidden layer neural networks, $\max \beta_0(f)$ is clearly $\Omega(h_0)$. Examining $\ell = 2, 3, 4$ in Figure 6(top), we conjecture[6] that as $\ell$ increases

$$\max_{f \in F_A} \beta_0(f) \in \Omega(h_0^\ell),$$

---

[6]We note that $\beta_0(f)$ does not depend on $\beta_1(\mathcal{D})$ in the experiment for datasets with $\beta_1(\mathcal{D}) > \beta_0(\mathcal{D})$ were not generated.
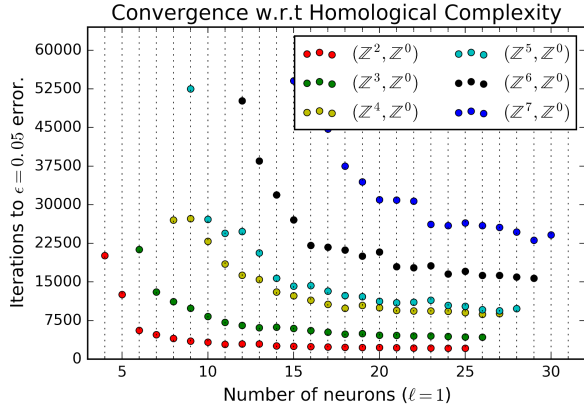
*Figure 5.* A scatter plot of the number of iterations required for single-layer architectures of varying hidden dimension to converge to 5% misclassification error. The colors of each point denote the topological complexity of the data on which the networks were trained. Note the emergence of monotonic bands. Multilayer plots given in the appendix look similar.

by application of Theorem 3.1 to the expressivity estimates. Therefore

$$h_{phase} \geq C \sqrt[\ell]{\beta_0(\mathcal{D})}. \tag{3.1}$$

Likewise, in Figure 6(bottom), each horizontal matrix gives the probability of expressing $\beta_1(\mathcal{D}) \in \{1, 2\}$ for each layer. This result indicates that higher order homology is extremely difficult to learn in the single-layer case, and as $\ell \to \infty$, $\max_{f \in f_A} \beta_1(f) \to n = 2$, the input dimension.

The importance of the foregoing empirical characterization for architecture selection is two-fold. First, by analyzing the effects of homological complexity on the optimization of different architectures, we were able to conjecture probabilistic bounds on the *learnable* homological capacity of neural networks. Thus, predictions of minimal architectures using those bounds are sufficient enough to learn data homology up to homeomorphism. Second, the analysis of individual Betti numbers enables data-first architecture selection using persistent homology.

## 4. Topological Architecture Selection

We have thus far demonstrated the discriminatory power of homological complexity in determining the expressivity of architectures. However, for homological complexity to have any practical use in architecture selection, it must be computable on real data, and more generally real data must have non-trivial homology. In the following section we present a method for relating persistent homology of any dataset to a minimally expressive architecture predicted by the foregoing empirical characterization, and then we experimentally validate our method on several datasets.

Topological architecture selection is comprised of three steps: given a dataset, compute its persistent homology;
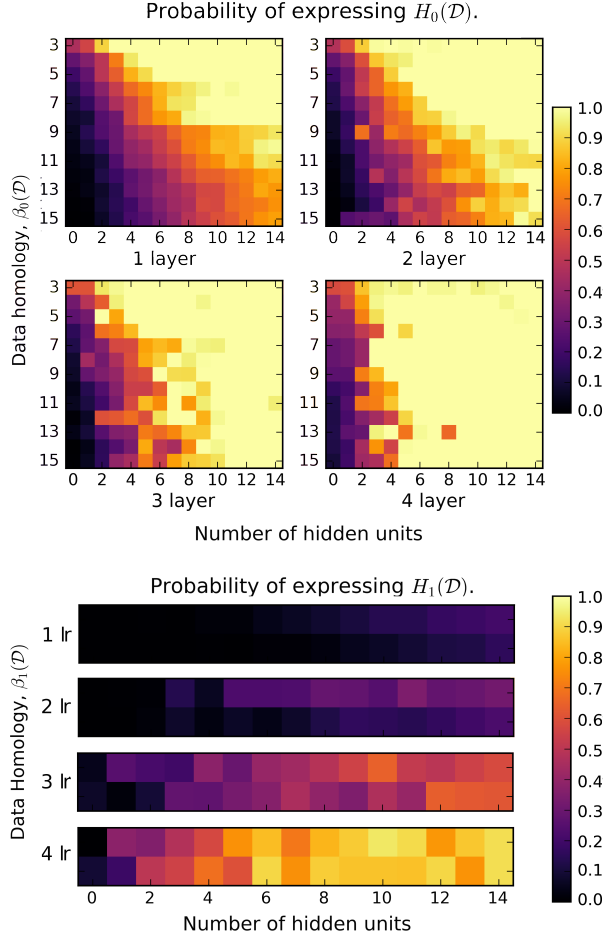




*Figure 6.* A table of estimated probabilities of different neural architectures to express certain homological features of the data after training. Top: the probabilities of express homologies with increasing $\beta_0$ as a function of layers and neurons. Bottom: The probabilities of expressing $\beta_1 \in \{1, 2\}$ as a function of layers and neurons.

determine an appropriate scale at which to accept topological features as pertinent to the learning problem; and infer a lower-bound on $h_{phase}$ from the topological features at or above the decided scale.

