
Assessing Neural Network Representations During Training Using Noise-Resilient Diffusion Spectral Entropy

Danqi Liao*
Yale University

Chen Liu*
Yale University

Benjamin W. Christensen
Yale University

Alexander Tong
Université de Montréal
Mila – Quebec AI Institute

Guillaume Huguet
Université de Montréal
Mila – Quebec AI Institute

Guy Wolf
Université de Montréal
Mila – Quebec AI Institute

Maximilian Nickel
Meta AI Research (FAIR)

Ian Adelstein
Yale University

Smita Krishnaswamy ✉
Yale University
Meta AI Research (FAIR)

*Equal contribution. Order of co-first authors determined by coin toss.

✉ Please direct correspondence to smita.krishnaswamy@yale.edu

Abstract

Entropy and mutual information in neural networks provide rich information on the learning process, but they have proven difficult to compute reliably in high dimensions. Indeed, in noisy and high-dimensional data, traditional estimates in ambient dimensions approach a fixed entropy and are prohibitively hard to compute. To address these issues, we leverage data geometry to access the underlying manifold and reliably compute these information-theoretic measures. Specifically, we define diffusion spectral entropy (DSE) in neural representations of a dataset as well as diffusion spectral mutual information (DSMI) between different variables representing data. First, we show that they form noise-resistant measures of intrinsic dimensionality and relationship strength in high-dimensional simulated data that outperform classic Shannon entropy, nonparametric estimation, and mutual information neural estimation (MINE). We then study the evolution of representations in classification networks

with supervised learning, self-supervision, or overfitting. We observe that (1) DSE of neural representations increases during training; (2) DSMI with the class label increases during generalizable learning but stays stagnant during overfitting; (3) DSMI with the input signal shows differing trends: on MNIST it increases, while on CIFAR-10 and STL-10 it decreases. Finally, we show that DSE can be used to guide better network initialization and that DSMI can be used to predict downstream classification accuracy across 962 models on ImageNet. The official implementation is available at <https://github.com/ChenLiu-1996/DiffusionSpectralEntropy>.

1 Introduction

Deep neural networks have emerged as a major breakthrough in data science, mainly because of their ability to learn increasingly meaningful representations of data. Neural networks function by transforming data through a series of nonlinear operations, so that each layer learns a new representation of the data. Although the representation vectors reside in high-dimensional spaces, they are in fact located in lower-dimensional manifolds [Fefferman et al., 2016]. Assessing the properties of this manifold therefore is key to better understanding the neural network.

Extended from a non-archival workshop version (Topology, Algebra, and Geometry in Machine Learning) presented at 40th International Conference on Machine Learning, Honolulu, Hawaii, USA. [Liao et al., 2023].

The information bottleneck principle [Tishby and Zaslavsky, 2015] proposed a framework to measure mutual information between layers within a neural network. While this sheds light on the inner workings of neural networks, these quantities are historically difficult to compute [Paninski, 2003] and often fail when the data dimension increases [Gao et al., 2015, Belghazi et al., 2018], which we refer to as the curse of dimensionality. Specifically, as the dimension D increases, the number of bins in a histogram used to calculate the probabilities increases exponentially in D , that is, b^D where b is the number of bins in each dimension. Indeed, recent investigations along this line are often restricted to small toy models: e.g., in [Saxe et al., 2018, Saxe et al., 2019], most experiments were performed on a 7-layer network of width 12-10-7-5-4-3-2.

Here, we use a powerful manifold learning paradigm, diffusion geometry, to present an alternative approach to quantifying entropy and mutual information. Diffusion geometry is known for its ability to access the underlying manifold of data representations [Coifman and Lafon, 2006], a property that is highly suitable for neural networks. With diffusion geometry, we can separate the noise dimensions from the signal dimensions and compute the entropy within the intrinsic dimensions of the manifold. By defining entropy on the eigenspectrum of the diffusion matrix, the problems incurred by the prohibitively high spatial dimension are mitigated. Further, by powering the matrix, we are able to low-pass filter the eigendimensions of the data, i.e., count noise dimensions with small eigenvalues proportionally less than dimensions with larger eigenvalues that explain more variance in diffusion coordinates.

A key contribution of this work is introducing diffusion spectral entropy (DSE), or spectral entropy of the diffusion operator, as a robust quantifier of the intrinsic information measure of the representation of data in the presence of noise. Furthermore, we extend the diffusion spectral entropy to a diffusion spectral mutual information (DSMI) in order to ascertain the information the embedding manifold has on the class labels or the raw input data from the dataset.

Our main contributions include the following.

- Introducing *diffusion spectral entropy (DSE)*, i.e., entropy of the diffusion operator spectrum as a measure of information in a data representation. Defining *diffusion spectral mutual information (DSMI)* to assess the relationships between different layers of information in a neural network.
- Demonstrating that DSE and DSMI remain descriptive on very high-dimensional data and hence are suitable to modern-sized neural networks.

- Utilizing DSE and DSMI to assess the evolution of neural representations during training.
- Demonstrating two ways to utilize DSE and DSMI. Specifically, using DSE to guide network initialization and using DSMI to predict downstream performance without full evaluation.

2 Background

2.1 Manifold learning and diffusion geometry

A useful assumption in representation learning is that high-dimensional data, commonly used in deep learning, originate 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 are collected. Conversely, given a dataset $X = \{x_1, \dots, x_N\} \subset \mathbb{R}^n$ of high-dimensional observations, manifold learning methods assume that the data points originate from a sampling $Z = \{z_i\}_{i=1}^N \in \mathcal{M}^d$ of the underlying manifold through $x_i = \mathbf{f}(z_i)$, $i = 1, \dots, n$ and aim to learn a low-dimensional intrinsic representation that approximates the manifold geometry of \mathcal{M}^d .

A paradigm that has emerged in manifold learning in recent years is diffusion geometry [Coifman and Lafon, 2006, Van Dijk et al., 2018, Burkhardt et al., 2019, Huguet et al., 2022]. It models data based on transition or random walk probabilities through the data and has been shown to be noise-tolerant. It has produced methods such as tSNE [Van der Maaten and Hinton, 2008], PHATE [Moon et al., 2019] and diffusion maps [Coifman and Lafon, 2006]. The geometry of the manifold can be learned with data diffusion by first computing the local similarities defined via a kernel \mathcal{K} . A popular choice is a Gaussian kernel $e^{-\|z_1 - z_2\|^2 / \sigma}$, where $\sigma > 0$ is interpreted as a user-configurable neighborhood size. To construct a diffusion geometry that is robust to sampling density variations, one would typically use an anisotropic kernel

$$\mathcal{K}(z_1, z_2) = \frac{\mathcal{G}(z_1, z_2)}{\|\mathcal{G}(z_1, \cdot)\|_1^\alpha \|\mathcal{G}(z_2, \cdot)\|_1^\alpha}, \text{ where} \quad (1)$$

$$\mathcal{G}(z_1, z_2) = e^{-\frac{\|z_1 - z_2\|^2}{\sigma}}$$

as proposed in [Coifman and Lafon, 2006]. $0 \leq \alpha \leq 1$ controls the separation of geometry from density.

Next, the similarities encoded by \mathcal{K} are normalized to define the 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

$$\mathbf{P}_{i,j} = p(z_i, z_j) = \frac{\mathcal{K}(z_1, z_2)}{\|\mathcal{K}(z_1, \cdot)\|_1} \quad (2)$$

which describes a Markovian diffusion process on the intrinsic geometry of the data.

2.2 Entropy and mutual information

Entropy Entropy 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 a spread-out probability mass, 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 follows.

$$H(X) = \mathbb{E}[-\log p(X)] = - \sum_{x \in X} p(x) \log p(x) \quad (3)$$

The von Neumann entropy [von Neumann, 2018] from the quantum information domain extends the entropy measure to the quantum mechanics context and operates on density matrices. If a density matrix ρ has a set of eigenvalues $\{\eta_i\}$, the von Neumann entropy is defined as $H(\rho) = -\text{tr}(\rho \log \rho) = -\sum_i \eta_i \log \eta_i$.

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 that 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 eigenvalues of a graph adjacency matrix and have been used in biology and other fields to compare graphs [Su et al., 2022, de Siqueira Santos et al., 2016, Takahashi et al., 2012, Merbis and de Domenico, 2023, Villafañe-Delgado and Aviyente, 2016].

Mutual information Mutual information is defined as a function of entropy. There are many equivalent formulations of mutual information that are equiv-

alent. The most useful formulation here is as the difference between the (unconditional) entropy of a variable and the entropy of a variable conditioned on the other variable. The conditional entropy is computed as a weighted sum over discrete categories.

$$\begin{aligned} I(X; Y) &= H(X) - H(X|Y) \\ &= H(X) - \sum_i p(Y = y_i) H(X|Y = y_i) \quad (4) \end{aligned}$$

2.3 The information bottleneck theory

The information bottleneck theory [Tishby and Zaslavsky, 2015] views a neural network with sequentially placed layers as a multistage information compressor. If we denote the neural representation of any layer as Z , the input signals as X and the ground truth labels as Y , the theory concludes that $I(Z; X)$ will decrease while $I(Z; Y)$ will increase during training. Additionally, the theory claims that the optimal representation can be achieved by minimizing $I(Z; X) - \beta I(Z; Y)$.

3 Related Work

In efforts to quantify information in neural networks, in [Tishby and Zaslavsky, 2015], the authors 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 method is the curse of dimensionality in the binning process that makes it impractical to analyze layers with more than a dozen neurons, which lag modern deep neural networks by orders of magnitude (see Supplement A).

Researchers proposed using kernel density estimators, Kraskov estimators, as well as other nonparametric or parametric estimators [Saxe et al., 2018, Kolchinsky and Tracey, 2017, Kraskov et al., 2004, Fraser and Swinney, 1986, Darbellay and Vajda, 1999, Suzuki et al., 2008, Kwak and Choi, 2002, Moon et al., 1995, Gao et al., 2015] for mutual information estimation, yet these methods require assumptions on the distributions of hidden layer activation or do not generalize well on high-dimensional data [Belghazi et al., 2018].

Among them, the nonparametric estimation toolbox [Ver Steeg, 2000, Kraskov et al., 2004] locally estimates the log probability density at each sample point and averages them for a holistic estimate. Mutual information neural estimation (MINE) [Belghazi et al., 2018] uses a neural network, optimized by gradient descent, to approximate a lower bound of mutual information. In contrast, our proposed method does not assume a specific distribution on hidden-layer activa-

tions and remains descriptive in high dimensions.

4 Methods

The main difficulty in entropy and mutual information quantification in neural networks is the prohibitively high dimensionality. Binning and quantization fail quickly when we increase the dimension of the data (Supplement A). We leverage concepts from diffusion geometry to access the underlying manifold and compute these quantities reliably at high dimensions.

4.1 Our definitions using diffusion geometry

While spectral entropy has often been used to measure 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. In particular, we utilize the anisotropically normalized diffusion operator from Eqn 2.

We define the symmetric matrix K by organizing the anisotropic kernel \mathcal{K} in Eqn 1 with $\alpha = \frac{1}{2}$

$$K_{i,j} = \mathcal{K}(z_i, z_j) \quad (5)$$

The row stochastic matrix \mathbf{P} computed using Eqn 2 with this K 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 4.1. We define *diffusion spectral entropy (DSE)* as an entropy of the eigenvalues of the diffusion operator \mathbf{P}_X calculated on a dataset X where $x \in X$ is a multidimensional vector $[x_1, x_2 \dots x_d]^T$:

$$S_D(\mathbf{P}_X, t) := - \sum_i \alpha_{i,t} \log(\alpha_{i,t}) \quad (6)$$

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 according to its probability of transition to every other data point if a random walk on the data is made. Thus, if a data point is disconnected or far away from others, then a random walk starting at the data point is likely to remain at the data point. In this setting, the 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 viewed as representations of the data, this is also a measure of the intrinsic dimensionality of the dataset. Note that the parameter t that parameterizes the entropy also gives us the capability to separate the 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 . In fact, \mathbf{P}^t has identical eigenvectors and powered eigenvalues as \mathbf{P} , which achieves a low-pass filtering of data values on the affinity graph [Van Dijk et al., 2018]. We note that a similar measure was used in the supplement of [Moon et al., 2019] to select parameters, but the diffusion spectral entropy was not defined or discussed there.

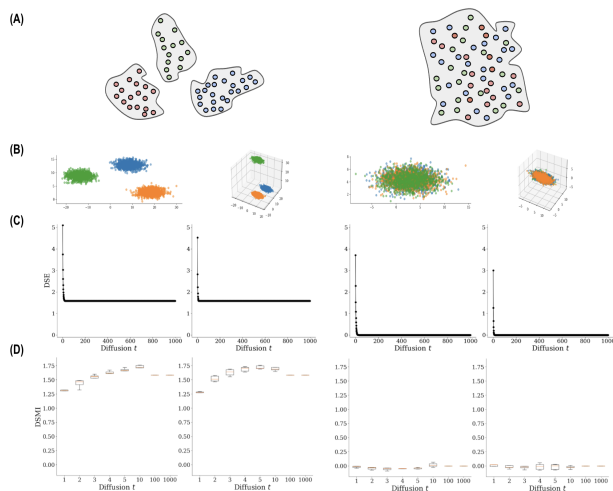


Figure 1: **Intuition for DSE and DSMI.** X : data, Y : class label. Colors represent classes. (A) sketch, (B) simulation, (C) DSE, (D) DSMI. In all panels, **Left:** X forms k clusters and Y corresponds to clusters. $S_D(\mathbf{P}_X, t) \approx \log(k)$ at sufficiently large t and $I_D(X; Y) > 0$. **Right:** X forms a single cluster and Y is randomly distributed over X . $S_D(\mathbf{P}_X, t) \approx 0$ at sufficiently large t and $I_D(X; Y) \approx 0$.

An important remark is that, since DSE is computed on the eigenspectrum of the diffusion geometric representation of the data manifold, diffusion spectral entropy essentially counts information in terms of the number of eigendirections with non-trivial eigenvalues. In other words, *DSE is more sensitive to the underlying dimension and structures (e.g., number of branches or clusters) than to the spread or noise in the data itself, which is contracted to the manifold by raising the diffusion operator to the power of t .* For an illustration of this phenomenon, see Figure 1(C), where as t increases, the DSE converges to the log of number of

distinct structures in the data. This is a highly valuable trait in the analysis of neural network representations, which often occurs on very noisy data. Indeed, neural representations that classify well create clear separations of the data into different eigendirections.

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 4.2. We define *Diffusion Spectral Mutual Information (DSMI)* as the difference between unconditional and conditional diffusion spectral entropy, similar to the definition in Eqn 4:

$$I_D(X; Y) = S_D(\mathbf{P}_X, t) - \sum_{y_i \in Y} p(Y = y_i) S_D(\mathbf{P}_{X|Y=y_i}, t) \quad (7)$$

$\mathbf{P}_{X|Y=y_i}$ is the transition matrix computed on the subset of X that has class label Y . To avoid the numeric issues involved in comparing spectra of different sizes of matrices, we also computed $S_D(\mathbf{P}_X, t)$ after subsampling X to the size of each class in Y . Since uniform subsampling maintains distributions, the sampled entropy would be the same as the total entropy as shown in our experiments (see Supplement H).

