

Assessing Neural Network Representations During Training Using Data Diffusion Spectra

Danqi Liao^{*1} Chen Liu^{*1} Alexander Tong² Guillaume Huguet² Guy Wolf² Maximilian Nickel³
Ian Adelstein¹ Smita Krishnaswamy^{1,3}

Abstract

Here we present information theoretic measures based on the data diffusion operator as characterisations of the representations learned by neural networks. Specifically, we define diffusion spectral entropy (DSE), i.e., entropy of the diffusion operator computed on the neural representation of a dataset as well as diffusion spectral mutual information (DSMI), which assesses the relationship between different sets of variables representing data. First, we show that these definitions form robust measures of intrinsic dimensionality and relationship strength respectively on toy data, outperforming binned Shannon entropy in terms of accuracy. Then we study the evolution of representations within classification networks and networks with self-supervised losses. In both cases, we see that generalizable training results in decrease in DSE over epochs — starting from a random initialization. We also see that there is an increase in DSMI with the class label over time. On the other hand, training with corrupt labels results in a maintenance or increase in entropy and near-zero DSMI with labels. We also assess DSMI with the input and observe differing trends. On MNIST it grows until plateaus, whereas on CIFAR it increases and then decreases. Overall results show that these measures can elucidate characteristics of network performance as well as data complexity. Code is available at <https://github.com/ChenLiu-1996/DiffusionSpectralEntropy>.

^{*}Equal contribution. Order of co-first authors determined by coin toss. ¹Yale University ²The Quebec AI Institute and Université de Montréal ³The FAIR Team, Meta AI. Correspondence to: Smita Krishnaswamy <smita.krishnaswamy@yale.edu>.

1. Introduction

Deep neural networks have emerged as a major breakthrough in data science largely because of their ability to learn increasingly meaningful representations of data. In fact neural networks function by transforming data via a series of non-linear operations such that each layer learns a new representation of the data. While the representations vectors reside in high dimensional spaces, they in fact lie on a lower dimensional manifold (Fefferman et al., 2016). Assessing the properties of this embedding manifold therefore is key to better understanding the neural network.

Here, we use a powerful manifold learning paradigm — diffusion geometry — to study the representations learned by neural networks. Diffusion geometry involves learning a data diffusion operator, which is a type of Markovian transition matrix describing relationships in the data. A key contribution of this work is in introducing diffusion spectral entropy (DSE), or spectral entropy of the diffusion operator as a robust quantifier of the intrinsic information measure of data representation despite the presence of noise. Further, we extend diffusion spectral entropy to a diffusion spectral mutual information (DSMI) in order to ascertain the information the embedding manifold has on the output labels or the raw input data of the dataset.

Our key contributions are:

- Establishing diffusion geometry as a tool for studying neural network representations.
- Introducing *diffusion spectral entropy*, i.e. Shannon entropy of the spectrum of the diffusion operator as a measure of the information content in a representation of the data.
- Defining *diffusion spectral mutual information* and providing a method of its computation for assessing relationships between different layers of information in a neural network.
- Utilizing both methods to assess the evolution of representations in neural networks over training. In specific we quantify DSE of neural representations as well as DSMI between neural representations and output labels or input data over training on different datasets.

2. Methods

2.1. Diffusion Geometric Quantifications of Manifold Characteristics

Here we define and motivate quantities that we use to characterize neural network representations. More information on the background can be found in supplementary materials A and B. In the supplements, we (1) introduce manifold learning and diffusion geometry, (2) explain the construction of diffusion operator from data graphs, and (3) discuss the definitions of entropy and mutual information.

Diffusion spectral entropy for data While spectral entropy has often been used for measuring entropy on graphs, it has not been used often to compute the entropy of data. Here we make a particular choice to compute a data-centric affinity matrix and then a spectral entropy from that matrix. We utilize the anisotropically normalized diffusion operator from the diffusion maps formulation (Coifman & Lafon, 2006) (see Eqn 6). This operator can be written as a symmetric anisotropic normalization of the Gaussian kernel $\mathcal{G}(z_i, z_j)$ by its diagonal degree matrix.

We define the symmetric matrix

$$K_{i,j} = \mathcal{K}(z_i, z_j) \quad (1)$$

with \mathcal{K} being the anisotropic kernel — an intermediate result during the diffusion maps computation (see Eqn 5). The row stochastic matrix

$$\mathbf{P} = D^{-1}K \quad (2)$$

with $D_{i,i} = \sum_j K_{i,j}$ is our diffusion matrix/operator.

When we compute the diffusion operator on the dataset X , we utilize the additional notation \mathbf{P}_X . We define diffusion spectral entropy, with respect to a particular value of diffusion time t as follows.

Definition 2.1. We define *Diffusion Spectral Entropy (DSE)* as an entropy of the eigenvalues of the diffusion operator \mathbf{P}_X computed on a dataset X where $x \in X$ is a multidimensional vector $x_1, x_2 \dots x_d$:

$$S_D(\mathbf{P}_X, t) := - \sum_i \alpha_{i,t} \log(\alpha_{i,t}), \quad (3)$$

where $\alpha_{i,t} := \frac{|\lambda_i^t|}{\sum_j |\lambda_j^t|}$, and $\{\lambda_i\}$ are the eigenvalues of the diffusion matrix \mathbf{P}_X .

In the matrix \mathbf{P} each data point is encoded based on its transition probability to every other data point if one takes a random walk on the data. Thus if a data point is disconnected or far away from others then it is likely that a random

walk starting at the data point would remain at the data point. In this setting eigenvectors of the diffusion operator are paths through the data that are stable states of the transition operator. Thus the entropy of the transition operator can be measured over the eigenbasis diffusion operator. Since rows of this matrix can also be thought to be representations of the data (based on their relationships to other points), this is also a measure of intrinsic dimensionality of the dataset. Note that the parameter t that parameterizes the entropy also gives us the capability of separating noise from the true entropy of the signal. As the value of t increases, the eigenspectrum shifts towards the low frequency eigenvectors (which move slowly over the graph) because the eigenvalues $|\lambda_i| < 1$ diminish at a rate inversely proportional to their value when they are raised to a power t . Indeed raising \mathbf{P}_K^T results in identical eigenvectors and powered eigenvalues, which achieves a low-pass filtering of the data values over the created affinity graph (Van Dijk et al., 2018), which to tune computation of entropy. We note that a similar measure was used in the supplementary material of (Moon et al., 2019) to select parameters, but diffusion spectral entropy was not formally defined or discussed there.

Some properties of S_D are discussed and proved in supplementary materials C.

Diffusion spectral mutual information We further extend the diffusion spectral entropy to define mutual information for understanding the information that some variables of a data representation have on others, for example the information that neurons in a hidden layer have about the primary output.

Definition 2.2. We define the *Diffusion Spectral Mutual Information (DSMI)* as the difference between conditional and unconditional diffusion spectral entropy

$$I_D(X; Y) = S_D(\mathbf{P}_X, t) - \sum_{y_i \in Y} p(Y = y_i) S_D(\mathbf{P}_{X|Y=y_i}, t). \quad (4)$$

The conditioned transition matrix $\mathbf{P}_{X|Y=y_i}$ is the transition matrix computed on the subset of X that has output label Y . To avoid numeric issues that are involved in comparing spectra of different sizes of matrices, we also compute $S_D(\mathbf{P}_X, t)$ using subsamples of X the average size of the classes of Y . Since uniform subsampling maintains distributions, the sampled entropy would be the same as the total entropy as shown in our experiments (see supplementary materials G).

In addition, empirically we show that DSMI can reflect the relationships between data points and their class labels, and that these degrade with corruption.

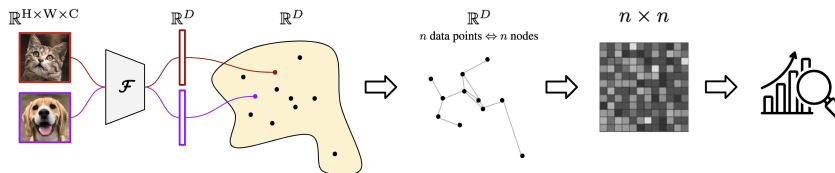


Figure 1. **Data processing to obtain the diffusion matrix of an embedding manifold.** Each data point x_i (in this case, an image) from the validation set is embedded by the network \mathcal{F} as a vector z on a D -dimensional embedding manifold (light yellow canvas). This yields a point cloud with n points, which can be converted into a graph based on local proximity. We can compute its diffusion matrix which allows for further analysis.

2.2. Efficiently Computing Aforementioned Quantities

Details can be found in supplementary materials E.