The extraction of static homology from persistence homology, while aesthetically valid (Carlsson et al., 2008), is ill-posed in many cases. For the purposes of architecture selection, however, exact reconstruction of the original homology is not necessary. Given a persistence diagram, $D_p$ containing the births $b_i$ and deaths $d_i$ of features in $H_p(\mathcal{D})$, let $\epsilon$ be given and consider all $\alpha_i = (b_i, d_i)$ such that when $|d_i - b_i| > \epsilon$ we assume $\alpha$ to be a topological component of the real space. Then the resulting homologies form a filtration $H_p^{\epsilon_1}(\mathcal{D}) \subset H_p^{\epsilon_2}$ for $\epsilon_1, \epsilon_2 \in \mathbb{R}^+$. If $\epsilon$ is chosen such that certain topologically noisy features (Fasy et al., 2014) are included in the estimate of $h_{phase}$, then at worst the architecture is overparameterized, but still learns. If $\epsilon$ is such that the estimated $h_{phase}$ is underrepresentitive of the
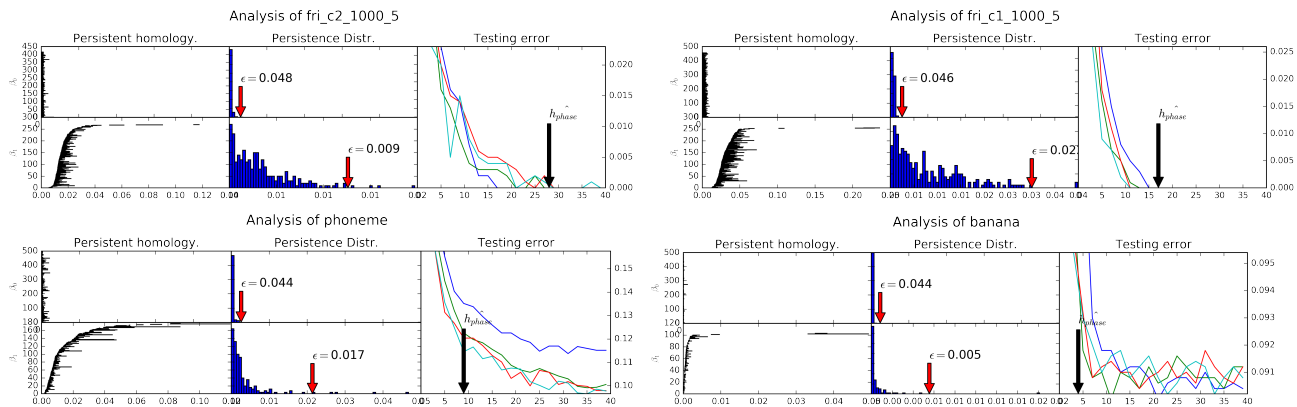
*Figure 7.* Topological architecture selection applied to four different datasets. The persistent homology, the histogram of topological lifespans, and the predicted $\hat{h}_{phase}$ are indicated. For each test error plot, the best performance of one (blue), two (green), three (red), and four (light blue) layer neural networks are given in terms of the number of hidden units on the first layer.

topological features of a space, then at worst, the architecture is underparameterized but potentially close to $h_{phase}$. As either case yields the solution or a plausibly useful seed for other algorithm selection algorithms (Zoph & Le, 2017; Feurer et al., 2015), we adopt the this static instantiation of persistence homology.

In order to select an architecture $(\ell, h_0)$ lower-bounding $h_{phase}$, we restrict the analysis to the case of $\ell = 1$ and regress a multilinear model training on pairs

$$(b_0, b_1) \mapsto \arg\min_m E^p_H(f_m, \mathcal{D}) \geq 1, \beta_*(\mathcal{D}) = (b_0, b_1)$$

over all $m$ hidden unit single layer neural networks $f_m$ and synthetic datasets of known homology $\mathcal{D}$ from our previous experiments. The resultant discretization of the model gives a lower-bound estimate after applying the bound from (3.1):

$$\hat{h}_{phase}(\beta_0, \beta_1) \geq \beta_1 C \sqrt[\ell]{(\beta_0)} + 2 \qquad (4.1)$$

Estimating this lower-bound is at the core of neural homology theory and is the subject of substantial future theoretical and empirical work.

