# Information Flow in Deep Neural Networks

A thesis submitted for the degree of

Doctor of Philosophy

by

Ravid Shwartz-Ziv

Submitted to the senate of the Hebrew University

June 2021

This work was carried out under the supervision of
Professor Naftali Tishby.

# Acknowledgments

First and foremost, I am extremely grateful to my supervisor and mentor, Professor Naftali Tishby, who passed away recently. His invaluable advice, continuous support, and patience during my Ph.D. study were so important. His immense knowledge and plentiful experience have encouraged me in all the time of my academic research and daily life. He was a remarkable man who gave me so much, and I learned a lot from him. An incredible scholar and a lovely person.

I would also like to thank Prof. Haim Sompolinsky and Dr. Alex Alemi for their support and guidance during my study. I wish to thank also my collaborators, Zoe Piran and Dr. Amichai Painsky. It was a great privilege to work with them.

Finally, I cannot begin to express my gratitude to my wife and children. Without their tremendous understanding and encouragement in the past few years, it would be impossible for me to complete my study. There were always there to support me and encourage me to succeed.

# Abstract

While deep neural networks have been immensely successful, a comprehensive theoretical understanding of how they work or how they are structured does not exist. Deep networks are often viewed as black boxes, where the interpretation of predictions and their reliability are still unclear. Today, understanding the groundbreaking performance of deep neural networks is one of the greatest challenges facing the scientific community. To use these algorithms more effectively and improve them, we need to understand their dynamic behavior and their ability to learn new representations.

This thesis addresses these issues by applying principles and techniques from information theory to deep learning models to increase our theoretical understanding and use it to design better algorithms. The main results and contributions of this thesis are structured in three parts, as detailed below.

Chapters 2 and 3 present our information-theoretic approach to deep learning models. As an explanation for deep learning systems, we propose using the Information Bottleneck (IB) theory. The novel paradigm for analyzing networks sheds light on their layered structure, generalization capabilities, and learning dynamics. Based on our analysis, we find that deep networks optimize each layer's mutual information on input and output variables, resulting in a trade-off between compression and prediction for each layer. Our analytical and numerical study of these networks demonstrated that the stochastic gradient descent (SGD) algorithm follows the IB trade-off principle by working in two phases: a fast empirical error minimization phase followed by a slow representation compression phase. These phases are distinguished by different signal-to-noise ratios (SNRs) for each layer. Moreover, we demonstrated that the SGD achieved this optimal bound due to the compression phase, derived a new Gaussian bound on the representation compression, and related it to the compression time. Furthermore, our results indicate that the network's layers converge to the IB theoretical bound, leading to a self-consistent relationship between the encoder and decoder distributions.

Chapter 4 deals with one of the most difficult problems of applying the IB to deep neural networks — estimating mutual information in high dimnesional space. Despite being an important quantity in data science, mutual information has historically posed

a computational challenge. Computing mutual information is only tractable for discrete variables or for a limited number of problems where probability distributions are known. To better estimate information-theoretic quantities and to investigate generalization signals, we research several frameworks and utilize recent theoretical developments, such as the neural tangent kernel (NTK) framework. In our study, we found that for infinite ensembles of infinitely wide neural networks, we could obtain tractable computations of many information-theoretic quantities and their bounds. Many quantities can be described in a closed-form solution by the network's kernels. By analyzing these derivations, we can learn the important information-theoretic quantities of the network and how compression, generalization, and the sample size are related.

Chapter 5 presents the dual Information Bottleneck (dualIB), a new information-theoretic framework. Despite the IB framework's advantages, it also has several drawbacks: The IB is completely non-parametric and operates only on the probability space. In addition, the IB formulation does not relate to the task of prediction over unseen patterns and assumes full access to the joint probability. Therefore, we developed the dualIB, which resolves some of the IB's drawbacks through a mere switch between terms in the distortion function. The dualIB can account for known features of the data and use them to make better predictions over unseen examples. We provide dualIB self-consistent equations, allowing us to obtain analytical solutions. A local stability analysis revealed the underlying structure of the critical points of the solutions, resulting in a full bifurcation diagram of the optimal pattern representations. We discovered several interesting properties of dualIB's objective. First, the dualIB retains its structure when expressed in a parametric form. It also optimizes the mean prediction error exponent, thereby improving prediction accuracy with respect to sample size. In addition to dualIB's analytic solutions, we provided a variational dualIB framework that optimizes the functional using deep neural networks. The framework enables a practical implementation of dualIB for real-world datasets. With it, we empirically evaluated its dynamics and validated the theoretical predictions in modern deep neural networks.

In conclusion, this thesis proposes a new information-theoretic perspective for studying deep neural networks that draws upon the correspondence between deep learning and the IB framework. Our unique perspective can provide a number of benefits, such as attaining a deeper understanding of deep neural networks, explaining their behavior, and improving their performance. At the same time, our study opens up new theoretical and practical research questions.

# Letter of Contribution

This dissertation includes four manuscripts that summarize Ravid Shwartz-Ziv's research under the supervision of Naftali Tishby. The work on this dissertation was also done in collaboration with Dr. Alex Alemi (Google), Dr. Amichai Painsky (Tel Aviv University) and Zoe Piran (The Hebrew University). They provided valuable guidance and contributed to the writing of the papers as co-authors. The lead author and primary contributor for three of these manuscripts is Ravid Shwartz-Ziv. Ravid Shwartz-Ziv and Zoe Piran each contributed equally to the manuscript "The Dual Information Bottleneck." For this manuscript, Ravid Shwartz-Ziv did the variational dual IB section and the numerical experiments. The manuscripts are listed below. One was published in a peer-reviewed venue, and three are yet to be published.

1. **"Opening the Black Box of Deep Neural Networks via Information"**.
   Co-author: Naftali Tishby. 2018.

2. **"Representation Compression and Generalization in Deep Neural Networks"**.
   Co-authors: Amichai Painsky and Naftali Tishby. 2019.

3. **"Information in Infinite Ensembles of Infinitely-Wide Neural Network"**.
   Co-author: Alex Alemi. Published in the Proceedings of the Symposium on Advances in Approximate Bayesian Inference, PMLR, 2020.

4. **"The Dual Information Bottleneck"**.
   Co-authors: Zoe Piran and Naftali Tishby. 2020.

# Table of Contents

# Introduction

Learning representations is at the core of many problems in computer vision, natural language processing, cognitive science, and machine learning (Bengio et al.; 2013). A complex data representation is required for classification and prediction since physical parameters, such as location, size, orientation, and intensity, are considered (Salakhutdinov et al.; 2013). However, it is unclear what constitutes a good representation and how it is related to learning and to the specific problem type.

By combining multiple transformations of simple neurons, deep neural networks (DNNs) can produce a more useful (and, in most cases, more abstract) representation. Due to their versatility and success across various domains, these systems have gained popularity over the past few years. The performance of DNNs demonstrates great improvements over conventional machine learning methods in various domains, including images, audio, and text (Devlin et al.; 2018; He et al.; 2016; Oord et al.; 2016). The latest deep learning models are more complex, and their architectures are becoming more complex with more parameters that need to be optimized. The ResNet-52 network, for example, contains about 23 million parameters optimized over millions of images.

However, the reasons for these performances are only partially understood from a theoretical perspective, and we only have a heuristic understanding of them. It is unclear why deep models perform so well on real-world data and what their key components are. In addition, current metrics do not provide insight into the internal structure of a network or the quality of its layers. As a result, even if the model is extremely accurate, it is difficult to use it as a basis for further scientific research. To make these algorithms more effective and improve them, we must understand their underlying dynamic behavior and how they learn representations.

In this thesis, we propose studying DNNs from the perspective of information theory. As an explanation for modern deep learning systems, we propose the information bottleneck (IB) theory. We hope to shed light on their layered structures, generalization capabilities, and learning dynamics through this innovative approach for analyzing DNNs.

To better understand DNNs, the first question is as follows: How can information theory in general and the IB framework, in particular, be used to better understand DNNs?

Shannon invented the information theory to determine the number of bits needed to transmit a message over a noisy channel. This theory has since been shown to be an invaluable measure of the influence between variables (Shannon; 1948). Given two random variables $X$ and $Y$, the mutual information between them measures the divergence of their joint probability distribution $P(x, y)$ from the product of their marginals $P(x)P(y)$ to determine how dependent or independent they are. The notion of mutual information, unlike correlation, can capture nonlinear statistical relationships between variables, strengthening our ability to analyze complex system dynamics (Kinney and Atwal; 2014). Although mutual information is an essential quantity in data science, it has historically been challenging to estimate (Paninski; 2003). Exact computations are tractable only for a limited number of problems with well-defined probability distributions (e.g., the exponential family). The calculation of mutual information is not possible for finite samples of data or general problems.

This leads us to the following research question: How can we calculate the mutual information for large-scale DNNs? To derive an exact calculation of information-theoretic quantities and to search for generalization signals, we examined several frameworks and utilized current theoretical developments, including the Neural Tangent Kernel (NTK) framework (Lee et al.; 2019). We obtained tractable calculations of information-theoretic quantities and their bounds for infinite ensembles of infinitely-wide neural networks. Our analysis revealed that the kernels described many quantities in a closed form. Furthermore, we found that the input's compression contributed to the generalization in this model family.

Although the IB framework has its advantages, it also has a few disadvantages, including an inability to preserve the structure of the data and suboptimal performance with finite data. The final research question concerns whether we can derive a new framework that can solve these problems and apply it to DNNs.

Therefore, we developed the dual IB (dualIB), which switches between the terms in the distortion function to solve some of the problems with the IB. A local stability analysis revealed the underlying structure and optimal pattern representations. We discovered that the dualIB retains its structure when expressed in a parametric form. Furthermore, it optimizes the mean prediction error exponent, improving the predictions' accuracy with respect to sample size. dualIB can be applied to real-world datasets using neural networks with the help of a variational framework. Using this framework, we evaluated the dynamics of the dualIB and validated the theoretical predictions.

First, we examine some of the information-theoretic principles upon which this thesis is based. The chapter concludes by highlighting the main contributions of the thesis.

## 1.1  Information theory

In 1948, Shannon's paper established the area of information theory (Shannon; 1948). Communication, which refers to sending information with the receiver's intent, is one of the main topics in information theory. Shannon's work served as the basis for quantifying the questions regarding this information. Let us review the results of information theory relevant to this thesis. In the present work, we use the following notation: Upper-case letters denote random variables (e.g., $X$), calligraphic letters denote their support (e.g., $\mathcal{X}$), and lower-case letters denote specific realizations (e.g., $x$).

### 1.1.1  Theoretical information quantities

As a formal concept, information is a function of several basic measures based solely on the probability distribution of random variables without any particular assumptions.

Consider $X \in \mathcal{X}$ and $Y \in \mathcal{Y}$ as random variables with the joint distribution function of $p(x, y)$. Our first definition of information is entropy, which takes into account the amount of uncertainty involved with a given distribution.

**Definition 1.1.1. Entropy**

The entropy of $X \sim p(x)$ is defined as

$$H(X) = - \sum_{x \in \mathcal{X}} p(x) \log(p(x))$$

Shannon axiomatically derived this definition by defining three intuitive properties of uncertainty measures (continuity, additivity, and monotonicity). Entropy is the only function that satisfies these requirements. Intuitively, entropy reflects the minimum description length of $X$ because it is (roughly) the minimal number of bits or binary questions needed to determine the exact value of $X$.

If the distribution is uniform, then the entropy reaches its maximum. If it is deterministic, then it is zero. As $H(X)$ is concave, the more $X$ strays away from being completely unpredictable, the more it is punished by $H(X)$.

Our next informational measure is the Kullback-Leibler (KL) divergence, also known as relative entropy, which measures the divergence between two distributions.

**Definition 1.1.2. Kullback Leibler Divergence**

The KL divergence between two distributions $p(x)$ and $q(x)$ is defined as

$$D[p||q] = \sum_{x \in \mathcal{X}} p(x) \log \left( \frac{p(x)}{q(x)} \right)$$

In the case where $p(x) = q(x)$ for some $x$, then it will not contribute to $D[p||q]$. If for all $x$, $p(x) = q(x)$, then $D(p||q) = 0$. When $p$ and $q$ have a low Kullback-Leibler distance, they are similar. Conversely, a high distance indicates they are dissimilar.

It can be shown that $D[p|q] \geq 0$ and that equality holds if and only if $p = q$ almost everywhere. Note that the KL divergence is not a metric because it is not symmetric and does not obey the triangle inequality. However, it is still useful to think of it as the natural "distance" between distributions for several reasons. First, notice that $D[p||q]$ is the expected log-likelihood ratio between $p(x)$ and $q(x)$. Thus, it controls the discriminability between these two distributions when $p$ is the true underlying distribution of $X$. Secondly, the KL divergence reflects the difference between the minimal description length of $X \sim p(x)$ and the description length if $q(x)$ is used instead of $p(x)$. To see this, notice that $D[p|q] = \sum_x p(x) log \frac{1}{q(x)} - H(X)$, and recall that $H(X)$ is the minimal description length. Third, Pinsker's inequality implies that the KL divergence upper bounds the $L1$ distance between $p$ and $q$. Therefore, it also bounds any $L_\rho$ distance for $\rho \geq 1$.

The last piece of information we will present is mutual information.

### Definition 1.1.3. Mutual Information

Let $(X, Y) \sim p(x, y)$, and let $p(x)$ and $p(y)$ be their marginal distributions, respectively. The mutual information between $X$ and $Y$ is defined as

$$I(X;Y) = \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} p(x, y) \log \frac{p(x, y)}{p(x)p(y)}$$

Mutual information is a useful method for measuring statistical dependence between variables and has many beneficial properties. For example, it remains invariant under bijective transformations. Its predecessor, called the transmission rate, was first introduced by Shannon in 1948 (Shannon; 1948) for a communication system. Mutual information can be axiomatically derived as a function satisfying several natural 'informativeness' conditions (Cover; 1999).

Mutual information could be understood by thinking of the $KL$ divergence between $p(x, y)$ and the hypothetical joint distribution if $X$ and $Y$ were independent.

$$I(X;Y) = D[p(x, y)||p(x)p(y)]$$

It thus follows that $I(X;Y) \geq 0$ and $I(X;Y) = 0$ if and only if $X$ and $Y$ are independent. Alternatively, mutual information can be interpreted as reducing uncertainty regarding one variable due to the knowledge of the other variable.

$$I(X;Y) = H(X) - H(X|Y) = H(Y) - H(Y|X)$$

Since $H(X|Y), H(Y|X) \geq 0$, it holds that $I(X;Y) \leq \min\{H(X), H(Y)\}$. If $H(X|Y) =$

0, that is, $X$ is completely known given $Y$, then $I(X;Y) = H(X)$. As a special case, it holds that $I(X;X) = H(X)$. In this sense, mutual information is more general than entropy, whereas entropy is sometimes referred to as self-information.

One of the important properties of mutual information is the data processing inequality (DPI), which implies that the information about $X$ contained in $Y$ cannot be increased by processing $Y$. Formally, we say that $(X, Y, Z)$ form a Markov chain if their joint distribution can be decomposed as $p(x, y, z) = p(x, y)p(z|y)$, $Z$ is a function of $Y$ and given $Y$, it is independent of $X$. Then, the DPI implies

$$I(X;Y) \geq I(X;Z) \tag{1.1}$$

for any 3 variables that form a Markov chain $X \to Y \to Z$.

Mutual information is also invariant with respect to invertible transformations:

$$I(X;Y) = I(\psi(X); \phi(Y))) \tag{1.2}$$

for any invertible functions $\phi$ and $\psi$.


## 1.2  Representation learning

Machine learning is a field of Artificial Intelligence (AI) that automatically learns and improves from experience. These models use an internal representation of the input based on prior observations or data records to make better decisions. Learning representations lie at the core of many computer vision problems, natural language processing, cognitive science, and machine learning (Bengio et al.; 2013). Even so, the question of what constitutes a good representation or the relationship between representation, learning process, and specific features of a problem remains unanswered.

David Marr (Marr; 2010) defined representation as a formal system for making explicit particular entities and types of information that an algorithm can use to process information. In this thesis, we focus on the second part – the idea that algorithms operate on representations. Representation learning refers to learning structures in the data. This learning makes it easier to extract useful information to classify and make predictions based on the data (Goodfellow et al.; 2016). The question of what a good representation is has many answers. For probabilistic models, a good representation often captures the posterior distribution of the underlying explanatory factors for the observed input (LeCun et al.; 2015). In their AI-tasks, Bengio and LeCun (Bengio and LeCun; 2007) introduced the concept of complex but highly structured dependencies. In most cases, the tasks require transforming a high-dimensional input structure into a low-dimensional output or learning low-level representations. Accordingly, most of the input's entropy is not relevant for the

output, and it is difficult to extract the input's relevant features (Tishby and Zaslavsky; 2015).

### 1.2.1 Large deviation theory and the Sanov theorem

Statisticians, information theorists, and machine learning scientists share concepts for defining and extracting the optimal representation from observations. How would we characterize the optimal representation of the random input variable $X$? In general, we would like to have a representation that achieves high values for some reward function. In addition, we aim to learn efficiently these representations by using an empirical sample of the (unknown) joint distribution, which will enable us to generalize to unknown datasets.

As we will see in the following section, the large deviation theory can be used to illustrate the idea that constraining information with the representation will produce a representation that is more likely to generalize. This approach is derived from the type analysis and the Sanov theorem (Cover; 1999). It is presented here for an independent and identically distributed (i.i.d.) process, although a Markov type analysis is also available (Csiszár et al.; 1987).

With the large deviation theory, we can understand the probability of rare events, i.e., events with an exponentially small probability. Type analysis is applied to analyze the most likely samples from a given probability distribution following a constraint. Sanov's theorem identifies the rate function for large deviations of the empirical measure of a sequence of i.i.d. random variables. The type definition and the Sanov Theorem for i.i.d. processes (Cover; 1999) are presented here with an adaptation of the annotation to our needs.

Let $T_1, T_2, \ldots$ be i.i.d. random variables with values in a finite set $A = \{\alpha_1, \ldots, \alpha_r\}$ and with a distribution $P$. Denote by $\mathcal{P}$ the set of probabilities on $\{\alpha_1, \ldots, \alpha_r\}$. Let $\hat{P}_n$ be the empirical distribution of $T$:

$$\hat{P}_n = \frac{1}{n} \sum_{k=1}^{n} \delta_{T_k}.$$

The law of large numbers states that $\hat{P}_n \to P$ almost surely. To define a rare event, we fix $E \subseteq \mathcal{P}$ that does not contain $P$. We are interested in the behavior of probabilities of the form $\mathcal{P}[\hat{P}_n \in E]$, as $n \to \infty$.

It is a trivial observation that each possible value $\{\alpha_1, \ldots, \alpha_r\}$ must appear an integer number of times among the samples $\{T_1, \ldots, T_n\}$. This implies, however, that the empirical measure $\hat{P}_n$ cannot take arbitrary values.

**Definition 1.2.1.** Type $P_t$ of the sequence $t = t_1, t_2, \ldots, t_n$ is the relative proportion of occurrences of each value. That is, $P_t(\alpha) = \frac{N(\alpha|t)}{n}$ where $N(\alpha|x)$ is the number of times $\alpha$ occurs in the sequence t.

By definition, each type contains only information about how often each value shows up in the sample, discarding the order in which they appear. The importance of it is that large numbers of sequences are associated with each type, and each type is associated with a single reward value $\sum_{\alpha \in A} P(\alpha)R(\alpha) = \bar{R}$. It is always the case that $\hat{P}_n \in \mathcal{P}_n$, where we define $\mathcal{P}_n$ as the set of types associated with sequences of samples with length $n$.

**Theorem 1.2.1.** *Sanov Let $t_1, t_2, \ldots t_{n-1}$ be i.i.d. sequence from $P(t)$. Let $E \subseteq \mathcal{P}_n$ be a set of types that is equal to the closure of its interior.*

$$\lim_{n \to \infty} \frac{1}{n} \log P^n[\hat{P}_n \in E] = -D_{KL}(P^\star|P^n)$$

*where $P^\star = arg \min_{Q \in E \cap \mathcal{P}_n} D_{KL}(Q|P^n)$.*

For our needs, we define the set $E$ to be all the types (and the sequences that are associated with them) that achieve a reward above the desired threshold $\theta$; i.e., $E = \sum \{P_{\alpha \in A} P(\alpha)R(\alpha) \geq \bar{R}\}$.

Sanov's theorem gives us an important insight: the most likely sequence to be selected as a type that is the closest from the $D_{KL}$ perspective to the distribution $P(t)$ (denoted as $P^\star$). This insight can be formalized as a trade-off, where on one side, we have the reward $\bar{R}$, and on the other, we have $D_{KL}$. The selected solution is governed by a parameter $\beta$, which is the Lagrange multiplier coefficient. Given the nature of the types, they can be visualized and represented as points of a simplex whose axes are the probabilities of each symbol to appear. In this space, the two crucial points are the maximal reward point (a deterministic point at the edge of the simplex) and $P$. The trade-off parameter $\beta$ defines a set of distributions that form a line between these points. These are known as geodesic lines. No two points along this line have the same $\beta$, and therefore they do not have the same $D_{KL}$ value.

*1.2.2  Minimal sufficient statistic*

An alternative definition of what constitutes a good representation is based on minimal sufficient statistics.

**Definition 1.2.2.** Let $(X, Y) \sim P(X, Y)$ . Let $T := t(X)$ , where $t$ is a deterministic function. We call T a sufficient statistic of $X$ for $Y$ if $Y - T - X$ forms a Markov chain.

Therefore, a sufficient statistic captures all the information about $Y$ that is available in $X$. The following theorem states this property:

**Theorem 1.2.2.** *Cover (1999) Let $T$ be a probabilistic function of $X$. Then, $T$ is a sufficient statistic for $Y$ if and only if (iff ) $I(T(X); Y) = I(X; Y)$*

7

As we can see, the sufficiency definition includes the trivial identity statistic $T = X$. Such statistics accomplish nothing since all they do is "copy" rather than "extract" important information. Therefore, it is necessary to prevent statistics from using observations in an inefficient manner.

To address this issue, the concept of minimal sufficient statistics was introduced:

**Definition 1.2.3.** (Minimal sufficient statistic (MSS)) A sufficient statistic $T$ is minimal if for any other sufficient statistic $S$, there exists a function $f$ such that $T = f(S)$ almost surely (a.s.).

MMS are the simplest sufficient statistics and induce the coarsest sufficient partition on $X$. MSS try to group the values of $X$ into as few partitions as possible without sacrificing any information. In addition, MSS can be shown to be the statistic with all the available information about $Y$ while retaining as little information about $X$ as possible. Generally, sufficient statistics are restricted, in the sense that their dimension always depends on the sample size, unless the data comes from an exponential family distribution (Koopman; 1936).

*1.2.3   The Information Bottleneck*

Since exact minimal sufficient statistics only exist for special distributions, Tishby et al. (2000) relaxed this optimization problem in two ways: (i) allowing the map to be stochastic, defined as an encoder $p(T|X)$, and (ii) allowing the map to capture only as much as possible of $I(X;Y)$, but not necessarily all of it. They introduced the Information Bottleneck (IB) as a principled approach to extract relevant information from observed signals related to a target. For a random variable $x \in X$, this framework finds the best trade-off between the accuracy and the complexity related to another random variable $y \in Y$ with a joint distribution. The IB has been used in several fields, including neuroscience (Buesing and Maass; 2010; Palmer et al.; 2015), slow feature analysis (Turner and Sahani; 2007), speech recognition (Hecht et al.; 2009) and deep learning (Tishby and Zaslavsky; 2015; Shwartz-Ziv and Tishby; 2017; Alemi et al.; 2016).

Let $X$ be an input random variable, $Y$ a target variable, and $p(x, y)$ their joint distribution. A representation $T$ is a stochastic function of $X$ defined by a mapping $p(t|x)$. This mapping can be viewed as transforming $X \sim P(X)$ into a representation of $T \sim P(T) := \int P_{T|X}(\cdot \mid x) dP_X(x)$ in the $\mathcal{T}$ space.

The triple $Y - X - T$ forms a Markov chain in that order w.r.t. the joint probability measure $P_{X,Y,T} = P_{X,Y} P_{T|X}$ and the mutual information terms $I(X;T)$ and $I(Y;T)$.

As part of the IB framework, we aim to find a representation $P(T \mid X)$ that extracts as much information as possible about $Y$ (high performance), while compressing $X$ maximally (keeping $I(X;T)$ small). We could also interpret it as extracting only the relevant information that $X$ contains about $Y$.