3. Results

In this study, we used the ResNet-50 (He et al., 2016) backbone, and assessed the penultimate layer of the network. As a convolutional network, its penultimate layer is different from most preceding layers as the representations learned do not contain explicit spatial dimensions. Specifically, the representations lie in \mathbb{R}^D rather than $\mathbb{R}^{h \times w \times c}$. This property allows us to interpret these vectors as points in a fixed-dimensional space, whereas for other layers we need to flatten the representations for analysis, which may raise concerns on mixing spatial information with channel information. Our evaluation framework can be easily adapted to many other vision backbones.

We trained the vision backbones under three conditions: (1) supervised learning, (2) contrastive learning (specifically, SimCLR (Chen et al., 2020)), and (3) purposeful overfitting. More details can be found in supplementary materials D.

To analyze the embedding manifold, we pass the entire validation set through the network and collect the embedding vectors. The key step is to convert this point cloud of n embedding vectors into a data graph, from which we can find the corresponding diffusion matrix. DSE and DSMI can be computed from eigenvalues of this diffusion matrix. This process is illustrated in Figure 1.

3.1. Toy Test Cases for DSE and DSMI

Results can be found in supplementary materials F.

3.2. Results on Neural Network Training Process

There are several key observations on the diffusion geometric quantities on the embedding manifold during neural network training. **Every figure in this section has an enlarged version in supplementary materials H.**

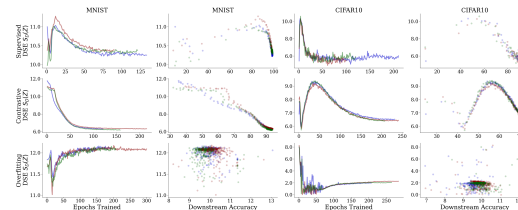


Figure 2. **Diffusion Spectral Entropy $S_D(Z)$ of embedding vectors Z .** Colors correspond to the three random seeds. t is empirically set to 1 for MNIST and 2 for CIFAR10. The same t setting is used in all subsequent figures.

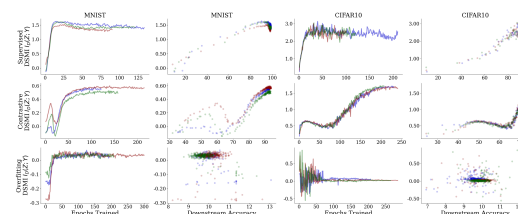


Figure 3. **Diffusion Spectral Mutual Information $I_D(Z; Y)$ between embedding vectors Z and the label classes Y .**

DSE We can observe from Figure 2 that DSE in proper learning (i.e., supervised or contrastive learning on correct labels) decreases as the model performs better on the downstream classification task. Meanwhile, this trend is completely absent when the model is forced to memorize random nonsense labels. This accords with the intuition that random initialization has high variance and entropy, while class labels have much lower entropy, and even organizations created by contrastive losses have lower "surprise" than random sampling.

DSMI with output In Figure 3, it can be observed that DSMI $I_D(Z; Y)$ consistently increases during proper learning. Under the same learning rate and scheduling, the DSMI climbs more slowly in contrastive learning compared to supervised learning, and it ends up at a lower terminal value. This may be attributed to the fact that contrastive learning lacks the direct supervision from explicit class labels. However since class labels relate to the data geometry, self-

supervised learning on the data alone still yields some mutual information with labels. In nonsense memorization, DSMI quickly converges to zero. This aligns well with the expectation, since a classifier that essentially performs random guessing has zero mutual information with the class label, whereas a functioning classifier corresponds to a positive mutual information.

Taken together the DSE and DSMI trends indicate that the representations coalesce to a less noisy and more streamlined form where they mainly contain information about the the output label during good training. Results on the classic Shannon version can be found in supplementary materials I.

DSMI with input We additionally show DSMI with the input signal in Figure 4. In nonsense memorization, $I_D(Z; X)$ stays close to zero just like $I_D(Z; Y)$, which ascertains the random projection tendencies. During proper learning, the information bottleneck theory would suggest $I_D(Z; X)$ shall decrease (Tishby & Zaslavsky, 2015) while counter-arguments have also been provided in (Saxe et al., 2019). Our results suggest they may both be correct, and the trend may depend on the nature of the dataset X . $I_D(Z; X)$ keeps increasing during learning on the MNIST dataset, whereas it begins to decrease after some point on the CIFAR10 dataset. This may be verified by future studies on more datasets.

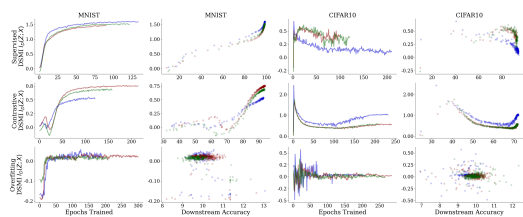


Figure 4. Diffusion Spectral Mutual Information $I_D(Z; X)$ between embedding vectors Z and the input X . Input X is flattened and spectral-clustered into same number of categories as the output Y for fair comparison.

Embedding visualization See supplementary materials J for how the visualized embeddings corroborate with the DSE trends assessed during training.

4. Conclusion

In conclusion, we introduced diffusion operator-based information theoretic measures as tools for assessing neural network representations. Specifically we proposed diffusion spectral entropy (DSE) for measuring information in embedding manifold. We further defined diffusion spectral mutual information (DSMI) and provided an efficient method of its computation, and proved bounds on its values in extreme as well as idealistic clustered cases. Through extensive simulation on toy datasets, we demonstrated diffusion spectral

entropy is a measure of intrinsic dimension on toy data, and spectral mutual information is a meaningful measure for mutual dependence between two variables. We also investigated the neural representation from the penultimate layer of ResNet-50 under supervised learning, contrastive learning and overfitting settings. We empirically showed that DSE during learning follows a general decreasing trend while such decreasing trend is not observed in overfitting setting. Further, we showed DSMI between a hidden layer and output increases during training and plateaus at some point whereas it quickly converges to zero in overfitting settings. We saw more complex data-dependent trends on DSMI with primary inputs. On the MNIST dataset DSMI with input increases until plateaus while on CIFAR it has non-monotonic trends characterized increase and later decrease.

See further discussions in supplementary materials K.

References

- Chen, T., Kornblith, S., Norouzi, M., and Hinton, G. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pp. 1597–1607. PMLR, 2020.
- Coifman, R. R. and Lafon, S. Diffusion maps. *Applied and computational harmonic analysis*, 21(1):5–30, 2006.
- Fefferman, C., Mitter, S., and Narayanan, H. Testing the manifold hypothesis. *Journal of the American Mathematical Society*, 29(4):983–1049, 2016.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Moon, K. R., van Dijk, D., Wang, Z., Gigante, S., Burkhardt, D. B., Chen, W. S., Yim, K., Elzen, A. v. d., Hirn, M. J., Coifman, R. R., et al. Visualizing structure and transitions in high-dimensional biological data. *Nature biotechnology*, 37(12):1482–1492, 2019.
- Saxe, A. M., Bansal, Y., Dapello, J., Advani, M., Kolchinsky, A., Tracey, B. D., and Cox, D. D. On the information bottleneck theory of deep learning. *Journal of Statistical Mechanics: Theory and Experiment*, 2019(12):124020, 2019.
- Tishby, N. and Zaslavsky, N. Deep learning and the information bottleneck principle. In *2015 IEEE information theory workshop (itw)*, pp. 1–5. IEEE, 2015.
- Van Dijk, D., Sharma, R., Nainys, J., Yim, K., Kathail, P., Carr, A. J., Burdziak, C., Moon, K. R., Chaffer, C. L., Pattabiraman, D., et al. Recovering gene interactions from single-cell data using data diffusion. *Cell*, 174(3):716–729, 2018.

Supplementary Materials

A. Background

Manifold learning and diffusion geometry A useful assumption in representation learning is that high dimensional data, which is commonly used in deep learning, originates from an intrinsic low dimensional manifold that is mapped via nonlinear functions to observable high dimensional measurements. This is commonly known as *the manifold assumption* (Fefferman et al., 2016). Let \mathcal{M}^d be a hidden d dimensional manifold that is only observable via a collection of $n \gg d$ nonlinear functions $f_1, \dots, f_n : \mathcal{M}^d \rightarrow \mathbb{R}$ that enable its immersion in a high dimensional ambient space as $F(\mathcal{M}^d) = \{\mathbf{f}(z) = (f_1(z), \dots, f_n(z))^T : z \in \mathcal{M}^d\} \subseteq \mathbb{R}^n$ from which data is collected. Conversely, given a dataset $X = \{x_1, \dots, x_N\} \subset \mathbb{R}^n$ of high dimensional observations, manifold learning methods assume data points originate from a sampling $Z = \{z_i\}_{i=1}^N \in \mathcal{M}^d$ of the underlying manifold via $x_i = \mathbf{f}(z_i)$, $i = 1, \dots, n$, and aim to learn a low dimensional intrinsic representation that approximates the manifold geometry of \mathcal{M}^d .