### 4.1. Results

In order to validate the approach, we applied topological architecture selection to several binary datasets from the OpenML dataset repository (Vanschoren et al., 2013): `fri_c`, `balance-scale`, `banana`, `phoneme`, and `delta_ailerons`. We compute persistent homology of each of the two labeled classes therein and accept topological features with lifespans greater than two standard deviations from the mean for each homological dimension. We then estimated a lowerbound on $h_{phase}$ for single hidden layer neural networks using 4.1. Finally we trained 100 neural networks for each architecture $(\ell, h_0)$ where $h_0 \in \{1, 3, \ldots, 99\}$ and $\ell \in \{1, \ldots, 4\}$. During training we

record the minimum average error for each $h_0$ and compare this with the estimate $\hat{h}_{phase}$. The results are summarized in 4.

These preliminary findings indicate that the empirically derived estimate of $h_{phase}$ provides a strong starting point for architecture selection, as the minimum error at $\hat{h}_{phase}$ is near zero in every training instance. Although the estimate is given only in terms of $0$ and $1$ dimensional homology of the data, it still performed well for higher dimensional datasets such as `phoneme` and `fri_c*`. In failure cases, choice of $\epsilon$ greatly affected the predicted $h_{phase}$, and thus it is imperative that more adaptive topological selection schemes be investigated.

While our analysis and characterization is given for the the decision regions of individual classes in a dataset, it is plausible that the true decision boundary is topologically simple despite the complexity of the classes. Although we did not directly characterize neural network by the topology of their decision boundaries, recent work by Varshney & Ramamurthy (2015) provides an exact method for computing the persistent homology of a decision boundary between any number of classes in a dataset. It is the subject of future work to provide a statistical foundation for (Varshney & Ramamurthy, 2015) and then reanalyze the homological capacity of neural networks in this context.

### 4.2. Homological Complexity of Real Data

In addition to testing topological measures of complexity in the setting of architecture selection, we verify that common machine learning benchmark datasets have non-trivial homology and the computation thereof is tractable.

**CIFAR-10.** We compute the persistent homology of several classes of CIFAR-10 using the Python library Dionysus.
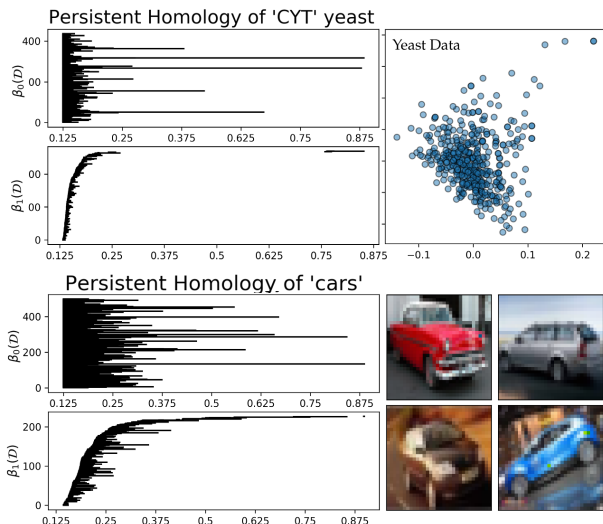
*Figure 8.* The persistent homology barcodes of classes in the CIFAR-10 Datasets; The barcode for the dimensions 0 and 1 for the 'cars' class along side different samples. Note how different orientations are shown.

Current algorithms for persistent homology do not deal well with high dimensional data, so we embed the entire dataset in $\mathbb{R}^3$ using local linear embedding (LLE; (Saul & Roweis, 2000)) with $K = 120$ neighbors. We note that an embedding of any dataset to a lower dimensional subspace with small enough error roughly preserves the static homology of the support of the data distribution distribution by Theorem 2.2. After embedding the dataset, we take a sample of 1000 points from example class 'car' and build a persistent filtration by constructing a Vietoris-Rips complex on the data. The resulting complex has 20833750 simplices and took 4.3 min. to generate. Finally, computation of the persistence diagram shown in Figure 8 took 8.4 min. locked to a single thread on a Intel Core i7 processor. The one-time cost of computing persistent homology could easily augment any neural architecture search.

Although we only give an analysis of dimension 2 topological features–and there is certainly higher dimensional homological information in CIFAR-10–the persistence barcode diagram is rich with different components in both $H_0(\mathcal{D})$ and $H_1(\mathcal{D})$. Intuitively, CIFAR contains pictures of cars rotated across a range of different orientations and this is exhibited in the homology. In particular, several holes are born and die in the range $\epsilon \in [0.15, 0.375]$ and one large loop from $\epsilon \in [0.625, 0.82]$.