The DPI implies $I(Y;T) \leq I(X;Y)$, so the compressed representation $T$ cannot convey more information than the original signal. As a result, there is a tradeoff between compressed representation and the preservation of relevant information about $Y$. The construction of an efficient representation variable is characterized by its encoder and decoder distributions, $P(T|X)$ and $P(Y|T)$, respectively. The efficient representation of $X$ means minimizing the complexity of the representation $I(T;X)$ while maximizing $I(T;Y)$. Formally, the IB optimization involves minimizing the following objective function:

$$\mathcal{F}[p_\beta(t \mid x); p_\beta(y \mid t)] = I(X;T) - \beta I(Y;T) , \tag{1.3}$$

where $\beta$ is a parameter that controls the trade-off between the complexity of $T$ and the amount of relevant information it preserves. Intuitively, we pass the information that $X$ contains about $Y$ through a "bottleneck" via the representation $T$. It can be shown that

$$I(T : Y) = I(X : Y) - \mathbb{E}_{x \sim p(x), t \sim p(t|x)} \left[ D \left[ p(y|x) || p(y|t) \right] \right] \tag{1.4}$$

IB representations can be found using the IB method, a variant of the Blatu Arimoto algorithm (Arimoto; 1972). For any $\beta$, the conditions for a stationary point of equation **??**, can be expressed via the following self-consistent equations (Tishby et al.; 2000):

$$\begin{cases} (i) & p_\beta(t \mid x) = \frac{p_\beta(t)}{Z_{t|\mathbf{x}}(x;\beta)} e^{-\beta D[p(y|x)||p_\beta(y|t)]} \\ (ii) & p_\beta(t) = \sum_x p_\beta(t \mid x)p(x) \\ (iii) & p_\beta(y \mid t) = \sum_x p(y \mid x)p_\beta(x \mid t) \end{cases} , \tag{1.5}$$

where $Z_{t|\mathbf{x}}(x;\beta)$ is a normalization term. If $X$, $Y$, and $T$ take values in finite sets, and $P(X,Y)$ is known, then alternating iterations of **??** locally converge to a solution, for any initial $P(T \mid X)$.

If we denote $I_X^\beta = I_\beta(T;X)$ and $I_Y^\beta = I_\beta(T;Y)$, the optimal information curve is defined as the optimal values of the trade-off $\left( I_X^\beta, I_Y^\beta \right)$ for some $\beta$. The information plane is the two-dimensional plane in which the IB curve resides. The equations in **??** are satisfied along the information curve, which separates the feasible and unfeasible regions of the information plane by a monotonic concave line.

### 1.3  Deep neural networks

In 1958, Frank Rosenblatt developed the perceptron algorithm, the first artificial neural network component (Rosenblatt; 1958). The system was designed to mimic the way the human brain processes visual data and identifies recognizable objects. It was extended to pattern recognition in the late 1980s.

Deep learning can perform hierarchical learning of the data by applying nonlinear

transformations, which distinguishes it from traditional neural networks. The data are cumulatively passed across multiple layers, which may be fully connected or partially connected. DNNs are multilayer structures constructed by processing units that are called neurons. Each neuron's activation involves the weighted summation of neuron inputs from the previous layer, followed by the transfer function's operation. Interconnected layers of these basic computing blocks are used to build complex deep learning architectures (Schmidhuber; 2014). Using this structure, DNNs can learn hierarchically sophisticated features directly from raw data without manually constructing them (Salakhutdinov et al.; 2013). Deep architectures are often more challenging to effectively train. They bring, however, two significant advantages: (1) they promote the reuse of features, and (2) they can potentially lead to progressively more abstract features at higher layers of representation, which hopefully make it easier to separate the explanatory factors in the data (Erhan et al.; 2009).

Although deep architectures have long existed, the term "deep learning" was first used in 2006 by Hinton et al. (2006). This work showed that a multilayer feedforward neural network could be more efficient by applying pretraining of one layer at a time and considering each layer as an unsupervised Restricted Boltzmann Machine (RBM) by using supervised backpropagation for finetuning. In 2007, Bengio et al. (2007) developed the Stack AutoEncoder (SAE), which comprises a deep architecture of many AutoeEcoders (AEs).

In 2012, Krizhevsky et al. (2012) proposed the AlexNet architecture, which won the ImageNet challenge. This was a significant breakthrough in artificial neural networks. They proposed a deep convolutional neural network with nine layers and implemented it over GPUs for the first time. Several extensions of the vanilla AlexNet network have been developed since then, including deeper convolutional nets as proposed by Simonyan and Zisserman (2014) and residual connections as demonstrated by Ren et al. (2015). DNNs have demonstrated their ability to improve state-of-the-art results in a wide range of machine learning tasks over the past few years. In many areas, from visual object recognition to speech recognition and genomics to drug discovery, they work well and enhance state-of-the-art results dramatically (Graves et al.; 2013; Zhang and LeCun; 2015; Hinton et al.; 2012; He et al.; 2015; LeCun et al.; 2015).

## 1.4 Information Bottleneck and DNNs

Even though DNN has been highly successful, little is known about its reasons for success, and no underlying principles have driven its development. DNNs are often considered black boxes, where the interpretation of predictions and reliability are still open questions. Additionally, their internal structure and the optimization process are still not fully understood. To use these algorithms more efficiently and improve them, we need to understand their dynamic behavior and their ability to learn representations.

In the literature, two areas of work involve DNNs and the IB. One uses the IB concept to analyze DNNs, while the other uses the IB to improve the learning algorithm for DNNs. The rest of this section is divided into these categories.

### 1.4.1  Information Bottleneck as optimization objective

Recently, the IB framework was explored as an objective for deep learning. This concept was achieved by optimizing the IB Lagrangian using a variational bound (Alemi et al.; 2016; **?**).

The variational information bottleneck (VIB) approach presented in (Alemi et al.; 2016) used DNNs to parameterize the IB model. A variational approximation of the objective is parametrized using DNN, and an efficient training algorithm is suggested for obtaining a stochastic network that maps inputs to randomized representations.

Given a DNN, denote its output representation as $T$, a randomized mapping operating on the input feature $X$, where the corresponding label is $Y$. The encoding of $X$ into $T$ is defined through a conditional probability distribution, which is parametrized as $p_{T|X}^{\theta}$. The VIB optimization objective is

$$\mathcal{L}_{\beta}^{(\mathsf{VIB})}\left(\theta\right) := \max_{\theta \in \Theta} I\big(T^{(\theta)}; Y\big) - \beta I\big(X; T^{(\theta)}\big). \tag{1.6}$$

However, since the data distributions $P_{X,Y}$ and $P_{T|X}$ are unknown, a direct optimization of Equation **??** is intractable. To overcome it, Alemi et al. (2016) suggested to lower bound Equation **??** in a form that we can optimize:

$$\mathbb{E}_{P_{Y,T}^{(\theta)}}\left[\log Q_{Y|T}^{(\phi)}\left(Y\Big|T^{(\theta)}\right)\right] - \beta \mathsf{D}_{\mathsf{KL}}\left(P_{T|X}^{(\theta)}\Big\|P_{T}^{(\theta)}\right), \tag{18}$$

In this case, we parametrize the decoder $Q_{(Y|T}^{phi}$ by a DNN. However, the main difficulty of this optimization is intractable marginal distribution $P_{T}^{\theta}$. To circumvent this, we take $p_{T|X}$ and $p_{T}$ distributions with a closed-form KL solution. In this case, we treat the network output's encoder as the parameters of $P_{T|X}$.

Using the reparametrization trick from Kingma and Welling (2013) and replacing $P(X, Y)$ with its empirical proxy, the loss function of IB can be approximated as

$$\hat{\mathcal{L}}_{\beta}^{(\mathsf{VIB})}\left(\theta, \phi, \mathcal{D}_{n}\right) := \frac{1}{n}\sum_{i=1}^{n}\mathbb{E}\left[-\log Q_{Y|T}^{(\phi)}\left(y_{n}\big|f\left(x_{n}, T\right)\right)\right]$$
$$+ \beta \mathsf{D}_{\mathsf{KL}}\left(P_{T|X}^{(\theta)}\left(\cdot|x_{n}\right)\Big\|R_{T}\left(\cdot\right)\right). \tag{19}$$

$T$ is an auxiliary noise variable, and the expectation is w.r.t to its law. Note that if $t$ represents the output of a DNN, the first term is cross-entropy, which is derived from

the loss function common in deep learning. The second term acts as a regularization term that penalizes the dependence of $X$ on $T$, thus encouraging compression. An unbiased estimate of the true variational lower bound can be obtained through standard stochastic gradient-based methods by calculating the estimator's gradient.

Since Alemi et al. (2016), various extensions have been developed, demonstrating promising attributes (Strouse and Schwab; 2017; Elad et al.; 2019; **?**). Recently, the conditional entropy bottleneck (CEB) (Fischer and Alemi; 2020) has been proposed. The CEB provides variational optimizing bounds on $I(Y;T)$ and $I(X;T)$ using a variational decoder $q(y|x)$, variational conditional marginal, $q(t|y)$, and a variational encoder, $p(t|x)$, all implemented by DNNs. Alemi et al. (2016) also showed that the variational auto-encoder (VAE) (Kingma and Welling; 2013) can be considered as an estimation for a special case of IB when $Y = X$, $\beta = 1$ and the prior distribution function is fixed.

### 1.4.2  Information Bottleneck theory for deep learning

Tishby and Zaslavsky (2015) proposed a theoretical framework to analyze DNNs based on the principle of the IB. They formulated the ultimate goal of the network as a trade-off between compression and prediction. An optimal point on the *information curve* exists for each layer where this trade-off can be addressed effectively. The layer's network structure forms a Markov chain, where each layer processes inputs from the previous layer. Due to DPI, any loss of information about $Y$ in one layer is not recoverable in higher layers.

Formally, define $t^i \in T_i$ as the compressed representation of $x$ in the i-th layer and $\hat{y}$ as the network's output. $T_i$ is uniquely mapped to a single point in the information-plane with coordinates $(I(X;T_i), I(T_i;Y))$. We map a network with $k$ layers to $K$ monotonic connected points in the plane. For any $k \geq j$, it holds that

$$
\begin{aligned}
H(X) \geq I(X;T_j) \geq I(X;T_k) \geq I\left(X;\hat{Y}\right) \\
I(Y;X) \geq I(Y;T_j) \geq I(Y;T_k) \geq I\left(Y;\hat{Y}\right) \\
I(X;T_j) \geq I(Y;T_j),
\end{aligned}
\tag{1.7}
$$

and the equality in the second line is achieved IFF each layer is a sufficient statistic of its input. Using this framework, each layer should extract a compact representation while preserving the relevant information.

By successively decreasing $\beta$, we shift the network from low representations and construct higher and more abstract ones. On the one hand, $I(Y;T_i)$ measures how much of the predictive features in $X$ for $Y$ is captured by the layer and can view it as an upper bound of the layer's quality. On the other hand, $I(X;T_i)$ can be interpreted as the complexity of the layer. As a result, we can assess the performance of DNNs not just in terms of output, as we do when we use other measures of error when evaluating DNNs.

Based on the theoretical IB limit and the limitations imposed by the DPI on the information flow between layers, we can get a good sense of each layer's optimality in the network. With each successive layer, the IB distortion level increases, but it also compresses the inputs, hopefully removing only non-relevant information (Tishby and Zaslavsky; 2015).
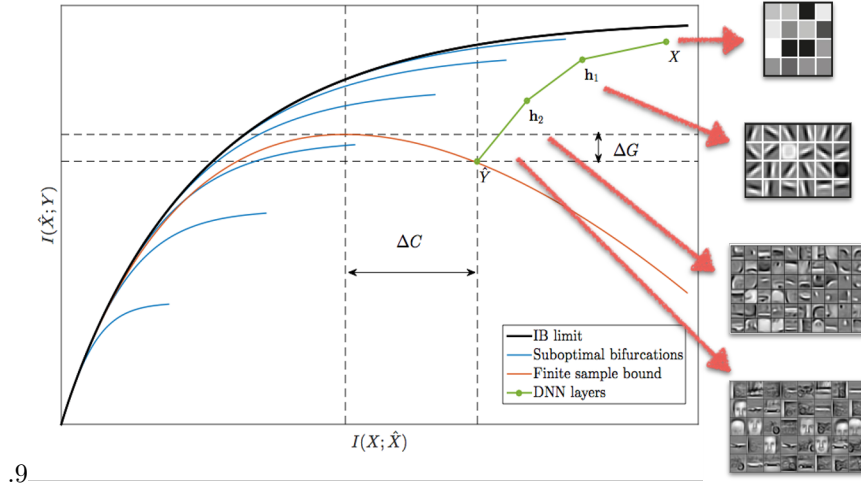


.9

**Figure 1.1: A qualitative information plane from Tishby and Zaslavsky (2015), a hypothesized path of the layers in a typical DNN (green line) on the training data. The black line is the optimal achievable IB limit, and the blue lines are sub-optimal IB bifurcations. The red line corresponds to the upper bound on the out-of-sample IB distortion when training from a finite sample. We want to shift the green DNN layers closer to the optimal curve to obtain lower complexity and better generalization by shifting the last layer to the maximum of the red curve.**

*1.4.3 Information bounds on the generalization gap*

Based on the statistical learning theory, models with many parameters tend to overfit by modeling learned data too accurately, which reduces their ability to generalize to new data (Boucheron et al.; 2005). In practice, however, we see that DNNs have a minimal generalization gap between training and testing. Recently, there has been much interest in understanding implicit regularization. Researchers' findings indicate that network size is not the most important factor involved in the process of learning multilayer feedforward networks and that some unknown factors are important (Neyshabur et al.; 2014). Moreover, conventional statistical learning theories, such as Rademacher complexity, VC-dimension, and uniform stability (Vapnik; 1998; Bartlett and Mendelson; 2002; Bousquet and Elisseeff; 2002), do not explain all of the unexpected results of numerical experiments. According to Zhang et al. (2016), regularization plays a unique role in deep learning that differs from empirical risk minimization. However, convincing experiments have indicated that

the generalization gap can be reduced even without explicit regularization. In recent years, several works have demonstrated that the mutual information between the training inputs and the inferred parameters provides a concise bound on the generalization gap (Xu and Raginsky; 2017; Pensia et al.; 2018; Negrea et al.; 2019; Asadi et al.; 2018; Russo and Zou; 2016; Steinke and Zakynthinou; 2020; Achille et al.; 2019). Achille and Soatto (2018a) investigated how using an IB objective on network parameters (rather than the representations) could avoid overfitting while enforcing invariant representations.

## 1.5  Infinitely-wide neural networks

Neural tangent kernel (NTK) is a powerful theoretical tool for modeling neural networks. Lee et al. (2019) showed that if neural network training is modeled as gradient flow, the training trajectory can be modeled by ordinary differential equations (ODEs). Equations like these represent a finite-width tangent kernel encoded by the network architecture and its current time-dependent weights. Moreover, they showed that if one scales the learning rate per layer appropriately ("NTK parametrization") and lets the width tend to infinity, this kernel converges to the infinite-width NTK. By being independent of weights and staying constant throughout training, this model simplifies the ODE in this limit. They showed that for $l2$ loss, the predictor at convergence would be produced by a kernel regression using an infinite-width NTK. Importantly, the NTK depends only on the architecture of the network and is not learned.

Infinitely-wide neural networks behave as they are linear in their parameters (Lee et al.; 2019):

$$z(x, \theta) = z_0(x) + \frac{\partial z_0}{\partial \theta}(\theta - \theta_0) \quad z_0(x) \equiv z(x, \theta_0). \tag{1.8}$$

This makes them particularly analytically tractable. An infinitely-wide neural network, trained by gradient flow to minimize the squared loss, admits a closed form expression for the evolution of its predictions as a function of time:

$$z(x, \tau) = z_0(x) - \Theta(x, \mathcal{X})\Theta^{-1} \left( I - e^{-\tau\Theta} \right) (z_0(\mathcal{X}) - \mathcal{Y}). \tag{1.9}$$

Here, $z$ denotes the output of our neural network acting on the input $x$. $\tau$ is a dimensionless representation of the time of our training process. $\mathcal{X}$ denotes the whole training set of examples, with their targets $\mathcal{Y}$; $z_0(x) \equiv z(x, \tau = 0)$ denotes the neural networks output at initialization. The evolution is governed by $\Theta$ (the NTK). For a finite width network, the NTK corresponds to $JJ^T$, the neural network gradients' gram matrix. As the network's width increases to infinity, this kernel converges in probability to a fixed value. Tractable ways to calculate the exact infinite-width kernel for broad classes of neural networks are available (Lee et al.; 2019). The shorthand $\Theta$ denotes the kernel function evaluated on the train data ($\Theta \equiv \Theta(\mathcal{X}, \mathcal{X})$).

Observe that infinitely-wide networks trained with gradient flow and squared loss behave as affine transformations of their initial predictions. As a result, for an infinite ensemble of such networks, if the initial weight configurations are drawn from a Gaussian distribution, the law of large numbers guarantees the distribution of the output conditioned on the input is Gaussian. As the evolution is an affine transformation of the initial predictions, the predictions remain Gaussian throughout.

$$p(z|x) \sim \mathcal{N}(\mu(x,\tau), \Sigma(x,\tau)) \tag{1.10}$$

$$\mu(x,\tau) = \Theta(x,\mathcal{X})\Theta^{-1}\left(I - e^{-\tau\Theta}\right)\mathcal{Y} \tag{1.11}$$

$$\Sigma(x,\tau) = \mathcal{K}(x,x) + \Theta(x,\mathcal{X})\Theta^{-1}\left(I - e^{-\tau\Theta}\right)\left(\mathcal{K}\Theta^{-1}\left(I - e^{-\tau\Theta}\right)\Theta(\mathcal{X},x) - 2\mathcal{K}(\mathcal{X},x)\right). \tag{1.12}$$

Here, $\mathcal{K}$ denotes another kernel, the *neural network gaussian process* kernel (NNGP). For a finite width network, the NNGP corresponds to the expected gram matrix of the output-$\mathbb{E}\left[zz^T\right]$. In the infinite width limit, this concentrates to a fixed value. Just as for the NTK, the NNGP can be tractably computed (Lee et al.; 2019) and should be considered only a function of the neural network architecture. These results, which give us a conditional posterior distribution for each time step, enable us to create a powerful model family to investigate DNNs. These tractable distributions can be used to derive many intractable information-theoretic quantities.

### 1.6   The impact of the work

In summary, the IB theory offers a new perspective on DNNs and their capability to learn meaningful representations. This theory suggests studying the system by grouping each layer in the network into two information pairs, one with input and another with output: $(I(X;T_\ell), I(T_\ell;Y))$.

Following this study, which provided a theoretical foundation for understanding deep learning, several studies have offered further explorations of DNNs using information theory tools. We now discuss some of these works and whether they support or challenge our claims. A partial list includes Achille and Soatto (2018a); Saxe et al. (2019); Yu et al. (2020); Cheng et al. (2018); Goldfeld et al. (2018); Wickstrøm et al. (2019); Amjad and Geiger (2019); Goldfeld et al. (2020); Cvitkovic and Koliander (2019).

Following this study, several researchers have analyzed the information plane using different estimation mechanisms in various DNN datasets, architectures, and activation functions. They found conflicting results: The authors of Saxe et al. (2019) did not observe compression in DNNs with ReLU activation functions. However, according to Chelombiev et al. (2019), compression can occur earlier in training or later in training, depending on how the DNN parameters are initialized. Goldfeld et al. (2018) developed a noisy DNNs framework with a rigorous estimator for $I(X;T)$. Using this estimator, they observed

the input's compression in various models. The authors related the noisy DNN to an information-theoretic communication problem, demonstrating that compression is driven by the progressive clustering of inputs belonging to the same class. The study clarified the geometric effects of mutual information compression during training. Other methods that have been proposed for estimating mutual information in DNNs are generative decoder networks (Darlow and Storkey; 2020; Nash et al.; 2018), the mutual information neural estimator (MINE) (Elad et al.; 2019), ensemble dependency graph estimator (EDGE) (Noshad and Hero III; 2018), adaptive approaches for density estimation (Chelombiev et al.; 2019), and noisy sounding entropy estimator (Goldfeld et al.; 2018), and more. For a detailed review of these works, see Geiger (2020).

In another line of research, our idea was used to analyze networks using information with the weights (Achille and Soatto; 2018a; Achille et al.; 2018; Achille and Soatto; 2018b). It has been demonstrated that flat minima, which have better generalization properties, bound the information with the weights, which bounds the information with the activations. In Achille et al. (2018), they used Fisher information on the weights to illustrate the two stages of learning that DNNs go through: increasing the information followed by progressively decreasing the information. By compressing the information with the weights, the network also compresses the information with the activations.

As shown in Achille and Soatto (2018b), minimizing a stochastic network using an approximation of the compression term as a regularizer is equivalent to minimizing cross-entropy over deterministic DNNs with multiplicative noise (information dropout). In addition, they found that Bernoulli noise is a special case that leads to the dropout method. Elad et al. (2019) shows that the binned information can be interpreted as a weight decay penalty, which is typical of DNN training.

Furthermore, many competing information objectives have been proposed for training DNNs based on the IB principle, which can be difficult to compute for high dimensions. Kirsch et al. (2020) examined these quantities and compared, unified, and related them to surrogate objectives that are more easily optimized. The unifying view provided insights into IB training limitations and demonstrated how to avoid the pathological behavior of IB objectives. In addition, they discussed how simple objectives, based on this intuition, could capture many desirable features of IB algorithms while also scaling to complex deep learning problems. Furthermore, they explored how applying practical constraints to a neural network's expressivity can provide new insights into measuring compression and possibly improving regularization.

### 1.7 Overview and main contributions

The purpose of our thesis is to provide an explanation of deep learning by using an information-theoretic framework, supported by empirical evidence, and overcoming some

of the shortcomings of IB. We have outlined the major results and contributions of the thesis in three sections below:

- Part 1 – **Opening the Black Box of Deep Neural Networks** – The first contribution of this thesis is a presentation of an information-based theory for DNNs. Combining two papers shows that DNNs can learn to optimize the mutual information that each layer preserves on the input and output variables, resulting in a trade-off between compression and prediction. We present a theoretical and numerical analysis of DNNs in the information plane and explain how the SGD algorithm follows the information bottleneck trade-off principle. SGD achieves this optimal bound by compressing each layer to a maximum conditional entropy state subject to the constraints of the labels' information. Moreover, we find that the network's training is characterized by a rapid increase in the mutual information between the layers and the target label, followed by a slower decrease in the mutual information between the layers and the input variable. By introducing a new generalization error bound, which is exponential in the input representation compression, we propose a novel analytic bound on the mutual information between successive layers in the network. Combining these with the empirical case study enables a comprehensive understanding of optimization dynamics, SGD training properties, and deep architectures' computational benefits.

- Part 2 – **Information in Infinite Ensembles of Infinitely-wide neural networks** – As previously mentioned, it is often challenging to measure information-theoretic quantities. We developed tractable computations for a wide range of information-theoretic quantities using the NTK framework. These results are used to investigate the critical quantities that correlate with generalization, dynamics, and optimality during the learning process in high-dimensional networks.

- Part 3 – **The Dual Information Bottleneck** – As mentioned before, the IB framework has several known drawbacks, including that it is not optimal for training with finite samples and that it does not preserve the problem structure. We present the dualIB, a framework for addressing some of these shortcomings. By switching the order between the representation decoder and data in the IB's distortion function, the optimal decoding becomes the geometric mean of input points instead of the arithmetic mean. Whenever the data can be modeled in a parametric form, this conversion preserves its structure and, thus, the original properties of the data. Furthermore, the prediction accuracy is improved by optimizing the mean prediction error exponent for each data size. By analyzing the framework's structure, we show how to solve this new representation learning formulation. A variational formulation of the dualIB for DNNs, (VdualIB) is presented to address large-scale problems. We examine its properties using real-world datasets and compare them to the original IB.

***Bibliography***

Achille, A., Paolini, G. and Soatto, S. (2019). Where is the information in a deep neural network?, *arXiv preprint arXiv:1905.12213* .

Achille, A., Rovere, M. and Soatto, S. (2018). Critical learning periods in deep networks, *International Conference on Learning Representations.*

Achille, A. and Soatto, S. (2018a). Emergence of invariance and disentanglement in deep representations, *The Journal of Machine Learning Research* **19**(1): 1947–1980.

Achille, A. and Soatto, S. (2018b). Information dropout: Learning optimal representations through noisy computation, *IEEE transactions on pattern analysis and machine intelligence* **40**(12): 2897–2905.

Alemi, A. A., Fischer, I., Dillon, J. V. and Murphy, K. (2016). Deep variational information bottleneck, *arXiv preprint arXiv:1612.00410* .

Amjad, R. A. and Geiger, B. C. (2019). Learning representations for neural network-based classification using the information bottleneck principle, *IEEE transactions on pattern analysis and machine intelligence* .

Arimoto, S. (1972). An algorithm for computing the capacity of arbitrary discrete memoryless channels, *IEEE Transactions on Information Theory* **18**(1): 14–20.