Figure 1 provides some insight into DSE and DSMI. In cases where the dataset X forms k clusters and the class label coincides exactly with well-separated clusters, DSE approaches $\log(k)$ with increasing t , because it counts the number of structures of the manifold instead of the data variance in the ambient space (panel (C) left). For DSMI, based on Proposition 4.4, the unconditional entropy with multiple clusters shall be greater than the conditional entropy within a single cluster, so DSMI will be positive (panel (D) left). In the second case, where X forms a single cluster and class labels are evenly distributed across the manifold, DSE approaches $\log(1) = 0$ with increasing t as there is only one structure (panel (C) right). For DSMI, the conditional entropy $S_D(X|Y)$ and unconditional entropy $S_D(X)$ would be similar. As a result, DSMI shall be close to zero or even slightly negative for numerical reasons but, nonetheless, indicates low MI (panel (D) right). See Supplement G and H for additional results.

4.2 Properties and propositions

Here we discuss some properties of DSE and DSMI, with **proofs in Supplement B**. First, we provide the lower bound and the upper bound of S_D when $t \rightarrow \infty$, and we explain the conditions when they are reached.

Proposition 4.1. S_D achieves a minimal entropy of 0 when the diffusion operator defines an ergodic Markov chain, and is in steady state (as $t \rightarrow \infty$).

Note that this also implies that if all data points are very similar and all points can easily transition to any other point, then it has minimal entropy.

Proposition 4.2. As $t \rightarrow \infty$, $S_D(\mathbf{P}_X, t)$ on data with k well-separated clusters is $\log(k)$.

This shows that S_D will reach its maximum value when the points are spread out very far apart.

Next, we examine the expected value of S_D .

Proposition 4.3. Let $X \in \mathbb{R}^{n \times d}$ be a dataset of n independent and identically distributed multivariate Gaussian vectors in \mathbb{R}^d , where $x_i \sim \mathcal{N}(0, I_d)$. Then, using K as defined in Eqn 1 with $\alpha = 1/2$,

$$\begin{aligned} & \mathbb{E}[S_D(\mathbf{P}_X, t = 1)] \\ & \lesssim \log\left(\frac{n}{1-\beta}\right) - \left(\frac{1}{n} + \left(\frac{n-1}{n}\right)\beta\right) \log\left(1 + \frac{\beta n}{1-\beta}\right) \\ & \text{where } \beta = \left(1 + \frac{4}{\sigma}\right)^{-\frac{d}{2}} \end{aligned}$$

This establishes a theoretical upper bound¹ on the DSE at any given layer. Since independence is not necessarily assumed in general, the DSE for any dataset in practice should be less than this bound. Furthermore, this bound intuitively reinforces that for large d , $\beta \approx 0$, so $\mathbb{E}[S_D(\mathbf{P}_X, t = 1)] \approx \log(n)$.

Our last proposition on DSE relates DSE of a single cluster to that of multiple clusters. This has implications for classifier training.

Proposition 4.4. Take n to be arbitrarily large. Let $X \in \mathbb{R}^{n \times d}$ be a matrix of i.i.d. random values $x_{ij} \sim f$. Let $Y \in \mathbb{R}^{n \times d}$ be a matrix of i.i.d. random values $y_{ij} \sim f$, but in $k \in \mathbb{N}$ distinct clusters such that when the anisotropic probability matrix is computed for $\alpha = 1/2$, the probability of diffusion between points of different clusters is arbitrarily small. Then, using β as defined in Proposition 4.3, the approximate upper bound on DSE increases by $\beta \log(k)$.

Recall the learning process of a classification network. As the model learns, it separates latent representations into distinct clusters throughout the process. Proposition 4.4 elucidates that the expected upper bound of S_D will increase during the training of a classifier.

¹It is also worth noting that the above bound remains valid for any distribution where $\|x_i - x_j\|_2^2 \sim f$ for $i \neq j$, where f has the real-valued moment generating function $M(t)$. Then, we simply set $\beta = M(-\frac{1}{\sigma})$.

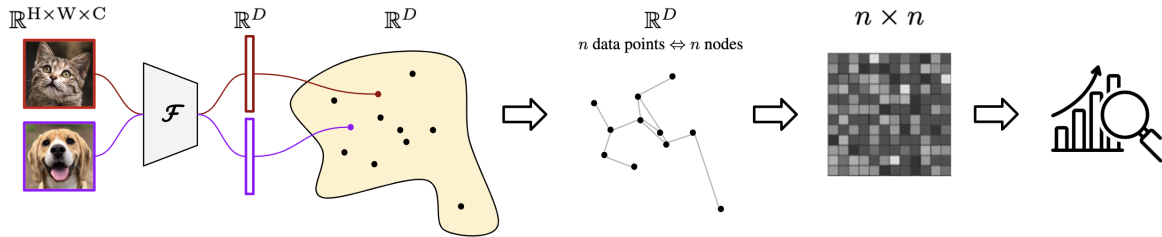


Figure 2: **Data processing to obtain the diffusion matrix of an embedding manifold.** Each data point is embedded by the network as a vector on a D -dimensional embedding manifold. This yields a point cloud that can be converted into a graph. We can compute its diffusion matrix, which allows for further analysis.

DSMI has a minimal value of 0, achieved when the random variables are mutually independent. On the other hand, its maximal value can be bounded as follows.

Proposition 4.5. *As $t \rightarrow \infty$ in a hidden layer X with k well-separated clusters and class labels Y perfectly indicating clusters, we will have $I_D(X; Y) = \log(k)$.*

4.3 Efficiently computing DSE and DSMI

To compute DSE and DSMI on neural network representations, we pass the training data X to the desired layer L of a neural network. At this layer, we collect the activations of each neuron into a vector $L(x_i) = [L_1(x_i), \dots, L_n(x_i)]$, where L_j is the j -th neuron of the L -th layer. $L(x_i)$ is a n -dimensional representation. If the activation is a multidimensional tensor, we can flatten it into a vector. We then compute $\mathbf{P}_{L(X)}$ and proceed to compute DSE and DSMI. This process is illustrated in Figure 2.

The calculation of eigenvalues via complete eigen-

decomposition is known to have a time complexity $\mathcal{O}(n^3)$. However, we take advantage of two characteristics to provide a faster method: 1) we do not need the eigenvectors, 2) \mathbf{P} has the same eigenvalues as the real symmetric matrix K defined in Eqn 5 as shown in [Coifman and Lafon, 2006]. Thus, we can compute the DSE efficiently using QR decomposition since we only need the eigenvalues of a real symmetric matrix.

5 Results

5.1 Toy test cases for DSE and DSMI

To demonstrate the behavior of DSE, we first performed several simulations as shown in Figure 3(A). The left panel indicates that, for an arbitrary real symmetric matrix, the closer it gets to an identity matrix, the higher its DSE — since an $n \times n$ identity matrix represents n distinct clusters/states. The two other panels show that DSE in general increases as the intrinsic dimensionality of the data manifold increases.

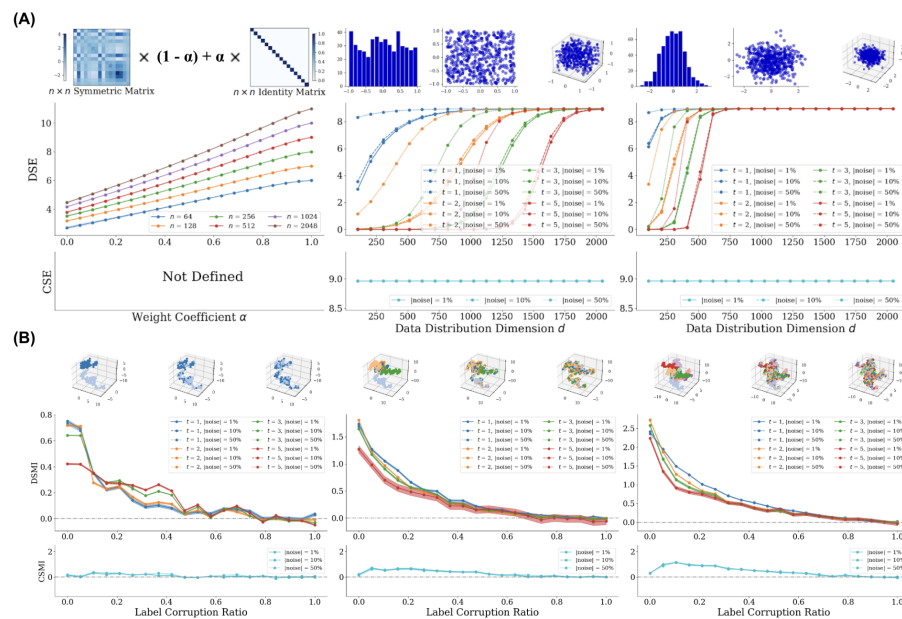


Figure 3: **Diffusion Spectral Entropy (DSE) and Diffusion Spectral Mutual Information (DSMI) on toy data.** (A) DSE increases as intrinsic dimension grows, while classic Shannon entropy (CSE) saturates to $\log(n) = \log_2(500) = 8.966$ due to curse of dimensionality. (B) When two random variables are dependent, DSMI negatively correlates with the level of data corruption, while classic Shannon mutual information (CSMI) does not capture this trend. See full description in Supplement C.

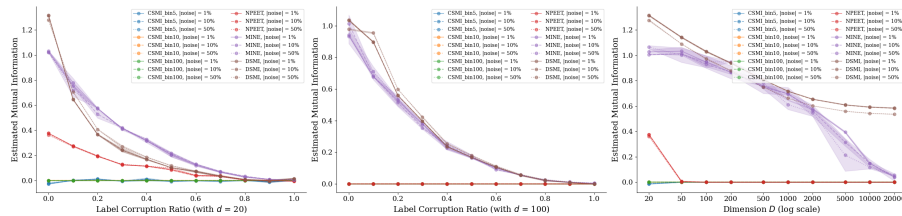


Figure 4: Mutual information estimation on Gaussian blobs.

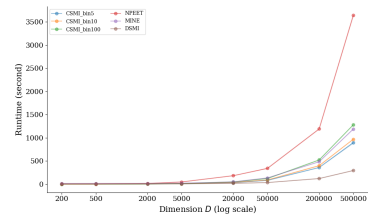


Figure 5: Runtime comparison.

To demonstrate the behavior of DSME, we performed similar simulations as shown in Figure 3(B). In these toy test cases, we gradually corrupt the class label on the simulated tree branches. Just as expected, DSME starts high when the association between the data point and the class label is high (i.e., when each branch is assigned a distinct label), and drops to zero as the corruption gradually brings the label to complete noise.

We also computed the classic Shannon entropy (CSE) and classic Shannon mutual information (CSMI) for comparison using the method from [Tishby and Zaslavsky, 2015]. Our proposed method outperforms the classic Shannon version in high-dimensional spaces. In Figure 3(A), DSE consistently captures entropy trends while CSE saturates to $\log(n)$. In Figure 3(B), DSME follows the expected trend, while CSMI does not.

5.2 DSME at very high dimensions

In this section, we show how DSME better scales to high dimensions compared to existing methods. We ran mutual information estimators (DSME, CSMI, a nonparametric entropy and mutual information estimation toolbox (NPEET) [Ver Steeg, 2000, Kraskov et al., 2004] and mutual information neural estimation (MINE) [Belghazi et al., 2018]) on labeled Gaussian blobs of various dimensions. The left and middle panels of Figure 4 show that all methods generally obey the expected behavior as they show positive mutual information when the label is clean and gradually degrade to zero when the label is fully corrupted. The right panel indicates that CSMI, NPEET, and MINE fail as the dimension gets large, while DSME remains significant.

Moreover, DSME scales better in runtime (Figure 5). This is expected since the time complexity for DSME is $\mathcal{O}(n^2 \cdot D + n^3)$ due to the cost of pairwise distance computation and eigendecomposition (without considering optimization), while that for CSMI is $\mathcal{O}(2^D)$ due to the cost of computing exponentially many data bins.

5.3 Evolution of DSE and DSME during training

To analyze DSE and DSME on real data, we trained common vision networks and assessed the penultimate layer of each network at the end of each epoch. To cover the variety of model architectures, we experimented with 3 convolutional neural networks and 3 vision transformers. See supplementary D for details.

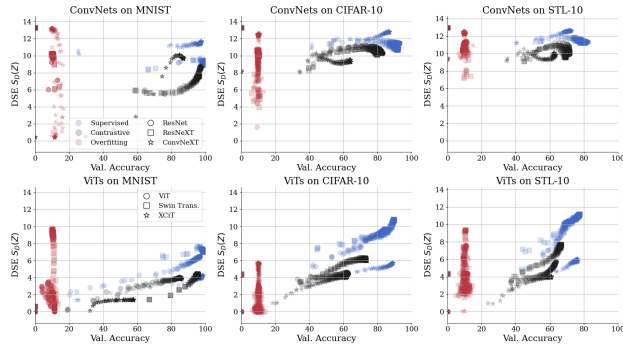
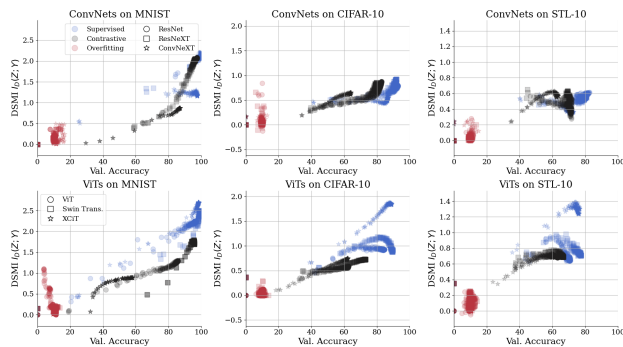
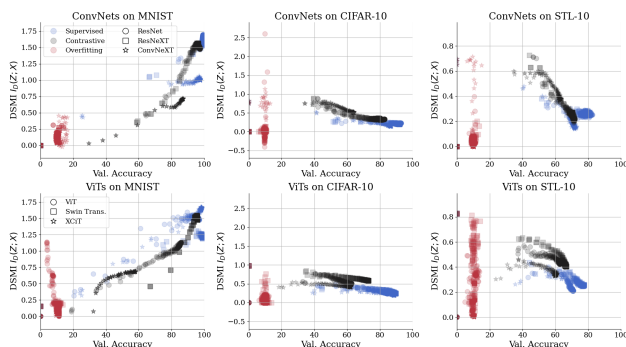
We specifically investigated the penultimate layer because it is usually believed that the representations in this layer reflect the learning of the entire network. As a reference, it is common practice to directly use the penultimate layer to evaluate the overall representation power of the network, e.g., linear probing [Chen et al., 2020]. Although our experiments mainly focus on this layer, our evaluation framework can be easily adapted to other layers in the network.

We trained the vision backbones under three conditions: supervised learning, contrastive learning, and intentional overfitting on randomized nonsense labels, repeated under 3 random seeds. More details can be found in Supplement D. For all of the DSE and DSME experiments below, results of the classic Shannon version can be found in Supplement A, and the raw results for each network are included in Supplement E.

DSE of a hidden layer We can observe from Figure 6 that DSE in proper learning (i.e. supervised or contrastive learning on correct labels) increases as the model performs better on the classification task. When the model is forced to memorize random nonsense labels, the entropy increases similarly even though the classification performance is stagnant.

DSME of a hidden layer with the class label

In Figure 7, it can be seen that DSME $I_D(Z; Y)$ increases consistently during proper learning. In many cases, DSME climbs more slowly in contrastive learning compared to supervised learning and ends up at a lower terminal value. This may be attributed to the fact that contrastive learning lacks direct supervision from explicit class labels. However, since class labels relate to the data geometry, self-supervised learning


 Figure 6: DSE $S_D(Z)$ of representation Z .