**UCI Datasets.** We further compute the homology of three low dimensional UCI datasets and attempt to assert the of non-trivial , $h_{phase}$. Specifically, we compute the persistent homology of the majority classes in the Yeast Protein Localization Sites, UCI Ecoli Protein Localization Sites, and HTRU2 datasets. For these datasets no dimensional-

ity reduction was used. In Figure 8(left), the persistence barcode exhibits two separate significant loops (holes) at $\epsilon \in [0.19, 0.31]$ and $\epsilon \in [0.76, 0.85]$, as well as two major connected components in $\beta_0(\mathcal{D})$. The Other persistence diagrams are relegated to the appendix.

## 5. Related Work

We will place this work in the context of deep learning theory as it relates to expressivity. Since the seminal work of Cybenko (1989) which established standard universal approximation results for neural networks, many researchers have attempted to understand the expressivity of certain neural architectures. Pascanu et al. (2013) and MacKay (2003) provided the first analysis relating the depth and width of architectures to the complexity of the sublevel sets they can express. Motivated therefrom, Bianchini et al. (2014) expressed this theme in the language of Pfefferian functions, thereby bounding the sum of Betti numbers expressed by sublevel sets. Finally Guss (2016) gave an account of how topological assumptions on the input data lead to optimally expressive architectures. In parallel, Eldan & Shamir (2016) presented the first analytical minimality result in expressivity theory; that is, the authors show that there are simple functions that cannot be expressed by two layer neural networks without exponential dependence on input dimension. This work spurred the work of Poole et al. (2016), Raghu et al. (2016) which reframed expressivity in a differential geometric lens.

Our work presents the first method to derive expressivity results empirically. Our topological viewpoint sits dually with its differential geometric counterpart, and in conjunction with the work of (Poole et al., 2016) and (Bianchini et al., 2014). This duality implies that when topological expression is not possible, exponential differential expressivity allows networks to bypass homological constraints at the cost of adversarial sets. Furthermore, our work opens a practical connection between the foregoing theory on neural expressivity and architecture selection, with the potential to substantially improve neural architecture search (Zoph & Le, 2017) by directly computing the capacities of different architectures.

## 6. Conclusion

Architectural power is closely related to the algebraic topology of decision regions. In this work we distilled neural network expressivity into an empirical question of the generalization capabilities of architectures with respect to the homological complexity of learning problems. This view allowed us to provide an empirical method for developing tighter characterizations on the the capacity of different architectures in addition to a principled approach to guiding

architecture selection by computation of persistent homology on real data.

There are several potential avenues of future research in using homological complexity to better understand neural architectures. First, a full characterization of neural networks with convolutional linearities and state-of-the-art topologies is a crucial next step. Our empirical results suggest that there are exact formulas describing the of power of neural networks to express decision boundaries with certain properties. Future theoretical work in determining these forms would significantly increase the efficiency and power of neural architecture search, constraining the search space by the persistent homology of the data. Additionally, we intend on studying how the topological complexity of data changes as it is propagated through deeper architectures.

## 7. Acknowledgements

## References

Betti, E. Il nuovo cimento. *Series*, 2:7, 1872.

Bianchini, Monica et al. On the complexity of shallow and deep neural network classifiers. In *ESANN*, 2014.

Bredon, Glen E. *Topology and geometry*, volume 139. Springer Science & Business Media, 2013.

Carlsson, Gunnar, Ishkhanov, Tigran, De Silva, Vin, and Zomorodian, Afra. On the local behavior of spaces of natural images. *International journal of computer vision*, 76(1):1–12, 2008.

Cybenko, George. Approximation by superpositions of a sigmoidal function. *Mathematics of Control, Signals, and Systems (MCSS)*, 2(4):303–314, 1989.

Daniely, Amit, Frostig, Roy, and Singer, Yoram. Toward deeper understanding of neural networks: The power of initialization and a dual view on expressivity. In *Advances In Neural Information Processing Systems*, pp. 2253–2261, 2016.

Dey, TK, Edelsbrunner, H, and Guha, S. Computational topology, invited paper in advances in discrete and computational geometry, eds. b. chazelle, je goodmann and r. pollack. *Contemporary Mathematics, AMS*, 1998.

Eldan, Ronen and Shamir, Ohad. The power of depth for feedforward neural networks. In *Conference on Learning Theory*, pp. 907–940, 2016.

Fasy, Brittany Terese, Lecci, Fabrizio, Rinaldo, Alessandro, Wasserman, Larry, Balakrishnan, Sivaraman, Singh, Aarti, et al. Confidence sets for persistence diagrams. *The Annals of Statistics*, 42(6):2301–2339, 2014.

Fernando, Chrisantha, Banarse, Dylan, Blundell, Charles, Zwols, Yori, Ha, David, Rusu, Andrei A., Pritzel, Alexander, and Wierstra, Daan. Pathnet: Evolution channels gradient descent in super neural networks. *CoRR*, abs/1701.08734, 2017. URL http://arxiv.org/abs/1701.08734.