Asadi, A., Abbe, E. and Verdu, S. (2018). Chaining mutual information and tightening generalization bounds, *in* S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi and R. Garnett (eds), *Advances in Neural Information Processing Systems 31*, Curran Associates, Inc., pp. 7234–7243.
**URL:** *http://papers.nips.cc/paper/7954-chaining-mutual-information-and-tightening-generalization-bounds.pdf*

Bartlett, P. L. and Mendelson, S. (2002). Rademacher and gaussian complexities: Risk bounds and structural results, *Journal of Machine Learning Research* **3**(Nov): 463–482.

Bengio, Y., Courville, A. and Vincent, P. (2013). Representation learning: A review and new perspectives, *IEEE transactions on pattern analysis and machine intelligence* **35**(8): 1798–1828.

Bengio, Y., Lamblin, P., Popovici, D. and Larochelle, H. (2007). Greedy layer-wise training of deep networks, *Advances in neural information processing systems*, pp. 153–160.

Bengio, Y. and LeCun, Y. (2007). Scaling learning algorithms towards ai ,in l. bottou, o. chapelle, d. decoste, and j. weston, editors,, *Large Scale Kernel Machines,MIT Press.* .

Boucheron, S., Bousquet, O. and Lugosi, G. (2005). Theory of classification: A survey of some recent advances, *ESAIM: probability and statistics* **9**: 323–375.

Bousquet, O. and Elisseeff, A. (2002). Stability and generalization, *Journal of machine learning research* **2**(Mar): 499–526.

Buesing, L. and Maass, W. (2010). A spiking neuron as information bottleneck, *Neural computation* **22**(8): 1961–1992.

Chelombiev, I., Houghton, C. and O'Donnell, C. (2019). Adaptive estimators show information compression in deep neural networks, *ICLR* .

Cheng, H., Lian, D., Gao, S. and Geng, Y. (2018). Evaluating capability of deep neural networks for image classification via information plane, *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 168–182.

Cover, T. M. (1999). *Elements of information theory*, John Wiley & Sons.

Csiszár, I., Cover, T. and Choi, B.-S. (1987). Conditional limit theorems under markov conditioning, *IEEE Transactions on Information Theory* **33**(6): 788–801.

Cvitkovic, M. and Koliander, G. (2019). Minimal achievable sufficient statistic learning, *arXiv preprint arXiv:1905.07822* .

Darlow, L. N. and Storkey, A. (2020). What information does a resnet compress?, *arXiv preprint arXiv:2003.06254* .

Devlin, J., Chang, M.-W., Lee, K. and Toutanova, K. (2018). Bert: Pre-training of deep bidirectional transformers for language understanding, *arXiv preprint arXiv:1810.04805* .

Elad, A., Haviv, D., Blau, Y. and Michaeli, T. (2019). Direct validation of the information bottleneck principle for deep nets, *Proceedings of the IEEE International Conference on Computer Vision Workshops*.

Erhan, D., Manzagol, P.-A., Bengio, Y., Bengio, S. and Vincent, P. (2009). The difficulty of training deep architectures and the effect of unsupervised pre-training., *AISTATS*, Vol. 5, pp. 153–160.

Fischer, I. and Alemi, A. A. (2020). Ceb improves model robustness, *arXiv preprint arXiv:2002.05380* .

Geiger, B. C. (2020). On information plane analyses of neural network classifiers–a review, *arXiv preprint arXiv:2003.09671* .

Goldfeld, Z., Berg, E. v. d., Greenewald, K., Melnyk, I., Nguyen, N., Kingsbury, B. and Polyanskiy, Y. (2018). Estimating information flow in deep neural networks, *arXiv preprint arXiv:1810.05728* .

Goldfeld, Z., Greenewald, K., Niles-Weed, J. and Polyanskiy, Y. (2020). Convergence of smoothed empirical measures with applications to entropy estimation, *IEEE Transactions on Information Theory* .

Goodfellow, I., Bengio, Y. and Courville, A. (2016). Deep learning. Book in preparation for MIT Press.
**URL:** *http://www.deeplearningbook.org*

Graves, A., Mohamed, A.-r. and Hinton, G. (2013). Speech recognition with deep recurrent neural networks, *Acoustics, speech and signal processing (icassp), 2013 ieee international conference on*, IEEE, pp. 6645–6649.

He, K., Zhang, X., Ren, S. and Sun, J. (2015). Deep residual learning for image recognition, *CoRR* **abs/1512.03385**.

He, K., Zhang, X., Ren, S. and Sun, J. (2016). Deep residual learning for image recognition, *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 770–778.

Hecht, R. M., Noor, E. and Tishby, N. (2009). Speaker recognition by gaussian information bottleneck, *Tenth Annual Conference of the International Speech Communication Association*.

Hinton, G. E., Osindero, S. and Teh, Y.-W. (2006). A fast learning algorithm for deep belief nets, *Neural computation* **18**(7): 1527–1554.

Hinton, G. E., Srivastava, N., Krizhevsky, A., Sutskever, I. and Salakhutdinov, R. (2012). Improving neural networks by preventing co-adaptation of feature detectors, *CoRR* **abs/1207.0580**.
**URL:** *http://arxiv.org/abs/1207.0580*

Kingma, D. P. and Welling, M. (2013). Auto-encoding variational bayes, *arXiv preprint arXiv:1312.6114* .

Kinney, J. B. and Atwal, G. S. (2014). Equitability, mutual information, and the maximal information coefficient, *Proceedings of the National Academy of Sciences* **111**(9): 3354–3359.

Kirsch, A., Lyle, C. and Gal, Y. (2020). Unpacking information bottlenecks: Unifying information-theoretic objectives in deep learning, *arXiv preprint arXiv:2003.12537* .

Koopman, B. O. (1936). On distributions admitting a sufficient statistic, *Transactions of the American Mathematical society* **39**(3): 399–409.

Krizhevsky, A., Sutskever, I. and Hinton, G. E. (2012). Imagenet classification with deep convolutional neural networks, *Advances in neural information processing systems*, pp. 1097–1105.

LeCun, Y., Bengio, Y. and Hinton, G. (2015). Deep learning, *Nature* .

Lee, J., Xiao, L., Schoenholz, S., Bahri, Y., Novak, R., Sohl-Dickstein, J. and Pennington, J. (2019). Wide neural networks of any depth evolve as linear models under gradient descent, *Advances in neural information processing systems*, pp. 8572–8583.

Marr, D. (2010). *Vision: A Computational Investigation Into the Human Representation and Processing of Visual Information*, London.
**URL:** *https://books.google.co.il/books?id=EehUQwAACAAJ*

Nash, C., Kushman, N. and Williams, C. K. (2018). Inverting supervised representations with autoregressive neural density models, *arXiv preprint arXiv:1806.00400* .

Negrea, J., Haghifam, M., Dziugaite, G. K., Khisti, A. and Roy, D. M. (2019). Information-theoretic generalization bounds for sgld via data-dependent estimates, *in* H. Wallach, H. Larochelle, A. Beygelzimer, F. d Alché-Buc, E. Fox and R. Garnett (eds), *Advances in Neural Information Processing Systems 32*, Curran Associates, Inc., pp. 11015–11025.
**URL:** *http://papers.nips.cc/paper/9282-information-theoretic-generalization-bounds-for-sgld-via-data-dependent-estimates.pdf*

Neyshabur, B., Tomioka, R. and Srebro, N. (2014). In search of the real inductive bias: On the role of implicit regularization in deep learning, *arXiv preprint arXiv:1412.6614* .

Noshad, M. and Hero III, A. O. (2018). Scalable mutual information estimation using dependence graphs, *arXiv preprint arXiv:1801.09125* .

Oord, A. v. d., Dieleman, S., Zen, H., Simonyan, K., Vinyals, O., Graves, A., Kalchbrenner, N., Senior, A. and Kavukcuoglu, K. (2016). Wavenet: A generative model for raw audio, *arXiv preprint arXiv:1609.03499* .

Palmer, S. E., Marre, O., Berry, M. J. and Bialek, W. (2015). Predictive information in a sensory population, *Proceedings of the National Academy of Sciences* **112**(22): 6908–6913.

Paninski, L. (2003). Estimation of entropy and mutual information, *Neural computation* **15**(6): 1191–1253.

Pensia, A., Jog, V. and Loh, P.-L. (2018). Generalization error bounds for noisy, iterative algorithms, *2018 IEEE International Symposium on Information Theory (ISIT)*, IEEE, pp. 546–550.

Ren, S., He, K., Girshick, R. and Sun, J. (2015). Faster r-cnn: Towards real-time object detection with region proposal networks, *Advances in neural information processing systems*, pp. 91–99.

Rosenblatt, F. (1958). The perceptron: a probabilistic model for information storage and organization in the brain., *Psychological review* **65**(6): 386.

Russo, D. and Zou, J. (2016). Controlling bias in adaptive data analysis using information theory, *Artificial Intelligence and Statistics*, pp. 1232–1240.

Salakhutdinov, R., Tenenbaum, J. B. and Torralba, A. (2013). Learning with hierarchical-deep models, *IEEE transactions on pattern analysis and machine intelligence* **35**(8): 1958–1971.

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.

Schmidhuber, J. (2014). Deep learning in neural networks: An overview, *CoRR* **abs/1404.7828**.

Shannon, C. E. (1948). A mathematical theory of communication, *The Bell system technical journal* **27**(3): 379–423.

Shwartz-Ziv, R. and Tishby, N. (2017). Opening the black box of deep neural networks via information, *arXiv preprint arXiv:1703.00810* .

Simonyan, K. and Zisserman, A. (2014). Very deep convolutional networks for large-scale image recognition, *arXiv preprint arXiv:1409.1556* .

Steinke, T. and Zakynthinou, L. (2020). Reasoning about generalization via conditional mutual information, *arXiv preprint arXiv:2001.09122* .

Strouse, D. and Schwab, D. J. (2017). The deterministic information bottleneck, *Neural computation* **29**(6): 1611–1630.

Tishby, N., Pereira, F. C. and Bialek, W. (2000). The information bottleneck method, *arXiv preprint physics/0004057* .

Tishby, N. and Zaslavsky, N. (2015). Deep learning and the information bottleneck principle, *Information Theory Workshop (ITW), 2015 IEEE*, IEEE, pp. 1–5.

Turner, R. and Sahani, M. (2007). A maximum-likelihood interpretation for slow feature analysis, *Neural computation* **19**(4): 1022–1038.

Vapnik, V. N. (1998). Statistical learning theory, *Wiley* .

Wickstrøm, K., Løkse, S., Kampffmeyer, M., Yu, S., Principe, J. and Jenssen, R. (2019). Information plane analysis of deep neural networks via matrix-based renyi's entropy and tensor kernels, *arXiv preprint arXiv:1909.11396* .

Xu, A. and Raginsky, M. (2017). Information-theoretic analysis of generalization capability of learning algorithms, *Advances in Neural Information Processing Systems*, pp. 2524–2533.

Yu, S., Wickstrøm, K., Jenssen, R. and Principe, J. C. (2020). Understanding convolutional neural networks with information theory: An initial exploration, *IEEE Transactions on Neural Networks and Learning Systems* .

Zhang, C., Bengio, S., Hardt, M., Recht, B. and Vinyals, O. (2016). Understanding deep learning requires rethinking generalization, *arXiv preprint arXiv:1611.03530* .

Zhang, X. and LeCun, Y. (2015). Text understanding from scratch, *CoRR* **abs/1502.01710**.
   **URL:** *http://arxiv.org/abs/1502.01710*

# Opening the Black Box of Deep Neural Networks via Information

**Unpublished**

Ravid Shwartz-Ziv and Naftali Tishby (2018)

# Opening the Black Box of Deep Neural Networks via Information

Ravid Shwartz-Ziv [1] Naftali Tishby[1,2]

[1] The Edmond and Lilly Safra Center for Brain Sciences, The Hebrew University,
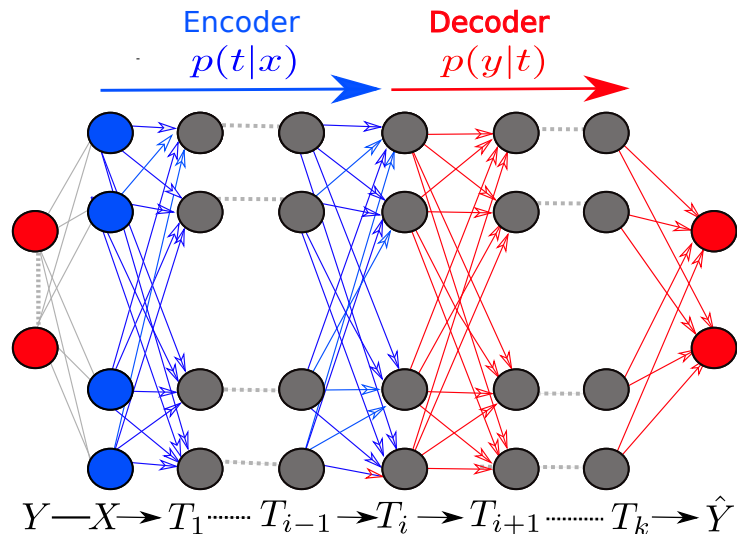Jerusalem, Israel.
[2] School of Computer Science and Engineering,
The Hebrew University,
Jerusalem, Israel.

## Abstract

Despite their great success, there is still no comprehensive theoretical understanding of learning with Deep Neural Networks (DNNs) or their inner organization. Previous work suggested analyzing DNNs in the *Information-Plane*; the plane of the mutual information values that each layer preserves on the input and output variables. They suggested that the network's goal is to optimize the Information Bottleneck (IB) trade-off between compression and prediction, successively, for each layer. In this work, we demonstrate the effectiveness of the information-plane visualization of DNNs. We first show that the stochastic gradient descent (SGD) epochs have two distinct phases: fast empirical error minimization followed by slow representation compression for each layer. We then argue that the DNN layers end up very close to the IB theoretical bound and present a new argument for the hidden layers' computational benefit.

## Introduction

DNNs heralded a new era in predictive modeling and machine learning. Their ability to learn and generalize has set a new bar on performance, compared to state-of-the-art methods. This improvement is evident across almost every application domain (Graves

25

$$Y \mathbin{\!-\!} X \!\rightarrow\! T_1 \cdots\!\cdots T_{i-1} \!\rightarrow\! T_i \!\rightarrow\! T_{i+1} \cdots\!\cdots T_k \!\rightarrow\! \hat{Y}$$

**Figure 1.2: The DNN layers form a Markov chain of successive internal representations of the input. Any representation, $T$, is defined through an encoder, $P(T|X)$, and a decoder $P(\hat{Y}|T)$, and can be quantified by its *information-plane* coordinates: $I(X;T)$ and $I(T;Y)$. The IB bound characterizes the optimal representations, which maximally compress the input $X$, for a given mutual information on the desired output $Y$.**

et al.; 2013; Zhang and LeCun; 2015; Hinton et al.; 2012; He et al.; 2015; LeCun et al.; 2015).

However, despite their success, there is still little understanding of their internal organization or optimization process. They are often seen as mysterious "black boxes" (Alain and Bengio; 2016).

The authors in Tishby and Zaslavsky (2015) pointed out that layered neural networks' representations of input layers form a Markov chain. They suggested studying these in the *information-plane* - The plane of the mutual information values of any other variable with the input variable $X$ and the desired output variable $Y$ (Figure 1.2). The rationale for this analysis was based on the mutual information's invariance to invertible re-parameterization and the data processing inequalities (DPI) (Cover and Thomas; 2006) along the Markov chain of the layers. Moreover, they suggested that optimized DNNs layers should approach the Information Bottleneck (IB) bound (Tishby et al.; 1999) of the optimally achievable representations of the input $X$.

In this paper, we extend their work and demonstrate the effectiveness of visualizing DNNs in the information-plane to understand better the training dynamics, learning processes, and internal representations of deep learning.

Our analysis reveals that the SGD optimization, commonly used in deep networks, has two different and distinct phases: empirical error minimization (ERM) and representation

compression. Each of these phases displays a very different signal-to-noise ratio of the stochastic gradients. During the ERM phase, gradient norms are much larger than their stochastic fluctuations, leading to a sharp increase in the mutual information on the label variable $Y$. At the compression phase, the errors fluctuate much more than their means, causing the weights to change much like Weiner processes, or random diffusion, with a small influence of error gradients. In this phase, the input variable $X$ undergoes slow representation compression or reduction of the mutual information.

In our experiments, most of the optimization iterations are spent on compressing the internal representations under the training error constraint. This compression occurs by the SGD without any other explicit regularization or sparsity imposed. We suspect that this helps avoid overfitting in deep networks. This observation also suggests that many (exponential in the number of weights) different randomized networks have essentially optimal performance. Hence, the interpretation of a single neuron (or weight) in the layers is practically meaningless.

Then we show that the optimized layers lay on or relatively close to the optimal IB bound for large enough training samples, resulting in a self-consistent relationship between each layer's encoder and decoder distributions (Figure 1.2). The hidden layers in the optimized model converge along particular lines in the information-plane and move up as the training sample size increases. In addition to this, the diffusive nature of the SGD dynamics explains the computational benefit of the hidden layers.

### Information Theory of Deep Learning

In supervised learning, we are interested in good representations, $T(X)$, of the input patterns $x \in X$, which enable good predictions of the label $y \in Y$. Moreover, we want to efficiently learn such representations from an empirical sample of the (unknown) joint distribution $P(X, Y)$, in a manner that provides good generalization.

DNNs generate a Markov chain of such representations, the hidden layers, by minimizing the empirical error over the network's weights, layer by layer. This optimization occurs via SGD, using a noisy estimate of the empirical error gradient at each weight through the backpropagation algorithm (Rumelhart et al.; 1986).

Our first important insight is to treat the whole layer $T$, as a single random variable, characterized by its encoder, $P(T|X)$, and decoder, $P(Y|T)$ distributions. As we are only interested in the information flowing through the network, invertible transformations of the representations that preserve information generate equivalent representations even if the individual neurons encode entirely different input features. For this reason, we quantify the representations by two numbers, or order parameters, that are invariant to any invertible re-parameterization of $T$, the mutual information of $T$ with the input $X$, $I(X; T)$, and the desired output $Y$, $I(T; Y)$.

Next, the layers' quality is quantified by comparing them to the information-theoretic optimal representation - the IB representations. Furthermore, we explore how the SGD can be used to obtain these optimal representations in DNNs.

*Mutual information*

Given any two random variables, $X$ and $Y$, with a joint distribution $p(x, y)$, their mutual information is defined as:

$$I(X; Y) = D_{KL}[p(x, y) || p(x)p(y)] \tag{1.13}$$

$$= \sum_{x \in X, y \in Y} p(x, y) \log \left( \frac{p(x, y)}{p(x)p(y)} \right) \tag{1.14}$$

$$= H(X) - H(X|Y) , \tag{1.15}$$

where $D_{KL}[p||q]$ is the Kullback-Liebler divergence of the distributions $p$ and $q$, and $H(X)$ and $H(X|Y)$ are the entropy and conditional entropy of $X$ and $Y$, respectively.

The mutual information quantifies the number of relevant bits that the input variable $X$ contains about the label $Y$, on average. The optimal learning problem can be cast as the construction of an *optimal encoder* of that relevant information via an efficient representation - a minimal sufficient statistic of $X$ with respect to $Y$. A minimal sufficient statistic can enable the *decoding* of the relevant information with the smallest number of binary questions (on average); i.e., an optimal code. The connection between mutual information and minimal sufficient statistics is discussed in the next section.

Among the characteristics of mutual information, two stand out particularly in relation to DNNs. The first is its invariance to invertible transformations:

$$I(X; Y) = I(\psi(X); \phi(Y))) \tag{1.16}$$

for any invertible functions $\phi$ and $\psi$.

The second is the DPI – (Cover and Thomas; 2006) For any three variables that form a Markov chain $X \to Y \to Z$,

$$I(X; Y) \geq I(X; Z) .$$

*The information-plane*

Any representation variable, $T$, defined as a (possibly stochastic) map of the input $X$, is characterized by its joint distributions with $X$ and $Y$, or by its encoder and decoder distributions, $P(T|X)$ and $P(Y|T)$, respectively. Given $P(X; Y)$, $T$ is uniquely mapped to a point in the information-plane with coordinates $(I(X; T), I(T; Y))$. When applied to the Markov chain of a K-layers DNN, with $T_i$ denoting the $i^{th}$ hidden layer as a single

multivariate variable (Figure 1.2), the layers are mapped to $K$ monotonic connected points in the plane. Therefore, a unique *information path* which satisfies the following DPI chains:

$$I(X;Y) \geq I(T_1;Y) \geq I(T_2;Y) \geq ...I(T_k;Y) \geq I(\hat{Y};Y)$$
$$H(X) \geq I(X;T_1) \geq I(X;T_2) \geq ...I(X;T_k) \geq I(X;\hat{Y}).$$

Layers related by invertible re-parametrization appear at the same point, meaning that each information path in the plane corresponds to a multitude of different DNNs, possibly corresponding to very different architectures.

*The Information Bottleneck optimal bound*

What characterizes the optimal representations of $X$ w.r.t. $Y$? The classical notion of minimal sufficient statistics provides good candidates for optimal representations. Sufficient statistics, in our context, are maps or partitions of $X$, $S(X)$, capturing all the information that $X$ has on $Y$. Namely, $I(S(X);Y) = I(X;Y)$ (Cover and Thomas; 2006).

Minimal sufficient statistics, $T(X)$, are the simplest sufficient statistics and induce the coarsest sufficient partition on $X$. Thus, they are functions of any other sufficient statistic. A simple way of formulating this is through the Markov chain: $Y \rightarrow X \rightarrow S(X) \rightarrow T(X)$, which should hold for a minimal sufficient statistics $T(X)$ with any other sufficient statistics $S(X)$. Using DPI, we can cast it into a constrained optimization problem:

$$T(X) = \arg \min_{S(X):I(S(X);Y)=I(X;Y)} I(S(X);X) \ .$$

Since exact minimal sufficient statistics only exist for special distributions (i.e., exponential families), Tishby et al. (1999) relaxed this optimization problem by first allowing the map to be stochastic, defined as an encoder $P(T|X)$, and then, by enabling the map to capture *as much as possible* of $I(X;Y)$, not necessarily all of it.

This leads to the IB trade off (Tishby et al.; 1999), which provides a computational framework for finding approximate minimal sufficient statistics, or the optimal trade-off between compression of $X$ and prediction of $Y$. In that sense, efficient representations are approximate minimal sufficient statistics.

If we define $t \in T$ as the compressed representations of $x \in X$, the representation of $x$ is now defined by the mapping $p(t|x)$. This IB trade-off is formulated by the following optimization problem, carried independently for the distributions, $p(t|x), p(t), p(y|t)$, with the Markov chain: $Y \rightarrow X \rightarrow T$,

$$\min_{p(t|x),p(y|t),p(t)} \{I(X;T) - \beta I(T;Y)\} \ . \tag{1.17}$$

The Lagrange multiplier $\beta$ determines the level of relevant information captured by the representation $T$, $I(T;Y)$, which is bounding the error in the label prediction from this representation. The local optimal solution to this problem is given by three self-consistent equations, the IB equations:

$$\begin{cases} p\left(t|x\right) = \frac{p(t)}{Z(x;\beta)} \exp\left(-\beta D_{KL}\left[p\left(y|x\right) \parallel p\left(y|t\right)\right]\right) \\ p\left(t\right) = \sum_x p\left(t|x\right) p\left(x\right) \\ p\left(y|t\right) = \sum_x p\left(y|x\right) p\left(x|t\right) \ , \end{cases} \quad (1.18)$$

where $Z\left(x;\beta\right)$ is the normalization function. These equations must be satisfied along the *information curve* which is a monotonic concave line that separates the achievable and unachievable regions in the information-plane. For smooth $P(X,Y)$ distributions; i.e., when $Y$ is not a completely deterministic function of $X$, the information curve is strictly concave with a unique slope, $\beta^{-1}$, at every point, and a finite slope at the origin. In these cases, $\beta$ determines a single point on the information curve with specified *encoder*, $P^\beta(T|X)$, and *decoder*, $P^\beta(Y|T)$, distributions that are related through Eq.(**??**).

*Visualizing DNNs in the information-plane*

As proposed by Tishby and Zaslavsky (2015), we study the *information paths* of DNNs in the information-plane. It is possible to do so when the underlying distribution, $P(X,Y)$, is known and the encoder and decoder distributions, $P(T|X)$ and $P(Y|T)$, can be derived directly. In "real-world" problems, these distributions and mutual information values should be estimated from samples or other modeling assumptions. We use two order parameters, $I(T;X)$ and $I(T;Y)$, to visualize and compare different network architectures in terms of their efficiency in preserving the relevant information in $P(X;Y)$.