 Figure 7: DSMI $I_D(Z; Y)$ between representation Z and the class label Y .

 Figure 8: DSMI $I_D(Z; X)$ between Z and input X .

on the data alone still yields some mutual information with labels. In nonsense memorization, DSMI quickly converges to around 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 positive mutual information.

DSMI of a hidden layer with the input We also show DSMI with the input signal in Figure 8. During proper learning, the information bottleneck theory would suggest $I_D(Z; X)$ should decrease [Tishby

and Zaslavsky, 2015] while counterarguments have also been provided in [Saxe et al., 2018]. Our results suggest that both may 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, while it mostly decreases after an initial increase on the CIFAR-10 and STL-10 datasets. We suspect that all models are essentially overfitting the MNIST data, and this may be verified by future studies. In nonsense memorization, $I_D(Z; X)$ rises to a significant level in most cases. This indicates that even when they overfit random labels, the neural representations still contain information about the input signals.

5.4 Guiding network initialization

We sought to assess the effects of network initialization in terms of DSE. We were motivated by two observations (see Figure S6 in Supplement E): (1) the initial DSEs for different models are not always the same despite using the same method for random initialization; (2) if DSE starts low, it grows monotonically; if DSE starts high, it first decreases and then increases.

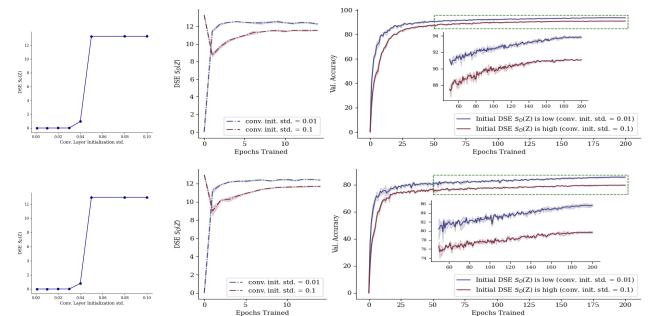


Figure 9: Initializing the network at a low DSE allows for faster convergence and better final performance. Top: CIFAR-10, bottom: STL-10. Shaded areas indicate standard deviation across 3 random seeds.

We found that if we initialize the convolutional layers with weights $\sim \mathcal{N}(0, \sigma)$, DSE $S_D(Z)$ is affected by σ (Figure 9 left). We then trained ResNet models with networks initialized at high ($\approx \log(n)$) versus low (≈ 0) DSE by setting $\sigma = 0.1$ and $\sigma = 0.01$, respectively. The training history suggests that initializing the network at a lower $S_D(Z)$ can improve the convergence speed and final performance (Figure 9). We believe this is because the high initial DSE from random initialization corresponds to an undesirable high-entropy state, which the network needs to get away from (causing the DSE decrease) before it migrates to the desirable high-entropy state (causing the DSE increase). More details can be found in Supplement F.

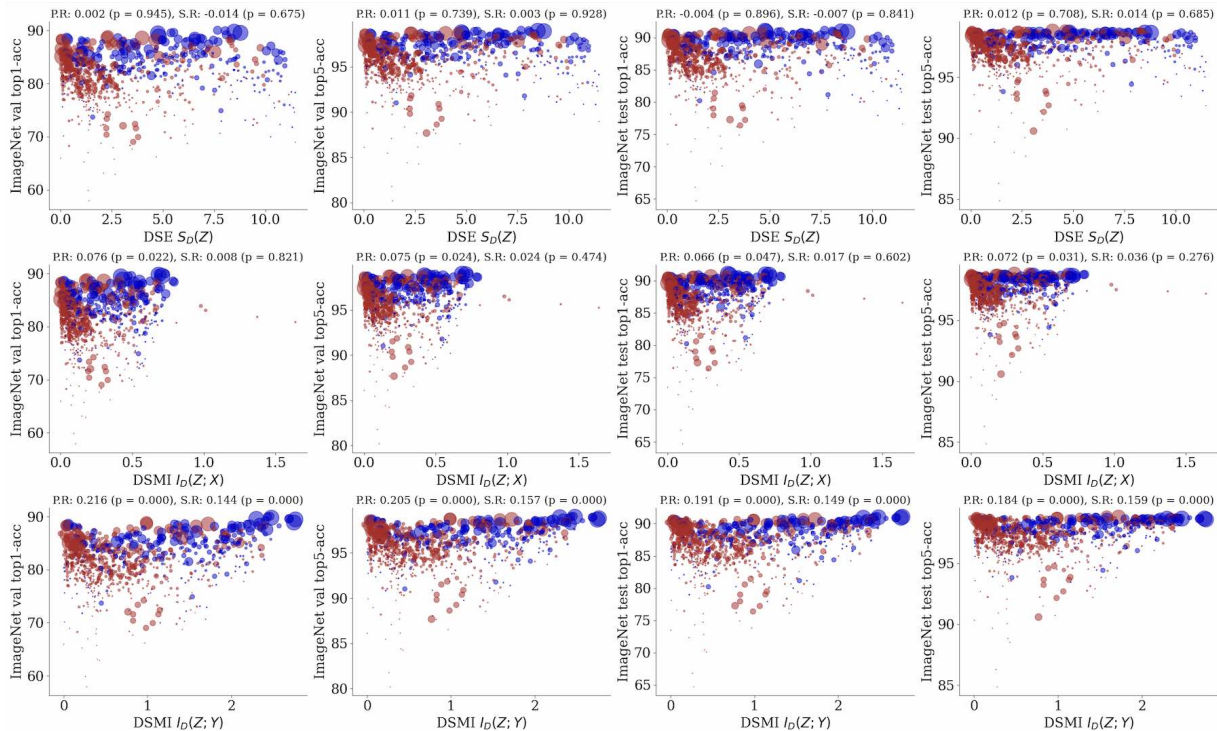


Figure 10: **Correlation analysis between DSE $S_D(Z)$, DSMI $I_D(Z; X)$, $I_D(Z; Y)$ and ImageNet accuracy evaluated on 962 pre-trained models.** Red circles are ConvNets and blue circles are ViTs. Circle sizes indicate model sizes. $I_D(Z; Y)$ shows a strong positive correlation ($p < 0.001$). P.R: Pearson correlation coefficient, S.R: Spearman correlation coefficient.

5.5 ImageNet cross-model correlation

In addition to evaluating DSE and DSMI *along the training process of the same model*, we investigated whether they can serve as meaningful indicators *across different models*. To that end, we computed the correlations between these diffusion geometric quantities on a representative subset of data and the downstream classification accuracy across 962 publicly available models pre-trained on ImageNet. Details of this experiment are described in Supplement D. As illustrated in Figure 10, $I_D(Z; Y)$ has a strong positive correlation with classification accuracy (with $p < 0.001$ consistently under four conditions). On the other hand, $S_D(Z; X)$ and $I_D(Z; X)$ do not show strong statistical significance. This indicates a trend that a higher $I_D(Z; Y)$ corresponds to a better downstream performance, which may facilitate the selection of promising models without evaluating on the entire dataset.

6 Conclusion

In this work, we define diffusion spectral entropy (DSE) and diffusion spectral mutual information (DSMI). Through simulation on toy datasets, we showed that DSE and DSMI are meaningful and noise-

resilient measures that can scale better to very high dimensions compared to existing methods in both reliability and runtime. We then analyzed the neural representations of six vision backbones in supervised learning, contrastive learning, and overfitting settings. We observed that DSE during learning increases under all training conditions. Furthermore, we showed that DSMI between a hidden layer and class labels increases during learning, whereas it converges to zero when overfitting random labels. We saw more complex trends in DSMI with the primary input: it increases in the MNIST dataset, while it decreases in CIFAR-10 and STL-10. Finally, we show that DSE can be used to guide network initialization and that DSMI is strongly correlated with downstream performance.

Some future directions include further investigating the effect of network initialization, exploring DSE and DSMI as regularizations for supervised learning, using DSE to regularize self-supervised learning, using DSMI to improve adversarially trained generative models, and many other exciting opportunities. Researchers can further extend this framework to data from other systems, in addition to neural networks, to understand how neural networks process information similarly or differently from other systems, such as brain networks.

References

- [Belghazi et al., 2018] Belghazi, M. I., Baratin, A., Rajeshwar, S., Ozair, S., Bengio, Y., Courville, A., and Hjelm, D. (2018). Mutual information neural estimation. In *International conference on machine learning*, pages 531–540. PMLR.
- [Burkhardt et al., 2019] 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. (2019). Enhancing experimental signals in single-cell rna-sequencing data using graph signal processing. *Nature Biotechnology*.
- [Chen et al., 2020] Chen, T., Kornblith, S., Norouzi, M., and Hinton, G. (2020). A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR.
- [Coifman and Lafon, 2006] Coifman, R. R. and Lafon, S. (2006). Diffusion maps. *Applied and computational harmonic analysis*, 21(1):5–30.
- [Darbellay and Vajda, 1999] Darbellay, G. A. and Vajda, I. (1999). Estimation of the information by an adaptive partitioning of the observation space. *IEEE Transactions on Information Theory*, 45(4):1315–1321.
- [de Siqueira Santos et al., 2016] de Siqueira Santos, S., Takahashi, D. Y., Sato, J. R., Ferreira, C. E., and Fujita, A. (2016). Statistical methods in graphs: parameter estimation, model selection, and hypothesis test. *Mathematical Foundations and Applications of Graph Entropy*, 6:183–202.
- [Fefferman et al., 2016] Fefferman, C., Mitter, S., and Narayanan, H. (2016). Testing the manifold hypothesis. *Journal of the American Mathematical Society*, 29(4):983–1049.
- [Fraser and Swinney, 1986] Fraser, A. M. and Swinney, H. L. (1986). Independent coordinates for strange attractors from mutual information. *Physical review A*, 33(2):1134.
- [Gao et al., 2015] Gao, S., Ver Steeg, G., and Galstyan, A. (2015). Efficient estimation of mutual information for strongly dependent variables. In *Artificial intelligence and statistics*, pages 277–286. PMLR.
- [Huguet et al., 2022] Huguet, G., Tong, A., Rieck, B., Huang, J., Kuchroo, M., Hirn, M., Wolf, G., and Krishnaswamy, S. (2022). Time-inhomogeneous diffusion geometry and topology. *arXiv preprint arXiv:2203.14860*.
- [Kolchinsky and Tracey, 2017] Kolchinsky, A. and Tracey, B. D. (2017). Estimating mixture entropy with pairwise distances. *Entropy*, 19(7):361.
- [Kraskov et al., 2004] Kraskov, A., Stögbauer, H., and Grassberger, P. (2004). Estimating mutual information. *Physical review E*, 69(6):066138.
- [Kwak and Choi, 2002] Kwak, N. and Choi, C.-H. (2002). Input feature selection by mutual information based on parzen window. *IEEE transactions on pattern analysis and machine intelligence*, 24(12):1667–1671.
- [Liao et al., 2023] Liao, D., Liu, C., Christensen, B., Tong, A., Huguet, G., Wolf, G., Nickel, M., Adelstein, I., and Krishnaswamy, S. (2023). Assessing neural network representations during training using data diffusion spectra. In *ICML 2023 Workshop on Topology, Algebra and Geometry in Machine Learning (TAG-ML)*.
- [Merbis and de Domenico, 2023] Merbis, W. and de Domenico, M. (2023). Complex information dynamics of epidemic spreading in low-dimensional networks.
- [Moon et al., 2019] 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. (2019). Visualizing structure and transitions in high-dimensional biological data. *Nature biotechnology*, 37(12):1482–1492.
- [Moon et al., 1995] Moon, Y.-I., Rajagopalan, B., and Lall, U. (1995). Estimation of mutual information using kernel density estimators. *Physical Review E*, 52(3):2318.
- [Paninski, 2003] Paninski, L. (2003). Estimation of entropy and mutual information. *Neural computation*, 15(6):1191–1253.
- [Saxe et al., 2018] Saxe, A. M., Bansal, Y., Dapello, J., Advani, M., Kolchinsky, A., Tracey, B. D., and Cox, D. D. (2018). On the information bottleneck theory of deep learning. In *International Conference on Learning Representations*.
- [Saxe et al., 2019] Saxe, A. M., Bansal, Y., Dapello, J., Advani, M., Kolchinsky, A., Tracey, B. D., and Cox, D. D. (2019). On the information bottleneck theory of deep learning. *Journal of Statistical Mechanics: Theory and Experiment*, 2019(12):124020.
- [Su et al., 2022] Su, H., Chen, D., Pan, G.-J., and Zeng, Z. (2022). Identification of network topology variations based on spectral entropy. *IEEE Transactions on Cybernetics*, 52(10):10468–10478.

- [Suzuki et al., 2008] Suzuki, T., Sugiyama, M., Sese, J., and Kanamori, T. (2008). Approximating mutual information by maximum likelihood density ratio estimation. In *New challenges for feature selection in data mining and knowledge discovery*, pages 5–20. PMLR.
- [Takahashi et al., 2012] Takahashi, D. Y., Sato, J. R., Ferreira, C. E., and Fujita, A. (2012). Discriminating different classes of biological networks by analyzing the graphs spectra distribution. *PLoS one*, 7(12):e49949.
- [Tishby and Zaslavsky, 2015] Tishby, N. and Zaslavsky, N. (2015). Deep learning and the information bottleneck principle. In *2015 IEEE Information Theory Workshop (ITW)*, pages 1–5. IEEE.
- [Van der Maaten and Hinton, 2008] Van der Maaten, L. and Hinton, G. (2008). Visualizing data using t-sne. *Journal of machine learning research*, 9(11).
- [Van Dijk et al., 2018] 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. (2018). Recovering gene interactions from single-cell data using data diffusion. *Cell*, 174(3):716–729.
- [Ver Steeg, 2000] Ver Steeg, G. (2000). Non-parametric entropy estimation toolbox (npeet). *Non-parametric entropy estimation toolbox (NPEET)*.
- [Villafañe-Delgado and Aviyente, 2016] Villafañe-Delgado, M. and Aviyente, S. (2016). 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)*, pages 1137–1141.
- [von Neumann, 2018] von Neumann, J. (2018). *Mathematical foundations of quantum mechanics: New edition*, volume 53. Princeton university press.

Supplementary Material

A	Limitations of the Classic Shannon Entropy and Mutual Information	12
A.1	Computing the Classic Shannon Entropy and Mutual Information	12
A.2	Limitations of Classic Shannon Entropy and Mutual Information	13
B	Propositions with Proofs on DSE and DSMI	15
C	Full Caption of the Toy Data Figure	21
D	Experimental Details	22
D.1	Three conditions for neural network training process	22
D.2	Computing DSMI with input	22
D.3	Details on Imagenet cross-model correlation	23
D.4	Logarithm base	23
E	All Raw Results	24
F	DSE Initialization Experiments	26
G	Additional DSE Results: Intuition and MAGIC visualization	27
G.1	DSE vs. Diffusion t in the Intuition Figure	27
G.2	MAGIC Visualization for DSE vs. Diffusion t	27
G.3	Ablation Study	28
H	Additional DSMI Results: Intuition and Subsampling Robustness	29
H.1	DSMI vs. Diffusion t in the Intuition Figure	29
H.2	Subsampling Technique for DSMI: Subsampling Robustness	30

A Limitations of the Classic Shannon Entropy and Mutual Information

A.1 Computing the Classic Shannon Entropy and Mutual Information

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

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 3. The most common way (e.g., in [Tishby and 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 the vectors along each dimension into b different bins. For example, if $b = 10$, 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 $[1, b]$. Each possible quantized vector can be referred to as a “bucket”. It can be easily noticed that the number of buckets is equal to b^D . After counting the number of vectors in each bucket, we can

estimate the distribution of probability density on the buckets. Finally, we compute the classic Shannon entropy using Eqn 3. Classic Shannon mutual information can be computed in a similar manner, using Eqn 4 and the aforementioned entropy computation.