Feurer, Matthias, Springenberg, Jost Tobias, and Hutter, Frank. Initializing bayesian hyperparameter optimization via meta-learning. In *AAAI*, pp. 1128–1135, 2015.

Goodfellow, Ian, Bengio, Yoshua, and Courville, Aaron. *Deep Learning*. MIT Press, 2016. http://www.deeplearningbook.org.

Guss, William H. Deep function machines: Generalized neural networks for topological layer expression. *arXiv preprint arXiv:1612.04799*, 2016.

He, Kaiming, Zhang, Xiangyu, Ren, Shaoqing, and Sun, Jian. Deep residual learning for image recognition. *CoRR*, abs/1512.03385, 2015. URL http://arxiv.org/abs/1512.03385.

Hinton, Geoffrey, Deng, Li, Yu, Dong, Dahl, George E, Mohamed, Abdel-rahman, Jaitly, Navdeep, Senior, Andrew, Vanhoucke, Vincent, Nguyen, Patrick, Sainath, Tara N, et al. Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups. *IEEE Signal Processing Magazine*, 29(6):82–97, 2012.

Kingma, Diederik P and Ba, Jimmy. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

Krizhevsky, Alex, Sutskever, Ilya, and Hinton, Geoffrey E. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pp. 1097–1105, 2012.

MacKay, David JC. *Information theory, inference and learning algorithms*. Cambridge university press, 2003.

Mnih, Volodymyr, Kavukcuoglu, Koray, Silver, David, Graves, Alex, Antonoglou, Ioannis, Wierstra, Daan, and Riedmiller, Martin. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.

Nair, Vinod and Hinton, Geoffrey E. Rectified linear units improve restricted boltzmann machines. In *Proceedings of the 27th international conference on machine learning (ICML-10)*, pp. 807–814, 2010.

Pascanu, Razvan, Montufar, Guido, and Bengio, Yoshua. On the number of response regions of deep feed forward networks with piece-wise linear activations. *arXiv preprint arXiv:1312.6098*, 2013.

Poole, Ben, Lahiri, Subhaneil, Raghu, Maithreyi, Sohl-Dickstein, Jascha, and Ganguli, Surya. Exponential expressivity in deep neural networks through transient chaos. In *Advances In Neural Information Processing Systems*, pp. 3360–3368, 2016.

Raghu, Maithra, Poole, Ben, Kleinberg, Jon, Ganguli, Surya, and Sohl-Dickstein, Jascha. On the expressive power of deep neural networks. *arXiv preprint arXiv:1606.05336*, 2016.

Saul, Lawrence K and Roweis, Sam T. An introduction to locally linear embedding. *unpublished. Available at: http://www. cs. toronto. edu/˜ roweis/lle/publications. html*, 2000.

Saxena, Shreyas and Verbeek, Jakob. Convolutional neural fabrics. *CoRR*, abs/1606.02492, 2016. URL http://arxiv.org/abs/1606.02492.

Simonyan, Karen and Zisserman, Andrew. Very deep convolutional networks for large-scale image recognition. *CoRR*, abs/1409.1556, 2014. URL http://arxiv.org/abs/1409.1556.

Smith, Samuel L, Kindermans, Pieter-Jan, and Le, Quoc V. Don't decay the learning rate, increase the batch size. *arXiv preprint arXiv:1711.00489*, 2017.

Szegedy, Christian, Liu, Wei, Jia, Yangqing, Sermanet, Pierre, Reed, Scott E., Anguelov, Dragomir, Erhan, Dumitru, Vanhoucke, Vincent, and Rabinovich, Andrew. Going deeper with convolutions. *CoRR*, abs/1409.4842, 2014. URL http://arxiv.org/abs/1409.4842.

Topaz, Chad M, Ziegelmeier, Lori, and Halverson, Tom. Topological data analysis of biological aggregation models. *PloS one*, 10(5):e0126383, 2015.

Vanschoren, Joaquin, van Rijn, Jan N., Bischl, Bernd, and Torgo, Luis. Openml: Networked science in machine learning. *SIGKDD Explorations*, 15(2):49–60, 2013. doi: 10.1145/2641190.2641198. URL http://doi.acm.org/10.1145/2641190.2641198.

Varshney, Kush R and Ramamurthy, Karthikeyan Natesan. Persistent topology of decision boundaries. In *Acoustics, Speech and Signal Processing (ICASSP), 2015 IEEE International Conference on*, pp. 3931–3935. IEEE, 2015.

Wu, Yonghui, Schuster, Mike, Chen, Zhifeng, Le, Quoc V, Norouzi, Mohammad, Macherey, Wolfgang, Krikun, Maxim, Cao, Yuan, Gao, Qin, Macherey, Klaus, et al. Google's neural machine translation system: Bridging the gap between human and machine translation. *arXiv preprint arXiv:1609.08144*, 2016.