By visualizing the paths of different networks in the information-plane, we explore the following fundamental issues:

1. The SGD layer's dynamics in the information-plane.

2. The effect of the training sample size on the layers.

3. The benefit of the hidden layers.

4. The final location in the information-plane of the hidden layers.

5. The optimality of the layers' representation.

**Experiments**

This section is organized as follows; First, we describe the experimental settings. Next, we discuss the results of our experiments, including the dynamics of the optimization process in the information-plane, the stochastic gradients, and the computational benefit of the hidden layers. In the last part, we present the layers' optimality and their evolution with the training size.

*Experimental setup*

The numerical studies in this paper explore fully connected feed-forward neural networks with no further architecture constraints. The activation function of all hidden layers is the hyperbolic tangent function. For the final layer, we use a sigmoidal function. The networks train using SGD and the cross-entropy loss function (with no explicit regularization). Unless otherwise noted, the DNNs use up to seven fully connected hidden layers, with widths of $12 - 10 - 7 - 5 - 4 - 3 - 2$ neurons (see Figure 1.5a). In our results below, layer one is the hidden layer closest to the input and the highest is the output layer.

To simplify our analysis, the tasks were chosen as binary decision rules, which are invariant under $O(3)$ rotations of the sphere. We define the input $X$ to be 12 binary inputs representing uniformly distributed points on a $2D$ sphere (for other, non-symmetric rules, see Supplementary material). The 4096 different patterns of the input variable $X$ are divided into 64 disjoint orbits of the rotation group. These orbits form a minimal sufficient statistics for spherically symmetric rules (Kazhdan et al.; 2003).

Given a function $f$ defined on the sphere, we can obtain a rotation and reflection invariant representation $\psi(f)$ by computing the spherical harmonic decomposition of the function.

$$f(\theta, \phi) = \sum_{l \geq 0} \sum_{m=-l}^{l} a_l^m Y_l^m (\theta, \phi)$$

and calculate the energy ($L_2$ norms) of the frequency components.

$$\psi(f) = \sqrt{\|a_0\|^2 + \|a_1\|^2 \ldots} \tag{1.19}$$

We refer the reader to Kazhdan et al. (2003) for a good exposition on the above. To generate the input-output distribution, $P(X, Y)$, we consider spherically symmetric real valued function of the pattern $f(x)$ minus a threshold, $\theta$, and apply a sigmoidal function, $\Omega(u) = 1/(1 + \exp(-\gamma u))$:

$$p(y = 1|x) = \Omega(\psi(f_x) - \theta)) . \tag{1.20}$$

The threshold $\theta$ was selected such that $p(y = 1) = \sum_x p(y = 1|x)p(x) \approx 0.5$, with uniform

$p(x)$. The sigmoidal gain $\gamma$ was high enough to keep the mutual information $I(X;Y) \approx 0.99$ bits.

*Estimating the mutual information of the layers*

As mentioned above, we model each layer in the network as a single variable $T_i, 1 \leq i \leq K$, and calculate its mutual information with the input and the labels.

To calculate the mutual information of the networks' layers with the input and output variables, we binned the neuron's *arctan* output activation into 30 equal intervals between $-1$ and $1$. These discretized values for each neuron in the layer, $t \in T_i$, are used to directly calculate the joint distributions, over the 4096 equally likely input patterns $x \in X$, $P(T_i, X)$ and $P(T_i, Y) = \sum_x P(x, Y)P(T_i|x)$, using the Markov chain $Y \to X \to T_i$ for every hidden layer. Using the above discrete joint distributions, we calculate the decoder mutual information, $I(Y; T_i)$, and the encoder mutual information, $I(T_i; X)$, for each hidden layer. Note that $I(T_i; Y)$ is calculate with the full data distribution; thus, it corresponds to the generalization error.
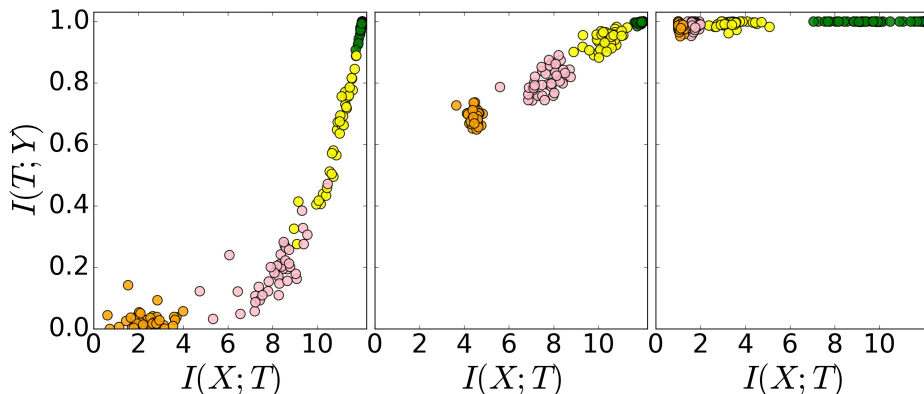
We repeat these calculations with 50 different randomized initialization of the network's weights and different random selections of the training samples, randomly distributed according to the rule $P(X, Y)$ in Eq.(**??**).

*The dynamics of the optimization process*

To understand the network's SGD optimization dynamics, we evaluate and visualize pairs of $I(X; T_i)$ and $I(T_i; Y)$ for each layer along the learning process. We repeat this process for 50 different randomized initializations with different randomized training samples. Figure 1.3 depicts the layers (each one in a different color) of all the 50 networks, trained with a randomized 85 percents of the input patterns, in the information-plane, at various times during the training.

As can be seen, the deeper layers of the randomly-initialized network do not preserve the relevant information at the beginning of the optimization, and $I(T; Y)$ decreases sharply along the path. During the SGD optimization, the layers first increase $I(T; Y)$, then significantly decrease $I(X; T)$ – compressing the representation. One striking observation is that the different randomized networks' layers follow similar paths through the optimization process and eventually converge to nearby points in the information-plane. The randomized networks are thus averaged over, and the average layer trajectory plot is shown in Figure 1.4.

The graph on the left shows those same trajectories when trained with 85 percent of all patterns, whereas on the right, the same trajectories when trained with only 5 percent. Note that the mutual information is calculated with the full rule distribution. Thus,

**Figure 1.3: Snapshots of layers (different colors) of 50 randomized networks during the training process in the information-plane: left - At the initialization time; center - After 400 epochs; right - After 9000 epochs.**

$I(T;Y)$ corresponds to the test or generalization error. In each case, the two optimization phases are visible. During the fast - ERM phase, which lasts a few hundred epochs, the layers increase the information on the labels (increase $I(T;Y)$). However, the layers' information on the input, $I(X;T)$, decreases in the second phase of training, resulting in loss of irrelevant information until converging (yellow points). We call this phase the *representation compression* phase. Because of the cross-entropy loss approximate $I(T;Y)$ up to a constant (see Shwartz-Ziv et al. (2018)), the increase of $I(T;Y)$ in the ERM phase is expected. However, the compression phase is surprising. There is no explicit regularization that could simplify the representations, such as $L1$ regularization, and there is no sparsification or norm reduction of the weights (see supplementary).

While both the small (5%) and large (85%) training sample sizes exhibit the same ERM phase, the compression phase drastically reduces the layers' label information for a small number of examples. However, with large sample sizes, the layers' label information increases. Those results seem very much to be caused by overfitting the small sample noise, a problem that can be eliminated by early stopping methods (Larochelle et al.; 2009). This overfitting is mainly the consequence of the compression phase. It simplifies the layers' representations and loses relevant information. Understanding what determines the convergence points of the layers in the information-plane for different training data sizes is an interesting theoretical goal.

*The two phases of SGD optimization*

By examining the stochastic gradients' behavior along the epochs, a better understanding of the ERM and the representation-compression phases can be gained. For each layer, we look at the mean and standard deviations of the weights' stochastic gradients (in the sample batches). Namely, $M_i = \left\| \langle \frac{\partial L}{\partial W_k} \rangle \right\|_F$ and $S_i = \left\| \text{STD} \left( \frac{\partial L}{\partial W_k} \right) \right\|_F$, where $\langle \cdot \rangle$ and $\text{STD}(\cdot)$ are
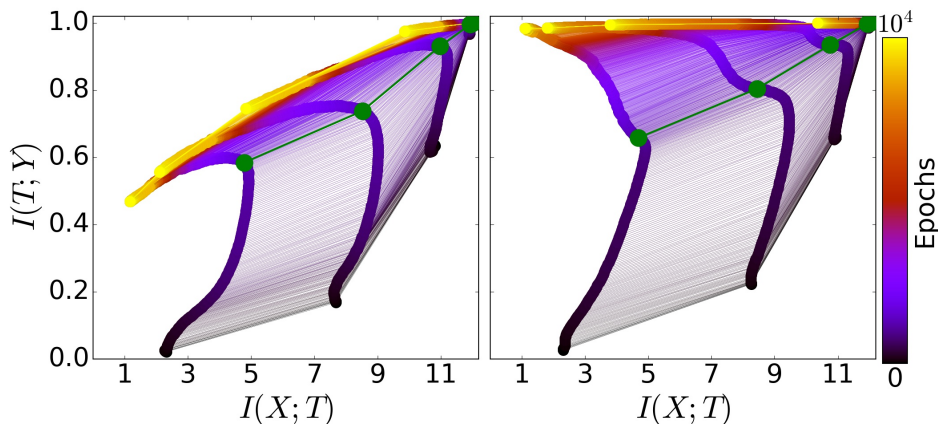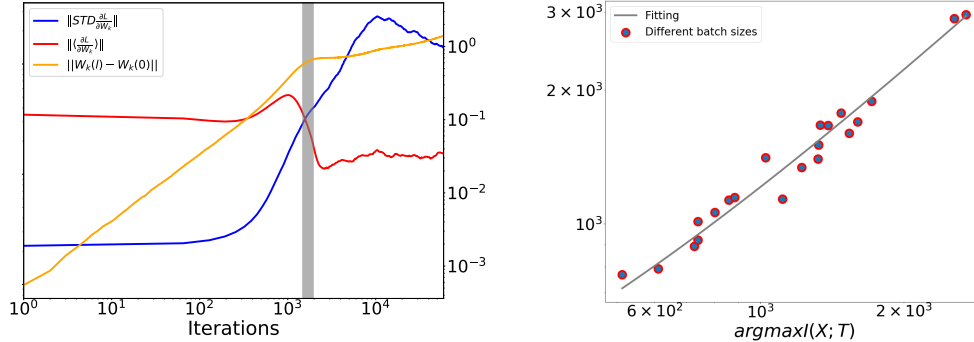
**Figure 1.4: The evolution of layers during the learning process - the information-plane. Left - Paths of layers trained with only 5% of the patterns. Right - Paths of the layers, trained with 85% of the input patterns. The line color indicates the number of epochs. There are 6 points for each path's different hidden layers, averaged over all 50 randomized networks. The layer with the smallest $I(X;T)$ is the output layer, and the one with the highest $I(X;T)$ is the first hidden layer. The green paths correspond to the SGD phase transition grey line on Figure 1.5a**

the mean and the element-wise standard deviation respectively, and $\|\cdot\|_F$ denotes the Frobenius norm. Figure 1.5a shows the mean and standard deviation of the weights' gradients for one layer in our network, as a function of the iterations (in log-log scale). As can be seen, there is a transition between two distinct phases (the vertical line). The phases are defined by the ratio between the means of the gradients and their standard deviations. We refer to this ratio as the signal-to-noise ratio (SNR). The first phase is an *drift phase*, when the gradient means are much larger than the standard deviations, indicating relatively small gradient stochasticity (high SNR). Gradient means are very small compared to batch-to-batch fluctuations (low SNR) in the second phase, which we call the *diffusion phase*. This transition is generally expected when the empirical error becomes saturated, and SGD is dominated by its fluctuations (Bertsekas; 2011). We claim that these distinct SGD phases (grey line in Figure 1.5a), correspond and explain the ERM and compression phases we observe in the information plane (marked green paths in Figure 1.4).

This dynamic phase transition occurs in the same number of iterations as in the layers' trajectories in the information plane. To relate the transition phase in the information plane to the gradients, we examine the mini-batch size effect on them. For each mini-batch size, we find both the starting point of the information compression and the gradient phase transition (the iteration where the derivative of the SNR is maximal). This is given in Figure 1.5b. The $X$-axis is the iteration where compression is started, and $Y$-axis is the iteration at which a phase transition occurs in the gradient. There is a linear trend between

**(a)** The change of weights, the means and standard deviations of the gradients for one layer, during the training (log-log scale). The grey line marks the transition between the *drift* to the *diffusion*

**(b)** The transition point of the SNR ($Y$-axis) versus the beginning of the information compression ($X$-axis), for different mini-batch sizes
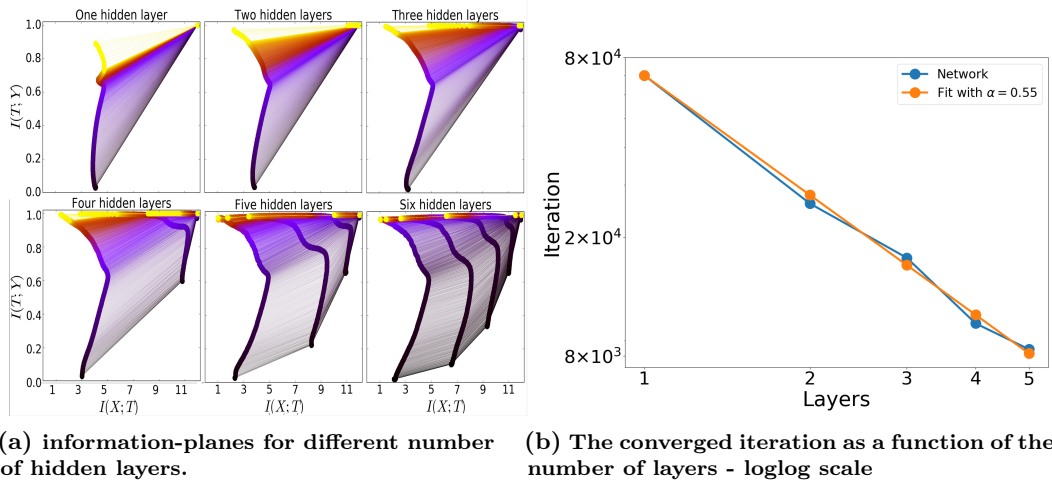
Figure 1.5: The Two Phases of SGD Optimization

the two.

In the drift phase, $I(T; Y)$ increases for every layer, as it rapidly reduces the empirical error. Conversely, the diffusion phase adds random noise to the weights, which evolve like Wiener processes under the training error or label information constraint. A Focker-Planck equation can describe such diffusion processes (see, e.g., Risken (1989)), whose stationary distribution maximizes the entropy of the distribution of the weight under the training error constraint. That, in turn, maximizes the conditional entropy, $H(X|T_i)$, or minimizes the mutual information $I(X; T_i) = H(X) - H(X|T_i)$, because the input entropy, $H(X)$ is constant. This entropy maximization by additive noise, also known as stochastic relaxation, is constrained by the empirical error or equivalently (for small errors) by $I(T; Y)$.

It remains unclear why different hidden layers converge to different points in the information-plane. Our analysis suggests that different layers have different noise levels in the gradients during the compression phase, explaining why they end up in different maximum entropy distributions. However, the gradient noises appear to vary and eventually decrease when the layers converge, suggesting that the convergence points are related to the *critical slowing down* of stochastic relaxation near phase transitions on the IB curve.

An interesting outcome of *compression by diffusion* is the randomized nature of the network's final weights. The correlations between the in-weights of different neurons in the same layer, which converge to essentially the same point in the plane, are very small. This indicates that there are many different networks with essentially optimal performance, and attempts to interpret single weights or even single neurons in such networks are meaningless.
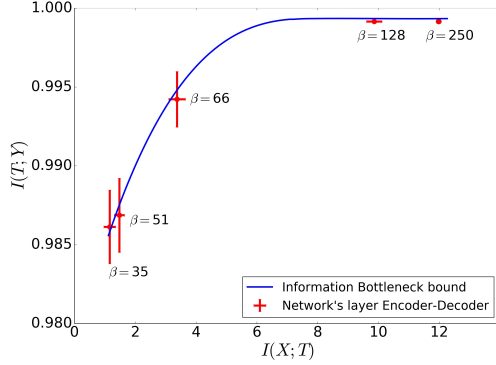
(a) information-planes for different number of hidden layers.

(b) The converged iteration as a function of the number of layers - loglog scale

Figure 1.6: The computational benefit of the layers

*The benefit of the hidden layers*

What is the benefit of the hidden layers? This question is one of the fundamentals of deep learning. For this reason, we train networks with a different number of layers $(1-5)$ and determine the minimum iteration for which the network converges. As before, we repeat each training 50 times with randomized initial weights and training samples. Figure 1.6a shows the information-plane paths for these six architectures during the training epochs, while Figure 1.6b shows the iteration number that the different networks reached 98% accuracy.

There are several important results from this experiment:

1. *The addition of hidden layers dramatically reduces the number of training iterations required for good generalization.* As more layers are added, the convergence time decreases by a factor of $k^{\frac{1}{\alpha}}$ where $alpha = 0.55$.

2. *When starting from previously compressed layers, each compression phase will be shorter.* By comparing the time to good generalization with four hidden layers and five hidden layers, we can see that convergence with four layers is much slower than with five or six hidden layers, where it takes half the time to reach the endpoints.

3. *The compression rate is faster for the deeper (narrower) layers, which are located closer to the output.* While in the drift phase, the lower layers move first (as a result of DPI), in the diffusion phase, the top layers compress first and pull the lower ones after them. Adding more layers seems to accelerate the compression.

**Figure 1.7: The DNN layers are fixed-points of the IB equations. The error bars represent standard error measures for $N = 50$. In each line, there are 5 points for the different layers. For each point, $\beta$ is the optimal value tof that layer.**

*The IB optimality of the layers*

Finally, to quantify the IB optimality of the layers, we test whether the converged layers satisfy the encoder-decoder relations of Eq. (**??**), for some value of the Lagrange multiplier $\beta$. With the encoder and decoder distributions based on quantized values of layer neurons, $p_i\left(t|x\right)$ and $p_i\left(y|t\right)$, respectively, we compute the information values for each converged layer.

To test the IB optimality of the layers, we calculate the optimal IB encoder, $p_{i,\beta}^{IB}\left(t|x\right)$ using the $i^{th}$ layer decoder, $p_i\left(y|t\right)$, through Eq.(**??**). This can be done for any value of $\beta$, with the known $P(X,Y)$. Then, we find the optimal $\beta_i$ for each layer, by minimizing the average KL divergence between the IB and the layer's encoders,

$$\beta_i^{\star} = \arg\min_{\beta} \mathbb{E}_x D_{KL}\left[p_i\left(t|x\right)||p_{\beta}^{IB}\left(t|x\right)\right] .$$

In figure 1.7 we see the information plotted over the layers together with the IB curve (blue line). The five empirical layers (trained with SGD) lie remarkably close to the theoretical IB limit. Furthermore, the slope of the curve, $\beta^{-1}$, matches their estimated optimal $\beta_i^{\star}$.

Hence, the DNN layers' encoder-decoder distributions satisfy the IB self-consistent equations within our numerical precision, decreasing $\beta$ as we go deeper in the network. The error bars are calculated over 50 randomized networks. As predicted by the IB equations, near the information curve $\Delta I_Y \sim \beta^{-1}\Delta I_X$.

*The evolution of the layers with training size*

Another problem in machine learning, which we just briefly discuss in this paper, is the importance of the training data size (Cho et al.; 2015). It is useful to visualize the hidden

layers' converged locations for different training data sizes in the information-plane (Figure 1.8).



**Figure 1.8: The effects of different training data sizes on the layers in the information-plane. Each line represents a converged network that had different training data sizes.**

As before, we train the network using six hidden layers, with different sample sizes ranging from 3 percent to 85 percent of the data. As expected, the layer's true label information, $I(T;Y)$ pushed up and gets closer to the theoretical IB bound on the rule distribution with increasing training size.

Despite the randomizations, the converged layers for different training sizes line up on a smooth line with remarkable regularity for each layer. We claim that the layers converge to specific points on the finite sample information curves, which can be calculated using the IB self-consistent equations (Eq. (**??**)), with the decoder replaced by the empirical distribution. This finite sample IB bound also explains the bounding shape on the left of Figure 1.4. Since the IB information curves are convex for any distribution, the layers converge to a convex curve in the plane even with very small samples.

However, the layers' training size effect is different for $I(T;Y)$ and $I(X;T)$. The training size hardly changes the information in the lower layers since even random weights keep most of the mutual information on both $X$ and $Y$. However, for deeper layers, the network learns to preserve the relevant information about $Y$ and compress the irrelevant information in $X$. More details on $X$ are relevant for $Y$ with a larger training set, and $I(X;T)$ increases in the middle layers.

### Related Work

Recently, the IB theory of deep learning has received much attention, including criticism of its rationale. Follow-up works attempt to address several points:

*Information compression*

Saxe et al. (2019) constructed several experiments to explore and refute the IB interpretation. They claim that the compression phase is an artifact of the quantization used to approximate $I(X;T)$ and the activation function (saturating nonlinearity). They observed no compression for the nonsaturating ReLU activation function.

Getting a reasonable estimate of mutual information between joint distributions in high-dimensional contexts is one of the main challenges in applying information-theoretic measures to real-world data. This problem has been extensively studied over the years (e.g., Paninski (2003)), showing that there is no "efficient" solution when the dimension of the problem is large and known approximations do not scale well with dimension and sample size (Gao et al.; 2015).

Several recent works have attempted to develop novel and efficient methods for estimating mutual information for DNNs. One line of works uses a generative decoder network (PixelCNN++) to estimate a lower bound on the mutual information (Darlow and Storkey; 2020; Nash et al.; 2018). In Darlow and Storkey (2020), the authors observed compression in the hidden layers during learning on the ImageNet dataset, using both ResNet and autoencoder architectures. This work confirms our two-stage learning for both classification and autoencoding tasks, characterized by (1) an initial short increase phase (2) following by a longer decrease in the mutual information with the input. Further, they observe that when the information is maximally compressed, the input images' class-irrelevant features are discarded; Namely, conditionally generated samples vary more while retaining information relevant to classification. Nash et al. (2018) reproduced our work using a similar approach. They observed a compression phase when using a convolution network with ReLU activation on the MNIST dataset. A two-phase behavior including a compression phase on the MNIST dataset with both fully-connected networks and convolutional networks was also reported by Noshad & Alfred (Noshad and Hero III; 2018). The authors introduced a new mutual information estimator, the ensemble dependency graph estimator (EDGE), which combines locality-sensitive hashing with dependency graphs and ensemble bias-reduction methods. The authors of Chelombiev et al. (2019) proposed adaptive approaches to estimating mutual information. These adaptive approaches compared the behavior of different activation functions and observed compression in DNNs with nonsaturating activation functions. It was observed that unlike saturating activation functions, compression does not always happen, and it is sensitive to initialization. They also observed that DNNs with L2 regularization strongly compress information.

In Goldfeld et al. (2018), the authors propose a new theoretical noisy entropy estimator to estimate mutual information. With both ReLU and linear activations, they observed the compression phase on our dataset and the MNIST dataset and related this compression behavior to geometric clustering. They also found that the behavior of $I(X;T)$ is strongly influenced by the "binning size," which is used for estimating the mutual information.

In Elad et al. (2018), the authors utilized the mutual information neural estimator (MINE) (Belghazi, Baratin, Rajeshwar, Ozair, Bengio, Hjelm and Courville; 2018), which estimates the KL divergence through the maximization of the dual representation of Donsker & Varadhan (Donsker and Varadhan; 1975). They showed that for the MNIST dataset with ReLU activation, the compression phase did not appear for "vanilla" cross-entropy training but did appear for training with weight decay regularization.

The authors in Achille et al. (2017) used Fisher information on the weights to demonstrate a two-phase learning process, involving an initial short increase, followed by a longer phase of decreasing information. Their paper described the reduction in information in the weights as implying a reduction in information in the activation. An additional confirmation on the two-phase behavior of DNNs can found in Li and Yuan (2017), who investigated shallow neural networks with residual connections and normal input distribution and showed that the SGD has two phases; (1) search and (2) convergence. Additionally, in Dieuleveut et al. (2017) they presented transient and stationary phases by looking at the inner product between successive mini-batch gradients in the network.

*Generalization and compression*

Another claim ofSaxe et al. (2019) is that the generalization does not require compression. They constructed a deep linear network toy example to illustrate generalization without compression.

The authors in Chelombiev et al. (2019) found that generalization accuracy was positively correlated with the degree of compression of the last layer. The connection of the generalization to compression is also discussed in Shwartz-Ziv et al. (2018). They showed that the generalization error depends exponentially on $I(X;T)$, once $I(T;X)$ becomes smaller than $\log 2n$ - the query sample complexity. Moreover, They showed that M bits of compression of $X$ are equivalent to an exponential factor of $2^M$ training examples. Furthermore, Achille and Soatto (2017) proved that flat minima, which have better generalization properties, bound information with the weights, and the information in the weights bound information in the activations.