A paradigm that has emerged as useful in manifold learning in recent years is diffusion geometry (Coifman & Lafon, 2006; Moon et al., 2019; Van Dijk et al., 2018; Burkhart et al., 2019; Huguet et al., 2022). Diffusion geometry seeks to describe data points based on random-walk probabilities to one another. This has been seen to be a noise-tolerant and adaptive way of representing data whose dimensionality reductions have yielded methods such as PHATE (Moon et al., 2019) and diffusion maps (Coifman & Lafon, 2006).

Diffusion maps begin with a kernel \mathcal{K} , often a Gaussian kernel $\exp(-\|z_1 - z_2\|^2/\sigma)$, where $\sigma > 0$ is interpreted as a user-configurable neighborhood size. Such a kernel transforms distances between data points to similarities or affinities. However, such neighborhoods encode sampling density information together with local geometric information. To construct a diffusion geometry that is robust to sampling density variations we may use an anisotropic kernel

$$\mathcal{K}(z_1, z_2) = \frac{\mathcal{G}(z_1, z_2)}{\|\mathcal{G}(z_1, \cdot)\|_1^\alpha \|\mathcal{G}(z_2, \cdot)\|_1^\alpha} \quad \mathcal{G}(z_1, z_2) = e^{-\frac{\|z_1 - z_2\|^2}{\sigma}} \quad (5)$$

as proposed in (Coifman & Lafon, 2006), where $0 \leq \alpha \leq 1$ controls the separation of geometry from density, with $\alpha = 0$ yielding the classic Gaussian kernel, and $\alpha = 1$ completely removing density and providing a geometric equivalent to uniform sampling of the underlying manifold. Next, the similarities encoded by \mathcal{K} are normalized to define transition probabilities $p(z_1, z_2) = \frac{\mathcal{K}(z_1, z_2)}{\|\mathcal{K}(z_1, \cdot)\|_1}$ that are organized in an $n \times n$ row stochastic matrix

$$\mathbf{P}_{i,j} = p(z_i, z_j) \quad (6)$$

that describes a Markovian diffusion process over the intrinsic geometry of the data. Finally, a diffusion map (Coifman & Lafon, 2006) is defined by taking the eigenvalues $1 = |\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_N|$ and corresponding eigenvectors $\{\phi_j\}_{j=1}^N$ of \mathbf{P} , and mapping each data point $x_i \in X$ to an N dimensional vector $\Phi_t(x_i) = [\lambda_1^t \phi_1(x_i), \dots, \lambda_N^t \phi_N(x_i)]^T$, where t represents a diffusion-time, i.e., number of transitions considered in the diffusion process. In general, as t increases, most of the eigenvalues λ_j^t , $j = 1, \dots, N$, become negligible, and thus truncated diffusion map coordinates can be used for dimensionality reduction (Coifman & Lafon, 2006). For example, PHATE involves computing a symmetric divergence between the rows of the diffusion operator and embedding this with multidimensional scaling.

Entropy and mutual information Entropy, a basic quantity in information theory, quantifies the amount of uncertainty or “surprise” when given the value of a random variable. If the variable is distributed with a distribution that has probability mass that is spread out, such as a uniform distribution, then the entropy is high. On the other extreme, if there is no uncertainty in the quantity of the variable, i.e., it is deterministic then the entropy is 0. The Shannon entropy is computed as below.

$$H(X) = \mathbb{E}[-\log p(X)] = - \sum_{x \in X} p(x) \log p(x) \quad (7)$$

The von Neumann entropy (von Neumann, 2018) from quantum information extends the entropy measure to the quantum mechanics domain, and in particular it operates on density matrices. If a density matrix ρ has a set of eigenvalues $\{\eta_i\}$, the von Neumann entropy is defined as

$$H(\rho) = -\text{tr}(\rho \ln \rho) = -\sum_i \eta_i \log \eta_i \quad (8)$$

Here von Neumann entropy is considered to be an extension of Gibbs entropy, which is a measure of the spread of a distribution on the microstates of a classical system. Classical systems can only exist in pure states (or standard basis states). However, quantum systems can exist in superposition states, and depending on the distribution of superposition states the stable or ground states can be redefined as the eigenfunctions of a density operator which describes the probabilities of superpositions.

This notion has subsequently been extended to graph spectra in several works. These methods generally compute the entropy of normalized eigenvalue of a graph adjacency matrix and have been used variously in biology and other fields to compare graphs (Su et al., 2022; de Siqueira Santos et al., 2016; Takahashi et al., 2012; Merbis & de Domenico, 2023; Villafañe-Delgado & Aviyente, 2016).

Mutual information is defined as a function of entropy. There are many alternative formulations of mutual information that are equivalent. The most useful formulation here is as a difference between the (unconditional) entropy of a variable, and the entropy of a variable when conditioned on the value of another variable.

$$I(X; Y) = H(X) - H(X|Y) = H(X) - \sum_i p(Y = y_i) H(X|Y = y_i) \quad (9)$$

Here the conditional entropy $H(X|Y)$ is given as a weighted sum over values of Y in the discrete case, and thus it is computed as in the form on the right hand of Eqn 9.

B. Related Works

Prior works attempted to study neural networks during training by visualizing the neural representation (Gigante et al., 2019) or the loss landscape (Li et al., 2018). These works provided some qualitative ways to analyze neural networks during training but did not offer quantification. The information bottleneck theory for deep learning (Tishby & Zaslavsky, 2015) introduced a framework for quantifying information content in neural networks during training. They binned the vectors along each feature dimension to form a probability distribution and computed the Shannon entropy and mutual information. The main limitation of their proposed method is the curse of dimensionality in the binning process that renders it impractical to analyze layers with more than a dozen neurons — which is ubiquitous in modern deep neural networks (see supplementary material I). Follow-up work (Saxe et al., 2019) used kernel density estimation (Kolchinsky & Tracey, 2017) and Kraskov estimator (Kraskov et al., 2004) for approximating mutual information, yet both estimation methods require specific assumptions on the distributions of hidden layer activation. Our proposed method does not assume specific distribution on hidden layer activations and also avoids the binning problem since it operates on the eigenspectrum rather than the set of embedding vectors.

C. Propositions on Diffusion Spectral Entropy and Diffusion Spectral Mutual Information

The propositions below establish some bounds on minimal and maximal values of DSE. In addition, they provide intuition on the definition of DSMI. Note that taking $t \rightarrow \infty$ allows us to talk about the major structures in the dataset.

Proposition C.1. S_D achieves minimal entropy of 0 when the diffusion operator defines an ergodic Markov chain, and is in steady state (as $t \rightarrow \infty$).

Proof. The eigenvalues of a finite ergodic Markov chain have the form $1 = |\lambda_1| > |\lambda_2| \geq \dots \geq |\lambda_N| \geq 0$, we see that $\lambda_i^t = 0$ as $t \rightarrow \infty \forall i > 1$. Thus the resultant entropy is $1 \log(1) + \sum_{i>1} 0 = 0$, proving the proposition. \square

Note that this also implies that if all data points are very similar, i.e., have equal probability of transitioning to any other point, then it has minimal entropy. This is because such a distribution is a steady state distribution for the diffusion matrix shown in Eqn 2.

Note that steady state distribution for ergodic \mathbf{P} (which is the case of our Kernel-based definition) is characterized by having the same transition probabilities from every starting state. Therefore the rows of the matrix \mathbf{P}^t are identical, and there is no distinguishability between points.

Next we discuss when this entropy will reach its maximal value. This happens effectively when points are all spread very far apart.

Additionally, we can show that there exists a $T \in \mathbb{N}$, such that the diffusion spectral entropy $S_D(\mathbf{P}_X, t)$ decreases for all $t > T$, implying that as the data is denoised the entropy decreases.

Proposition C.2. *Assuming that \mathbf{P}_X has at least one eigenvalue $\lambda_j \in (0, 1)$, then, there exist a $T \in \mathbb{N}$ such that for all $t > T$*

$$\frac{\partial}{\partial t} S_D(\mathbf{P}_X, t) < 0.$$

Proof. We recall that the eigenvalues of \mathbf{P}_X can be ordered as $1 = |\lambda_0| \geq |\lambda_1| \geq \dots \geq |\lambda_n| \geq 0$. Therefore $\sum_j |\lambda_j^t| \geq 1$. By computation we have:

$$\begin{aligned} \frac{\partial}{\partial t} S_D(\mathbf{P}_X, t) &= \frac{\partial}{\partial t} \left[- \sum_{\lambda} \frac{|\lambda^t|}{\sum_{\lambda} |\lambda^t|} \log \frac{|\lambda^t|}{\sum_{\lambda} |\lambda^t|} \right] \\ &= - \sum_i \frac{|\lambda_i^t|}{\left(\sum_j |\lambda_j^t|\right)^2} \left(\sum_j |\lambda_j^t| [\log |\lambda_j| - \log |\lambda_i|] \right) \left(\log \sum_j |\lambda_j^t| - t \log |\lambda_i| - 1 \right) \end{aligned}$$

The first term is always positive. The asymptotic behavior of the second term is positive. This can be seen by noting that taking the limit of $t \rightarrow \infty$, the only summands that do not converge to zero are when $\lambda_j = 1$. Since these are positive, the sum in the second term is positive in the limit of $t \rightarrow \infty$, and furthermore, since it is continuous, it remains positive for sufficiently big t .