A.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 scales exponentially with respect to the feature dimension D** . Take ResNet-50 as an 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 result in unique bucketing which leads to the maximum entropy, this metric has very limited expressiveness. This phenomenon, which we call the curse of dimensionality in the CSE and CSMI computation, is illustrated in Figure S1.

This phenomenon can be seen in both toy data (Figures 6 and 7) and real data (Figures S2, S3 and S4).

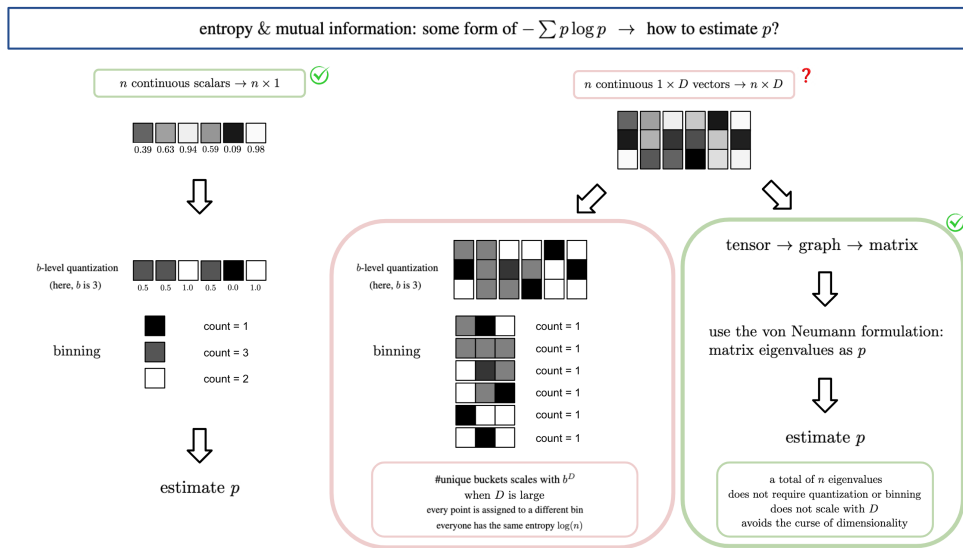


Figure S1: Curse of dimensionality in Classic Shannon Entropy / Mutual Information.

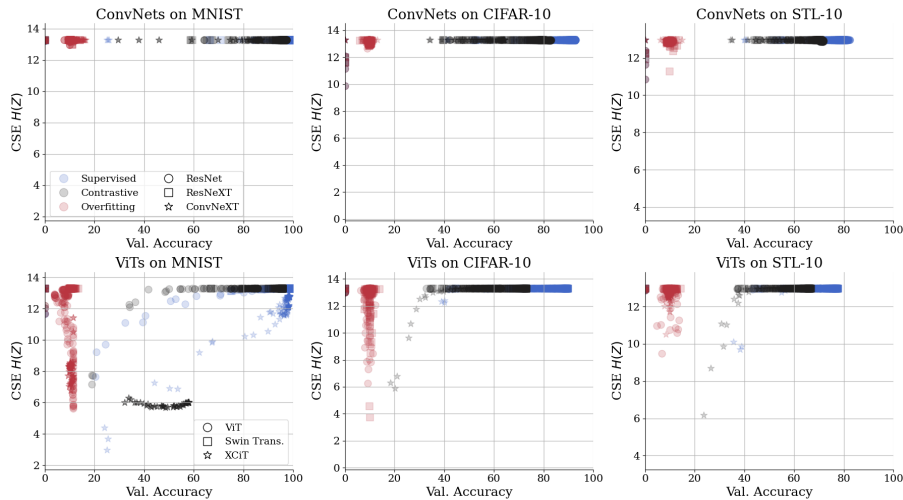


Figure S2: **Classic Shannon entropy version for Figure 6.** The embedding vectors are frequently allocated to unique buckets, which leads to the maximum possible entropy of $-\log_2 \frac{1}{10000} = 13.288$ for 10,000 data points.

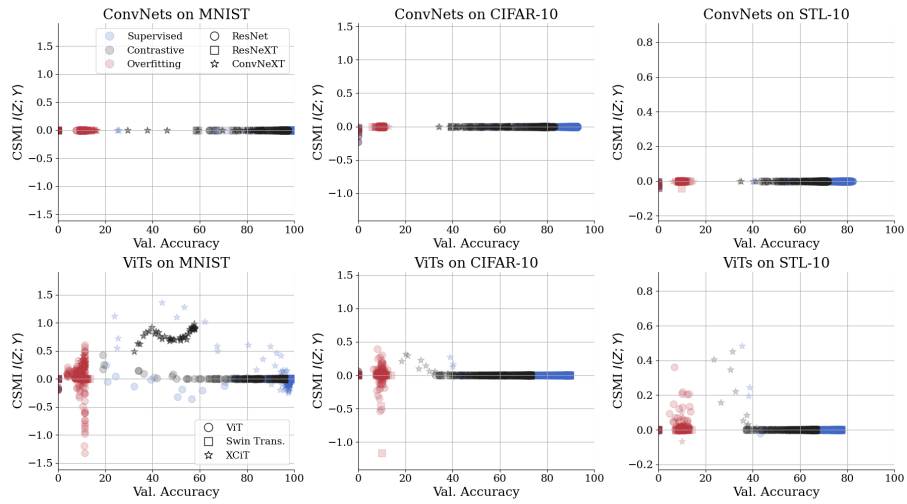


Figure S3: Classic Shannon mutual information version for Figure 7.

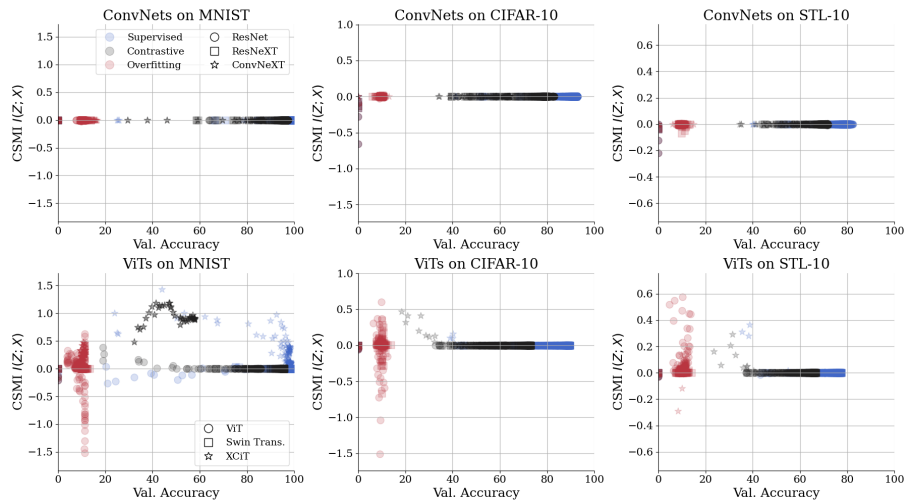


Figure S4: Classic Shannon mutual information version for Figure 8.

B Propositions with Proofs on DSE and DSMI

The following propositions establish some bounds on the minimal and maximal values of Diffusion Spectral Entropy (DSE). In addition, they provide insight on the definition of Diffusion Spectral Mutual Information (DSMI). Note that taking $t \rightarrow \infty$ allows us to talk about the major structures in the dataset.

Proposition 4.1 S_D achieves a 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

Proposition 4.2 As $t \rightarrow \infty$, $S_D(\mathbf{P}_X, t)$ on data with k well-separated, connected 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 $\sum_k \frac{1}{k} \log(k) = \log(k)$ completing the proof. \square

Lemma B.1. Let K be the anisotropic kernel matrix, where for data points x_i and x_j , $K_{ij} = \frac{G(x_i, x_j)}{\|G(x_i, \cdot)\|_2^\alpha \|G(x_j, \cdot)\|_2^\alpha}$, and $G(x_i, x_j) = \exp(-\|x_i - x_j\|_2^2 / \sigma)$. Then, the anisotropic kernel matrix K is a positive semidefinite matrix.

Proof. Note that $K = D^{-\alpha} G D^{-\alpha}$, where G is the Gaussian kernel matrix and $D_{ii} = \sum_{j=1}^n G_{ij}$ for all $i \in [n]$, and $D_{ij} = 0$ for $i \neq j$. It is well known that the Gaussian kernel matrix G is positive semidefinite, so $x^T G x \geq 0 \forall x \in \mathbb{R}^n \setminus \{0\}$. Let $\tilde{x} = D^{-\alpha} x \in \mathbb{R}^n \setminus \{0\}$, and $\tilde{x}^T K \tilde{x} = (x^T D^{-\alpha}) G (D^{-\alpha} x) \geq 0 \forall \tilde{x} \in \mathbb{R}^n \setminus \{0\}$ since \tilde{x} is a linear reparameterization of x . \square

Lemma B.2. Given n data points, let

$$\begin{aligned} \phi_m(\lambda_1, \dots, \lambda_m) = & -(n - m + 1) \frac{\lambda_m}{(n - m + 1)\lambda_m + \sum_{j=1}^{m-1} \lambda_j} \log \left(\frac{\lambda_m}{(n - m + 1)\lambda_m + \sum_{j=1}^{m-1} \lambda_j} \right) \\ & - \sum_{i=1}^{m-1} \frac{\lambda_i}{(n - m + 1)\lambda_m + \sum_{j=1}^{m-1} \lambda_j} \log \left(\frac{\lambda_i}{(n - m + 1)\lambda_m + \sum_{j=1}^{m-1} \lambda_j} \right) \end{aligned}$$

be a function of the eigenvalues of the anisotropic diffusion matrix K for $m \in \{2, \dots, n\}$, where the eigenvalues are ordered such that $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$. For all $m \in \{2, n-1\}$, $\phi_m(\lambda_1, \dots, \lambda_m) \geq \phi_{m+1}(\lambda_1, \dots, \lambda_m, \lambda_{m+1})$.

Proof. By Lemma B.1, K has all non-negative eigenvalues. Staying in convention with DSE, let $0 \log \left(\frac{0}{0} \right) = 0$.

First, we observe the following:

$$\phi_m(\lambda_1, \dots, \lambda_m) = \phi_{m+1}(\lambda_1, \dots, \lambda_m, \lambda_m)$$

By the Fundamental Theorem of Calculus,

$$\int_{\lambda_{m+1}}^{\lambda_m} \frac{\partial \phi_{m+1}}{\partial \lambda_{m+1}} d\lambda_{m+1} = \phi_m - \phi_{m+1}$$

Thus, it suffices to show that on $[\lambda_m, \lambda_{m+1}]$, $\int_{\lambda_{m+1}}^{\lambda_m} \frac{\partial \phi_{m+1}}{\partial \lambda_{m+1}} d\lambda_{m+1} \geq 0$. In pursuit of tidiness, let $\phi_{m+1} = f + g$, where $f(\lambda_1, \dots, \lambda_{m+1}) = (n - m) \frac{\lambda_{m+1}}{(n - m)\lambda_{m+1} + \sum_{j=1}^m \lambda_j} \log \left(\frac{\lambda_{m+1}}{(n - m)\lambda_{m+1} + \sum_{j=1}^m \lambda_j} \right)$, and g represents the summation term.

For f ,

$$\begin{aligned}
 \frac{\partial f}{\partial \lambda_{m+1}} &= -(n-m) \left(\log \left(\frac{\lambda_{m+1}}{(n-m)\lambda_{m+1} + \sum_{j=1}^m \lambda_j} \right) + 1 \right) \\
 &\quad \left[\frac{-(n-m)\lambda_{m+1}}{\left((n-m)\lambda_{m+1} + \sum_{j=1}^m \lambda_j \right)^2} + \frac{1}{(n-m)\lambda_{m+1} + \sum_{j=1}^m \lambda_j} \right] \\
 &= -(n-m) \left(\log \left(\frac{\lambda_{m+1}}{(n-m)\lambda_{m+1} + \sum_{j=1}^m \lambda_j} \right) + 1 \right) \left[\frac{-(n-m)\lambda_{m+1} + (n-m)\lambda_{m+1} + \sum_{j=1}^m \lambda_j}{\left((n-m)\lambda_{m+1} + \sum_{j=1}^m \lambda_j \right)^2} \right] \\
 &= -(n-m) \left(\log \left(\frac{\lambda_{m+1}}{(n-m)\lambda_{m+1} + \sum_{j=1}^m \lambda_j} \right) + 1 \right) \left[\frac{\sum_{j=1}^m \lambda_j}{\left((n-m)\lambda_{m+1} + \sum_{j=1}^m \lambda_j \right)^2} \right]
 \end{aligned}$$

For g ,

$$\begin{aligned}
 \frac{\partial g}{\partial \lambda_{m+1}} &= -\sum_{i=1}^m \left(\log \left(\frac{\lambda_i}{(n-m)\lambda_{m+1} + \sum_{j=1}^m \lambda_j} \right) + 1 \right) \left[\frac{-(n-m)\lambda_i}{\left((n-m)\lambda_{m+1} + \sum_{j=1}^m \lambda_j \right)^2} \right] \\
 &= \sum_{i=1}^m \left(\log \left(\frac{\lambda_i}{(n-m)\lambda_{m+1} + \sum_{j=1}^m \lambda_j} \right) + 1 \right) \left[\frac{(n-m)\lambda_i}{\left((n-m)\lambda_{m+1} + \sum_{j=1}^m \lambda_j \right)^2} \right]
 \end{aligned}$$

Together, letting $S = (n-m)\lambda_{m+1} + \sum_{j=1}^m \lambda_j$,

$$\begin{aligned}
 \frac{\partial \phi_{m+1}}{\partial \lambda_{m+1}} &= \frac{\partial g}{\partial \lambda_{m+1}} + \frac{\partial f}{\partial \lambda_{m+1}} \\
 &= \frac{1}{S^2} \left((n-m) \sum_{i=1}^m \lambda_i \log \left(\frac{\lambda_i}{S} \right) + (n-m) \sum_{i=1}^m \lambda_i - (n-m) \sum_{i=1}^m \lambda_i \log \left(\frac{\lambda_{m+1}}{S} \right) - (n-m) \sum_{i=1}^m \lambda_i \right) \\
 &= \frac{1}{S^2} \left((n-m) \sum_{i=1}^m \lambda_i \log \left(\frac{\lambda_i}{S} \right) - (n-m) \sum_{i=1}^m \lambda_i \log \left(\frac{\lambda_{m+1}}{S} \right) \right) \\
 &= \frac{n-m}{S^2} \sum_{i=1}^m \lambda_i \log \left(\frac{\lambda_i}{\lambda_{m+1}} \right)
 \end{aligned}$$

We conclude by observing that λ_{m+1} is the smallest eigenvector, so every log term must be non-negative. \square

Here, note that $\phi_n = S_D(\mathbf{P}_X, t = 1)$.