Zomorodian, Afra and Carlsson, Gunnar. Computing persistent homology. *Discrete & Computational Geometry*, 33(2):249–274, 2005.

Zoph, Barret and Le, Quoc V. Neural architecture search with reinforcement learning. 2017. URL https://arxiv.org/abs/1611.01578.

# A. Proofs, Conjectures, and Formal Definitions

## A.1. Homology

Homology is naturally described using the language of category theory. Let $Top^2$ denote the category of topological spaces and $Ab$ the category of abelian groups.

**Definition A.1** (Homology Theory, (Bredon, 2013)). *A homology theory on the on $Top^2$ is a function $H : Top^2 \to Ab$ assigning to each pair $(X, A)$ of spaces a graded (abelian) group $\{H_p(X, A)\}$, and to each map $f : (X, A) \to (Y, B)$, homomorphisms $f_* : H_p(X, A) \to H_p(Y, B)$, together with a natural transformation of functors $\partial_* : H_p(X, A) \to H_{p-1}(X, A)$, called the connecting homomorphism (where we use $H_*(A)$ to denote $H_*(A, \emptyset)$) such that the following five axioms are satisfied.*

1. *If $f \simeq g : (X, A) \to (Y, B)$ then $f_* = g_* : H_*(X, A) \to H_*(Y, B)$.*

2. *For the inclusions $i : A \to X$ and $j : X \to (X, A)$ the sequence sequence of inclusions and connecting homomorphisms are exact.*

3. *Given the pair $(X, A)$ and an open set $U \subset X$ such that $cl(U) \subset int(A)$ then the inclusion $k : (X - U, A - U) \to (X, A)$ induces an isomorphism $k_* : H_*(X - U, A - U) \to H_*(X, A)$*

4. *For a one point space $P$, $H_i(P) = 0$ for all $i \neq 0$.*

5. *For a topological sum $X = +_\alpha X_\alpha$ the homomorphism*

$$\bigoplus (i_\alpha)_* : \bigoplus H_n(X_\alpha) \to H_n(X)$$

*is an isomorphism, where $i_\alpha : X_\alpha \to X$ is the inclusion.*

For related definitions and requisite notions we refer the reader to (Bredon, 2013).

## A.2. Proof of Theorem 3.1

**Theorem A.2.** *Let $X$ be a topological space and $X^+$ be some open subspace. If $\mathcal{F} \subset 2^X$ such that $f \in \mathcal{F}$ implies $H_S(f) \neq H(X^+)$, then for all $f \in \mathcal{F}$ there exists $A \subset X$ so that $f(A \cap X^+) = \{0\}$ and $f(A \cap (X \setminus X^+)) = \{1\}$.*

*Proof.* Suppose the for the sake of contraiction that for all $f \in \mathcal{F}$, $H_S(f) \neq H(X^+)$ and yet there exists an $f$ such that for all $A \subset X$, there exists an $x \in A$ such that $f(x) = 1$. Then take $\mathcal{A} = \{x\}_{x \in X}$, and note that $f$ maps each singleton into its proper partition on $X$. We have that for any open subset of $V \subset X^+$, $f(V) = \{1\}$, and for any closed subset $W \subset X \setminus X^+$, $f(W) = \{0\}$. Therefore