Examining the third term we have that for sufficiently large t the $-t \log |\lambda_i|$ term dominates, giving a positive third term, and has positive derivative.

since all three terms are positive for sufficiently large t , $\frac{\partial}{\partial t} S_D(\mathbf{P}_X, t) < 0$ completing the proof. \square

If instead of using a Gaussian kernel for computation of \mathbf{P} , we use a k-nearest-neighbor or other thresholded kernel we can make the following statement.

Proposition C.3. *As $t \rightarrow \infty$, $S_D(\mathbf{P}_X, t)$ on data with k well-separated clusters is $\log k$.*

Proof. If the data has k well separated clusters then \mathbf{P}_X has eigenvalues of the form $1 = |\lambda_1| = |\lambda_2| = \dots = |\lambda_k| > |\lambda_{k+1}| \geq \dots \geq |\lambda_N| \geq 0$. In other words the multiplicity of 1 eigenvalues of \mathbf{P}_X corresponds to the number of connected components in the underlying graph, which here is k , all other eigenvalues are strictly less than 1 and greater than or equal to 0. Therefore as $t \rightarrow \infty$ only these eigenvalues remain and the resultant DSE is $\sum_k 1/k \log k = \log k$ completing the proof. \square

Corollary 1. *S_D achieves maximal entropy in a matrix where each point only transitions to itself, and the entropy here will be $\log(n)$*

Proof. In this case the transition matrix corresponds to the identity matrix, hence, each of the n eigenvalues is 1. Thus the diffusion spectral entropy is the uniform distribution on n states $-\sum_i^n (1/n) \log(1/n) = \log(n)$, which maximizes the entropy. \square

Proposition C.4. *As $t \rightarrow \infty$ on a hidden layer X with k well-separated clusters, and output labels perfectly coinciding with clusters, we will have $I_D(X; Y) = \log k$, i.e., DDMI between hidden layer and output layer being positive.*

Proof. Based on proposition 3.3, the entropy of k well-separated clustered data with $t \rightarrow \infty$ is $\log k$, thus $S_D(\mathbf{P}_x, t) = \log(k)$. However for each label y $S_D(\mathbf{P}_x, t) = \log(1) = 0$. Thus $I_D(X; Y) = \log k$. \square

D. Experimental Details

D.1. Three Training Conditions on Real Data

Supervised Learning We trained ResNet-50 (He et al., 2016) models end-to-end using an AdamW optimizer (Loshchilov & Hutter, 2017) at an initial learning rate 1e-5 for MNIST or 1e-3 for CIFAR10. Learning rate is modulated by a Cosine Annealing Scheduler (Loshchilov & Hutter, 2016) with a linear warmup for the first 10 epochs. Early stopping kicks in if the validation accuracy no longer increases for 15 epochs. Experiments are repeated over 3 random seeds. At the end of each epoch, we pass the entire validation set (10,000 images for MNIST or CIFAR10) through the model and collect the \mathbb{R}^{2048} representation vectors as the outputs of the penultimate layer for further analysis.

Contrastive Learning For the contrastive learning experiments, we followed the SimCLR (Chen et al., 2020) paradigm: we create two augmentations of the same image and asks the model to embed them closer on the embedding manifold, while encouraging bigger separation between them and the other images in the same batch. The standard training procedure of contrastive learning first trains the backbone, with the classifier (the final, linear layer) detached, for some epochs and then either performs a linear probing or fine-tuning. In either case, a linear classifier layer is attached and trained for some more epochs. The weights of the backbone are frozen in the former case versus learnable in the latter.

Since we need to assess how well the model is learning at each epoch, we instead performs linear probing by the end of each epoch. Specifically, we freeze the backbone weights, attach a re-initialized linear classifier layer, and train the classifier for 10 epochs with the training set. Then, we record the end-to-end validation accuracy as well as collect the embedding vectors on the validation set, similar to the supervised learning case. Finally, we unfreeze the backbone weights for the next epoch of SimCLR training.

For fair comparison, the training details are otherwise the same as the supervised learning case, including learning rate and scheduling.

Purposeful Overfitting In purposeful overfitting, we train the model in the same manner as the supervised learning case, except that the data labels are randomly permuted. In this way, the models are forced to learn nonsense labels. To better overfit, we reduced the extent of data augmentation during training, since data augmentation is proven effective to mitigate overfitting. We also increased the early-stop patience from 15 to 30, with the triggering metric being the train-validation accuracy divergence instead of validation accuracy.

For fair comparison, the training details are otherwise the same as the supervised learning case, including learning rate and scheduling.

D.2. Computing DSMI with input

We compute $I_D(Z; X)$ in the same fashion as we compute $I_D(Z; Y)$. By a simple change of variables, Eqn 4 can be rewritten as:

$$\begin{aligned} I_D(Z; X) &= S_D(Z) - S_D(Z|X) \\ &= S_D(\mathbf{P}_Z, t) - \sum_{x_i \in X} p(X = x_i) S_D(\mathbf{P}_Z | X = x_i, t) \end{aligned}$$

Compared to DSMI between the neural representation and the output, the DSMI with the input is slightly more complicated because the input signals do not fall in discrete categories, i.e., x_i are not naturally defined. To that end, for the set of n input images, we flatten them respectively and perform spectral clustering. For fair comparison with $I_D(Z; Y)$, we cluster these vectors into the same number of categories as the number of classes in Y . The remaining process is the same as how we compute DSMI with the output.

D.3. Logarithm Base

It can be verified that the choice of logarithm base does not affect the validity of our formulations. By convention, we choose base 2 (\log_2) in our implementation. Hence the unit of entropy will be bits.

E. Efficiently Computing Diffusion Geometric Quantities on Neural Networks

To compute diffusion geometric quantities on neural network representations, we pass the training data X to the desired layer L of a neural network. At this layer, we collect activations of each neuron into a vector $L(x_i) = [L_1(x_i), L_2(x_i) \dots L_n(x_i)]$, where L_j is the j -th neuron of the L -th layer. $L(x_i)$ is high dimensional representation of data point $x_i \in X$. In case the activation is a multi-dimensional tensor with spatial and channel information, we can flatten it into a vector. We then compute $\mathbf{P}_{L(X)}$ and proceed to compute $S_D(\mathbf{P}_{L(X)}, t)$. We choose t to be consistent over an experiment, i.e., different evaluations of the same network, but tuned to the level of noise in the data. A rule of thumb is, after powering eigenvalues to t , there should be approximately 1 percent of eigenvalues that remain larger than 0.01. To assess the evolution of the representation, we compute this quantity after every epoch of training.

Computing eigenvalues via full eigendecomposition is known to have time complexity $O(n^3)$. However, we take advantage of two characteristics of our situation to provide a faster method: 1) We only need eigenvalues and not eigenvectors for spectral entropy computation, 2) \mathbf{P} has the same eigenvalues as K defined in Eqn 1 as discussed in (Coifman & Lafon, 2006). Thus, we can compute the DSE fast using QR decomposition since we actually compute the eigenvalues of a real symmetric matrix. For further speedup, one can use the Chebyshev moments to approximate the eigenvalues, and the implementation is also provided in our codebase.