Lemma B.3. *Define ϕ as in Lemma B.2. Then, the function $\phi_2(\lambda_1, \lambda_2)$ is concave for the anisotropic diffusion matrix.*

Proof. Since $\lambda_1 \equiv 1$ by construction, it suffices to show that $\frac{\partial^2 \phi_2}{\partial \lambda_2^2}(1, \lambda_2) \leq 0$ for all $\lambda_2 \leq 1$.

$$\begin{aligned}
 \frac{\partial \phi_2}{\partial \lambda_2} &= - \left(\log \left(\frac{1}{1 + (n-1)\lambda_2} \right) + 1 \right) \frac{-1(n-1)}{(1 + (n-1)\lambda_2)^2} \\
 &\quad - (n-1) \left(\log \left(\frac{\lambda_2}{1 + (n-1)\lambda_2} \right) + 1 \right) \left(\frac{-\lambda_2(n-1)}{(1 + (n-1)\lambda_2)^2} + \frac{1}{1 + (n-1)\lambda_2} \right) \\
 &= \frac{(n-1)}{(1 + (n-1)\lambda_2)^2} \left(\log \left(\frac{1}{1 + (n-1)\lambda_2} \right) + 1 \right) \\
 &\quad + \frac{(n-1)}{(1 + (n-1)\lambda_2)^2} \left(\log \left(\frac{\lambda_2}{1 + (n-1)\lambda_2} \right) + 1 \right) (\lambda_2(n-1) - (1 + (n-1)\lambda_2)) \\
 &= \frac{(n-1)}{(1 + (n-1)\lambda_2)^2} \left(\log \left(\frac{1}{1 + (n-1)\lambda_2} \right) + 1 - \log \left(\frac{\lambda_2}{1 + (n-1)\lambda_2} \right) - 1 \right) \\
 &= \frac{(n-1)}{(1 + (n-1)\lambda_2)^2} \left(\log \left(\frac{1}{\lambda_2} \right) \right) \\
 &= - \frac{(n-1)}{(1 + (n-1)\lambda_2)^2} \log(\lambda_2)
 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial^2 \phi_2}{\partial \lambda_2^2} &= -(-2)(n-1) \frac{(n-1)}{(1 + (n-1)\lambda_2)^3} \log(\lambda_2) - \frac{(n-1)}{\lambda_2(1 + (n-1)\lambda_2)^2} \\
 &= 2(n-1) \frac{(n-1)}{(1 + (n-1)\lambda_2)^3} \log(\lambda_2) - \frac{(n-1)(1 + (n-1)\lambda_2)}{\lambda_2(1 + (n-1)\lambda_2)^3} \\
 &= \frac{n-1}{(1 + (n-1)\lambda_2)^3} \left(2(n-1) \log(\lambda_2) - \frac{(1 + (n-1)\lambda_2)}{\lambda_2} \right) \\
 &= \frac{n-1}{(1 + (n-1)\lambda_2)^3} \left(2(n-1) \log(\lambda_2) - \left(\frac{1}{\lambda_2} + (n-1) \right) \right)
 \end{aligned}$$

Concavity requires $\frac{1}{\lambda_2} + (n-1) \geq 2(n-1) \log(\lambda_2)$. Here, we observe that $\lambda_2 \leq 1$, and $\frac{1}{\lambda_2} + (n-1) > 0 \geq 2(n-1) \log(\lambda_2)$ completes the proof. \square

Proposition 4.3 Let $X \in \mathbb{R}^{n \times d}$ be a dataset of n independent and identically distributed multivariate Gaussian vectors in \mathbb{R}^d , where $x_i \sim \mathcal{N}(0, I_d)$. Then, using K as defined in Eqn 1 with $\alpha = 1/2$,

$$\begin{aligned}
 &\mathbb{E}[S_D(\mathbf{P}_X, t = 1)] \\
 &\lesssim \log\left(\frac{n}{1-\beta}\right) - \left(\frac{1}{n} + \left(\frac{n-1}{n}\right)\beta\right) \log\left(1 + \frac{\beta n}{1-\beta}\right) \\
 &\quad \text{where } \beta = \left(1 + \frac{4}{\sigma}\right)^{-\frac{d}{2}}
 \end{aligned} \tag{8}$$

Proof. We prove the above theorem in four steps:

1. Expectation is linear.
2. One can approximate the expected random anisotropic kernel matrix.
3. Based on this, one can derive the expected eigenvalues of the expected matrix.
4. Use Jensen's Inequality with the bounded DSE.

First, eigenvalues λ are defined such that for matrix $K \in \mathbb{R}^{n \times n}$, $Kx = \lambda x$ for all vectors $x \in \mathbb{R}^n$. By the linearity of expectation, $\mathbb{E}[K]x = \mathbb{E}[\lambda]x$. Thus, computing $\mathbb{E}[K]$ leads to a solution for $\mathbb{E}[\lambda]$, which can be used to upper bound the DSE.

Second, $\mathbb{E}[K] = \mathbb{E}[D^{-1/2}GD^{-1/2}]$. Here, we approximate the expectation by $\mathbb{E}[K] \approx (D | \mathbb{E}[G])^{-1/2} \mathbb{E}[G] (D | \mathbb{E}[G])^{-1/2}$, as confirmed empirically. Computations based on $\mathbb{E}[G]$ are much more approachable, given the nature of the problem.

For diagonal entries of G , $\mathbb{E}[G_{ii}] = \mathbb{E}[\exp(-\|x_i - x_i\|_2^2/\sigma)] = 1$. For off-diagonal entries of G , we proceed inductively. First, $x_i - x_j \sim \mathcal{N}(0, 2I_d)$. Let $z_{ij} \equiv \|x_i - x_j\|_2^2$. Since $\frac{z_{ij}}{2} \sim \chi_d^2$, one can verify that the chi-squared distribution moment generating function is $\mathbb{E}[e^{tX}] = (1 - 2t)^{-d/2}$ for $t < 1/2$. Thus, $\mathbb{E}[\exp(-\|x_i - x_j\|_2^2/\sigma)] = \mathbb{E}[e^{(-1/\sigma)z_{ij}}] = (1 + \frac{1}{\sigma})^{-d/2}$.

For $\mathbb{E}[K]$,

$$\begin{aligned} \mathbb{E}[K_{ii}] &\approx \frac{1}{1 + (n-1)(1 + \frac{1}{\sigma})^{-d/2}} \\ \mathbb{E}[K_{ij}] &\approx \frac{(1 + \frac{1}{\sigma})^{-d/2}}{1 + (n-1)(1 + \frac{1}{\sigma})^{-d/2}} \quad \forall i \neq j \end{aligned}$$

Third, observe that our approximate $\mathbb{E}[K]$ may be decomposed as follows:

$$\mathbb{E}[K] \approx \mathbb{1}\mathbb{1}^T \left(\frac{(1 + \frac{1}{\sigma})^{-d/2}}{1 + (n-1)(1 + \frac{1}{\sigma})^{-d/2}} \right) + I_n \left(\frac{1 - (1 + \frac{1}{\sigma})^{-d/2}}{1 + (n-1)(1 + \frac{1}{\sigma})^{-d/2}} \right)$$

One easily verifies that the first matrix has eigenvalues $\lambda_1 = n \left(\frac{(1 + \frac{1}{\sigma})^{-d/2}}{1 + (n-1)(1 + \frac{1}{\sigma})^{-d/2}} \right)$ and $\lambda_2 = \dots = \lambda_n = 0$, and the second matrix has eigenvalues $\lambda_1 = \dots = \lambda_n = \frac{1 - (1 + \frac{1}{\sigma})^{-d/2}}{1 + (n-1)(1 + \frac{1}{\sigma})^{-d/2}}$. By construction, the eigenvalues are additive, so $\mathbb{E}[K]$ has $\mathbb{E}[\lambda_1] = 1$, and $\mathbb{E}[\lambda_2] = \dots = \mathbb{E}[\lambda_n] = \frac{1 - (1 + \frac{1}{\sigma})^{-d/2}}{1 + (n-1)(1 + \frac{1}{\sigma})^{-d/2}}$.

Fourth, we refer to Lemma B.2 in claiming that $\phi_2(\lambda_1, \lambda_2) \geq \phi_n(\lambda_1, \lambda_2, \dots, \lambda_n) = S_D(\mathbf{P}_X, t = 1)$. By the concavity established in Lemma B.3, Jensen's Inequality implies that $\mathbb{E}[\phi_2(\lambda_1, \lambda_2)] \leq \phi_2(\mathbb{E}[\lambda_1], \mathbb{E}[\lambda_2])$. Simply put,

$$\mathbb{E}[S_D(\mathbf{P}_X, t = 1)] \leq \phi_2(\mathbb{E}[\lambda_1], \mathbb{E}[\lambda_2])$$

Finally, letting $\beta = (1 + \frac{4}{\sigma})^{-d/2}$,

$$\begin{aligned}
 \phi_2(\mathbb{E}[\lambda_1], \mathbb{E}[\lambda_2]) &= -\frac{(n-1)\mathbb{E}[\lambda_2]}{(n-1)\mathbb{E}[\lambda_2] + \mathbb{E}[\lambda_1]} \log\left(\frac{\mathbb{E}[\lambda_2]}{(n-1)\mathbb{E}[\lambda_2] + \mathbb{E}[\lambda_1]}\right) \\
 &\quad - \frac{\mathbb{E}[\lambda_1]}{(n-1)\mathbb{E}[\lambda_2] + \mathbb{E}[\lambda_1]} \log\left(\frac{\mathbb{E}[\lambda_1]}{(n-1)\mathbb{E}[\lambda_2] + \mathbb{E}[\lambda_1]}\right) \\
 &= -\frac{(n-1)\frac{1-\beta}{1+(n-1)\beta}}{(n-1)\frac{1-\beta}{1+(n-1)\beta} + 1} \log\left(\frac{\frac{1-\beta}{1+(n-1)\beta}}{(n-1)\frac{1-\beta}{1+(n-1)\beta} + 1}\right) \\
 &\quad - \frac{1}{(n-1)\frac{1-\beta}{1+(n-1)\beta} + 1} \log\left(\frac{1}{(n-1)\frac{1-\beta}{1+(n-1)\beta} + 1}\right) \\
 &= -(n-1)\frac{(1-\beta)}{(n-1)(1-\beta) + 1 + (n-1)\beta} \log\left(\frac{(1-\beta)}{(n-1)(1-\beta) + 1 + (n-1)\beta}\right) \\
 &\quad - \frac{1 + (n-1)\beta}{(n-1)(1-\beta) + 1 + (n-1)\beta} \log\left(\frac{1 + (n-1)\beta}{(n-1)(1-\beta) + 1 + (n-1)\beta}\right) \\
 &= -(n-1)\frac{(1-\beta)}{(n-1) + 1} \log\left(\frac{(1-\beta)}{(n-1) + 1}\right) - \frac{1 + (n-1)\beta}{(n-1) + 1} \log\left(\frac{1 + (n-1)\beta}{(n-1) + 1}\right) \\
 &= -(n-1)\frac{(1-\beta)}{n} \log\left(\frac{(1-\beta)}{n}\right) - \frac{1 + (n-1)\beta}{n} \log\left(\frac{1 + (n-1)\beta}{n}\right) \\
 &= \frac{n-1-n\beta+\beta}{n} \log\left(\frac{n}{1-\beta}\right) + \frac{1+n\beta-\beta}{n} \log\left(\frac{n}{1+(n-1)\beta}\right) \\
 &= \frac{n}{n} \log\left(\frac{n}{1-\beta}\right) + \frac{1+n\beta-\beta}{n} \left(\log\left(\frac{n}{1+(n-1)\beta}\right) - \log\left(\frac{n}{1-\beta}\right)\right) \\
 &= \log\left(\frac{n}{1-\beta}\right) + \frac{1+n\beta-\beta}{n} \left(\log\left(\frac{1-\beta}{1+(n-1)\beta}\right)\right) \\
 &= \log\left(\frac{n}{1-\beta}\right) + \left(\frac{1}{n} + \frac{n-1}{n}\beta\right) \log\left(\frac{1-\beta}{1+(n-1)\beta}\right) \\
 &= \log\left(\frac{n}{1-\beta}\right) - \left(\frac{1}{n} + \frac{n-1}{n}\beta\right) \log\left(\frac{1-\beta+n\beta}{1-\beta}\right) \\
 &= \log\left(\frac{n}{1-\beta}\right) - \left(\frac{1}{n} + \frac{n-1}{n}\beta\right) \log\left(1 + \frac{n\beta}{1-\beta}\right)
 \end{aligned}$$

□

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. □

Lemma B.4. Let f be a probability distribution, and let $X \in \mathbb{R}^{kn \times d}$ be a dataset of $n \in \mathbb{N}$ independent and identically distributed d -dimensional points $x_i \sim f$. Let these points lie in $k \in \mathbb{N}$ distinct clusters, where the diffusion probability between points of different clusters is arbitrarily small. Let $Y \in \mathbb{R}^{n \times d}$ be a dataset of n i.i.d. points $y_i \sim f$ in one cluster. Then, the expected DSE $\mathbb{E}[S_D(\mathbf{P}_X, t)]$ may be related to the expected DSE $\mathbb{E}[S_D(\mathbf{P}_Y, t)]$ by the following equality:

$$\mathbb{E}[S_D(\mathbf{P}_X, t)] = \mathbb{E}[S_D(\mathbf{P}_Y, t)] + \log(k)$$

Proof. Let $K \in \mathbb{R}^{n \times n}$ be the anisotropic kernel matrix for Y , and let $\lambda_1, \dots, \lambda_n$ denote its ordered eigenvalues. Let $\tilde{K} \in \mathbb{R}^{kn \times kn}$ be the anisotropic kernel matrix for X , where $\tilde{\lambda}_1, \dots, \tilde{\lambda}_{kn}$ are its ordered eigenvalues. Note

that as a block diagonal matrix, $\tilde{\lambda}_{ki+1} = \dots \tilde{\lambda}_{ki+k}$ in law for $i \in [0, n-1]$. Since each block is composed of n data points distributed as f , $\lambda_{i+1} = \lambda_{ki+1}$ in law. Thus, we can reformulate $\mathbb{E}[S_D(\mathbf{P}_X, t)]$ as follows:

$$\begin{aligned}
 \mathbb{E}[S_D(\mathbf{P}_X, t)] &= \mathbb{E} \left[- \sum_{i=1}^{kn} \left(\frac{|\tilde{\lambda}_i^t|}{\sum_{i=1}^{kn} |\tilde{\lambda}_j^t|} \right) \log \left(\frac{|\tilde{\lambda}_i^t|}{\sum_{i=1}^{kn} |\tilde{\lambda}_j^t|} \right) \right] \\
 &= \mathbb{E} \left[-k \cdot \sum_{i=1}^n \left(\frac{|\lambda_i^t|}{k \cdot \sum_{i=1}^n |\lambda_j^t|} \right) \log \left(\frac{|\lambda_i^t|}{k \cdot \sum_{i=1}^n |\lambda_j^t|} \right) \right] \\
 &= \mathbb{E} \left[-\frac{k}{k} \cdot \sum_{i=1}^n \left(\frac{|\lambda_i^t|}{\sum_{i=1}^n |\lambda_j^t|} \right) \log \left(\frac{|\lambda_i^t|}{k \cdot \sum_{i=1}^n |\lambda_j^t|} \right) \right] \\
 &= \mathbb{E} \left[- \sum_{i=1}^n \left(\frac{|\lambda_i^t|}{\sum_{i=1}^n |\lambda_j^t|} \right) \left(\log \left(\frac{|\lambda_i^t|}{\sum_{i=1}^n |\lambda_j^t|} \right) - \log(k) \right) \right] \\
 &= \mathbb{E} \left[- \sum_{i=1}^n \left(\frac{|\lambda_i^t|}{\sum_{i=1}^n |\lambda_j^t|} \right) \log \left(\frac{|\lambda_i^t|}{\sum_{i=1}^n |\lambda_j^t|} \right) \right] - \log(k) \cdot \mathbb{E} \left[- \sum_{i=1}^n \left(\frac{|\lambda_i^t|}{\sum_{i=1}^n |\lambda_j^t|} \right) \right] \\
 &= \mathbb{E}[S_D(\mathbf{P}_Y, t)] + \log(k)
 \end{aligned}$$

□

Proposition 4.4 Take n to be arbitrarily large. Let $X \in \mathbb{R}^{n \times d}$ be a matrix of i.i.d. random values $x_{ij} \sim f$. Let $Y \in \mathbb{R}^{n \times d}$ be a matrix of i.i.d. random values $y_{ij} \sim f$, but in $k \in \mathbb{N}$ distinct clusters such that when the anisotropic probability matrix is computed for $\alpha = 1/2$, the probability of diffusion between points of different clusters is arbitrarily small. Then, using β as defined in Proposition 4.3, the approximate upper bound on DSE increases by $\beta \log(k)$.