$X^+ = \bigcup_{A \in \tau_{X^+ \cap X}} A \subset supp(f)$ as the subspace topology $\tau_{X^+ \cap X} = \tau_{X^+} \cap \tau_X$ where $\tau_{X^+} = \{A \in \tau_X \mid A \subset X^+\}$ and $\tau_X$ denotes the topology of $X$. Likewise, $int(X^-) \subset X \setminus supp(F)$ under the same logic. Therefore $supp(f)$ has the exact same topology as $X^+$ and so by Theorem 2.2 $H(X^+) = H(supp(f))$ but this is a contradiction. This completes the proof. $\square$
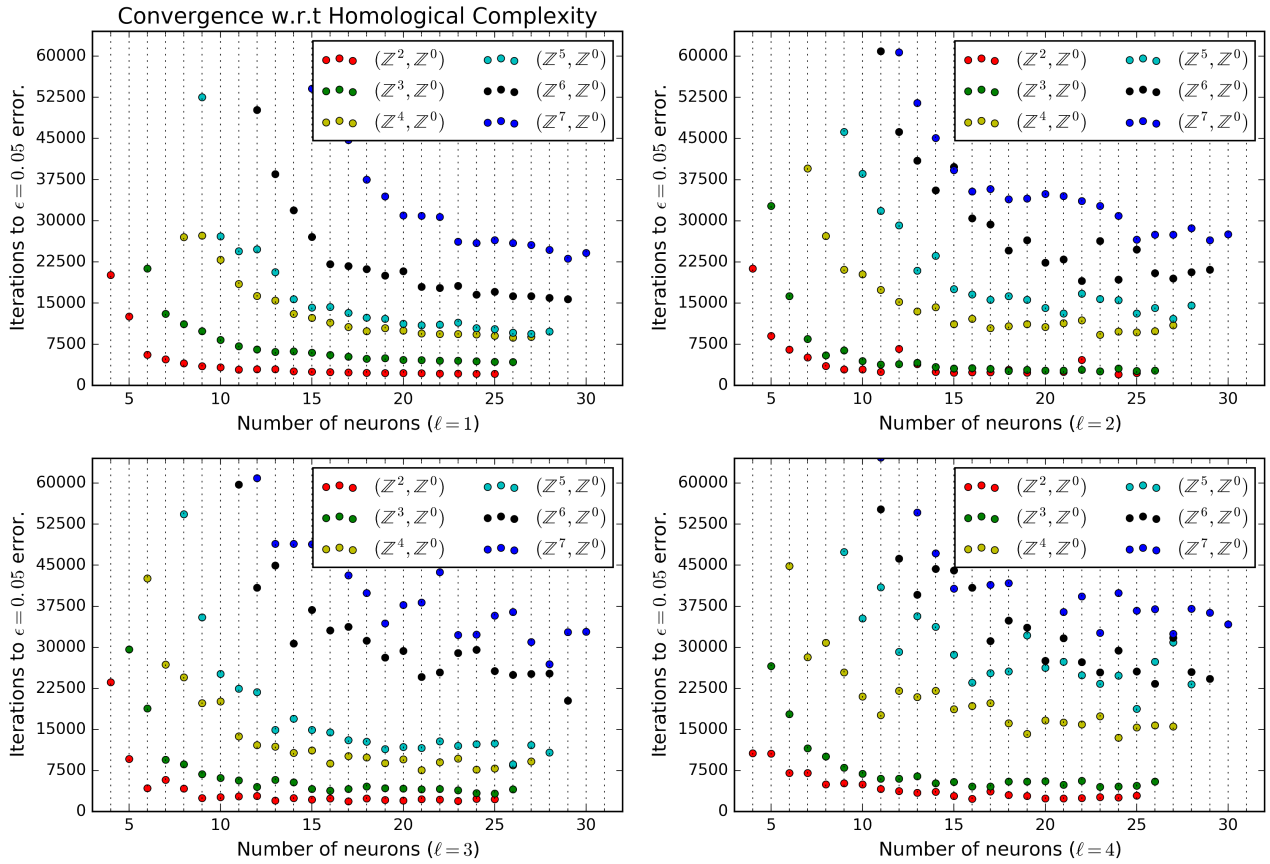
# B. Additional Figures

See next page.

*Figure 9.* A scatter plot of the number of iterations required for architectures of varying hidden dimension and number of layers to converge to 5% misclassification error. The colors of each point denote the topological complexity of the data on which the networks were trained. Note the emergence of convergence bands.
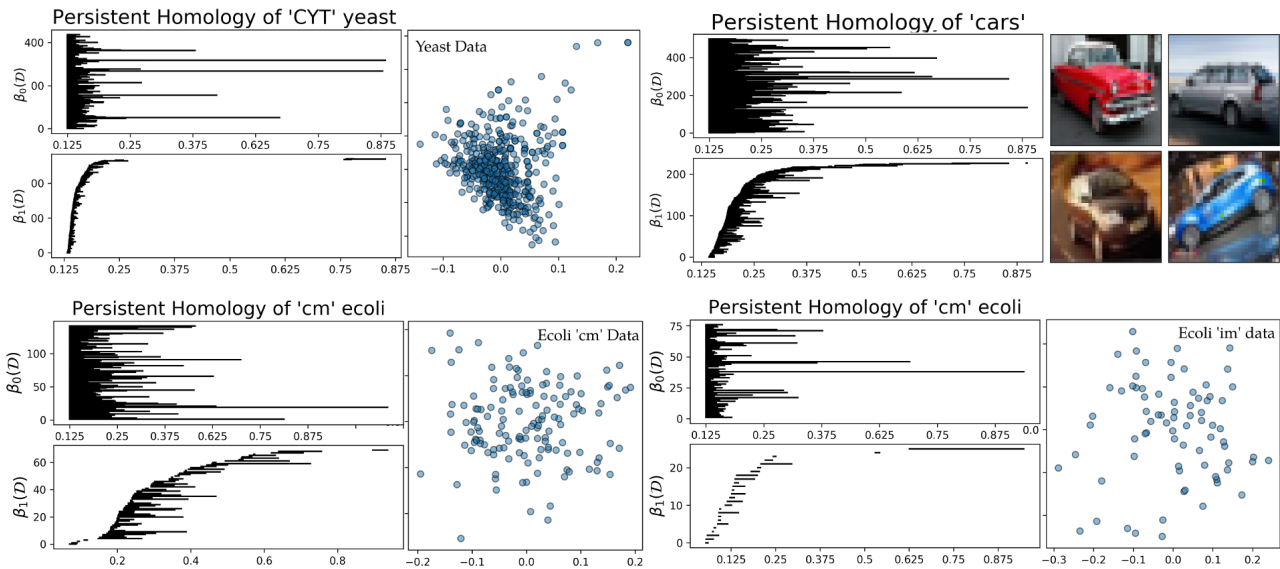