F. DSE and DSMI on Toy Datasets

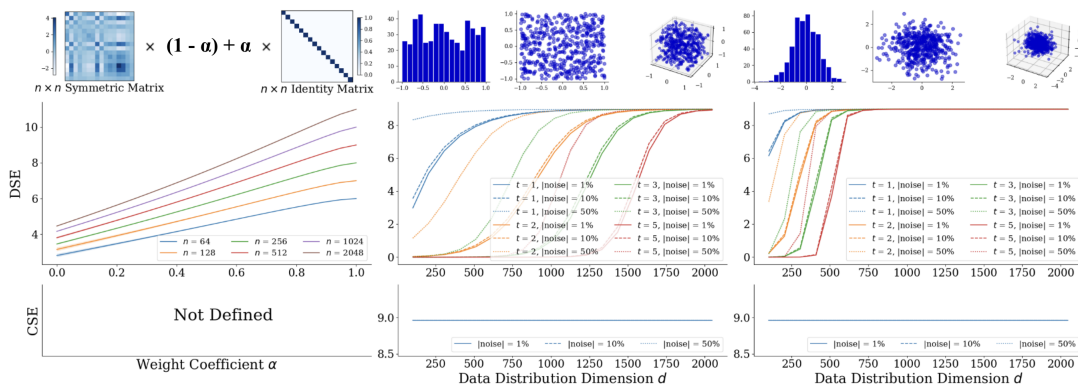


Figure S1. Diffusion Spectral Entropy (DSE) and Classic Shannon Entropy (CSE) on toy data. **Left:** Weighted sum of a random $n \times n$ symmetric positive definite matrix (to simulate a diffusion matrix) and the $n \times n$ identity matrix. In theory, the identity matrix shall have the highest entropy at each respective n . **Mid:** d -dimensional $U[-1, 1]$ inside a 2048-dimensional space. **Right:** d -dimensional $\mathcal{N}(0, I)$ inside a 2048-dimensional space. Shaded areas indicate standard deviation from 5 independent runs. For the latter two distributions, additive noise is injected to the coordinates, and schematics for $d = \{1, 2, 3\}$ are illustrated on top. Note that in the latter two case for DSE we compute the diffusion matrix of the data manifold prior to entropy evaluation, whereas in the first case we skip that step as the matrix is already provided.

To demonstrate the behavior of DSE, we first run several simulations as shown in Figure S1. The left panel indicates that, for an arbitrary real symmetric matrix, the closer it gets to an identity matrix, the higher its DSE — since identity matrix is a steady state diffusion matrix. The two other panels show that, DSE in general increases as the intrinsic dimensionality of the data manifold increases.

To demonstrate the behavior of DSMI, we run similar simulations as shown in Figure S2. In these toy test cases, we gradually corrupt the class label, and just as expected, DSMI starts high when the association between the data point and class label is high, and drops to zero as the label corruption increases. The parameter t should be tuned to the level of noise in the data, otherwise there is a risk of not quantifying entropy in signal.

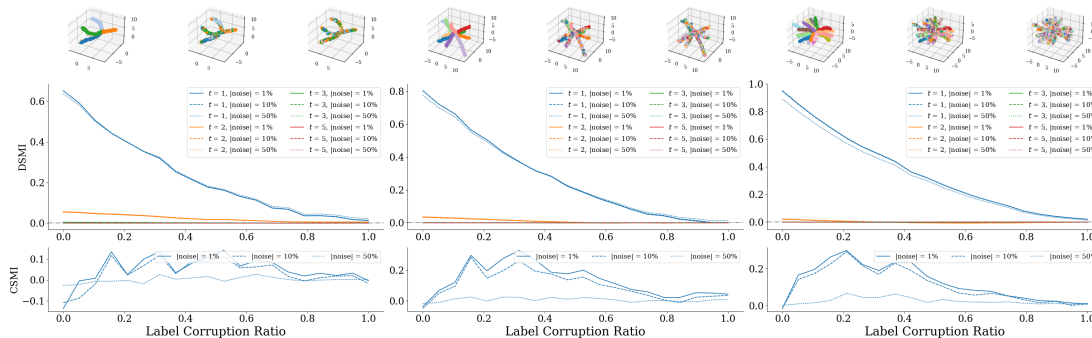


Figure S2. **Diffusion Spectral Mutual Information (DSMI) and Classic Shannon Mutual Information (CSMI) on toy data.** DSMI $I_D(Z; Y)$ and CSMI are computed on synthetic, 30-dimensional trees with $\{5, 10, 20\}$ branches (Left, Mid, Right). Along the x-axes, an increasing percentage of labels are corrupted, with 3-dimensional schematics demonstrating corruption ratios $\{0, 0.5, 1\}$ displayed on top. At full corruption, $I_D(Z; Y)$ converges to zero, as the embedding vectors contain no information on the labels. Noise injection and repeated experiment is the same as in Figure S1.

We also computed the classic Shannon entropy and mutual information for comparison using the method from (Tishby & Zaslavsky, 2015). Our proposed method outperforms the classic Shannon version in high-dimensional spaces. In Figure S1, DSE consistently captures the entropy trends while CSE saturates to the maximal value. In Figure S2, DSMI is considerably more robust to noise compared to the classic CSMI — when the noise level is at 50% the amplitude of the signal, CSMI is significantly dampened while DSMI still remains close to the noiseless counterpart.

G. DSMI Results

G.1. DSMI Intuition

Figure S3 provides some intuition for DSMI. In cases where the class label coincides exactly with well-separated clusters, based on Proposition C.3 the overall data will have entropy $\log K$ as $t \rightarrow \infty$. However, in the computation of the conditional entropy is done on single-clusters, the entropy should converge to 0, so the DSMI will be positive. In the second case, where class labels are spread throughout the manifold the conditioned and unconditioned entropy would be similar, based on their similar spread. Thus in such cases DSMI can be zero, or occasionally also slightly negative — but nevertheless indicates low MI. See supplementary material for more results and discussion of DSMI.

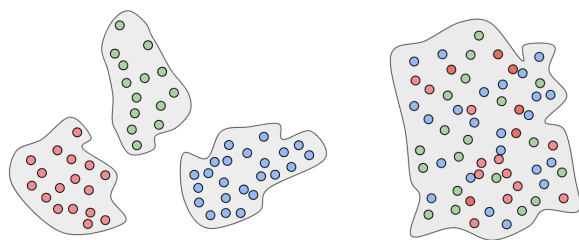


Figure S3. **Intuition for mutual information.** Class labels are colored. **Left:** Y is on separated clusters. $S_D(X|Y) < S_D(X)$ and $I_D(X; Y)$ is positive. **Right:** Y is close to a uniform sub-sampling of X . $I_D(X; Y)$ is around 0, but can have negative values in case $S_D(X|Y)$ is larger than $S_D(X)$ due to numeric reasons.

G.2. DSMI Intuition, Quantified

Figure S4 illustrates DSMI for the cases in the intuition figure (Figure S3). It can be seen that DSMI is positive on nicely separated clusters, while close to zero on well-mixed clusters.

G.3. Subsampling Technique for DSMI: DSE Subsampling Robustness

As mentioned in Definition 2.2, we use a subsampling technique to compute *Diffusion Spectral Mutual Information (DSMI)*. We hereby justify this technique by showing that Diffusion Spectral Entropy (DSE) is robust to subsampling (Figure S5). Meanwhile, we also need to point out that the subsampling robustness is constrained by the DSE upper bound: DSE is

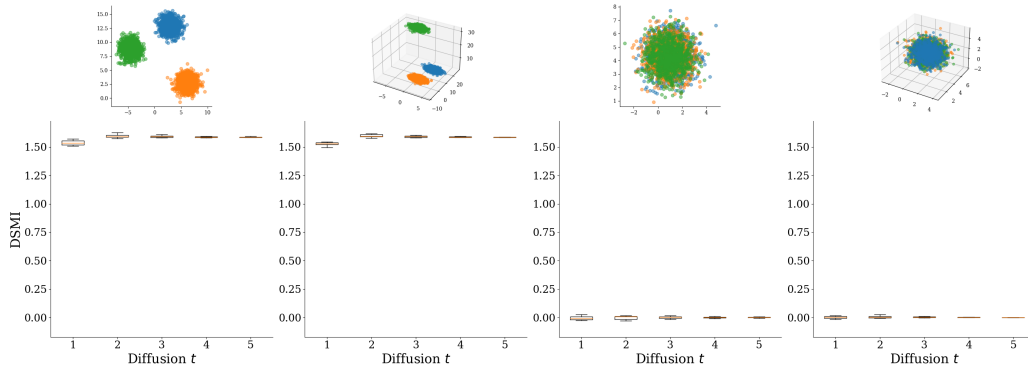


Figure S4. **DSMI simulations for Figure S3.** **Left:** Y is on separated clusters. Simulated with 3 blobs in 2 or 3 dimensions. $S_D(X|Y) < S_D(X)$ and $I_D(X; Y)$ is positive. **Right:** Y is close to a uniform sub-sampling of X . Simulated with 1 blob in 2 or 3 dimensions. $I_D(X; Y)$ is around 0, but can have negative values in case $S_D(X|Y)$ is larger than $S_D(X)$ due to numeric reasons.