Proof. We take n to be large such that the upper bound in Prop. 4.3 approaches $\log\left(\frac{n}{1-\beta}\right) - \beta \log\left(\frac{\beta n}{1-\beta}\right)$, for $\mathbb{E}[S_D(\mathbf{P}_X, t=1)]$. Let $Z \in \mathbb{R}^{n/k \times d}$ also contain i.i.d. random entries according to f . For $\mathbb{E}[S_D(\mathbf{P}_Z, t=1)]$, the aforementioned upper bound evaluates to $\log\left(\frac{n}{1-\beta}\right) - \beta \log\left(\frac{\beta n}{1-\beta}\right) - (1-\beta) \log(k)$ following simple manipulations.

By Lemma B.4, $\mathbb{E}[S_D(\mathbf{P}_Y, t=1)] = \mathbb{E}[S_D(\mathbf{P}_Z, t=1)] + \log(k)$.

$$\begin{aligned}
 \mathbb{E}[S_D(\mathbf{P}_Y, t=1)] &= \mathbb{E}[S_D(\mathbf{P}_Z, t=1)] + \log(k) \\
 &\leq \log\left(\frac{n}{1-\beta}\right) - \beta \log\left(\frac{\beta n}{1-\beta}\right) - (1-\beta) \log(k) + \log(k) \\
 &\leq \log\left(\frac{n}{1-\beta}\right) - \beta \log\left(\frac{\beta n}{1-\beta}\right) + \beta \log(k)
 \end{aligned}$$

□

Proposition 4.5 As $t \rightarrow \infty$ in a hidden layer X with k well-separated clusters and class labels Y perfectly indicating clusters, we will have $I_D(X; Y) = \log(k)$.

Proof. Based on proposition 4.2, 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 $I_D(X; Y) = \log(k)$. □

C Full Caption of the Toy Data Figure

Here we show Figure 3 with the full caption.

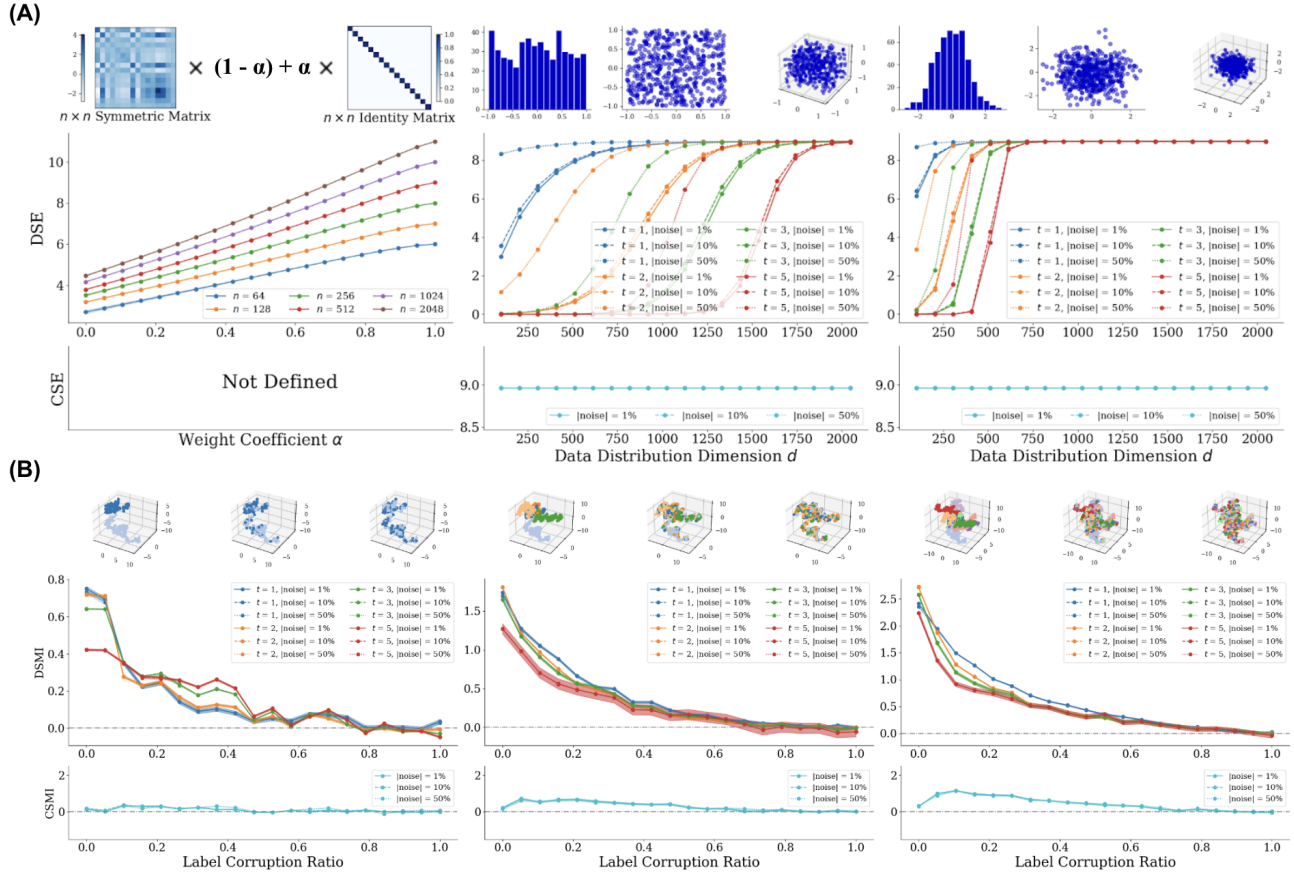


Figure S5: **Diffusion Spectral Entropy (DSE) and Diffusion Spectral Mutual Information (DSMI) on toy data.**

(A) DSE increases as intrinsic dimension grows, while classic Shannon entropy (CSE) saturates to $\log(n)$ due to curse of dimensionality. **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 the standard deviation from 5 independent runs. For the latter two distributions, additive noise is injected into the coordinates and schematics for $d = \{1, 2, 3\}$ are illustrated on top. The number of data points for the simulation is 500. Note that in the latter two cases for DSE we compute the diffusion matrix of the data manifold prior to entropy evaluation, whereas in the first case we skip that step because the matrix is already provided. In the latter two cases, CSE saturates to $\log(n) = \log_2(500) = 8.966$.

(B) When two random variables are dependent, DSMI negatively correlates with the level of data corruption, while classic Shannon mutual information (CSMI) does not capture this trend. DSMI $I_D(Z; Y)$ and CSMI are computed on synthetic, 20-dimensional trees with $\{2, 5, 10\}$ branches (**Left, Mid, Right**). The x-axis represents the level of label corruption, ranging from completely clean label (left) to fully corrupted label (right), with 3 dimensional schematics demonstrating corruption ratios $\{0.0, 0.5, 1.0\}$ displayed on top. At full corruption, $I_D(Z; Y)$ converges to zero, as the embedding vectors do not contain information on the labels.

D Experimental Details

D.1 Three conditions for neural network training process

All experiments were run on a single NVIDIA A100 GPU.

D.1.1 Supervised learning

We trained the 6 networks (3 ConvNets and 3 vision transformers as mentioned in Section 5.3) end-to-end. The ConvNets included ResNet [He et al., 2016], ResNeXT [Xie et al., 2017], and ConvNeXT [Liu et al., 2022]. The vision transformers included the original Vision Transformer (ViT) [Dosovitskiy et al., 2020], Swin Transformer [Liu et al., 2021], and Cross-Covariance Image Transformer (XCiT) [Ali et al., 2021].

We used an AdamW optimizer [Loshchilov and Hutter, 2017] with an initial learning rate of 1e-5 for MNIST [Lecun, 1998] or 1e-4 for CIFAR-10 [Krizhevsky et al., 2009] and STL-10 [Coates et al., 2011]. The learning rate is modulated by a Cosine Annealing Scheduler [Loshchilov and Hutter, 2016] with a linear warm-up for the first 10 epochs. The total duration of training is 200 epochs. At the end of each epoch, we pass the entire validation set through the model and collect the \mathbb{R}^D representation vectors as output from the penultimate layer for further analysis.

D.1.2 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 ask the model to embed them closer on the embedding manifold, while encouraging greater 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 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 perform linear probing by the end of each epoch. Specifically, we freeze the backbone weights, attach a reinitialized linear classifier layer, and train the classifier for 10 epochs with the training set. Then, we record the end-to-end validation accuracy and 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 in the supervised learning case, including the learning rate and scheduling.

D.1.3 Purposeful overfitting

In purposeful overfitting, we train the model in the same way as in 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 in mitigating overfitting. For fair comparison, the training details are otherwise the same as in the supervised learning case, including learning rate and scheduling.

D.2 Computing DSMI with input

We compute $I_D(Z; X)$ the same way as we compute $I_D(Z; Y)$. By a simple change of variables, Eqn 7 can be rewritten as:

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

Compared to DSMI between the neural representation and the output, the DSMI with input is slightly more complicated because the input signals do not fall into discrete categories, i.e. x_i are not naturally defined. To this end, for the set of n input images, we flatten them and perform spectral clustering. For a 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 output.

D.3 Details on Imagenet cross-model correlation

D.3.1 Pre-trained models

The pre-trained models are obtained from the Huggingface PyTorch Image Models repository [Wightman, 2019]. Specifically, we used the models listed in <https://github.com/huggingface/pytorch-image-models/blob/main/results/results-imagenet.csv>, which the owners shared under Apache License 2.0. There are originally 1002 models, but 40 models failed our standardized pipeline, likely due to the presence of auxiliary structures, leading to a final set of 962 models. All of these models are pre-trained on ImageNet [Deng et al., 2009].

D.3.2 ImageNet accuracy

The ImageNet validation accuracy (top-1 and top-5) is taken from <https://github.com/huggingface/pytorch-image-models/blob/main/results/results-imagenet.csv>. Test accuracy (top-1 and top-5) is taken from <https://github.com/huggingface/pytorch-image-models/blob/main/results/results-imagenet-real.csv>.

D.3.3 DSE and DSMI

The full ImageNet validation set contains 1000 classes each with 50 samples. For efficient computation, we measured DSE $S_D(Z)$, DSMI $I_D(Z; X)$, and $I_D(Z; Y)$ on the union of two commonly used subsets of ImageNet validation set: namely Imagenette and Imagewoof. The former contains 10 very distinct classes (tench, English springer, cassette player, chain saw, church, French horn, garbage truck, gas pump, golf ball, parachute), while the latter contains 10 very similar classes (Australian terrier, Border terrier, Samoyed, Beagle, Shih-Tzu, English foxhound, Rhodesian ridgeback, Dingo, Golden retriever, Old English sheepdog). We believe that the union of these two subsets forms a representative subset of ImageNet.

D.4 Logarithm base

It can be verified that the choice of a 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 All Raw Results

Here we show the raw results corresponding to the following figures: Figures 6, 7, 8, S2, S3, S4.

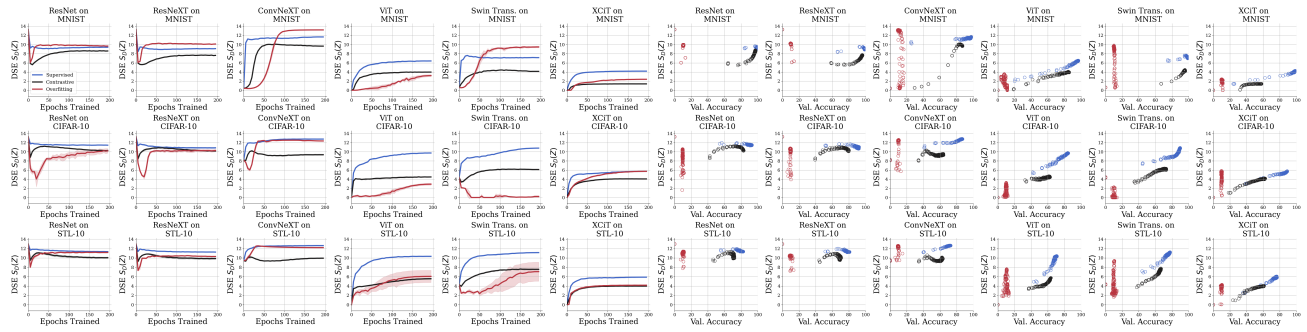


Figure S6: Raw results for Figure 6.

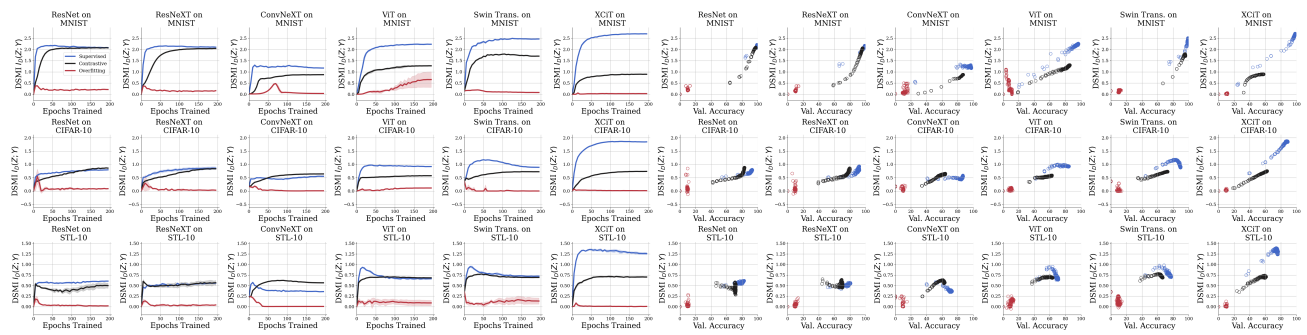


Figure S7: Raw results for Figure 7.