*Information in deterministic networks*

Several works (Saxe et al.; 2019; Amjad and Geiger; 2018; Goldfeld et al.; 2018), state that the term $I(X;T)$ in the IB functional is theoretically either infinite or a constant for deterministic DNNs with continuous input. Thus, they wonder about the meaning of measuring it.

This problem was addressed in the literature by (1) training stochastic DNNs, which ensure the information is finite, and by (2) using a tractable variational approximation of

the mutual information (Alemi et al.; 2016; Kolchinsky and Tracey; 2017; Chalk et al.; 2016; Achille and Soatto; 2018; Belghazi, Rajeswar, Baratin, Hjelm and Courville; 2018). These works inject stochasticity by adding noise (or quantizing the input), only for quantifying the mutual information between the input and the hidden layer and not into the DNN itself. Goldfeld et al. (2018) introduced an auxiliary (noisy) DNN framework and showed that it is a good proxy for the original (deterministic) DNN both regarding performance and the learned representations.

To add noise, we can inject it directly into the representation (the activation values) and get a noisy variable we can measure. Saxe et al. (2019) used a non-parametric KDE estimator outlined by Kolchinsky and Tracey (2017), which directly adds small Gaussian noise to the data. However, all the training processes had a fixed level of noise, which led to failure due to the huge variation in the activation values (See (Chelombiev et al.; 2019) for detailed explanation).

It is also possible to add noise by discretizing the continuous variables into bins and approximating the representation as a discrete variable. In Goldfeld et al. (2018), they claimed that this estimation injects noise that is not present in the actual network, which is highly sensitive to selecting bin size and does not track $I(X;T)$ for different choices of the bin size. They developed a noisy DNNs framework and a rigorous estimator for $I(X;T)$. Using this estimator, they observed compression in various models. By relating $I(X;T)$ in the noisy DNN to an information-theoretic communication problem, they showed that compression is driven by the progressive clustering of hidden representations of inputs from the same class. They also proved that the estimator of $I(X;T)$ using binning is a measure for clustering.

Binning saturating activation functions facilitates mutual information estimation since all hidden activity is bounded within a predetermined range. However, with nonsaturating functions, the estimation procedure's noise level must be adapted for every layer of the network each time. In Saxe et al. (2019), the activation values are binned using a single range for the whole training, which is terminated by the maximum value across all epochs and all layers. However, since the network at every epoch is different, by this binning procedure, the network's estimation mixes unrelated factors (Chelombiev et al.; 2019).

Achille and Soatto (2018) stated that minimizing a stochastic network with an approximate compression term from the IB functional as a regularizer is equivalent to minimizing cross-entropy over deterministic DNNs with multiplicative noise (information dropout). Moreover, they proved that the special case of a Bernoulli noise results in the dropout method. In Elad et al. (2018), the authors showed that the binned information could be interpreted as a weight decay penalty, which aligns with common practice in DNN training.

*Invertible convolutional neural networks*

The invertible DNN (Jacobsen et al.; 2018) is a network's architecture that achieves state-of-the-art performance. The authors claimed that in such networks, $I(X;T_k)$ is never discarded regardless of the network parameters, and it is possible to reconstruct the input from each layer. However, although reversible networks do not discard the input's irrelevant components, we hypothesize that these networks progressively separate the irrelevant components from the relevant, allowing the final classification mapping to discard this information.

### Discussion and Conclusions

Motivated by the IB framework, our numerical experiments demonstrate that the visualization of the layers in the information-plane reveals many - so far unknown, details about the inner working of DNNs. They reveal the distinct phases of the SGD optimization, drift, and diffusion, explaining the ERM and the representation compression trajectories of the layers' information. The stochasticity of SGD methods is usually motivated by escaping local minima of the training error. In this paper, we give it a new, perhaps much more important role: It generates highly efficient internal representations through *compression by diffusion*. This is consistent with other recent suggestions on the role of noise in DNNs (Achille and Soatto; 2018; Kadmon and Sompolinsky; 2016).

We also argue that SGD seems an overkill during the diffusion phase, which consumes most of the training epochs, and that much simpler optimization algorithms, such as Monte-Carlo relaxations (Geman and Geman; 1988), can be more efficient. However, the IB framework may provide even more. If the layers converge to the IB theoretical bounds, there is an analytic connection between the encoder and decoder distributions for each layer, which can be exploited during training. Combining the IB iterations with stochastic relaxation methods may significantly boost DNN training.

To conclude, our analysis suggests that SGD with DNNs is, in essence, learning algorithms that effectively find efficient representations approximate minimal sufficient statistics in the IB sense.

### Bibliography

Achille, A., Rovere, M. and Soatto, S. (2017). Critical learning periods in deep neural networks, *arXiv preprint arXiv:1711.08856* .

Achille, A. and Soatto, S. (2017). Emergence of invariance and disentangling in deep representations, *arXiv preprint arXiv:1706.01350* .

Achille, A. and Soatto, S. (2018). Information dropout: Learning optimal representations through noisy computation, *IEEE transactions on pattern analysis and machine intelligence* **40**(12): 2897–2905.

Alain, G. and Bengio, Y. (2016). Understanding intermediate layers using linear classifier probes.

Alemi, A. A., Fischer, I., Dillon, J. V. and Murphy, K. (2016). Deep variational information bottleneck, *arXiv preprint arXiv:1612.00410* .

Amjad, R. A. and Geiger, B. C. (2018). How (not) to train your neural network using the information bottleneck principle, *arXiv preprint arXiv:1802.09766* .

Belghazi, I., Rajeswar, S., Baratin, A., Hjelm, R. D. and Courville, A. (2018). Mine: mutual information neural estimation, *arXiv preprint arXiv:1801.04062* .

Belghazi, M. I., Baratin, A., Rajeshwar, S., Ozair, S., Bengio, Y., Hjelm, R. D. and Courville, A. C. (2018). Mutual information neural estimation, *ICML*.

Bertsekas, D. P. (2011). Incremental gradient, subgradient, and proximal methods for convex optimization: A survey, *Optimization for Machine Learning* **2010**(1-38): 3.

Chalk, M., Marre, O. and Tkacik, G. (2016). Relevant sparse codes with variational information bottleneck, *Advances in Neural Information Processing Systems*, pp. 1957–1965.

Chelombiev, I., Houghton, C. and O'Donnell, C. (2019). Adaptive estimators show information compression in deep neural networks, *ICLR* .

Cho, J., Lee, K., Shin, E., Choy, G. and Do, S. (2015). How much data is needed to train a medical image deep learning system to achieve necessary high accuracy?, *arXiv preprint arXiv:1511.06348* .

Cover, T. M. and Thomas, J. A. (2006). *Elements of Information Theory (Wiley Series in Telecommunications and Signal Processing)*, Wiley-Interscience.

Darlow, L. N. and Storkey, A. (2020). What information does a resnet compress?, *arXiv preprint arXiv:2003.06254* .

Dieuleveut, A., Durmus, A. and Bach, F. (2017). Bridging the gap between constant step size stochastic gradient descent and markov chains, *arXiv preprint arXiv:1707.06386* .

Donsker, M. D. and Varadhan, S. S. (1975). Asymptotic evaluation of certain markov process expectations for large time, i, *Communications on Pure and Applied Mathematics* **28**(1): 1–47.

Elad, A., Haviv, D., Blau, Y. and Michaeli, T. (2018). The effectiveness of layer-by-layer training using the information bottleneck principle.

Gao, S., Ver Steeg, G. and Galstyan, A. (2015). Efficient estimation of mutual information for strongly dependent variables, *Artificial Intelligence and Statistics*, pp. 277–286.

Geman, S. and Geman, D. (1988). Stochastic relaxation, gibbs distributions, and the bayesian restoration of images, neurocomputing: foundations of research.

Goldfeld, Z., van den Berg, E., Greenewald, K., Melnyk, I., Nguyen, N., Kingsbury, B. and Polyanskiy, Y. (2018). Estimating Information Flow in Neural Networks, *ArXiv e-prints* .

Graves, A., Mohamed, A.-r. and Hinton, G. (2013). Speech recognition with deep recurrent neural networks, *Acoustics, speech and signal processing (icassp), 2013 ieee international conference on*, IEEE, pp. 6645–6649.

He, K., Zhang, X., Ren, S. and Sun, J. (2015). Deep residual learning for image recognition, *CoRR* **abs/1512.03385**.

Hinton, G. E., Srivastava, N., Krizhevsky, A., Sutskever, I. and Salakhutdinov, R. (2012). Improving neural networks by preventing co-adaptation of feature detectors, *CoRR* **abs/1207.0580**.
**URL:** *http://arxiv.org/abs/1207.0580*

Jacobsen, J.-H., Smeulders, A. W. and Oyallon, E. (2018). i-revnet: Deep invertible networks, *International Conference on Learning Representations*.
**URL:** *https://openreview.net/forum?id=HJsjkMb0Z*

Kadmon, J. and Sompolinsky, H. (2016). Optimal architectures in a solvable model of deep networks, *NIPS*.

Kazhdan, M., Funkhouser, T. and Rusinkiewicz, S. (2003). Rotation invariant spherical harmonic representation of 3d shape descriptors, *Eurographics Symposium on Geometry Processing* .

Kolchinsky, A. and Tracey, B. D. (2017). Estimating mixture entropy with pairwise distances, *Entropy* **19**(7): 361.

Larochelle, H., Bengio, Y., Louradour, J. and Lamblin, P. (2009). Exploring strategies for training deep neural networks, *J. Mach. Learn. Res.* **10**: 1–40.

LeCun, Y., Bengio, Y. and Hinton, G. (2015). Deep learning, *Nature* .

Li, Y. and Yuan, Y. (2017). Convergence analysis of two-layer neural networks with relu activation, *Advances in Neural Information Processing Systems*, pp. 597–607.

Nash, C., Kushman, N. and Williams, C. K. (2018). Inverting supervised representations with autoregressive neural density models, *arXiv preprint arXiv:1806.00400* .

Noshad, M. and Hero III, A. O. (2018). Scalable mutual information estimation using dependence graphs, *arXiv preprint arXiv:1801.09125* .

Paninski, L. (2003). Estimation of entropy and mutual information, *Neural Comput.* **15**(6): 1191–1253.

Risken, H. (1989). *The Fokker-Planck Equation: Methods of Solution and Applications*, number isbn9780387504988, lccn=89004059 in *Springer series in synergetics*, Springer-Verlag.

Rumelhart, D. E., Hinton, G. E. and Williams, R. J. (1986). Learning representations by back-propagating errors, *nature* **323**(6088): 533.

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.

Shwartz-Ziv, R., Painsky, A. and Tishby, N. (2018). Representation compression and generalization in deep neural networks.

Tishby, N., Pereira, F. C. and Bialek, W. (1999). The information bottleneck method, *In Proceedings of the 37-th Annual Allerton Conference on Communication, Control and Computing* .

Tishby, N. and Zaslavsky, N. (2015). Deep learning and the information bottleneck principle, *Information Theory Workshop (ITW), 2015 IEEE*, IEEE, pp. 1–5.

Zhang, X. and LeCun, Y. (2015). Text understanding from scratch, *CoRR* **abs/1502.01710**.
    **URL:** *http://arxiv.org/abs/1502.01710*

# Representation Compression and Generalization in Deep Neural Networks

# Representation Compression and Generalization in Deep Neural Networks

Ravid Shwartz-Ziv [1] Amichai Painsky [2] Naftali Tishby[1,2]

[1] The Edmond and Lilly Safra Center for Brain Sciences, The Hebrew University,
Jerusalem, Israel.
[2] School of Computer Science and Engineering,
The Hebrew University,
Jerusalem, Israel.

## Abstract

Understanding the groundbreaking performance of Deep Neural Networks (DNNs) is one of the most significant challenges to the scientific community today. In this work, we introduce an information-theoretic viewpoint on the behavior of deep network optimization processes and their generalization abilities. Specifically, we study DNNs on the information plane, the plane of the MI between each layer with the input variable, and the desired label, during the training dynamics. We show that we can characterize the network's training by a rapid increase in the mutual information (MI) between the layers and the target label, followed by a longer decrease in the MI between the layers and the input variable.

Furthermore, we explicitly show that these two fundamental information-theoretic quantities govern the network's generalization error by introducing a new generalization gap bound that is exponential in the input representation compression. The analysis focuses on typical patterns of large-scale problems. For this purpose, we introduce a novel analytic bound on the MI between consecutive layers in the network. An important consequence of our analysis is a superlinear boost in training time with the number of non-degenerate hidden layers, demonstrating the hidden layers' computational benefit.

### *Introduction*

Deep Neural Networks (DNNs) heralded a new era in predictive modeling and machine learning. Their ability to learn and generalize has set a new bar on performance compared to state-of-the-art methods. This improvement is evident across almost every application domain, especially in areas involving complicated dependencies between the input variable and the target label (LeCun et al.; 2015). However, despite their great empirical success, there is still no comprehensive understanding of their optimization process and its relationship to their (remarkable) generalization abilities.

This work examines DNNs from an information-theoretic viewpoint. For this purpose, we utilize the Information Bottleneck (IB) principle (Tishby et al.; 1999). The IB is a computational framework for extracting the most compact yet informative representation of the input variable ($X$) with respect to a target label variable ($Y$). The IB bound defines the optimal tradeoff between representation complexity and its predictive power. Specifically, it is achieved by minimizing the mutual information (MI) between the representation and the input, subject to the level of MI between the representation and the target label.

Recent results (Shwartz-Ziv and Tishby; 2017), demonstrated that the layers of DNNs tend to converge to the IB optimal bound. The results pointed to a distinction between the two phases of the training process. The first phase is characterized by an increase in the MI with the label (i.e., fitting the training data), whereas in the second and most important phase, the training error slowly reduces the MI between the layers and the input (i.e., representation compression). These two phases appear to correspond to fast convergence to a flat minimum (drift) following a random walk, or diffusion, in the vicinity of the training error's flat minimum, as reported in other studies (e.g., (Zhang, Liao, Rakhlin, Miranda, Golowich and Poggio; 2018)).

These observations raised several interesting questions: (a) which properties of the SGD optimization cause these two training phases? (b) how can the diffusion phase improve the generalization performance? (c) can the representation compression explain the convergence of the layers to the optimal IB bound? (d) can this diffusion phase explain the benefit of many hidden layers?

In this work, we attempt to answer these questions. Specifically, we draw important connections between recent results inspired by statistical mechanics and information-theoretic principles. We show that the layers of a DNN indeed follow the behavior described by Shwartz-Ziv and Tishby (2017). We claim that the reason is the Stochastic Gradient Descent (SGD) optimization mechanism. We show that the first phase of the SGD is characterized by a rapid decrease in the training error, which corresponds to an increase in the MI with the labels. Then, the SGD behaves like a nonhomogeneous Brownian motion in the weights space, in the proximity of a flat error minimum. This nonhomogeneous diffusion corresponds to a decrease in MI between the layers and the input variable, in "directions" that are irrelevant to the target label.

One of the main challenges in applying information-theoretic measures to real-world data is a reasonable estimation of high-dimensional joint distributions. This problem has been extensively studied over the years (e.g., (Paninski; 2003)) and has led to the conclusion that there is no "efficient" solution when the dimension of the problem is large. Recently, several studies have focused on calculating the MI in DNNs using statistical mechanics. These methods have generated promising results in various special cases (Gabrié et al.; 2018).

In this work, we provide an analytic bound on the MI between consecutive layers, which is valid for any nonlinearity of the units and directly demonstrates the representation's compression during the diffusion phase. Specifically, we derive a Gaussian bound that only depends on the linear part of the layers. This bound gives a superlinear dependence of the layers' convergence time, which enables us to prove the superlinear computational benefit of the hidden layers. Furthermore, the Gaussian bound allows us to study the MI in DNNs in real-world data without estimating it directly.

*Preliminaries and Notations*

Let $X \in \mathcal{X}$ and $Y \in \mathcal{Y}$ be a pair of random variables of the input patterns and their target label (respectively). We consider the practical setting where $X$ and $Y$ are continuous random variables represented in a finite precision machine throughout this work. This means that both $X$ and $Y$ are practically binned (quantized) into a finite number of discrete values. Alternatively, $X, Y$, may be considered as continuous random variables measured in the presence of small independent additive (Gaussian) noise, corresponding to their numerical precision. We use these two interpretations interchangeably, at the limit of infinite precision, where the limit is applied at the final stage of our analysis.

We denote the joint probability of $X$ and $Y$ as $p(x, y)$, whereas their corresponding MI is defined as $I(X; Y) = D\left[p(y|x)||p(y)\right] = D\left[p(x|y)||p(x)\right]$. We use the standard notation $D[p||q]$ for the Kullback-Liebler (KL) divergence between the probability distributions $p$ and $q$. Let $f_{W^K}(x)$ denote a DNN, with $K$ hidden layers, where each layer consists of $d_k$ neurons, each with some activation function $\sigma_k(x)$, for $k = 1, \ldots, K$. We denote the values of the $k^{th}$ layer by the random vector $T_k$. The DNN mapping between two consecutive layers is defined as $T_k = \sigma_k (W_k T_{k-1})$, where $W_k$ is a $d_k \times d_{k-1}$ real weight matrix. Note that we consider both the weights, $W_k$ and the layer representations, $T_k$, as stochastic entities because they depend on the network's stochastic training rule and the random input pattern (as described in the next section). However, when the network weights are given, the weights are fixed realizations of the random training process (i.e., they are "quenched"). Note that given the weights, the layers form a Markov chain of successive internal representations of the input variable $X$: $Y \rightarrow X \rightarrow T_1 \rightarrow ... \rightarrow T_K$, and their MI values obey a chain of Data Processing Inequalities (DPI), as discussed by Shwartz-Ziv and Tishby (2017).

We denote the set of all $K$ layers weight matrices as $W^K = \{W_1, \ldots, W_K\}$. Let the *training sample*, $S^n = \{(x_1, y_1), \ldots, (x_n, y_n)\}$ be a collection of $n$ independent samples from $p(x, y)$. Let $\ell_{W^K}(x_i, y_i)$ be a (differentiable) loss function that measures the discrepancy between a prediction of the network $f_{W^K}(x_i)$ and the corresponding true target value $y_i$, for a given set of weights $W^K$. Then, the empirical error is defined as $\mathcal{L}_{W^K}(S^n) = \frac{1}{n} \sum_{i=1}^n \ell_{W^K}(x_i, y_i)$. The corresponding error gradients (with respect to the weights) are denoted as $\nabla_{W^K} \mathcal{L}_{W^K}(S^n)$.

*Deep Neural Networks*

DNNs opened a new era of machine learning capabilities. They dramatically improved predictive abilities as they outperform state-of-the-art methods in a wide variety of fields, ranging from visual object recognition, speech recognition to drug discovery, genomics, and automatic game playing (Graves et al.; 2013). DNNs are multilayer data structures, where each layer consists of multiple simple processing units called neurons (as an analogy to the structure of a biological brain). Each layer's neurons are defined as a linear function of the neurons of the previous layer, followed by a nonlinear activation function.

*Training the Network – the SGD Algorithm*

Training a DNN corresponds to setting the weights $W^K$ from a given set of samples $S^n$. It is typically done by minimizing the empirical error, which approximates the expected loss. The SGD algorithm is a standard optimization method for this purpose (Robbins and Monro; 1951).

Let $S^{(m)}$ be a random set of $m$ samples drawn (uniformly, with replacement) from $S^n$, where $m < n$. We refer to $S^{(m)}$ as a *mini-batch* of $S^n$. Define the corresponding empirical error and gradient of the minibatch as $\mathcal{L}_{W^K}(S^{(m)}) = \frac{1}{m} \sum_{\{x_i, y_i\} \in S^{(m)}} \ell_{W^K}(x_i, y_i)$ and $\nabla_{W^K} \mathcal{L}_{W^K}(S^{(m)}) = \frac{1}{m} \sum_{\{x_i, y_i\} \in S^{(m)}} \nabla_{W^K} \ell_{W^K}(x_i, y_i)$ respectively. Then, the SGD algorithm is defined by the update rule: $W^K(l) = W^K(l-1) - \eta \nabla_{W^K(l-1)} \mathcal{L}_{W^K(l-1)}(S^{(m)})$, where $W^K(l)$ are the weights after $l$ iterations of the SGD algorithm and $\eta \in \mathbb{R}_+$ is the learning rate.

*The Different Phases of SGD Optimization*

The SGD algorithm plays a key role in the astonishing performance of DNNs. As a result, it has been extensively studied in recent years, especially in the context of flexibility and generalization (Chee and Toulis; 2017). Here, we examine the SGD as a stochastic process that can be decomposed into two separate phases. This idea has been studied in several works (Murata; 1998; Jin et al.; 2017; Hardt et al.; 2015). Murata argued that stochastic iterative procedures are initiated at some starting state and then move through

a fast *transient phase* towards a *stationary phase*, where the distribution of the weights becomes time-independent. However, this may not be when the SGD induces non-isotropic state-dependent noise, as argued, for example, by Chaudhari and Soatto (2017).

In contrast, Shwartz-Ziv and Tishby (2017) described the transient phase of the SGD as having two very distinct dynamic phases. The first is a *drift* phase, where the means of the error gradients in every layer are large compared to their batch-to-batch fluctuations. This behavior indicates small variations in the gradient directions, or *high-SNR gradients*. In the second part of the transient phase, which they refer to as *diffusion*, the gradient means become significantly smaller than their batch-to-batch fluctuations – *low-SNR gradients*. The transition between the two phases occurs when the training error saturates and the weights growth is dominated by the gradient batch-to-batch fluctuations. Typically, most SGD updates are expended in the diffusion phase before reaching Murata's stationary phase. In this work, we rigorously argue that this diffusion phase causes the representation compression, the observed reduction in $I(T_k; X)$, for most hidden layers.

*Drift and Diffusion with SGD*

It is well known that the discrete-time SGD can be considered as an approximation of a continuous-time stochastic gradient flow if the discrete-time iteration parameter $l$ is replaced by a continuous parameter $\tau$. Li et al. (2015) showed that when the minibatch gradients are unbiased with bounded variance, the discrete-time SGD is an approximation of a continuous-time Langevin dynamics,

$$dW^K(\tau) = -\nabla_{W^K(\tau)} \mathcal{L}_{W^K(\tau)}(S_n) \, d\tau + \sqrt{2\beta^{-1} C\left(W^K(\tau)\right)} dB(\tau) \qquad (1.21)$$

where $C\left(W^K(\tau)\right)$ is the sample covariance matrix of the weights, $B(\tau)$ is a standard Brownian motion (Wiener process) and $\beta$ is the Langevin temperature constant. The first term in (**??**) is called the gradient flow or drift component, and the second term corresponds to random diffusion. Although this stochastic dynamics holds for the entire SGD training process, the first term dominates the process during the high SNR gradient phase, while the second term becomes dominant when the gradients are small, due to the saturation of the training error in the low SNR gradient phase. Hence, these two SGD phases are referred to as drift and diffusion.

The *mean $L_2$ displacement* (MSD) measures the Euclidean distance from a reference position over time, which is used to characterize a diffusion process. Normal diffusion processes are known to exhibit a power-law MSD in time, $\mathbb{E}\left[\left\|W^K(\tau) - W^K(0)\right\|_2\right] \sim \gamma t^\alpha$, where $t$ is the diffusion time, $\gamma$ is related to the diffusion coefficient, and $0 < \alpha \leq 0.5$ is the diffusion exponent. For a standard flat space diffusion, the MSD increases as a square root of time ($\alpha = 0.5$). Hu et al. (2017) showed (empirically) that the weights' MSD in a DNN trained with SGD, indeed behaves (asymptotically) like a normal diffusion, where the

diffusion coefficient $\gamma$ depends on the batch size and learning rate. In contrast, Hoffer et al. (2017) showed that the weights' MSD demonstrates a much slower logarithmic increase. This type of dynamics is also called "ultra-slow" diffusion.

### Information Plane Analysis

Following Tishby and Zaslavsky (2015) and Shwartz-Ziv and Tishby (2017), we study the layer representation dynamics in the two-dimensional plane $(I(X;T_k), I(T_k;Y))$. Specifically, for any input and target variables, $X, Y$, let $T \triangleq T(X)$ denote a representation, or an encoding (not necessarily deterministic), of $X$. Clearly, $T$ is fully characterized by its *encoder*, the conditional distribution $p(t|x)$. Similarly, let $p(y|t)$ denote any (possibly stochastic) *decoder* of $Y$ from $T$. Given a joint probability function $p(x, y)$, the *information plane* is defined the set of all possible pairs $I(X;T)$ and $I(T;Y)$ for any possible representation, $p(T|X)$.

It is evident that not all points on the plane are feasible (achievable), as there is a tradeoff between these quantities; the more we compress $X$ (reduce $I(X;T)$), the less information can be maintained about the target, $I(T;Y)$.