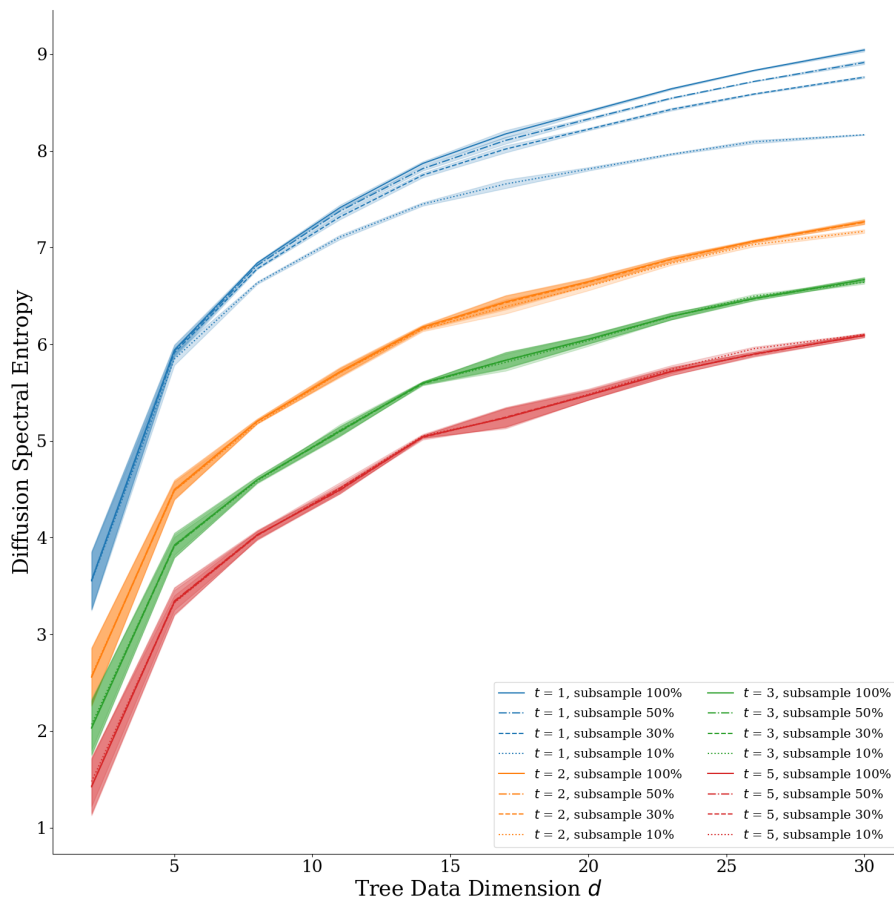


Figure S5. **Diffusion spectral entropy estimation is robust to subsampling.** Subsampled $S_D(Z)$ are computed on synthetic multi-dimensional trees with 3 branches.

capped by $\log_2 n$ where n is the number of data points. For example, if the entire population has a DSE of 8 while we subsample fewer than $2^8 = 256$ data points, the subsampled DSE will be an underestimation. In our studies, the subsampling ratio is about 10%, yielding significant number of data points per category. This technique shall be used with caution if it is applied to other studies with more extreme subsampling.

H. Enlarged Figures from Main Results

Here we display the enlarged versions of Figure 2, 3, and 4.

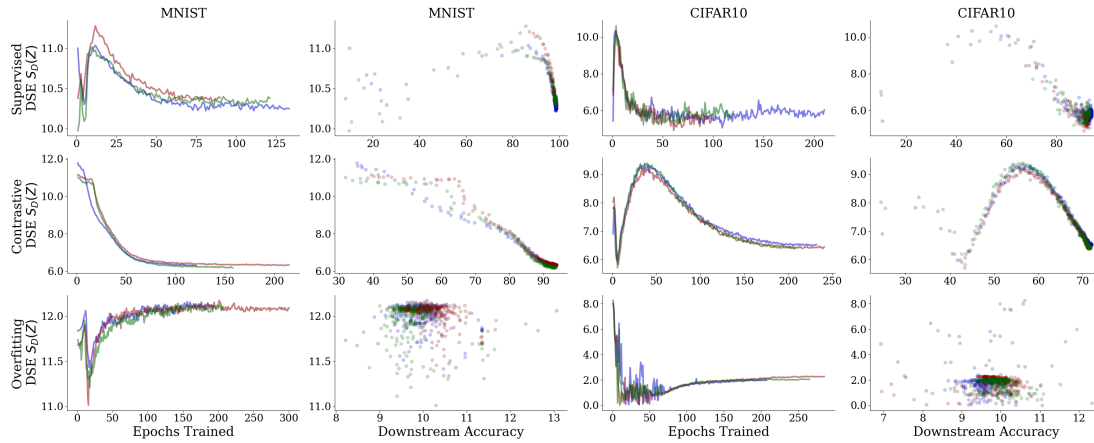


Figure S6. Enlarged version of Figure 2.

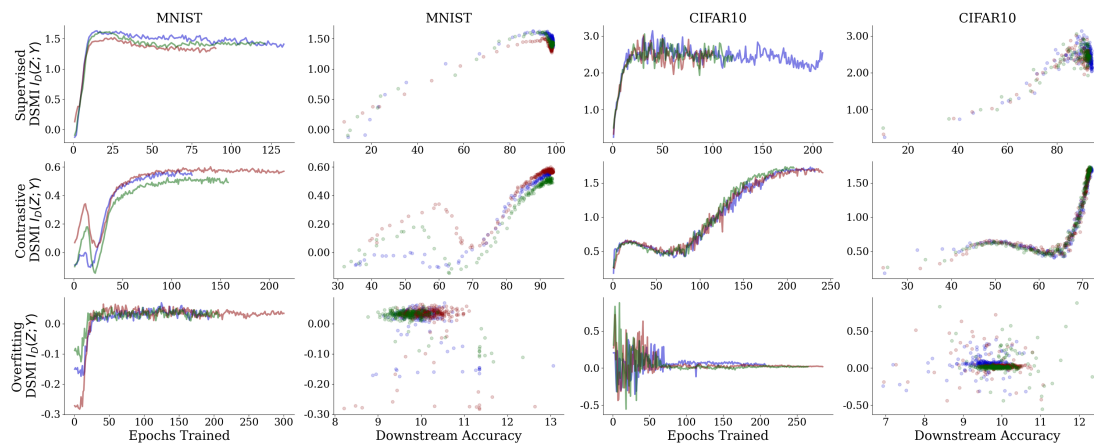


Figure S7. Enlarged version of Figure 3.

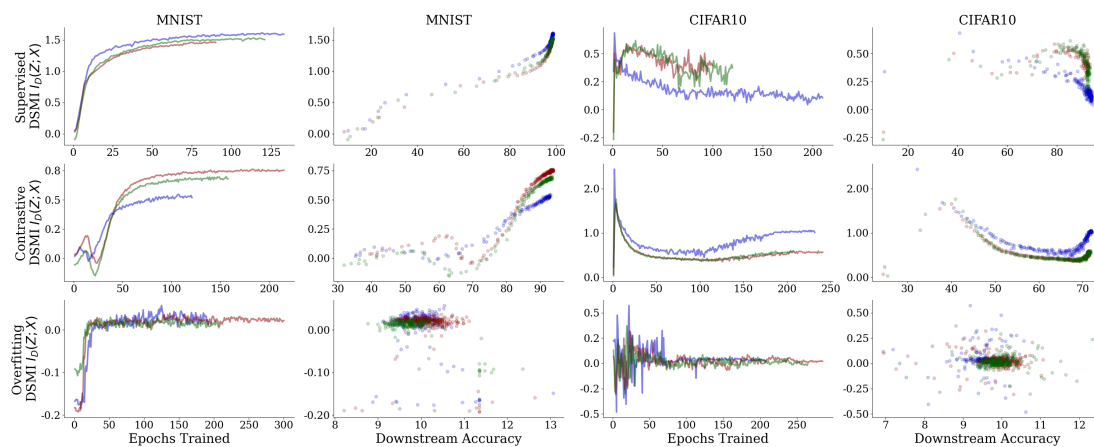


Figure S8. Enlarged version of Figure 4.

I. Limitations of the Classic Shannon Entropy and Mutual Information

I.1. Computing the Classic Shannon Entropy and Mutual Information

To compute the classic Shannon entropy of a variable Z , instead of computing the entropy over the eigenvalue spectrum of the diffusion matrix which gives $S_D(Z)$, we directly compute the entropy over the embedding vectors which gives $H(Z)$.

For n data points, the set of embedding vectors contains a total of n vectors in \mathbb{R}^D . We need to convert the n vectors to probability densities and use Eqn 7. The most common way (e.g., in (Tishby & Zaslavsky, 2015)) is to bin these vectors along each of the D feature dimensions.

Specifically, we will compute the global range of all n vectors along each feature dimension $i \in \{1, 2, \dots, D\}$, and normalize them to $[0, 1]$. Then, we can quantize all vectors along each dimension into k different bins. For example, if $k = 10$, values in $[0, 0.1)$ will be assigned to bin 1; values in $[0.1, 0.2)$ to bin 2, etc. As a result, each vector is converted to a quantized version, with each entry being an integer in $[1, k]$. Every possible quantized vector can be referred to as a “bucket”. It can be easily noticed that the number of buckets is equal to k^D . After counting the number of vectors in each bucket, we can estimate a probability density distribution over the buckets. Finally, we compute the classic Shannon entropy using Eqn 7.