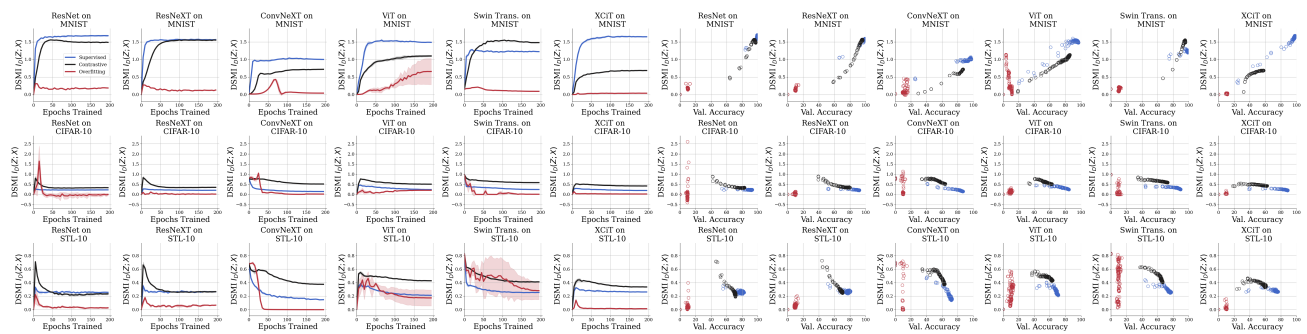


Figure S8: Raw results for Figure 8.

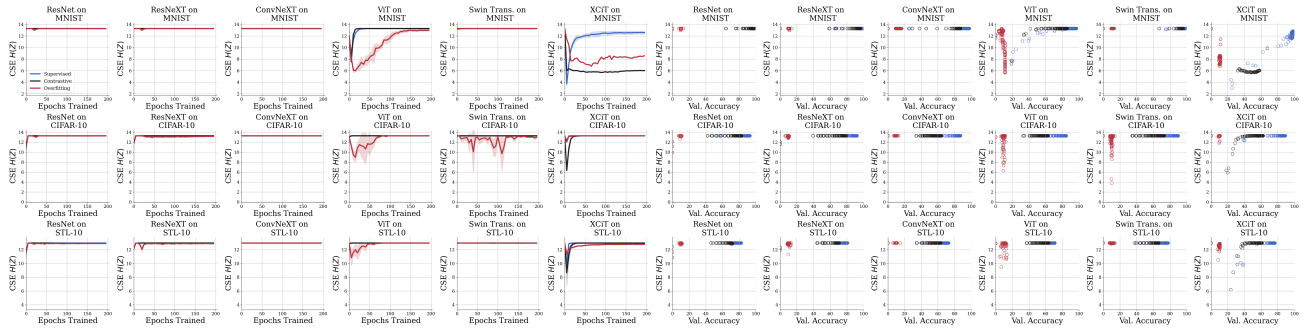


Figure S9: Raw results for Figure S2.

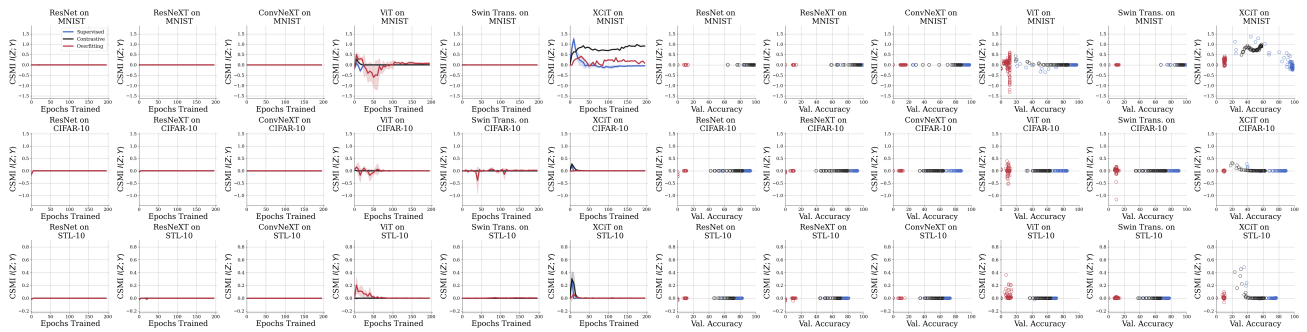


Figure S10: Raw results for Figure S3.

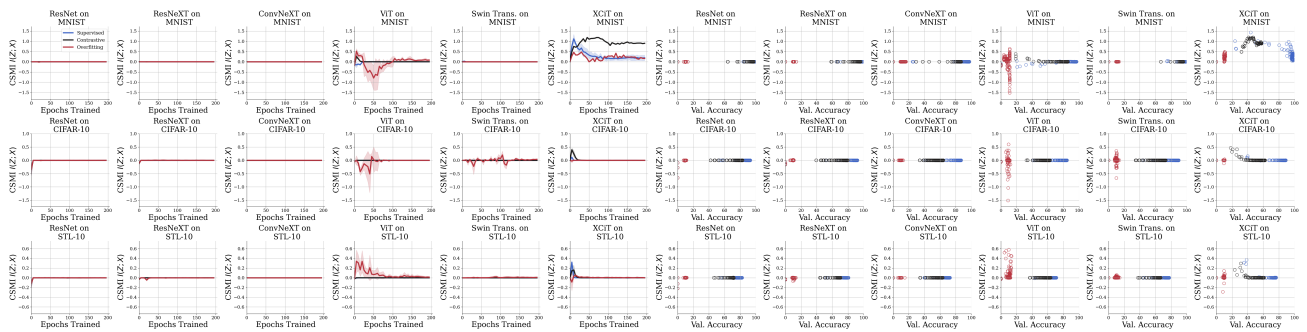


Figure S11: Raw results for Figure S4.

F DSE Initialization Experiments

For all experiments on the evolution of DSE and DSMI during neural network training, we initialized the model parameters using the same standard initialization process.

For convolutional layers (i.e. Conv and ConvTranspose), we used Kaiming normal initialization [He et al., 2015]. For linear layers (a.k.a. fully-connected layers), we initialized them with a normal distribution with mean 0 and standard deviation $1e-3$. For batch normalization weights, we initialized them with a constant value of 1.0. For all biases, we initialized them with a constant value of 0.

Despite the fact that we used the same initialization, the initial value of DSE is not always the same across different experiments. For example, in Figure S6, ResNet on MNIST began at a high DSE, while ViT on MNIST began at a low DSE. While this does not change the general increasing trend of DSE during training, we wonder: **Will initializing the network at a high DSE versus a low DSE affect the learning process?**

After further investigation, we found that we can control this behavior by modifying the initialization parameters on the convolutional layers. Specifically, we replaced the initialization method for the convolutional layers with a normal distribution with a mean of 0 and a tunable standard deviation. We observed a monotonic trend in the initial DSE as we adjusted the standard deviation (Figure S12). However, as the transition is abrupt, it is not easy to target an arbitrary DSE value.

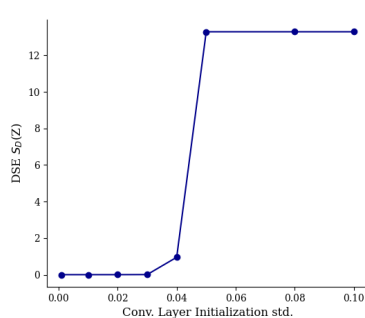


Figure S12: Initialization parameters affects the initial DSE (CIFAR-10).

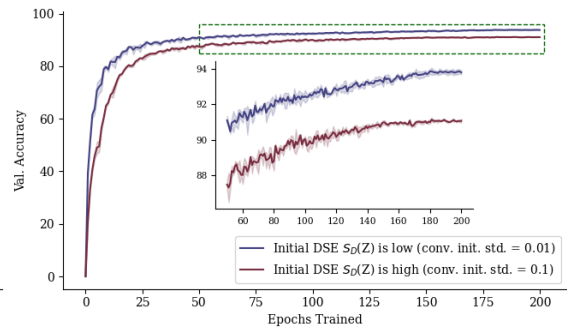
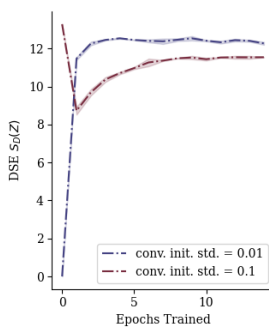


Figure S13: Initializing the network at a low DSE allows for faster convergence and better final performance (CIFAR-10).

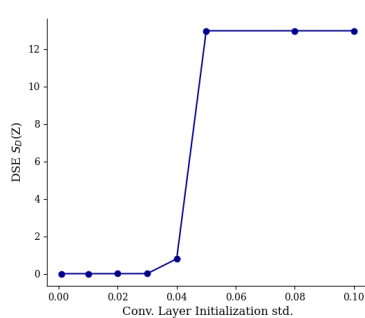


Figure S14: Initialization parameters affects the initial DSE (STL-10).

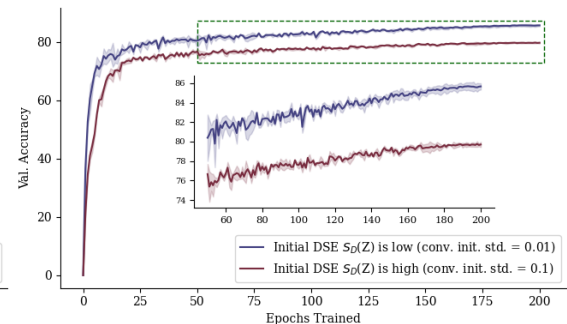
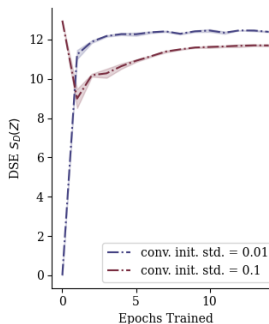


Figure S15: Initializing the network at a low DSE allows for faster convergence and better final performance (STL-10).

To answer the **question in bold** above, we trained two ResNets on the CIFAR-10 dataset, one starting with a low DSE (conv. init. std. = 0.01) and the other starting with a high DSE (conv. init. std. = 0.1). Experiments are repeated under 3 random seeds. Intuitively, we would expect the former to converge faster than the latter, because in the former case the embedding manifold only needs to monotonically increase DSE (Figure S13 left,

blue line), while in the latter case the DSE needs to first decrease from the peak and then increase (Figure S13 left, red line). Indeed, this intuition is supported by empirical results. In the right panel of Figure S13, it can be seen that, at least in our preliminary experiments, **initializing the network at a low DSE allows faster convergence and better final performance**.

We conducted the same experiments on STL-10 and observed the same results. The corresponding figures are Figure S14 and Figure S15.

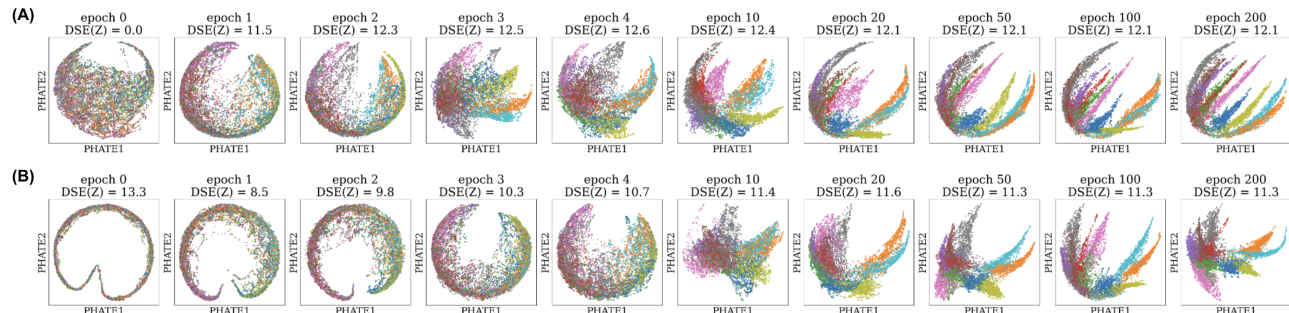


Figure S16: PHATE representation of the embedding spaces during training for low (panel A) and high (panel B) initial DSE. Colors represent ground truth class labels.

This observation can be further corroborated by the visualizations of the embedding space. When we initialize the network at a low DSE (Figure S16(A)), the PHATE [Moon et al., 2019] representation of the embedding space is close to a uniform blob. As the network undergoes more training, the embedding space becomes organized into distinct branches, each corresponding to a different class. This process is relatively smooth and monotonic.

On the other hand, when we initialize the network at a high DSE (Figure S16(B)), the PHATE representation almost looks like a curved ribbon. As training progresses, the embedding space first expanded to become more blob-like and then stretches out into branches. Compared to the other case above, which only has a branching phase, this process has an additional expansion phase before the branching phase, which may explain the slower convergence. Moreover, the final representations are not as clean and well-separated as the other case.

All of these results together suggest that different network initializations may affect the convergence rate and final performance, and that some properties of network initialization can be captured by DSE. However, since this is not the main focus of this paper, we leave it to future research to explore this topic in more depth.

G Additional DSE Results: Intuition and MAGIC visualization

G.1 DSE vs. Diffusion t in the Intuition Figure

Figure S17 illustrates DSE vs. diffusion time t for the cases in the intuition figure (Figure 1). It can be seen that DSE approaches $\log_2(3) \approx 1.585$ with 3 nicely separated blobs, while it approaches zero with one single blob.

G.2 MAGIC Visualization for DSE vs. Diffusion t

To emphasize that DSE reflects the structure of the data manifold instead of data variations in the ambient space, we employed MAGIC [Van Dijk et al., 2018] for visual representation. This visualization technique showcases the diffusion geometric representation of the data manifold across various diffusion t values. As evident in Figure S18, with an increase in diffusion t , a dataset characterized by three distinct blobs evolves into three separate ‘lines’. Conversely, a dataset featuring a singular blob transforms into a single ‘line.’ This visual representation using MAGIC underscores the inherent structures within the data manifold and it is evident that DSE aligns well with these visual interpretations.

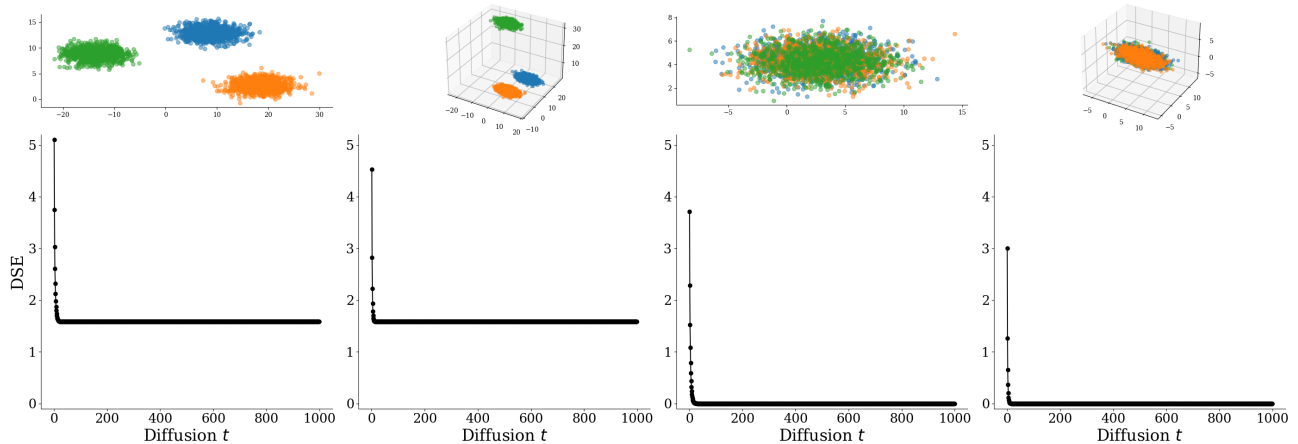


Figure S17: **DSE vs. Diffusion t for Figure 1.** **Left:** X are k separated clusters. Simulated with 3 blobs in 2 or 3 dimensions. $S_D(\mathbf{P}_X, t) \approx \log(k)$ when t increases **Right:** X is a single blob. Simulated with 1 blob in 2 or 3 dimensions. $S_D(\mathbf{P}_X, t) \approx 0$ with increasing t .