Our analysis is based on the fundamental role of these two MI quantities. We argue that for large-scale (high dimensional $X$) learning, for almost all (*typical*) input patterns, with mild assumptions (ergodic Markovian input patterns): (i) the MI values concentrate with the input dimension, (ii) the minimal sample complexity for a given generalization gap is controlled by $I(X;T)$, and (iii) the accuracy - the generalization error - is governed by $I(T;Y)$, with the Bayes optimal decoder representation.

We argue that these two MI quantities characterize the sample-size and the accuracy tradeoff of large-scale representation learning. For DNNs, this amounts to a dramatic reduction in the complexity of the analysis of the problem. We discuss these ideas in the following sections and prove the connection between the input representation compression, the generalization gap (the difference between training and generalization errors), and the minimal sample complexity (Theorem 1 below).

### Label Information and Generalization Error

The optimization of MI quantities is not a new concept in supervised or unsupervised learning (Deco and Obradovic; 2012; Linsker; 1988; Painsky et al.; 2016). This is not surprising, as it can be shown that $I(T;Y)$ corresponds to the irreducible error when minimizing the logarithmic loss (Painsky and Wornell; 2018b; Harremoes and Tishby; 2007). Here, we emphasize that $I(T;Y)$, for the optimal decoder of the representation $T$, governs all reasonable generalization errors (under the mild assumption that label $y$ is not entirely deterministic; $p(y|x)$ is in the interior of the simplex, $\Delta(Y)$, for all typical $x \in X$). First,

note that for the Markov chain $Y - X - T$, $I(T;Y) = I(X;Y) - \mathbb{E}_{X,T} D\left[p(y|x)||p(y|t)\right]$. By using the Pinsker inequality (Cover and Thomas; 2012), the variation distance between the optimal and the representation decoders can be bounded by their KL divergence,

$$D\left[p(y|x)||p(y|t)\right] \geq \frac{1}{2\ln 2}\left|p(y|x) - p(y|t)\right|_1^2. \tag{1.22}$$

Hence, by maximizing $I(T;Y)$ we minimize the expected *variation risk* between the representation decoder $p(y|t)$ and $p(y|x)$. For more similar bounds on the error measures see (Painsky and Wornell; 2018a).

*Representation Compression and Sample Complexity*

The *Minimum Description Length* (MDL) principle (Rissanen; 1978) suggests that the best representation for a given set of data is the one that leads to the minimal code length needed to represent the data. This idea has inspired the use of $I(X;T)$ as a regularization term in many learning problems (e.g., Chigirev and Bialek (2004); Painsky et al. (2018)). Here, we argue that $I(X;T)$ plays a much more fundamental role; we show that for large scale learning (high dimensional $X$) and typical input patterns, $I(X;T)$ controls the sample complexity of the problem, given a generalization error gap.

**Theorem 1** (Input Compression bound). *Let $X$ be a d-dimensional random variable that obeys an ergodic Markov random field probability distribution, asymptotically in d. Let $T \triangleq T(X)$ be a representation of $X$ and denote by $T_m = \{(t_1, y_1), \ldots, (t_m, y_m)\}$ an m-sample vector of $T$ and $Y$, generated with m independent samples of $x_i$, with $p(y|x_i)$ and $p(t|x_i)$. Assume that $p(x,y)$ is bounded away from 0 and 1 (strictly inside the simplex interior). Then, for large enough d, with probability $1 - \delta$, the typical expected squared generalization gap satisfies*

$$\left|\mathcal{L}\left(T_m\right) - \mathbb{E}_{T_m}\left[\mathcal{L}\left(T_m\right)\right]\right|^2 \leq \frac{2^{I(X;T)} + \log\frac{2}{\delta}}{2m}. \tag{1.23}$$

*where the typicality follows the standard Asympthotic Equipartition Property (AEP) (Cover and Thomas; 2012).*

The proof of this theorem is given in Appendix A. This theorem is also related to the bound proved by Shamir et al. (2010), with the typical representation cardinality, $|T(X)| \approx 2^{I(T;X)}$. The ergodic Markovian assumption is common in many large-scale learning problems. It means that $p(x) \approx \prod_{i=1:d} p(x_i|Pa(x_i))$, where $Pa(x_i)$ is a finite set of adjacent "parents" of $x_i$ in the $d$ dimensional pattern $X$.

The consequences of this input-compression bound are quite striking: the generalization error decreases exponentially with $I(X;T)$, once $I(T;X)$ becomes smaller than $\log 2m$ - the query sample-complexity. Moreover, it means that $M$ bits of representation compression,

beyond $\log 2m$, are equivalent to a factor of $2^M$ training examples. The tightest bound on the generalization bound is obtained for the most compressed representation or the last hidden layer of the DNN. The input-compression bound can yield a tighter and more realistic sample complexity than any of the worst-case PAC bounds with any reasonable estimate of the DNN class dimensionality, as typically, the final hidden layers are compressed to a few bits.

Nevertheless, two important caveats are in order. First, the layer representation in deep learning is learned from the training data; hence, the encoder, the partition of the typical patterns $X$, and the *effective "hypothesis class"*, depend on the training data. This can lead to considerable overfitting. Training with SGD avoids this potential overfitting because of the way the diffusion phase works. Second, for low $I(T;Y)$, there are exponentially (in $d$) many random encoders (or soft partitions of $X$) with the same value of $I(T;X)$. This seems to suggest that there is a missing exponential factor in our estimate of the hypothesis class cardinality. However, note that the vast majority (almost all) of these possible encoders are never encountered during a typical SGD optimization. Moreover, as $I(T;Y)$ increases, the number of such random encoders rapidly collapses to $O(1)$ when $I(T;Y)$ approaches the optimal IB limit, as we show next.

*The Information Bottleneck Limit*

As presented above, we are interested in the boundary of the achievable region in the information plane, or encoder-decoder pairs that minimize the sample complexity (minimize $I(X;T)$) and generalize well (maximize $I(T;Y)$).

These optimal encoder-decoder pairs are given precisely by the IB framework (Tishby et al.; 1999), which is formulated by the following optimization problem: $\min_{p(t|x)} I(X;T) - \beta I(T;Y)$, over all possible encoders-decoders pairs that satisfy the Markov condition $Y - X - T$. Here, $\beta$ is a positive Lagrange multiplier associated with the decoder information on $I(T;Y)$, which determines the representation's complexity.

The IB limit defines the set of optimal encoder-decoder pairs for the joint distribution $p(x,y)$. Furthermore, it characterizes the achievable region in the information plane, similar to Shannon's rate-distortion theory (Cover and Thomas; 2012). This analysis, also determines the optimal tradeoff between sample complexity and generalization error. The IB can only be solved analytically in exceptional cases (e.g., jointly Gaussian $X, Y$ (Chechik et al.; 2005)). In general, a (locally optimal) solution can be found by iterating the self-consistent equations, similar to the Blahut-Arimoto algorithm in rate-distortion theory (Tishby et al.; 1999). For general distributions, no efficient algorithm for solving the IB is known, though there are several approximation schemes (Chalk et al.; 2016; Painsky and Tishby; 2017).

The self-consistent equations are satisfied along the *information curve*. This monotonic curve separates between the achievable and non-achievable regions on the information plane.

Notice that for smooth joint distributions $p(x, y)$, the information curve is strictly concave with a unique slope, $\beta^{-1}$, at every point, and a finite slope at the origin. In these cases of interest (where $Y$ is not a deterministic function of $X$), every value of $\beta$ corresponds to a single point on the information curve with a corresponding optimal encoder-decoder pair.

### The Information Plane and SGD Dynamics for DNNs

By applying the DPI to the Markov chain of the DNN's layers, we obtain the following chains:

$$I(X; Y) \geq I(T_1; Y) \geq I(T_2; Y) \geq ...I(T_k; Y) \geq I(\hat{Y}; Y)$$
$$H(X) \geq I(X; T_1) \geq I(X; T_2) \geq ...I(X; T_k) \geq I(X; \hat{Y}).$$

where $\hat{Y}$ is the output of the network.

The pairs $(I(X; T_k), I(T_k, Y))$, for each SGD update, form a unique concentrated *information path* for each layer of a DNN, as demonstrated by Shwartz-Ziv and Tishby (2017).

As we are only interested in the information that flows through the network, invertible transformations of the representations that preserve information generate equivalent representations even if the individual neurons encode entirely different input features. Therefore, we quantify the representations by two numbers, or order parameters for each layer - the MI of $T$ with the input $X$ (the decoder's information) and with the desired output $Y$ (the information of the encoder). These quantities are invariant to any invertible re-parameterization of $T$.

For any fixed realization of the weights, the network is, in principle, a deterministic map. This does not imply that information is not lost between the layers; the layers' inherent finite precision, with possible saturation of the nonlinear activation function $\sigma_k$, can result in non-invertible mapping between the layers. Moreover, we argue below that for large networks, this mapping becomes effectively stochastic due to the diffusion phase of the SGD.

On the other hand, the paths of the layers in the information plane are invariant to invertible transformations of the representation $T_k$. Thus, the same paths are shared by very different weights and architectures and possibly different encoder-decoder pairs. This freedom is drastically reduced when the target information, $I(T_k, Y)$, increases, and the layers approach the IB limit. Minimizing the training error, together with standard uniform convergence arguments, clearly increases $I(T; Y)$. This raised the question what in the SGD dynamics can lead to the observed representation compression, which further improves the generalization? Moreover, can the SGD dynamics push the layer representations to the IB limit, as claimed in Shwartz-Ziv and Tishby (2017)?

We provide affirmative answers to both questions, using the properties of the drift and diffusion phases of the SGD dynamics.

*Representation Compression by Diffusion*

This section quantifies the roles of the drift and diffusion SGD phases and their influence on the MI between consecutive layers. Specifically, we show that the drift phase corresponds to an increase in information with the target label $I(T_k; Y)$, whereas the diffusion phase corresponds to representation compression or reduction of the $I(X; T_k)$. The representation compression is accompanied by further improvement in the generalization.

The general idea is as follows: the drift phase increases $I(T_k; Y)$ as it reduces the cross-entropy empirical error. On the other hand, the diffusion phase in high-dimensional weight space effectively adds an independent nonuniform random component to the weights, mainly in the directions that do not influence the loss - i.e., *irrelevant directions*. This results in a reduction of the SNR of the patterns' irrelevant features, which leads to a reduction in $I(X; T_k)$, or representation compression. We further argue that different layers filter out different irrelevant features, resulting in their convergence to different information plane locations.

*The SGD Compression Mechanism*

First, DPI implies that $I(X; T_k) \leq I(T_{k-1}; T_k)$. We focus on the second term during the diffusion phase and prove an asymptotic upper bound for $I(T_{k-1}; T_k)$, which reduces sub-linearly with the number of SGD updates.

For clarity, we describe the case where $T_k \in \mathbb{R}^{d_k}$ is a vector and $T_{k+1} \in \mathbb{R}$ is a scalar. The generalization to higher $d_{k+1}$ is straightforward. We examine the network during the diffusion phase, after $\tau$ iterations of the SGD beyond the drift-diffusion transition. For each layer, $k$, the weights matrix, $W^k(\tau)$ can be decomposed as follows,

$$W^k(\tau) = W^{k\star} + \delta W^k(\tau). \tag{1.24}$$

The first term, $W^{k\star}$, denotes the weights at the end of the drift phase ($\tau_0 = 0$) and remains constant with increasing $\tau$. As we assume that the weights converge to a (local, flat) optimum during the drift phase, $W^{k\star}$ is close to the weights at this local optimum. The second term, $\delta W^k(\tau)$, is the accumulated Brownian motion in $\tau$ steps due to the batch-to-batch fluctuations of the gradients near the optimum. For large $\tau$ we know that $\delta W^k(\tau) \sim \mathcal{N}(0, \tau C(W^k(\tau_0)))$ where $\tau_0$ is the time of the beginning of the diffusion phase. Note that at any given $\tau$, we can treat the weights as a fixed (quenched) realization, $w^k(\tau)$, of the random Brownian process $W^k(\tau)$. We can now model the mapping between the

layers $T_k$ and $T_{k+1}$ at that time as

$$T_{k+1} = \sigma_k \left( w^{*T} T_k + \delta w^k(\tau)^T T_k + Z \right) \tag{1.25}$$

where $w^* \in \mathbb{R}^{d_k}$ is the SGD's empirical minimizer, and $\delta w \in \mathbb{R}^{d_k}$ is a realization from a Gaussian vector $\delta w \sim \mathcal{N}(0, C_{\delta w})$, of the Brownian process. In addition, we consider $Z \sim \mathcal{N}(0, \sigma_z^2)$ to be the small Gaussian measurement noise, or quantization, independent of $\delta w^k$ and $T_k$. This standard additive noise allows us to treat all the random variables as continuous.

For simplicity, we assume that the $d_k$ components of $T_k$ have zero mean and are asymptotically independent for $d_k \to \infty$, and that $\lim_{d_k \to \infty} w^{*T} \delta w = 0$ almost surely.

**Proposition 2.** *Assume that the moments of $T_k$ are finite. Further assume that the components of $w^*$ and $\delta w(\tau)$ are in-general-positions, satisfying $\lim_{d_k \to \infty} \sum_{i=1}^{d_k} w_i^{*4} / \left( \sum_{i=1}^{d_k} w_i^{*2} \right)^2 = 0$ and $\lim_{d_k \to \infty} \sum_{i=1}^{d_k} \delta w_i^4 / \left( \sum_{i=1}^{d_k} \delta w_i^2 \right)^2 = 0$ almost surely. Then,*

$$\frac{1}{\sqrt{\sigma_{T_k}^2}} \left[ \frac{w^{*T} T_k}{||w^*||_2} \quad \frac{\delta w^T T_k}{||\delta w||_2} \right]^T \xrightarrow[d_k \to \infty]{\mathcal{D}} \mathcal{N}(0, I) \tag{1.26}$$

*almost surely, where $\sigma_{T_k}^2$ is the variance of the components of $T_k$.*

A proof for this CLT proposition is given in Appendix B.

Proposition 2 shows that under the standard conditions above, $w^{*T} T_k$ and $\delta w^T T_k$ are asymptotically jointly Gaussian and independent, almost surely. We stress that the components of $T_k$ do not have to be identically distributed to satisfy this property; Proposition 2 may be adjusted for this case with different normalization factors. Similarly, the i.i.d. assumption on $T_k$ can be relaxed to Markovian ergodic. It is easy to verify that Proposition 2 can be extended to the general case where $w^*, \delta w \in \mathbb{R}^{d_k \times d_{k+1}}$, under similar general position conditions, with almost sure orthogonality of $w^*$ and $\delta w$.

We can now bound the mutual information between $T_{k+1}$ and the linear projection of the previous layer $W^* T_k$, during the diffusion phase, for sufficiently high dimensions $d_k, d_{k+1}$, under the above conditions. Note that in this case, Equation **??** behaves like an additive Gaussian channel where $w^{*T} T_k$ is the signal and $\delta w^T T_k + Z$ is an independent additive Gaussian noise (i.e., independent of signal and normally distributed). Hence, for

sufficiently large $d_k$ and $d_{k+1}$, we can write

$$I(T_{k+1}; T_k | w^*) \leq I(T_{k+1}; w^{*T} T_k | w^*) \tag{1.27}$$

$$\leq I\left(w^{*T} T_k + \delta w^T T_k + Z; w^{*T} T_k | w^*\right)$$

$$= \frac{1}{2} \log \left( \frac{\left| \sigma_{T_k}^2 w^{*T} w^* + \sigma_{T_k}^2 \delta w^T \delta w + \sigma_z^2 I \right|}{\left| \sigma_{T_k}^2 \delta w^T \delta w + \sigma_z^2 I \right|} \right)$$

almost surely, where the first inequality is due to DPI for the Markov chain $T_k - w^{*T} T_k - T_{k+1}$. Finally, we apply an orthogonal eigenvalue decomposition to the multivariate Gaussian channel in Equation **??**. Let $\delta w^T \delta w = Q \Lambda Q^T$ where $Q Q^T = I$ and $\Lambda$ is a diagonal matrix whose diagonal elements are the corresponding eigenvalues, $\lambda_i$, of $\delta w^T \delta w$. Then, we have that

$$\left| \sigma_{T_k}^2 w^{*T} w^* + \sigma_{T_k}^2 \delta w^T \delta w + Z \right| = \sigma_{T_k}^2 |Q| \cdot |Q^T w^{*T} w^* Q + \Lambda + \frac{\sigma_z^2}{\sigma_{T_k}^2} Q^T Q| \cdot |Q^T| \tag{1.28}$$

$$= \sigma_{T_k}^2 |Q^T w^{*T} w^* Q + \Lambda + \frac{\sigma_z^2}{\sigma_{T_k}^2} I|$$

$$\leq \sigma_{T_k}^2 \prod_{i=1}^{d_{k+1}} \left( A_{ii} + \lambda_i + \frac{\sigma_z^2}{\sigma_{T_k}^2} \right)$$

where $A \triangleq Q^T W^{*T} W^* Q$. The last inequality is due to the Hadamard inequality. Plugging Equation **??** into Equation **??** yields that for sufficiently large $d_k$ and $d_{k+1}$,

$$I(T_{k+1}; T_k | w^*) \leq \frac{1}{2} \log \left( \frac{\prod_{i=1}^{d_{k+1}} \left( A_{ii} + \lambda_i + \frac{\sigma_z^2}{\sigma_{T_k}^2} \right)}{\prod_{i=1}^{d_{k+1}} \left( \lambda_i + \frac{\sigma_z^2}{\sigma_{T_k}^2} \right)} \right) \tag{1.29}$$

$$= \frac{1}{2} \sum_{i=1}^{d_{k+1}} \log \left( 1 + \frac{A_{ii}}{\lambda_i + \frac{\sigma_z^2}{\sigma_{T_k}^2}} \right) \xrightarrow[\sigma_z^2 \to 0]{} \frac{1}{2} \sum_{i=1}^{d_{k+1}} \log \left( 1 + \frac{A_{ii}}{\lambda_i} \right).$$

As previously established, $\delta w$ is a Brownian motion along the SGD iterations during the diffusion phase. This process is characterized by a low (and fixed) variance of the informative gradients (relevant dimensions), whereas the remaining irrelevant directions suffer from increasing variances as the diffusion proceeds (see, e.g. Sagun et al. (2017); Zhu et al. (2018); Jastrzebski et al. (2017)). In other words, we expect the "informative" $\lambda_i$ to remain fixed, while the irrelevant consistently grow as sub-linearly with time. Denote the set of "informativ" directions as $\Lambda^*$ and the set of "non-informative" as $\Lambda_{NI}$. Then our final limit, as the number of SGD steps grows, is

$$I(T_{k+1}; T_k | w^*) \leq \frac{1}{2} \sum_{\lambda_i^* \in \Lambda^*} \log \left( 1 + \frac{A_{ii}}{\lambda_i^*} \right).$$

Note that for real problems the distinction between informative and non-informative directions may not be that sharp and we can expect a gradual (exponential asymptotically) decrease of $A_{ii}$ with $i \to \infty$. Which directions are compressed and which are preserved depend on the required compression level. This is why different layers converge to different values of $I(T_k; X)$.

*Relation to Other Works*

The analysis above suggests that the SGD compresses during the diffusion phase in many directions of the gradients. We argue that these directions are the ones in which the gradients' variance is increasing (non-informative), whereas the information is preserved in the directions where the variance of the gradients remains small.

This statement is consistent with recent work on the statistical properties of gradients and generalization (Sagun et al.; 2017; Zhu et al.; 2018; Zhang, Saxe, Advani and Lee; 2018). These works showed that the gradients' covariance matrix is typically highly non-isotropic and that this is crucial for generalization by SGD. They suggested that the reason lies in the gradients' covariance matrix's proximity to the Hessian of the loss approximation. Furthermore, Zhang, Saxe, Advani and Lee (2018); Keskar et al. (2016); Jastrzebski et al. (2017) argued that SGD tends to converge to flat minima, which often results in a better generalization. Zhang, Saxe, Advani and Lee (2018) emphasized that SGD converges to flat minima values characterized by high entropy due to the non-isotropic nature of the gradients' covariance and its alignment with the error Hessian at the minima. In other words, the findings above suggest that non-isotropic gradients and Hessian typically characterize good generalization performance in orthogonal directions to the flat minimum of the training error objective.

**The Computational Benefit of the Hidden Layers**

Our Gaussian bound on the representation compression (Equation **??**) allows us to relate the convergence time of the layer representation information, $I(T_k; X)$, to the diffusion exponent $\alpha$, defined above. Denote the representation information at the diffusion time $\tau$ as $I(X; T_k)(\tau)$. It follows from Eqaution **??** that

$$I(X; T_k)(\tau) \leq C + \frac{1}{2} \sum_{\lambda_i \in \Lambda^{NI}} \log\left(1 + \frac{A_{ii}}{\lambda_i(\tau)}\right) \leq C + \frac{1}{2} \sum_{\lambda_i \in \Lambda^{NI}} \left(\frac{A_{ii}}{\lambda_i(\tau)}\right) \qquad (1.30)$$

where $C$ depends on the informative for this layer, but not on $\tau$.

Notice that $\lambda_i(\tau)$ are the singular values of the weights of a diffusion process, which

grow as $\tau^\alpha$, where $\alpha$ is the diffusion exponent. Hence, $\lambda_i(\tau) = \lambda_i^0 \cdot \tau^\alpha$. Therefore,

$$I(X; T_k)(\tau) \leq C + \frac{1}{\tau^\alpha} \sum_{\lambda_i \in \Lambda^{NI}} \left( \frac{A_{ii}}{\lambda_i^0} \right)$$

Inverting this relation, the time to compress the representation $T_k$ by $\Delta I(X; T_k) = \Delta I_k$ scales as: $\tau(\Delta I_k) \propto \left( \frac{-R}{\Delta I(X;T)} \right)^{\frac{1}{\alpha}}$, where $R = \frac{1}{2} \sum_{\lambda_i \in \Lambda^{NI}} \left( \frac{A_{ii}}{\lambda_i^0} \right)$. Note that $R$ depends solely on the problem, $f(x)$ or $p(y, x)$, and not on the architecture. The idea behind this argument is as follows - one can expand the function in any orthogonal basis (e.g. Fourier transform). The expansion coefficients determine both the dimensionality of the relevant/informative dimensions and the total trace of the irrelevant directions. Since these traces are invariant to the specific function basis, these traces remain the same when expanding the function in the network functions using the weights.

With $K$ hidden layers, each layer only needs to compress from the previous (compressed) layer, by $\Delta I_k$, and the total compression is $\Delta I_X = \sum_k \Delta I_k$. Under these assumptions, even if the layers compress one after the other, the total compression time can be broken down into $k$ smaller steps, as at

$$\left( \frac{R}{\sum_k \Delta I_k} \right)^{\frac{1}{\alpha}} \ll \sum_k \left( \frac{R}{\Delta I_k} \right)^{\frac{1}{\alpha}}$$

If the $\Delta I_k$ are similar, we obtain a super-linear boost in the computational time by a factor $K^{\frac{1}{\alpha}}$. Since $\alpha \leq 0.5$ this is at least a quadratic boost in $K$. For ultra-slow diffusion, we obtain an exponential boost (in $K$) in the convergence time to a good generalization.

### Experiments

We now illustrate our results in a series of experiments. We examine several different setups.

**MNIST dataset** – In the first experiment, we evaluate the MNIST handwritten digit recognition task (LeCun et al.; 1990). For this data set, we use a fully connected network with 5 hidden layers of width $500 - 250 - 100 - 50 - 20$, with an hyperbolic tangent (tanh) activation function. The relatively low dimension of the network and the bounded activation function allow us to empirically measure the MI in the network. The MI is estimate by binning the neurons' output into the interval $[-1, 1]$. The discretized values are then used to estimate the joint distributions and the corresponding MI, as described by Shwartz-Ziv and Tishby (2017).

Figure 1.9a depicts the norms of the weights, the signal-to-noise ratio (the ratio between the means of the gradients and their standard deviations), the compression rate $I(X; T)$ and the Gaussian upper bound on $I(X; T)$, as defined in Equation **??**. As expected, the two

**(a) The change of weights, the SNR of the gradients, the MI and the Gaussian bound during the training for one layer. In log-log scale**

**(b) The transition point of the SNR ($Y$-axis) versus the beginning of the information compression ($X$-axis), for different mini-batch sizes**

**Figure 1.9: MNIST data-set**

distinct phases correspond to the drift and diffusion phases. Furthermore, these two phases are evident by independently observing the SNR, the change of the weights $||W(l) - W(0)||$, the MI, and the upper bound. In the first phase, the weights grow almost linearly with the iterations, the SNR of the gradients is high, and there is almost no change in the MI. Then, after the transition point (that accrues almost at the same iteration for all the measures above), the weights behave as a diffusion process. In this phase, the SNR and MI decrease remarkably. In this phase, there is also a clear-cut reduction of the bound.