The classic Shannon mutual information can be computed in a similar manner, using Eqn 9 and the aforementioned entropy computation.

I.2. Limitations of Classic Shannon Entropy and Mutual Information

The major limitation lies in the binning process. It is a known problem that **the number of buckets scale exponentially with respect to the feature dimension D** . In our ResNet-50 example, the penultimate layer has $D = 2048$. It is overwhelmingly likely that all embedding vectors are assigned to different buckets, even if we use the minimal choice of 2 bins per feature dimension — which is already a very coarse-grained binning. When the majority of cases results in unique bucketing which leads to the maximum entropy, this metric has very limited expressiveness.

This can be seen in both toy data (Figure S1 and S2) and real data (Figure S9, S10 and S11).

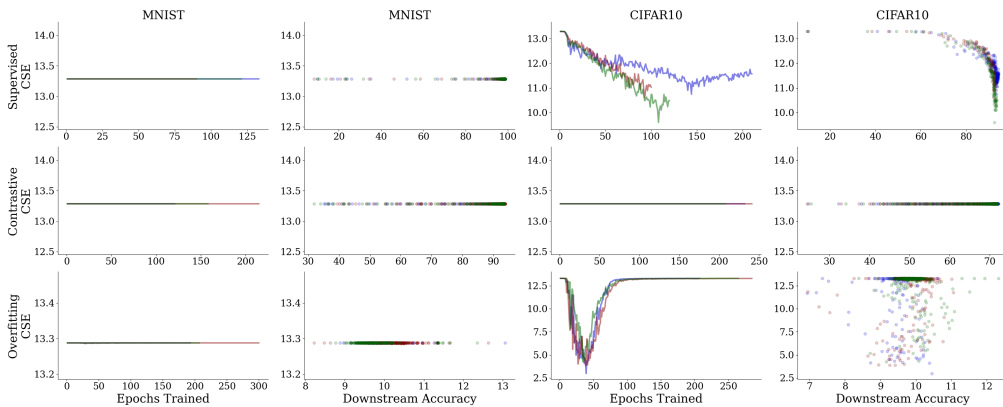


Figure S9. Classic Shannon entropy version for Figure 2. The embedding vectors are frequently allocated to unique buckets, which leads to the maximum possible entropy of $-\log_2 \frac{1}{10000} = 13.288$ for 10,000 data points.

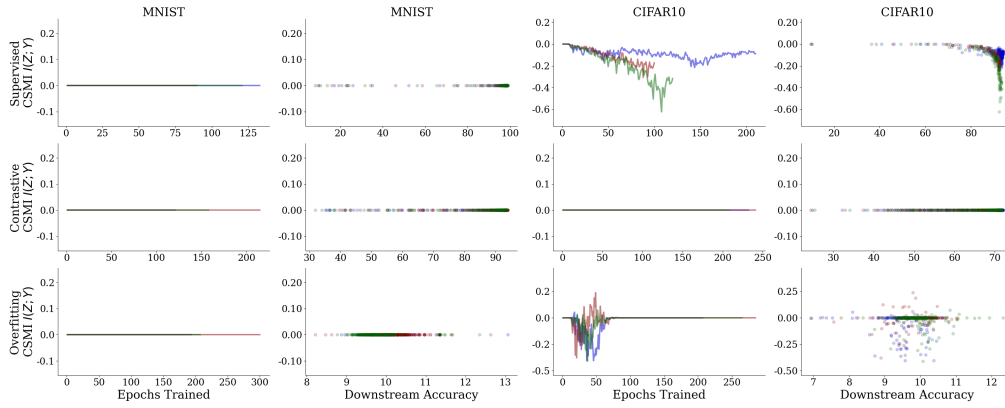


Figure S10. Classic Shannon mutual information version for Figure 3.

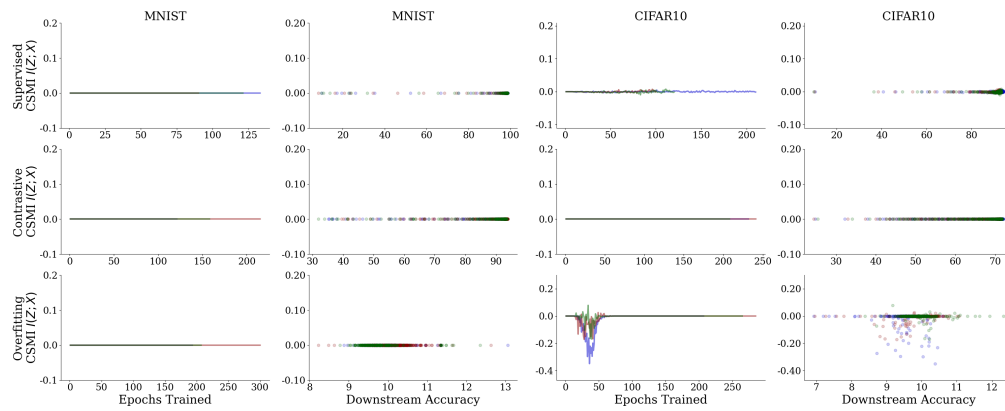


Figure S11. Classic Shannon mutual information version for Figure 4.

J. Embedding visualization

Lastly, we visualize the embedding vectors with PHATE (Moon et al., 2019), a visualization method based on the data diffusion operator. These visualizations nicely corroborate with the DSE trends assessed during training (Figure 2). In proper learning, the transition from 30% (crescent shape) to 80% (full blob) implies increase in DSE, whereas from 80% (full blob) to best accuracy (well-separated branches) implies DSE decrease — both trends match well with the empirical values of DSE in Figure 2. On the other hand, during overfitting on random labels, the dimensionality-reduced embedding vectors form a chaotic blob and remain unchanged throughout the process, which also confirms the empirical measurements of DSE.

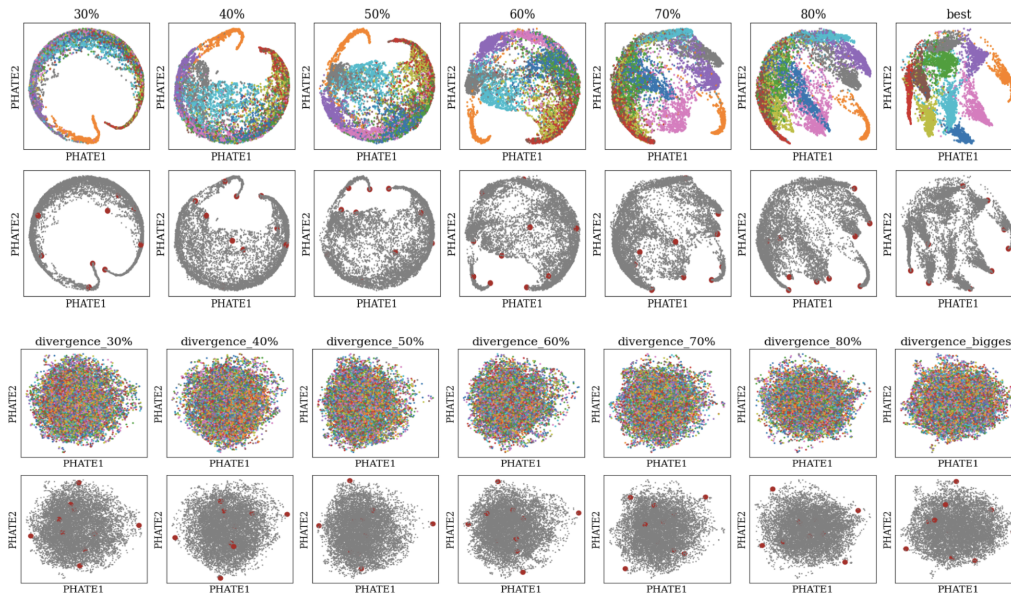


Figure S12. **Visualization of embedding vectors over the course of neural network training.** The \mathbb{R}^{2048} embedding vectors, representing images from MNIST validation set processed by various ResNet-50 models (up until the penultimate layer), are visualized in 2 dimensions. The top panel illustrate the training of a proper supervised learning, whereas the bottom panel illustrate the improper learning as the model overfits on random labels. Snapshots are taken from fixed intervals of validation accuracy (top) or train-validation accuracy divergence (bottom). In each panel, the top row is the PHATE-dimensionality-reduced data points colored by ground truth labels, while the bottom row is the same set of data points in gray-scale accompanied by the top 10 Laplacian extrema of the embedding manifold.