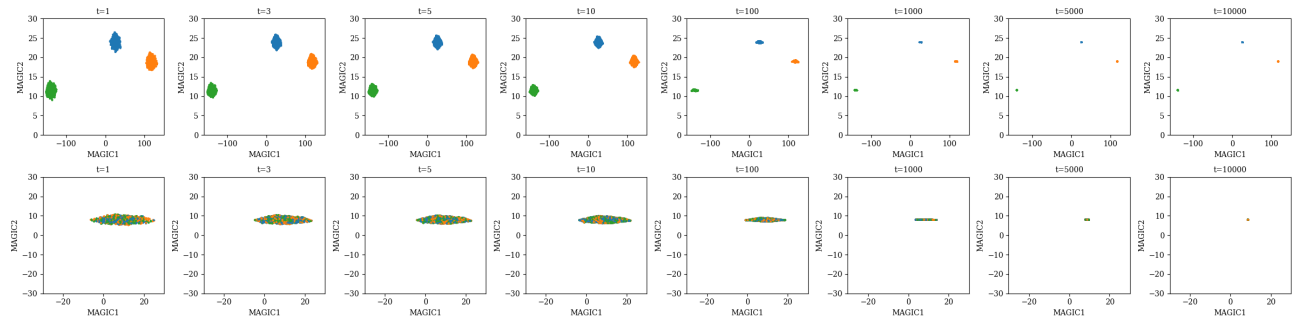


Figure S18: **MAGIC representation of simulated blobs data in 3 dimension.** **Top panel:** X are k separated clusters. **Bottom panel:** X is a single blob. As diffusion t increases, a dataset characterized by three distinct blobs evolves into three separate “lines.” Conversely, a dataset featuring a singular blob transforms into a single “line”.

G.3 Ablation Study

We conducted an ablation study on diffusion spectral entropy by experimenting with alternative methods for computing entropy. In addition to diffusion spectral entropy (DSE) and classic Shannon entropy (CSE), we investigate three alternatives: diffusion matrix entry entropy (DMEE) in which the classic Shannon entropy is computed on the entries in the diffusion matrix \mathbf{P} ; entropy with eigenvalues from the binary adjacency matrix of k -nearest-neighbor graph; entropy with eigenvalues from the weighted adjacency matrix \mathcal{G} constructed with the Gaussian kernel.

Table S1: Descriptions of methods being compared.

Method	Description
DSE	diffusion spectral entropy: proposed method (Eqn 6)
CSE	classic Shannon entropy (Eqn 3)
DMEE	diffusion matrix entry entropy (Shannon entropy over $\{\mathbf{P}_{i,j}\}$)
KNN	with eigenvalues from binary adjacency matrix of k -nearest-neighbor graph
Gaussian	with eigenvalues from the weighted adjacency matrix \mathcal{G} , where $\mathcal{G}(z_1, z_2) = e^{-\frac{\ z_1 - z_2\ ^2}{\sigma}}$

We visualized the behavior of different entropy methods against the intrinsic dimension of the data. We can see

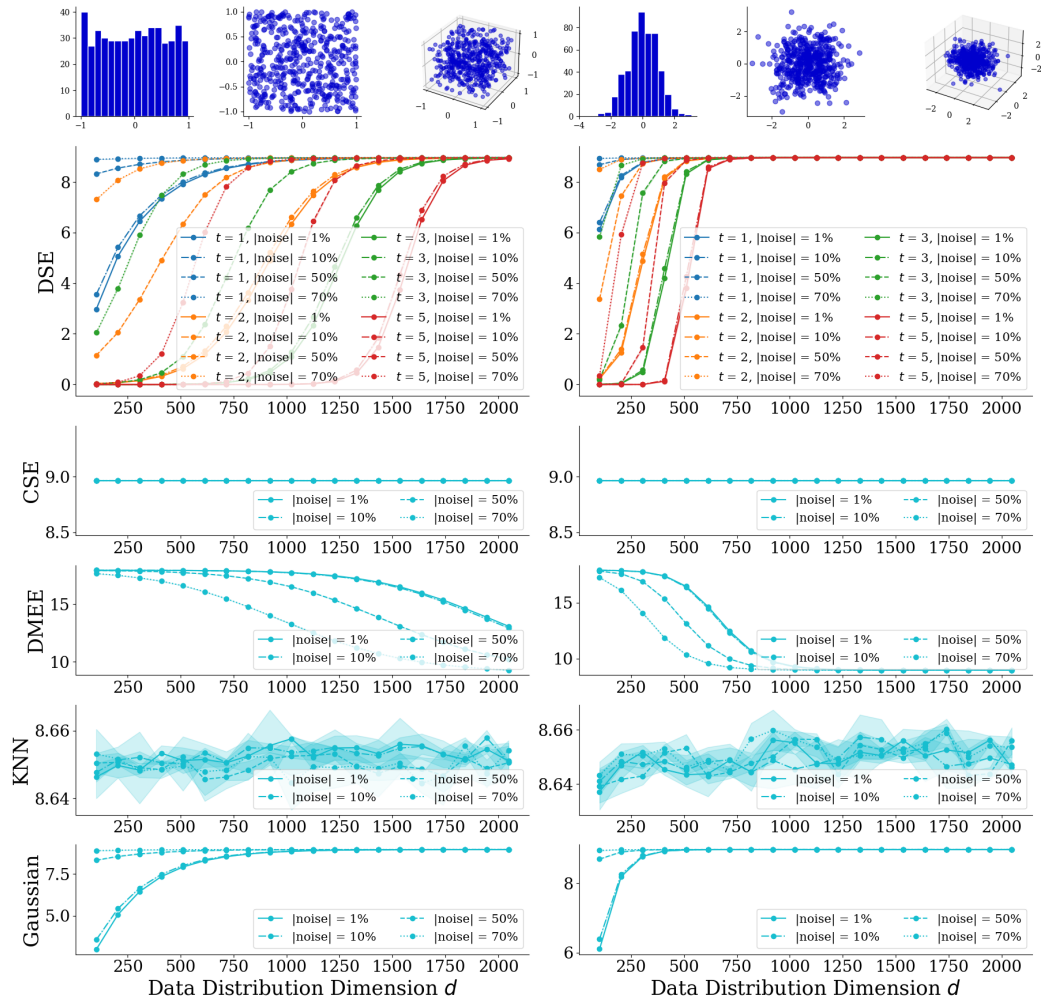


Figure S19: Ablation study on diffusion spectral entropy (DSE).

from Figure S19 that diffusion matrix entry entropy (DMEE) decreases as the intrinsic dimension of the data increases. The entropy computed on eigenvalues from the binary adjacency matrix of the k -nearest-neighbor graph (KNN method) jitters around a fixed value as the intrinsic data dimension increases. The entropy on the eigenvalues of the weighted adjacency matrix \mathcal{G} starts with a high entropy in a relatively small data dimension and continues to increase until it becomes stagnant. It plateaus in much smaller dimensions on both uniform and Gaussian distributed data. On data with a high percentage of noise (50%, 70%), we observe that entropy on eigenvalues of Gaussian adjacency matrix barely changes despite increasing data dimension, while diffusion spectral entropy (DSE) can remain resilient to noise and effectively reflect the increasing intrinsic dimension as the diffusion parameter t is adjusted.

H Additional DSMI Results: Intuition and Subsampling Robustness

H.1 DSMI vs. Diffusion t in the Intuition Figure

Figure S20 illustrates DSMI for the cases in the intuition figure (Figure 1). It can be seen that DSMI is positive on nicely separated clusters, while close to zero on well-mixed clusters. Moreover, when t is sufficiently high, DSMI approaches $\log(k)$, which in this case is $\log_2(3) \approx 1.585$.

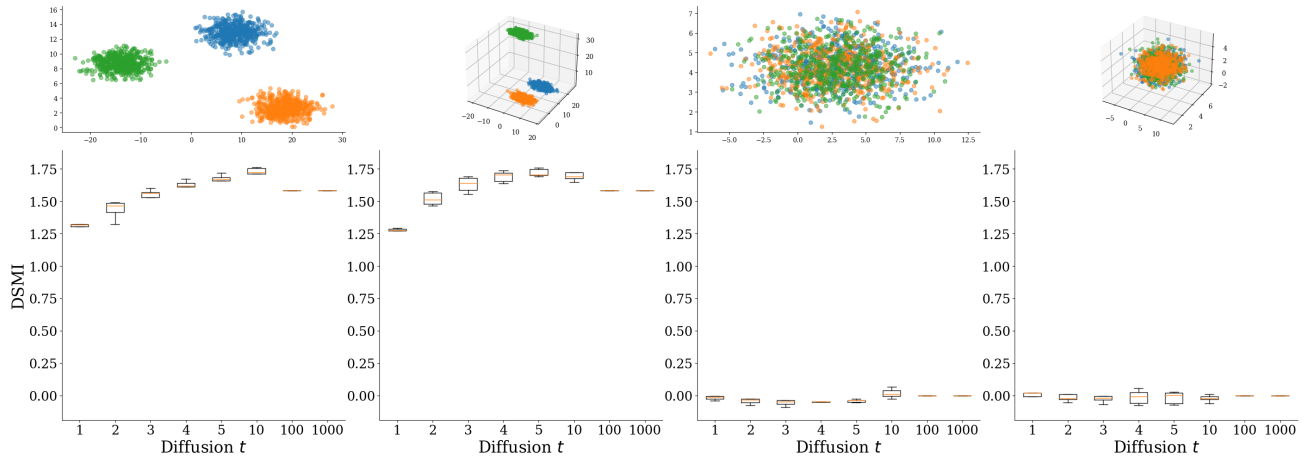


Figure S20: **DSMI simulations for Figure 1.** **Left:** Y is on separated clusters. Simulated with 3 blobs in 2 or 3 dimensions. $S_D(X|Y) < S_D(X)$ and $I_D(X; Y)$ are positive. **Right:** Y is close to a uniform sub-sampling of X . Simulated with 1 blob in 2 or 3 dimensions. $I_D(X; Y)$ is around 0, but can have negative values if $S_D(X|Y)$ is larger than $S_D(X)$ for numeric reasons.

H.2 Subsampling Technique for DSMI: Subsampling Robustness

As mentioned in **Definition 4.2**, we use a subsampling technique to compute *Diffusion Spectral Mutual Information (DSMI)*.

Since DSMI is computed as a conditional DSE, we hereby justify this technique by showing that Diffusion Spectral Entropy (DSE) is robust to subsampling (Figure S21). Four line styles of the same color, respectively, represent different subsampling ratio ranging from 100% (full dataset) to 10%. It can be observed that DSE can be reliably computed despite significant subsampling.

On the other hand, we also need to point out that the subsampling robustness is constrained by the DSE upper bound. DSE is 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 underestimate. In our studies, the subsampling ratio is about 10%, which always yields a 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.

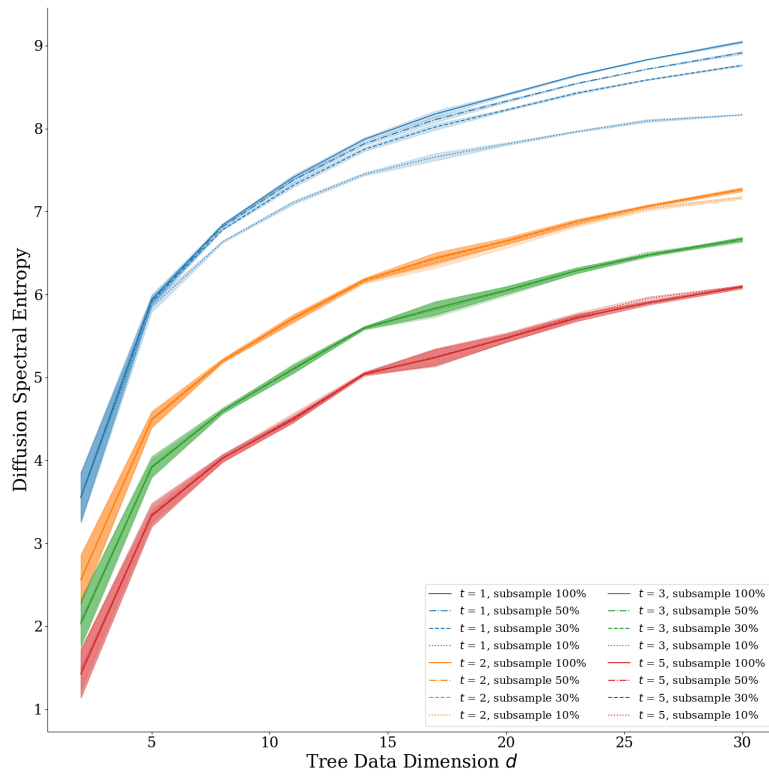


Figure S21: **Diffusion spectral entropy estimation is robust to subsampling.** Subsampled $S_D(Z)$ are computed on synthetic multi-dimensional trees.

Supplementary References

- [Ali et al., 2021] Ali, A., Touvron, H., Caron, M., Bojanowski, P., Douze, M., Joulin, A., Laptev, I., Neverova, N., Synnaeve, G., Verbeek, J., et al. (2021). Xcit: Cross-covariance image transformers. *Advances in neural information processing systems*, 34:20014–20027.
- [Chen et al., 2020] Chen, T., Kornblith, S., Norouzi, M., and Hinton, G. (2020). A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR.
- [Coates et al., 2011] Coates, A., Ng, A., and Lee, H. (2011). An analysis of single-layer networks in unsupervised feature learning. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pages 215–223. JMLR Workshop and Conference Proceedings.
- [Deng et al., 2009] Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. (2009). Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee.
- [Dosovitskiy et al., 2020] Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al. (2020). An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*.
- [He et al., 2015] He, K., Zhang, X., Ren, S., and Sun, J. (2015). Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, pages 1026–1034.
- [He et al., 2016] He, K., Zhang, X., Ren, S., and Sun, J. (2016). Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778.
- [Krizhevsky et al., 2009] Krizhevsky, A., Hinton, G., et al. (2009). Learning multiple layers of features from tiny images.
- [LeCun, 1998] LeCun, Y. (1998). The mnist database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>.
- [Liu et al., 2021] Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., and Guo, B. (2021). Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10012–10022.
- [Liu et al., 2022] Liu, Z., Mao, H., Wu, C.-Y., Feichtenhofer, C., Darrell, T., and Xie, S. (2022). A convnet for the 2020s. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11976–11986.
- [Loshchilov and Hutter, 2016] Loshchilov, I. and Hutter, F. (2016). Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983*.
- [Loshchilov and Hutter, 2017] Loshchilov, I. and Hutter, F. (2017). Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*.
- [Moon et al., 2019] 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. (2019). Visualizing structure and transitions in high-dimensional biological data. *Nature biotechnology*, 37(12):1482–1492.
- [Tishby and Zaslavsky, 2015] Tishby, N. and Zaslavsky, N. (2015). Deep learning and the information bottleneck principle. In *2015 IEEE information theory workshop (itw)*, pages 1–5. IEEE.
- [Wightman, 2019] Wightman, R. (2019). Pytorch image models. <https://github.com/rwightman/pytorch-image-models>.
- [Xie et al., 2017] Xie, S., Girshick, R., Dollár, P., Tu, Z., and He, K. (2017). Aggregated residual transformations for deep neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1492–1500.