**CIFAR-10 and CIFAR-100** – Next, we validate our theory on large-scale modern networks. In the second experiment, we consider two data sets, CIFAR-10 and CIFAR-100. Here, we train a ResNet-32 network, using a standard architecture (including ReLU activation functions as described in (He et al.; 2016). In this experiment, we do not estimate the MI directly due to the problem's large scale. Figure 1.10 shows the SNR of the gradients and the Gaussian bound for one layer in CIFAR-10 and CIFAR-100 on the ResNet-32 network averaged over 50 runs. Here, we observed similar behavior, as reported in the MNIST experiment. Specifically, there is a clear distinction between the two phases and a reduction of the MI bound along with the diffusion phase. Note that the same behavior was observed in most of the 32 layers in the network.

Recently, several attempts characterize the correspondence between the diffusion rate of the SGD and the size of the mini-batch (Hu et al. (2017); Hoffer et al. (2017)). In these articles, the authors claimed that a larger mini-batch size corresponds to a lower diffusion rate. Here, we examine the effect of the mini-batch size on the transition phase in the information plane. For each mini-batch size, we find both the starting point of the information compression and the gradient phase transition (the iteration where the derivative of the SNR is maximal). Figure 1.9b illustrates the results. The $X$-axis is the

**(a) CIFAR-10**          **(b) CIFAR-100**

**Figure 1.10:** Change in the SNR of the gradients and the Gaussian bound on the MI during the training of the network for one layer on ResNet-32, in log-log scale.



**(a) Symmetric dataset**          **(b) MNIST**

**Figure 1.11:** The computational benefit of the layers - The converged iteration as function of the number of layers in the network

iteration where the compression started, and the $Y$-axis is the iteration where the phase transition in the gradients accrued for different minibatch sizes. There is a clear linear trend between the two. This further justifies our suggested model since the two measures are strongly related.

NNext, we validate our results on the computational benefit of the layers. We train networks with different number of layers (1-5 layers) and examine the number of iterations a network takes to converge. Then, we find the $\alpha$ which fits the best trend $K^{\frac{1}{\alpha}}$, where $K$ is the number of layers. Figure 1.11 shows the results for two data-sets - MNIST and the symmetric dataset from Shwartz-Ziv and Tishby (2017). As our theory suggest, as we increase the number of layers, the convergence time decreases with a factor of $k^{\frac{1}{\alpha}}$ for different values of $\alpha$.

### Discussion and Conclusions

In this work, we study DNNs using information-theoretic principles. We describe the network's training process as two separate phases, as has been previously done by others. In the first phase (drift), we show that $I(T_k; Y)$ increases, corresponding to an improved generalization with ERM. In the second phase (diffusion), the representation information, $I(X; T_k)$ slowly decreases, while $I(T_K; Y)$ continues to increase. We rigorously prove that the representation compression is a direct consequence of the diffusion phase, independent of the nonlinearity of the activation function. We provide a new Gaussian bound on the representation compression and then relate the diffusion exponent to the compression time. One key outcome of this analysis is a new proof of the computational benefit of the hidden layers, where we show that they boost the overall convergence time of the network by at least a factor of $K^2$, where $K$ is the number of non-degenerate hidden layers. This boost can be exponential in the number of hidden layers if the diffusion is "ultra-slow", as recently reported.

### Bibliography

Billingsley, P. (2008). *Probability and measure*, John Wiley & Sons.

Chalk, M., Marre, O. and Tkacik, G. (2016). Relevant sparse codes with variational information bottleneck, *Advances in Neural Information Processing Systems*, pp. 1957–1965.

Chaudhari, P. and Soatto, S. (2017). Stochastic gradient descent performs variational inference, converges to limit cycles for deep networks, *CoRR* **abs/1710.11029**.
**URL:** *http://arxiv.org/abs/1710.11029*

Chechik, G., Globerson, A., Tishby, N. and Weiss, Y. (2005). Information bottleneck for gaussian variables, *Journal of machine learning research* **6**(Jan): 165–188.

Chee, J. and Toulis, P. (2017). Convergence diagnostics for stochastic gradient descent with constant step size, *arXiv preprint arXiv:1710.06382* .

Chigirev, D. V. and Bialek, W. (2004). Optimal manifold representation of data: an information theoretic approach, *Advances in Neural Information Processing Systems*, pp. 161–168.

Cover, T. M. and Thomas, J. A. (2012). *Elements of information theory*, John Wiley & Sons.

Deco, G. and Obradovic, D. (2012). *An information-theoretic approach to neural computing*, Springer Science & Business Media.

Gabrié, M., Manoel, A., Luneau, C., Barbier, J., Macris, N., Krzakala, F. and Zdeborová, L. (2018). Entropy and mutual information in models of deep neural networks, *arXiv preprint arXiv:1805.09785* .

Graves, A., Mohamed, A.-r. and Hinton, G. (2013). Speech recognition with deep recurrent neural networks, *Acoustics, speech and signal processing (icassp), 2013 ieee international conference on*, IEEE, pp. 6645–6649.

Hardt, M., Recht, B. and Singer, Y. (2015). Train faster, generalize better: Stability of stochastic gradient descent, *arXiv preprint arXiv:1509.01240* .

Harremoes, P. and Tishby, N. (2007). The information bottleneck revisited or how to choose a good distortion measure, *Information Theory, 2007. ISIT 2007. IEEE International Symposium on*, IEEE, pp. 566–570.

He, K., Zhang, X., Ren, S. and Sun, J. (2016). Deep residual learning for image recognition, *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 770–778.

Hoeffding, W. (1963). Probability inequalities for sums of bounded random variables, *Journal of the American statistical association* **58**(301): 13–30.

Hoffer, E., Hubara, I. and Soudry, D. (2017). Train longer, generalize better: closing the generalization gap in large batch training of neural networks, *Advances in Neural Information Processing Systems*, pp. 1731–1741.

Hu, W., Li, C. J., Li, L. and Liu, J.-G. (2017). On the diffusion approximation of nonconvex stochastic gradient descent, *arXiv preprint arXiv:1705.07562* .

Jastrzebski, S., Kenton, Z., Arpit, D., Ballas, N., Fischer, A., Bengio, Y. and Storkey, A. (2017). Three factors influencing minima in sgd, *arXiv preprint arXiv:1711.04623* .

Jin, C., Ge, R., Netrapalli, P., Kakade, S. M. and Jordan, M. I. (2017). How to escape saddle points efficiently, *arXiv preprint arXiv:1703.00887* .

Keskar, N. S., Mudigere, D., Nocedal, J., Smelyanskiy, M. and Tang, P. T. P. (2016). On large-batch training for deep learning: Generalization gap and sharp minima, *arXiv preprint arXiv:1609.04836* .

LeCun, Y., Bengio, Y. and Hinton, G. (2015). Deep learning, *Nature* .

LeCun, Y., Boser, B., Denker, J., Henderson, D., Howard, R., Hubbard, W. and Jackel, L. (1990). Handwritten digit recognition with a back-propagation network, *Advances in neural information processing systems 2, NIPS 1989*, Morgan Kaufmann Publishers, pp. 396–404.

Li, Q., Tai, C. et al. (2015). Stochastic modified equations and adaptive stochastic gradient algorithms, *arXiv preprint arXiv:1511.06251* .

Linsker, R. (1988). Self-organization in a perceptual network, *Computer* **21**(3): 105–117.

Murata, N. (1998). A statistical study of on-line learning, *Online Learning and Neural Networks. Cambridge University Press, Cambridge, UK* pp. 63–92.

Painsky, A., Feder, M. and Tishby, N. (2018). An information-theoretic framework for non-linear canonical correlation analysis, *arXiv preprint arXiv:1810.13259* .

Painsky, A., Rosset, S. and Feder, M. (2016). Generalized independent component analysis over finite alphabets, *IEEE Transactions on Information Theory* **62**(2): 1038–1053.

Painsky, A. and Tishby, N. (2017). Gaussian lower bound for the information bottleneck limit, *The Journal of Machine Learning Research* **18**(1): 7908–7936.

Painsky, A. and Wornell, G. W. (2018a). Bregman divergence bounds and the universality of the logarithmic loss, *arXiv preprint arXiv:1810.07014* .

Painsky, A. and Wornell, G. W. (2018b). On the universality of the logistic loss function, *arXiv preprint arXiv:1805.03804* .

Paninski, L. (2003). Estimation of entropy and mutual information, *Neural Comput.* **15**(6): 1191–1253.

Rissanen, J. (1978). Modeling by shortest data description, *Automatica* **14**(5): 465–471.

Robbins, H. and Monro, S. (1951). A stochastic approximation method, *The annals of mathematical statistics* pp. 400–407.

Sagun, L., Evci, U., Guney, V. U., Dauphin, Y. and Bottou, L. (2017). Empirical analysis of the hessian of over-parametrized neural networks, *arXiv preprint arXiv:1706.04454* .

Sauer, N. (1972). On the density of families of sets, *Journal of Combinatorial Theory, Series A* **13**(1): 145–147.

Shamir, O., Sabato, S. and Tishby, N. (2010). Learning and generalization with the information bottleneck, *Theoretical Computer Science* **411**(29): 2696 – 2711. Algorithmic Learning Theory (ALT 2008).
**URL:** *http://www.sciencedirect.com/science/article/pii/S030439751000201X*

Shelah, S. (1972). A combinatorial problem; stability and order for models and theories in infinitary languages, *Pacific Journal of Mathematics* **41**(1): 247–261.

Shwartz-Ziv, R. and Tishby, N. (2017). Opening the black box of deep neural networks via information, *arXiv preprint arXiv:1703.00810* .

Tishby, N., Pereira, F. C. and Bialek, W. (1999). The information bottleneck method, *In Proceedings of the 37-th Annual Allerton Conference on Communication, Control and Computing* .

Tishby, N. and Zaslavsky, N. (2015). Deep learning and the information bottleneck principle, *Information Theory Workshop (ITW), 2015 IEEE*, IEEE, pp. 1–5.

Vapnik, V. N. and Chervonenkis, A. Y. (1968). The uniform convergence of frequencies of the appearance of events to their probabilities, *Doklady Akademii Nauk*, Russian Academy of Sciences, pp. 781–783.

Zhang, C., Liao, Q., Rakhlin, A., Miranda, B., Golowich, N. and Poggio, T. A. (2018). Theory of deep learning iib: Optimization properties of SGD, *CoRR* **abs/1801.02254**. **URL:** *http://arxiv.org/abs/1801.02254*

Zhang, Y., Saxe, A. M., Advani, M. S. and Lee, A. A. (2018). Energy–entropy competition and the effectiveness of stochastic gradient descent in machine learning, *Molecular Physics* pp. 1–10.

Zhu, Z., Wu, J., Yu, B., Wu, L. and Ma, J. (2018). The Anisotropic Noise in Stochastic Gradient Descent: Its Behavior of Escaping from Minima and Regularization Effects, *ArXiv e-prints* .

*Appendix*

### Appendix A - proof of Theorem 1

We first revisit the well-known *Probably Approximately Correct* (PAC) bound. Let $\mathcal{H}$ be a finite set of hypotheses. Let $\ell_h(x_i, y_i)$ be a bounded loss function, for every $h \in \mathcal{H}$. For example, $\ell_h(x_i, y_i) = (y_i - h(x_i))^2$ is the squared loss while $\ell_h(x_i, y_i) = -y_i \log h(x_i)$ is the logarithmic loss (which may be treated as bounded, assuming that the underlying distribution is bounded away from zero and one). Let $\mathcal{L}_h(S_m) = \frac{1}{m} \sum_{i=1}^{m} \ell_h(x_i, y_i)$ be the empirical error. Hoeffding's inequality Hoeffding (1963) shows that for every $h \in \mathcal{H}$

$$\mathbb{P}\left[\left|\mathcal{L}_h(S_m) - \mathbb{E}_{S_m}[\mathcal{L}_h(S_m)]\right| \geq \epsilon\right] \leq 2 \exp\left(-2\epsilon^2 m\right). \tag{1.31}$$

Then, we can apply the union bound and conclude that

$$\mathbb{P}\left[\exists h \in \mathcal{H}\,\Big|\, \left|\mathcal{L}_h(S_m) - \mathbb{E}_{S_m}[\mathcal{L}_h(S_m)]\right| \geq \epsilon\right] \leq 2|\mathcal{H}| \exp\left(-2\epsilon^2 m\right).$$

We want to control the above probability with a confidence level of $\delta$. Therefore, we ask that $2|\mathcal{H}| \exp\left(-2\epsilon^2 m\right) \leq \delta$. This leads to a PAC bound, which states that for a fixed $m$ and for every $h \in \mathcal{H}$, we have with probability $1 - \delta$ that

$$\left|\mathcal{L}_h(S_n) - \mathbb{E}_{Sm}[\mathcal{L}_h(S_m)]\right|^2 \leq \frac{\log|\mathcal{H}| + \log\frac{2}{\delta}}{2m}. \tag{1.32}$$

Note that under the definitions stated above, we have that $|\mathcal{H}| \leq 2^{\mathcal{X}}$. However, the PAC bound above also holds for a infinite hypotheses class, where $\log|\mathcal{H}|$ is replaced with the VC dimension of the problem, with several additional constants (Vapnik and Chervonenkis; 1968; Shelah; 1972; Sauer; 1972).

Let us now assume that $X$ is a $d$-dimensional random vector that follows a Markov random field structure. As stated above, this means that $p(x_i) = \prod_i p(x_i|Pa(x_i))$ where $Pa(X_i)$ is a set of components in the vector $X$ that are adjacent to $X_i$. Assuming that the Markov random field is ergodic, we can define a *typical set* of realizations from $X$ as a set that satisfies the *Asymptotic Equipartition Property* (AEP) (Cover and Thomas; 2012). Therefore, for every $\epsilon > 0$, the probability of a sequence drawn from $X$ to be in the typical set $A_\epsilon$ is greater than $1 - \epsilon$ and $|A_\epsilon| \leq 2^{H(X)+\epsilon}$. Hence, if we only consider a typical realization of $X$ (as opposed to every possible realization), we have that asymptotically $|\mathcal{H}| \leq 2^{H(X)}$. Finally, let $T$ be a mapping of $X$. Then, $2^{H(X|T)}$ is the number of typical realizations of $X$ that are mapped to $T$. This means that the size of the typical set of $T$ is bounded from above by $2^{H(X)}/2^{H(X|T)} = 2^{I(X;T)}$. Plugging this into the PAC bound above yields that with probability $1 - \delta$, the typical squared generalization error of $T$, $\epsilon_T^2$

satisfies

$$\epsilon_T^2 \leq \frac{2^{I(X;T)} + \log \frac{2}{\delta}}{2m}. \tag{1.33}$$

### Appendix B - Proof of Proposition 2

We make the following technical assumptions:

1. $w^*$ and $\delta w$ satisfy $\lim_{d_k \to \infty} w^{*T} \delta w = 0$ almost surely.

2. The moments of $T_k$ are finite.

3. The components of $w^*$ and $\delta w$ are *in-general-positions*, satisfying

$$\lim_{d_k \to \infty} \sum_{i=1}^{d_k} w_i^{*4} / \left( \sum_{i=1}^{d_k} w_i^{*2} \right)^2 = 0$$

and

$$\lim_{d_k \to \infty} \sum_{i=1}^{d_k} \delta w_i^4 / \left( \sum_{i=1}^{d_k} \delta w_i^2 \right)^2 = 0$$

almost surely.

Consider a sequence of i.i.d. random variables, $\{X_i\}_{i=1}^{d}$ with zero mean and finite moments, $\mathbb{E}[X_i^r] < \infty$ for every $r \geq 1$.

Let $\{a_i\}_{i=1}^{d}$ be a sequence of constants. Denote $Y_i = a_i X_i$, so that $\{Y_i\}_{i=1}^{d}$ are independent with zero mean and $\text{Var}(Y_i) = a_i^2 \mathbb{E}[X^2]$. Let $S = \sum_{i=1}^{d} a_i X_i = \sum_{i=1}^{d} Y_i$ and denote $U_d^2 = \sum_{i=1}^{d} \text{Var}(Y_i) = \mathbb{E}[X^2] \sum_{i=1}^{d} a_i^2$.

The Lyapunov Central Limit Theorem (CLT) Billingsley (2008) states that if there exists some $\delta > 0$ for which

$$\lim_{d \to \infty} \frac{1}{U_d^{2+\delta}} \sum_{i=1}^{d} \mathbb{E}\left[|Y_i|^{2+\delta}\right] = 0 \tag{1.34}$$

then

$$\frac{1}{U_d} \sum_{i=1}^{d} Y_i \xrightarrow[d \to \infty]{\mathcal{D}} \mathcal{N}(0, 1). \tag{1.35}$$

Plugging $\delta = 2$ yields the following sufficient condition,

$$\lim_{d \to \infty} \frac{1}{U_d^4} \sum_{i=1}^{d} \mathbb{E}\left[Y_i^4\right] = \frac{\sum_{i=1}^{d} a_i^4}{\left( \sum_{i=1}^{d} a_i^2 \right)^2} \frac{\mathbb{E}[X^4]}{\mathbb{E}^2[X^2]} = 0 \tag{1.36}$$

Let us apply the Lyapunov CLT to our problem. Here, the components of $T_k$ are i.i.d. for sufficiently large $d_k$, with zero mean and finite $r^{th}$ moments for every $r \geq 1$. Furthermore, we assume that the components of $w^*$ and $\delta w$ are in-general-positions. This means that Lyapunov condition (**??**) is satisfied for both $w^{*T} T_k$ and $\delta w^T T_k$ almost surely, which means that

$$\frac{1}{\sqrt{\sigma_{T_k}^2} ||w^*||_2} w^{*T} T_k \xrightarrow[d_k \to \infty]{\mathcal{D}} \mathcal{N}(0,1) \tag{1.37}$$

and

$$\frac{1}{\sqrt{\sigma_{T_k}^2} ||\delta w||_2} \delta w^T T_k \xrightarrow[d_k \to \infty]{\mathcal{D}} \mathcal{N}(0,1). \tag{1.38}$$

almost surely, where $\sigma_{T_k}^2$ is the variance of the components of $T_k$.

Furthermore, for every pair of constants $a$ and $b$, the linear combination $(aw^* + b\delta w)^T T_k$ also satisfies Lyapunov's condition almost surely, which means that $w^{*T} T_k$ and $\delta w^T T_k$ are asymptotically jointly Gaussian, with

$$\mathbb{E}\left[ w^{*T} T_k \left( \delta w^T T_k \right)^T \right] = \sigma_{T_k}^2 w^{*T} \delta w \xrightarrow[d_k \to \infty]{} 0$$

almost surely. □

# Information in Infinite Ensembles of Infinitely-Wide Neural Networks

Ravid Shwartz-Ziv and Alexander A. Alemi.

# Information in Infinite Ensembles of Infinitely-Wide Neural Networks

Ravid Shwartz-Ziv [1] Alexander A. Alemi[2]

[1] The Edmond and Lilly Safra Center for Brain Sciences, The Hebrew University,
Jerusalem, Israel.
[2] Google Research
USA

## Abstract

In this work, we study the generalization properties of infinite ensembles of infinitely-wide neural networks. Amazingly, this model family admits tractable calculations for many information-theoretic quantities. We report analytical and empirical investigations in the search for signals that correlate with generalization.

## *Introduction*

According to the statistical learning theory (Boucheron et al.; 2005), models with many parameters tend to overfit by representing the learned data too accurately, therefore diminishing their ability to generalize to unseen data. However, in DNNs, we see that the 'generalization gap; i.e., the difference between 'training error' and 'test error' is minimal. One promising research direction is to view deep neural networks through the lens of information theory(Tishby and Zaslavsky; 2015). Abstractly, deep connections exist between the information a learning algorithm extracts and its generalization capabilities (Bassily et al.; 2017; Banerjee; 2006). Inspired by these general results, recent papers have attempted to measure information-theoretic quantities in ordinary deterministic neural networks (Shwartz-Ziv and Tishby; 2017; Achille and Soatto; 2017; Achille and Soatto; 2019).

Both practical and theoretical problems arise in the deterministic case (Amjad and Geiger; 2018; Saxe et al.; 2019; Kolchinsky et al.; 2018). These difficulties stem from

the fact that mutual information (MI) is reparameterization independent (Cover and Thomas; 2012).[1] One workaround is to make a network explicitly stochastic, either in its activations (Alemi et al.; 2016) or its weights (Achille and Soatto; 2017). Here we take an alternative approach, harnessing the stochasticity in our choice of initial parameters. We consider an *ensemble* of neural networks, all trained with the same training procedure and data. This generates an ensemble of predictions. In other words, we consider an infinite ensemble of neural networks that describes a distribution over the output space when we marginalize out our choice of initial parameters. Characterizing the generalization properties of the ensemble should characterize the generalization of individual draws from this ensemble. However, the challenge is in describing this ordinarily intractable distribution.

Infinitely-wide neural networks behave as if they are linear in their parameters (Lee et al.; 2019). Their evolution is fully described by the *neural tangent kernel* (NTK). The NTK is constant in time and can be tractably computed (Novak et al.; 2019) as a function of the network's architecture, e.g., the number and the structure of layers, nonlinearity, initial parameters' distributions, etc.

(Lee et al.; 2019) showed that the output of an infinite ensemble of infinitely-wide neural networks initialized with Gaussian weights and biases and trained with gradient flow to minimize a square loss is simply a conditional Gaussian distribution:

$$p(z|x) \sim \mathcal{N}(\mu(x,\tau), \Sigma(x,\tau)), \tag{1.39}$$

where $z$ is the output of the network and $x$ is its input. The mean $\mu(x,\tau)$ and covariance $\Sigma(x,\tau)$ functions can be computed (Novak et al.; 2019).

Recently, there has been much interest in understanding the importance of implicit regularization as a tool for explaining the generalization gap of DNNs. Numerical experiments demonstrate that network size may not be the main form of capacity control, and hence, some other unknown form plays a central role in learning multi-layer network (Neyshabur et al.; 2014, 2015). A line of works have derived interesting results proving that some information-theoretic quantities provide concise bounds on the generalization gap, hence behave as implicit regularization (Russo and Zou; 2019; Xu and Raginsky; 2017; Pensia et al.; 2018; Negrea et al.; 2019; Asadi et al.; 2018; Russo and Zou; 2016; Steinke and Zakynthinou; 2020; Achille and Soatto; 2019; Banerjee; 2006; Bassily et al.; 2017). However there is not direct evidence what are the most important quantities which determinate the generalization ability of the network.

In this work, the simple structure of the NTK allows us to bound several interesting information-theoretic quantities, including: the MI between the representation and the

---

[1] This implies that if we send a random variable through an invertible function, its MI with respect to any other variable remains unchanged.

targets, $I(Z; Y)$, the MI between the representation and the inputs after training, $I(Z; X|D)$, and the MI between the representations and the training set, conditioned on the input, $I(Z; D|X)$. We are also able to compute in closed form: the Fisher information metric, the distance the parameters move and the MI between the parameters and the data, $I(\Theta; D)$. All derivation of all these information-theoretic quantities allow us to explore what are the important factors in the generalization ability of DNNs.

### *Background*

#### *1.7.1   Neural Tangent Kernel*

In this section we describe fully-connected deep neural net architecture and its infinite width limit, and how training it with respect to the $l_2$ loss gives rise to a kernel regression problem involving the NTK. We denote by $f_t(\theta, x) \in \mathbb{R}$ the output of a neural network at time $t$ where $\theta \in \mathbb{R}^N$ is all the parameters in the network and $x \in \mathcal{R}^d$ is the input. Let $\ell(\hat{y}, y) : \mathbb{R} \times \mathbb{R} \to \mathbb{R}$ denote the loss function where the first argument is the prediction and the second argument the true label. Given a training dataset $(x_i, y_i)_{i=1}^N \subseteq \mathbb{R}^d \times \mathbb{R}$, consider training the neural network by minimizing the empirical loss over training data: $\mathcal{L} = \sum_{i=1}^n \ell(f_t(x_i, \theta), y)$. Let $\eta$ be the learning rate. Via continuous time gradient descent, the evolution of the parameters $\theta$ and the logits $f$ can be written as

**Lemma 1.7.1.**

$$\frac{d\theta_t}{dt} = -\eta \nabla_\theta f_t(\mathcal{X})^T \nabla_{f_t(\mathcal{X})} \mathcal{L}$$
$$\frac{df_t(\mathcal{X})}{dt} = \nabla_\theta f_t(\mathcal{X}) \frac{d\theta_t}{dt} = -\eta \hat{\Theta}_t (\mathcal{X}, \mathcal{X}) \nabla_{f_t(\mathcal{X})} \mathcal{L}$$

where $f_t(\mathcal{X}) = vec\left([f_t(x)] \, x \in \mathcal{X}\right)$ , the $k|D| \times 1$ vector of concatenated logits for all example, and $\hat{\Theta}_t \equiv \Theta_t(\hat{\mathcal{X}}, \mathcal{X})$ is the tangent kernel at time t - an $n \times n$ positive semidefinite matrix whose $(i, j)$-th entry is

$$\left\langle \frac{\partial f_t(x_i)}{\partial \theta}, \frac{\partial f_t(x_j)}{\partial \theta} \right\rangle$$

The shorthand $\Theta$ denotes the kernel function evaluated on the train data $(\Theta \equiv \Theta(\mathcal{X}, \mathcal{X}))$. For a finite width network, the NTK corresponds to $JJ^T$, the gram matrix of neural network gradients. As the width of a network increases to infinity, this kernel converges in probability to a fixed value. There exist tractable ways to calculate the exact infinite-width kernel for wide classes of neural networks (Novak et al.; 2019).