*Figure 10.* The topological persistence diagrams of several datasets. In each plot, the barcode diagram is given for 0 and 1 dimensional features along with a 2-dimensional embedding of each dataset. Note that we do not compute the homologies of the datasets in the embeeding with the exception of CIFAR.
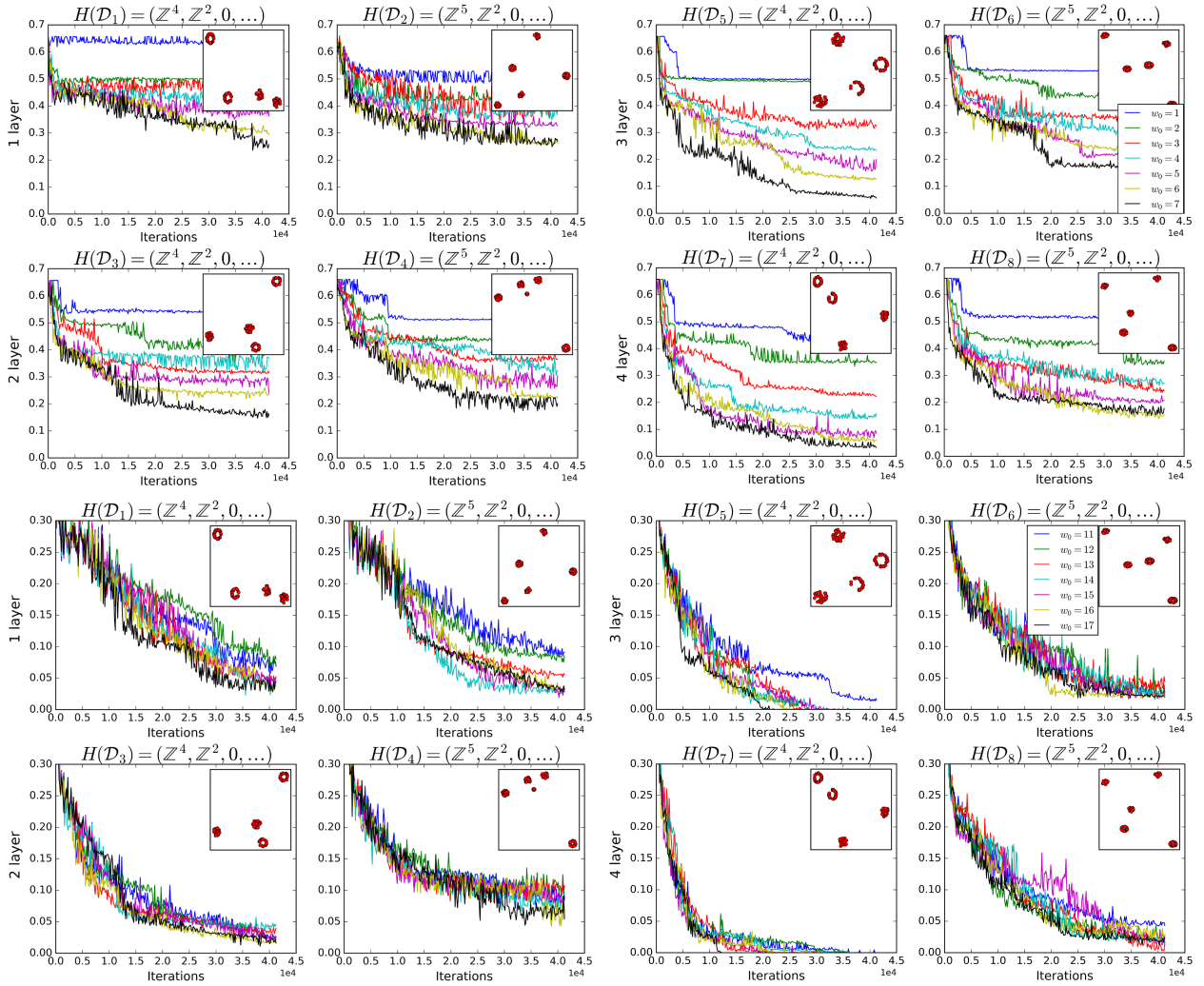
*Figure 11.* Additional topological phase transitions in low dimensional neural networks as the homological complexity of the data increases. The upper right corner of each plot is a dataset on which the neural networks of increasing first layer hidden dimension are trained.