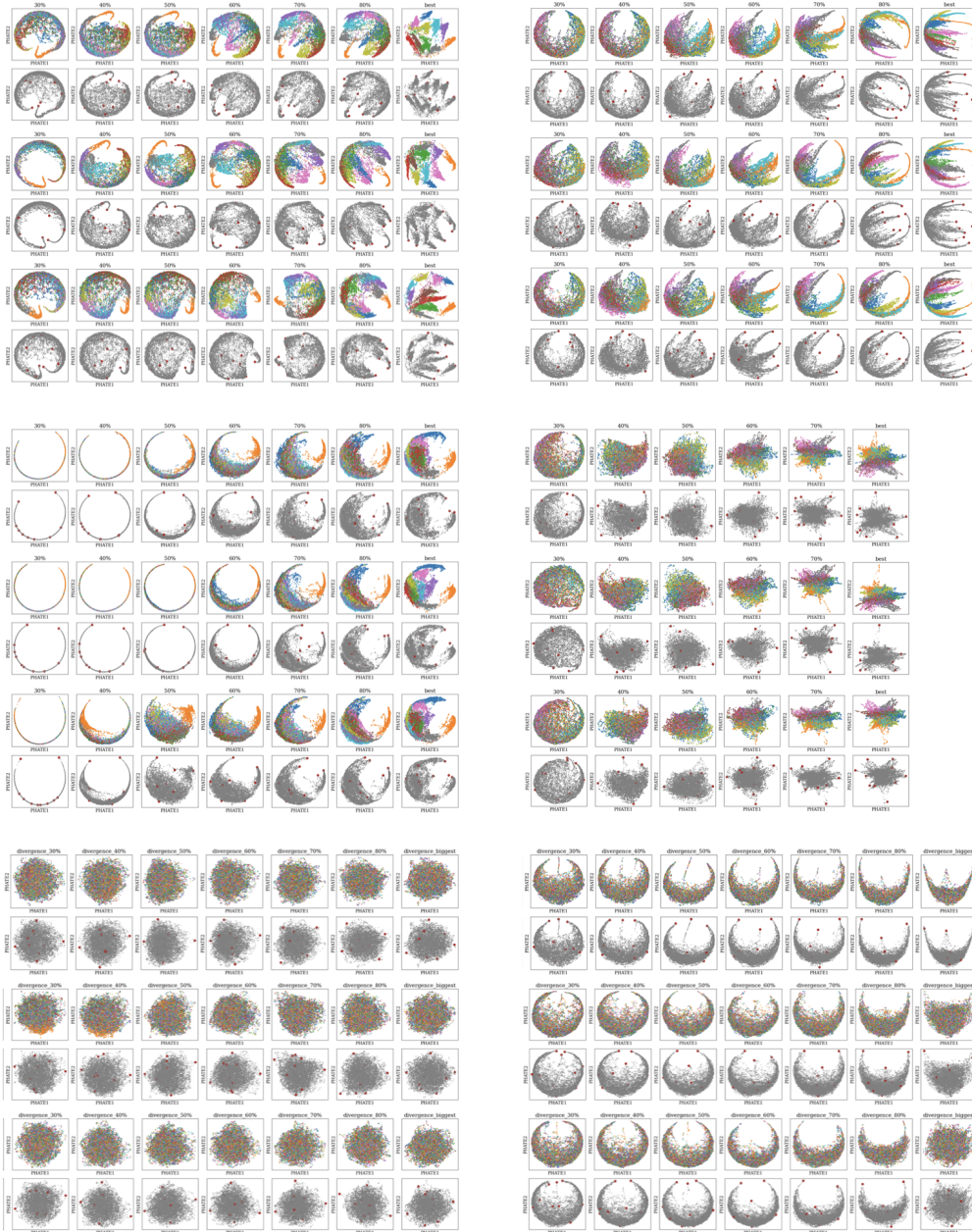


Figure S13. Extended results for Figure S12. **Top Left:** MNIST supervised learning. **Mid Left:** MNIST SimCLR contrastive learning. **Bottom Left:** MNIST overfitting on wrong labels. **Top Right:** CIFAR10 supervised learning. **Mid Right:** CIFAR10 SimCLR contrastive learning. **Bottom Right:** CIFAR10 overfitting on wrong labels.

K. Advantages and Limitations

A key advantage of diffusion spectral entropy over direct computation of Shannon entropy directly in embedded space is that it avoids artificially binning the vectors and it does not require assumptions on distributions of embeddings. Compared to classic binned Shannon entropy, our method performs with greater accuracy.

A limitation of this study is that we did not apply this framework to data from other systems, aside from neural networks to understand how neural networks process information similarly or differently from other systems (such as potentially even brain networks). Further, we did not study the effect initialization or training choices (such as momentum).

Supplementary References

- Burkhardt, D. B., Stanley III, J. S., Perdigoto, A. L., Gigante, S. A., Herold, K. C., Wolf, G., Giraldez, A. J., van Dijk, D., and Krishnaswamy, S. Enhancing experimental signals in single-cell rna-sequencing data using graph signal processing. *Nature Biotechnology*, 2019.
- Chen, T., Kornblith, S., Norouzi, M., and Hinton, G. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pp. 1597–1607. PMLR, 2020.
- Coifman, R. R. and Lafon, S. Diffusion maps. *Applied and computational harmonic analysis*, 21(1):5–30, 2006.
- de Siqueira Santos, S., Takahashi, D. Y., Sato, J. R., Ferreira, C. E., and Fujita, A. Statistical methods in graphs: parameter estimation, model selection, and hypothesis test. *Mathematical Foundations and Applications of Graph Entropy*, 6: 183–202, 2016.
- Fefferman, C., Mitter, S., and Narayanan, H. Testing the manifold hypothesis. *Journal of the American Mathematical Society*, 29(4):983–1049, 2016.
- Gigante, S., Charles, A. S., Krishnaswamy, S., and Mishne, G. Visualizing the phase of neural networks. *Advances in neural information processing systems*, 32, 2019.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Huguet, G., Tong, A., Rieck, B., Huang, J., Kuchroo, M., Hirn, M., Wolf, G., and Krishnaswamy, S. Time-inhomogeneous diffusion geometry and topology. *arXiv preprint arXiv:2203.14860*, 2022.
- Kolchinsky, A. and Tracey, B. D. Estimating mixture entropy with pairwise distances. *Entropy*, 19(7):361, 2017.
- Kraskov, A., Stögbauer, H., and Grassberger, P. Estimating mutual information. *Physical review E*, 69(6):066138, 2004.
- Li, H., Xu, Z., Taylor, G., Studer, C., and Goldstein, T. Visualizing the loss landscape of neural nets. *Advances in neural information processing systems*, 31, 2018.
- Loshchilov, I. and Hutter, F. Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983*, 2016.
- Loshchilov, I. and Hutter, F. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- Merbis, W. and de Domenico, M. Complex information dynamics of epidemic spreading in low-dimensional networks, 2023.
- Moon, K. R., van Dijk, D., Wang, Z., Gigante, S., Burkhardt, D. B., Chen, W. S., Yim, K., Elzen, A. v. d., Hirn, M. J., Coifman, R. R., et al. Visualizing structure and transitions in high-dimensional biological data. *Nature biotechnology*, 37(12):1482–1492, 2019.
- Saxe, A. M., Bansal, Y., Dapello, J., Advani, M., Kolchinsky, A., Tracey, B. D., and Cox, D. D. On the information bottleneck theory of deep learning. *Journal of Statistical Mechanics: Theory and Experiment*, 2019(12):124020, 2019.
- Su, H., Chen, D., Pan, G.-J., and Zeng, Z. Identification of network topology variations based on spectral entropy. *IEEE Transactions on Cybernetics*, 52(10):10468–10478, 2022. doi: 10.1109/TCYB.2021.3070080.

- Takahashi, D. Y., Sato, J. R., Ferreira, C. E., and Fujita, A. Discriminating different classes of biological networks by analyzing the graphs spectra distribution. *PloS one*, 7(12):e49949, 2012.
- Tishby, N. and Zaslavsky, N. Deep learning and the information bottleneck principle. In *2015 ieee information theory workshop (itw)*, pp. 1–5. IEEE, 2015.
- Van Dijk, D., Sharma, R., Nainys, J., Yim, K., Kathail, P., Carr, A. J., Burdziak, C., Moon, K. R., Chaffer, C. L., Pattabiraman, D., et al. Recovering gene interactions from single-cell data using data diffusion. *Cell*, 174(3):716–729, 2018.
- Villafañe-Delgado, M. and Aviyente, S. Graph information theoretic measures on functional connectivity networks based on graph-to-signal transform. In *2016 IEEE Global Conference on Signal and Information Processing (GlobalSIP)*, pp. 1137–1141, 2016. doi: 10.1109/GlobalSIP.2016.7906019.
- von Neumann, J. *Mathematical foundations of quantum mechanics: New edition*, volume 53. Princeton university press, 2018.