Infinitely-wide neural networks behave as though they were linear in their parameters (Lee et al.; 2019):

$$f_t(x) = f_0(x) + \frac{\partial f_0(X)}{\partial \theta}(\theta - \theta_0) \tag{1.40}$$

This makes them particularly analytically tractable. An infinitely-wide neural network, trained by gradient flow to minimize squared loss admits a closed form expression for evolution of its predictions as a function of time. For an arbitrary point $x$

$$f_t(x) = f_0(x) - \hat{\Theta}(x, \mathcal{X})\hat{\Theta}^{-1}\left(I - e^{-\tau\hat{\Theta}}\right)(f_0(\mathcal{X}) - \mathcal{Y}). \tag{1.41}$$

Notice that the behavior of infinitely-wide neural networks trained with gradient flow and squared loss is just a time-dependent affine transformation of their initial predictions. As such, if we now imagine forming an infinite ensemble of such networks as we vary their initial weight configurations, if those weights are sampled from a Gaussian distribution, the law of large numbers enforces that the distribution of outputs of the ensemble of networks at initialization is Gaussian, conditioned on its input. Since the evolution is an affine transformation of the initial predictions, the predictions remain Gaussian at all times.

**Corollary 1.7.1.1.** *For every test points $x \in X$ and $t \leq 0$, $z = f_t(x)$ converges in distribution as width goes to infinity to*

$$p(z|x) \sim \mathcal{N}(\mu(x,t), \Sigma(x,t))$$
$$\mu(x,t) = \Theta(x, \mathcal{X})\Theta^{-1}\left(I - e^{-t\Theta}\right)\mathcal{Y}$$
$$\Sigma(x,t) =$$
$$\mathcal{K}(x,x) + \Theta(x, \mathcal{X})\Theta^{-1}\left(I - e^{-t\Theta}\right)\left(\mathcal{K}\Theta^{-1}\left(I - e^{-\tau\Theta}\right)\Theta(\mathcal{X}, x) - 2\mathcal{K}(\mathcal{X}, x)\right)$$

For more details see Lee et al. (2019).

Here, $\mathcal{K}$ denotes the *neural network gaussian process* kernel (NNGP). For a finite width network, the NNGP corresponds to the expected gram matrix of the outputs: $\mathcal{K}^{i,j}(x, x') = \mathbb{E}\left[f_t^i(x)f_t^j(x')\right]$. In the infinite width limit, this concentrates on a fixed value. Just as for the NTK, the NNGP can be tractably computed (Novak et al.; 2019), and should be considered just a function of the neural network architecture. For this family of models, we would like to derive information-theoretic characterizations of their performance.

### 1.7.2   The Information Bottleneck Optimal Bound

What characterizes the optimal representations of $X$ w.r.t. $Y$? One of the candidate to this question is the classical notion of minimal sufficient statistics. Sufficient statistics, are

maps or partitions of $X$, $S(X)$, capturing all the information that $X$ has on $Y$. Namely,

$$I(S(X); Y) = I(X; Y)$$

(Cover and Thomas; 2012).

Minimal sufficient statistics, $T(X)$, are the simplest sufficient statistics and induce the coarsest sufficient partition on $X$. In other words, they are functions of any other sufficient statistic. A simple way of formulating this is through the Markov chain: $Y \rightarrow X \rightarrow S(X) \rightarrow T(X)$, which should hold for a minimal sufficient statistics $T(X)$ with any other sufficient statistics $S(X)$. Using DPI we can cast it into a constrained optimization problem:

$$T(X) = \arg \min_{S(X): I(S(X); Y) = I(X; Y)} I(S(X); X) \ .$$

Since exact minimal sufficient statistics only exist for very special distributions, (i.e., exponential families), (Tishby et al.; 1999) relaxed this optimization problem by first allowing the map to be stochastic, defined as an encoder $P(T|X)$, and then, by enabling the map to capture *as much as possible* of $I(X; Y)$, not necessarily all of it.

This leads to the *Information Bottleneck* (IB) trade off (Tishby et al.; 1999), which provides a computational framework for finding approximate minimal sufficient statistics, or the optimal trade off between compression of $X$ and prediction of $Y$.

If we define $t \in T$ as the compressed representations of $x \in X$, the representation of $x$ is now defined by the mapping $p(t|x)$. This IB trade off is formulated by the following optimization problem, carried independently for the distributions, $p(t|x), p(t), p(y|t)$, with the Markov chain: $Y \rightarrow X \rightarrow T$,

$$\min_{p(t|x), p(y|t), p(t)} \{ I(X; T) - \beta I(T; Y) \} \ . \tag{1.42}$$

The Lagrange multiplier $\beta$ determines the level of relevant information captured by the representation $T$, $I(T; Y)$.

If we denote $I_X^\beta = I_\beta(T; X)$ and $I_Y^\beta = I_\beta(T; Y)$ for some $\beta$, the optimal information curve is then defined as the optimal values of the trade-off $\left( I_X^\beta, I_Y^\beta \right)$ for each $\beta$. The two-dimensional plane in which the IB curve resides is coined as the information plane. The information curve is a monotonic concave line of optimal representations that separates the achievable and unachievable regions in the information-plane.

### 1.7.2.1 *Gaussian Information Bottleneck*

Generally speaking, solving the IB problem **??** for an arbitrary joint distribution is an hard task. Tishby et al. (1999) defined a set of self-consistent equations which formulate the

necessary conditions for the optimal solution of **??**. In general, these equations do not hold a tractable solution and are usually approximated by different means (Slonim; 2002).

A special exception is the Gaussian case, where $X$ and $Y$ are follow a jointly normal distribution. Namely -

$$p(x, y) = \mathcal{N} \left( 0, \begin{pmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{pmatrix} \right)$$

In this case, the Gaussian IB problem is analytically solved by linear projections to the canonical correlation vector space (Chechik et al.; 2005). In this case, the Gaussian IB problemis solved by a noisy linear projection, $T = AX + \epsilon$ and everything is governed by the eigenspectrum of:

$$\Sigma_{x|y}\Sigma_{xx}^{-1} = I - \Sigma_{xy}\Sigma_{yy}^{-1}\Sigma_{yx}\Sigma_{xx}^{-1}$$

### 1.7.3  Information Bounds on the Generalization Gap

Many works tried to understand the importance of implicit regularization in the small generalization gap of DNNs. Numerical experiments in (Neyshabur et al.; 2014, 2015) demonstrate that network size may not be the main factor to explain the capacity of the network, and hence, some other unknown form plays a central role in the learning. Further work in Zhang et al. (2016), found that in contrast with classical convex empirical risk minimization, regularization plays a rather different role in deep learning. From the theoretical viewpoint, regularization seems to be an indispensable component, while convincing experiments support the idea that the absence of explicit regularization does not necessarily induce poor generalization. A line of works have derived interesting results proving that the mutual information between the training inputs and the inferred parameters provides a concise bound on the gap, which crucially depends on a mapping of the training set into the network parameters, whose characterization is not an easy task (Russo and Zou; 2019; Xu and Raginsky; 2017; Pensia et al.; 2018; Negrea et al.; 2019; Asadi et al.; 2018; Russo and Zou; 2016; Steinke and Zakynthinou; 2020; Achille and Soatto; 2019). (Achille and Soatto; 2017) explored how the use of an IB objective on the network parameters may help avoid overfitting while enforcing invariant representations. On the other side, direct use of statistical learning theory, such as Rademacher complexity (Bartlett and Mendelson; 2002), VC-dimension (Vapnik; 1998), and uniform stability (Bousquet and Elisseeff; 2002) seem to be inadequate to explain the unexpected numerical observations on the generalization gap.

### Family of Trained Models

In the recent years there are two lines of works which investigate generalization in DNNs; (1) Some works try to use information quantities to train DNNs better, while others (2) try

to explain the current DNNs generalization ability. In this work, we would like to follow the second line of works.

Based on the NTK framework, we derive several information-theoretic quantities that are good candidates to explain the performance of DNNs. Following this, we empirically investigate these measures by varying the hyperparameters of the networks. Specifically, we check three types of hyperparameters relevant to the NTK framework –

- Network architecture (number of layers, activation's function, etc.).

- The initial noise of the weights and biases.

- The dataset (type, number of training examples, etc.) – We train the networks on three different datasets - MNIST, CIFAR-10, and a jointly Gaussian dataset where we can compare the network performance to that of an optimal analytical solution

The NTK framework is full batch training, so there is no effect of the learning rate or the learning algorithm on dynamics.

More train examples, for example, should in principle store more information about the data. However, the usual PAC-Bayes generalization bounds for information suggest that a high-performance network should be compressed(Banerjee; 2006).

### Derive Information Metrics

We can compute several information-theoretic quantities in the presence of a tractable form for representing an infinitely-wide ensemble of networks. As a result, we can shed light on previous attempts to explain generalization in neural networks and identify candidates for empirical investigations into quantities that can predict generalization.

#### Loss

To compute our ensemble's expected loss, we need to marginalize the stochasticity in the output of the network. Training with squared loss is equivalent to assuming a Gaussian observation model $p(y|z) \sim \mathcal{N}(0, 1)$. We can marginalize out our representation to obtain

$$q(y|x) = \int dz \, q(y|z)p(z|x) \sim \mathcal{N}(\mu(x, \tau), I + \Sigma(x, \tau)). \quad (1.43)$$

The expected log loss has contributions both from the square loss of the mean prediction, as well as a term which couples to the trace of the covariance:

$$\mathbb{E}\left[\log q(y|z)\right] = \frac{1}{2}\mathbb{E}\left[(y - z(x, \tau))^2\right] = \frac{1}{2}(y - \mu(x, \tau))^2 + \frac{1}{2}\operatorname{Tr}\Sigma(x, \tau) - \frac{k}{2}\log 2\pi \quad (1.44)$$

here $k$ is the dimensionality of $y$.

$H(Y|Z)$

The informativeness, or accuracy, of the representation is measured by $I(Z;Y)$, which is the amount of relevant information about $Y$ preserved by the representation. It measures how much of the predictive features in $X$ for $Y$ is captured by our model. Because we know that $I(X:Y) = H(Y) - H(Y|Z)$, we would like to estimating the conditional entropy. To calculate an lower bound, similarly to what we did before, we need to marginalize the stochasticity in the output by finding the observation model which gives us the lowest conditional entropy. However, in this case, we are not assuming a model with Gaussian's variance one, that could be sub-optimal.

$$
\mathbb{E}\left[\log q(y|z)\right] = \frac{1}{2}\mathbb{E}\left[(y - z(x,\tau))^2\right] =
$$
$$
\frac{1}{2}\left((y - \mu(x,\tau))\Sigma_r^{-1}(y - \mu(x,\tau)) + \frac{\operatorname{Tr}\Sigma(x,\tau) * d}{n}\right)
$$
$$
-\frac{k}{2}\log 2\pi\Sigma_r
$$

$I(Z;Y)$

While the MI between the network's output and the targets is intractable in general, we can obtain a tractable variational lower bound: (Poole et al.; 2019)

$$
I(Z;Y) = \mathbb{E}\left[\log\frac{p(y|z)}{p(y)}\right] \leq \mathbb{E}\left[\log\frac{q(y|z)}{p(y)}\right] = H(Y) + \mathbb{E}\left[\log q(y|z)\right] \tag{1.45}
$$

$I(Z;X|D)$

According to IB, the complexity of the representation is measured by $I(X;Z|D)$, which is roughly the number of bits that are required for representing the input ($X$) using the network's output ($Z$) conditioned on the dataset ($D$):

$$
I(Z;X|D) = \mathbb{E}\left[\log\frac{p(z|x,D)}{p(z|D)}\right]. \tag{1.46}
$$

This requires knowledge of the marginal distribution $p(z|D)$. Without knowledge of $p(x)$, this is in general intractable, but there exist simple tractable multi sample upper and lower bounds (Poole et al.; 2019):

$$
\frac{1}{N}\sum_i \log\frac{p(z_i|x_i,D)}{\frac{1}{N}\sum_j p(z_i|x_j,D)} \leq I(Z;X|D) \leq \frac{1}{N}\sum_i \log\frac{p(z_i|x_i,D)}{\frac{1}{N-1}\sum_{j\neq i} p(z_i|x_j,D)}. \tag{1.47}
$$

In this work, we show the minibatch lower bound estimates, which are upper bounded themselves by the log of the batch size.

$I(Z; D|X)$

We can also estimate a variational upper bound on the MI between the representation of our networks and the training dataset.

$$I(Z; D|X) = \mathbb{E}\left[\log \frac{p(z|x, D)}{p(z|x)}\right] \leq \mathbb{E}\left[\log \frac{p(z|x, D)}{p_0(z|x)}\right]. \tag{1.48}$$

Here, the MI we extract from the dataset involves the expected log ratio of our posterior distribution of outputs to the marginal over all possible datasets. Not knowing the data distribution, this is intractable in general, but we can variationally upper bound it with an approximate marginal. A natural candidate is the prior distribution of outputs, for which we have a tractable estimate.

*Fisher information*

It is usually assumed that the Fisher matrix approximates the Hessian spectrum, which can be used to estimate the objective function shape. In the literature on deep learning, it has been shown that eigenvalues close to zero locally form flat minima, leading to better generalization empirically (Keskar et al.; 2016; Liang et al.; 2019). Achille and Soatto (2019) connected flat minima (low Fisher information) to path stability of SGD (Hardt et al.; 2016) and information stability (Xu and Raginsky; 2017), showing that optimization algorithms that converge to flat minima and are path stable also satisfy a form of information stability, and hence generalization, by PAC-Bayes bound (Xu and Raginsky; 2017).

Infinitely-wide networks behave as though they were linear in their parameters with a fixed Jacobian. This leads to trivially flat information geometry. For squared loss, the true Fisher can be computed simply as $F = J^T J$ (Kunstner et al.; 2019). While the trace of the Fisher information has recently been proposed as an important quantity for controlling generalization in neural networks (Achille and Soatto; 2019), for infinitely-wide networks we can see that the trace of the Fisher is the same as the trace of the NTK, which is a constant and does not evolve with time

$$\operatorname{Tr} F = \operatorname{Tr} J^T J = \operatorname{Tr} JJ^T = \operatorname{Tr} \Theta$$

. In so much as infinite ensembles of infinitely-wide neural networks generalize, the degree to which they cannot be explained by the time evolution of the trace of the Fisher, given that the trace of the Fisher does not evolve.

*Parameter distance*

How much do the parameters of an infinitely-wide network change? Lee et al. (2019) emphasizes that the relative Frobenius norm change of the parameters throughout training vanish in the limit of infinite width. This is a justification for the linearization becoming more accurate as the network becomes wider. But is it thus fair to say the parameters are not changing? Instead of looking at the Frobenius norm we can investigate the *length* of the parameters path over the course of training. This reparameterization independent notion of distance utilizes the information geometric metric provided by the Fisher information:

$$L(\tau) = \int_0^\tau ds = \int_0^\tau d\tau \sqrt{\dot{\theta}_\alpha(\tau) g_{\alpha\beta} \dot{\theta}_\beta(\tau)} = \int_0^\tau d\tau \left\| \Theta e^{-\tau\Theta}(z_0(\mathcal{X}) - \mathcal{Y}) \right\| \qquad (1.49)$$

The length of the trajectory in parameter space is the integral of a norm of our residual at initialization projected along $\Theta e^{-\tau\Theta}$. This integral is both positive and finite even as $t \to \infty$. To get additional understanding into the structure of this term, we can consider its expectation over the ensemble, where we can use Jensen's inequality to bound the expectation of trajectory lengths. Since we know that at initialization $z_0(\mathcal{X}) \sim \mathcal{N}(0, \mathcal{K})$ we obtain further simplifications:

$$\mathbb{E}[L(\tau)]^2 \le \mathbb{E}[L^2(\tau)] = \int_0^\tau d\tau \, \mathbb{E}\left[ (z_0(\mathcal{X}) - \mathcal{Y})^T \Theta^2 e^{-2\tau\Theta}(z_0(\mathcal{X}) - \mathcal{Y}) \right] \qquad (1.50)$$

$$= \frac{1}{2} \mathbb{E}\left[ (z_0(\mathcal{X}) - \mathcal{Y})^T \Theta \left(1 - e^{-2\tau\Theta}\right)(z_0(\mathcal{X}) - \mathcal{Y}) \right] \qquad (1.51)$$

$$= \frac{1}{2} \left[ \mathrm{Tr}\left( \mathcal{K}\Theta \left(1 - e^{-2\tau\Theta}\right) \right) + \mathcal{Y}^T \Theta \left(1 - e^{-2\tau\Theta}\right) \mathcal{Y} \right]. \qquad (1.52)$$

$D_{\mathrm{KL}}\left[ p(\theta|D) \| p_0(\theta) \right]$

The MI between the parameters and the dataset $I(\theta; D)$ has been shown to control for overfitting (Bassily et al.; 2017). We can generate a variational upperbound on this quantity by consider the KL divergence between the posterior distribution of our parameters and the prior distribution $D_{\mathrm{KL}}\left[ p(\theta|D) \| p_0(\theta) \right]$, a quantity that itself has been shown to provide generalization bounds in PAC Bayes frameworks (Achille and Soatto; 2017). For our networks, the prior distribution is known and simple, but the posterior distribution can be quite rich. However, we can use the *instantaneous change of variables* formula (Chen et al.; 2018)

$$\log p(\theta_\tau) = \log p(\theta_0) - \int_0^\tau d\tau \, \mathrm{Tr}\left( \frac{\partial \dot{\theta}}{\partial \theta} \right), \qquad (1.53)$$

which gives us a value for the log likelihood the parameters of a trained model at any point in time in terms of its initial log likelihood and the integral of the trace of the kernel

governing its time evolution. For our infinitely-wide neural networks this is tractable:

$$I(\theta; D) \leq \mathbb{E}_{p(\theta_\tau)} \left[ \log p(\theta_\tau) - \log p_0(\theta_\tau) \right] \tag{1.54}$$

$$= \mathbb{E}_{p(\theta_0)} \left[ \log p(\theta_0) - \log p_0(\theta_\tau) - \int_0^\tau d\tau \; \mathrm{Tr} \left( \frac{\partial \dot{\theta}}{\partial \theta} \right) \right] \tag{1.55}$$

$$= \mathbb{E}_{p(\theta_0)} \left[ \log p(\theta_0) - \log p_0(\theta_0 + \Delta\theta_\tau) + \tau \, \mathrm{Tr}\,\Theta \right] \tag{1.56}$$

$$= \mathbb{E}_{p(\theta_0)} \left[ -\frac{1}{2}\theta_0^2 + \frac{1}{2}(\theta_0 + \Delta\theta_\tau)^2 \right] + \tau \, \mathrm{Tr}\,\Theta \tag{1.57}$$

$$= \mathbb{E}_{p(\theta_0)} \left[ \theta_0^T \Delta\theta_\tau + (\Delta\theta_\tau)^2 \right] + \tau \, \mathrm{Tr}\,\Theta \tag{1.58}$$

$$= \mathbb{E}_{p(\theta_0)} \left[ \theta_\tau^T \Delta\theta_\tau \right] + \tau \, \mathrm{Tr}\,\Theta \tag{1.59}$$

$$= \mathbb{E}_{p(\theta_0)} \left[ (z_0(\mathcal{X}) - \mathcal{Y})^T (I - e^{-\tau\Theta})^2 \Theta^{-1} (z_0(\mathcal{X}) - \mathcal{Y}) \right] + \tau \, \mathrm{Tr}\,\Theta \tag{1.60}$$

$$= \mathrm{Tr} \left( \mathcal{K}\Theta^{-1}(I - e^{-\tau\Theta})^2 \right) + \mathcal{Y}^T \Theta^{-1}(I - e^{-\tau\Theta})^2 \mathcal{Y} + \tau \, \mathrm{Tr}\,\Theta. \tag{1.61}$$

This tends to infinity as the time goes to infinity. This renders the usual PAC-Bayes style generalization bounds trivially vacuous for the generalization of infinitely wide neural networks at late times. Yet, infinite networks can generalize well (Arora et al.; 2019).

*WAIC*

Watanabe-Akaike Information Criterion (WAIC) first introduced in (Watanabe and Opper; 2010), and gives an asymptotically correct estimate of the gap between the training set and test set expectations. It is defined by the difference between Bayes' and Gibb's' errors;

$$WAIC = \sum_i \left( \log(\mathbb{E}\left[ p(y_i|\theta) \right]) - \mathbb{E}[[\log(p(y_i|\theta))]] \right)$$

Training with squared loss is equivalent to assuming a Gaussian observation model $p(y|z) \sim \mathcal{N}(z, I)$ and as we showed before that the expected log loss (Gibbs loss) has contributions both from the square loss of the mean prediction, as well as a term which couples to the trace of the covariance:

$$\mathbb{E}\left[ \log q(y|z) \right] = \tag{1.62}$$

$$-\frac{1}{2}(y - \mu_t)^2 + \frac{1}{2}Tr(\Sigma) \tag{1.63}$$

$$+\frac{k}{2} \log \left( (2\pi)|I + \Sigma_t| \right) \tag{1.64}$$

For the Bayes loss

$$\log \mathbb{E}\left[ q(y|z) \right] = -\frac{1}{2}(y - \mu_t)^T (I + \Sigma_t))^{-1} (y - \mu_t)$$

$$+\frac{k}{2} \log \left( (2\pi)|I + \Sigma_t| \right)$$

The Woodbury matrix identity tells us –

$$(A + UCV)^{-1} = A^{-1} - A^{-1}U \left(C^{-1} + VA^{-1}U\right)^{-1} VA^{-1}$$

and if $A = I$, and $C = I$, and $U = I$ –

$$(I + V)^{-1} = I - (I + V)^{-1} V$$

we get:

$$log \mathbb{E}\left[q(y|z)\right] = -\frac{1}{2} (y - \mu_t)^2 + \frac{1}{2} (y - \mu_t)(I + \Sigma_t))^{-1} \Sigma (y - \mu_t)$$
$$+ \frac{k}{2} \log\left((2\pi)|I + \Sigma_t|\right).$$

Combining all together

$$WAIC = \frac{1}{2} (y - \mu_t)(I + \Sigma_t))^{-1} \Sigma (y - \mu_t) + \frac{1}{2}Tr(\Sigma)$$

### Experiments

The Gaussian Information Bottleneck  (Chechik et al.; 2005) gives the optimal trade-off betweenn $I(Z; X)$ and $I(Z; Y)$ for jointly Gaussian data, where $X$ is the input, $Y$ is the label, and $Z$ is a stochastic representation, $p(z|x)$ of the input.
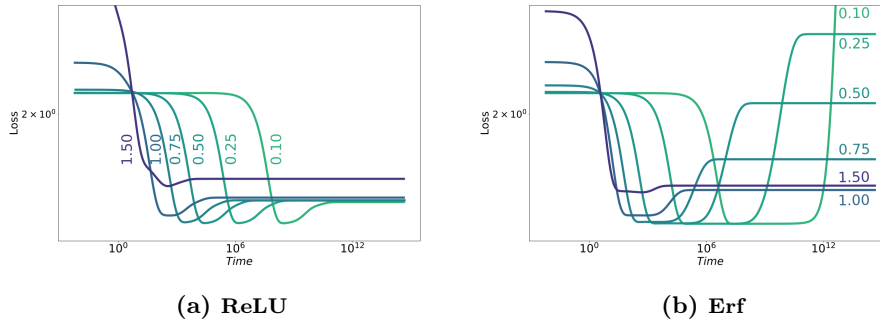
In the following, we fit infinite ensembles of infinitely-wide neural networks to jointly Gaussian data and calculate estimates of their mutual information. This allows us to determine how close these networks are to being optimal.

The Gaussian dataset we create has $|X| = 30$ and $|Y| = 1$ (for details, see Appendix A). We train a three-layer FC network with both RELU and ERF activation functions.

**??** shows the test set loss as a function of time for different choices of initial weight variance ($\sigma_w^2$). For both the RELU and ERF networks, at the highest $\sigma_w$ shown (darkest purple), the networks *underfit*. For lower initial weight variances, they all show signs of *overfitting* in the sense that the networks would benefit from early stopping. This overfitting is worse for the ERF nonlinearity, where we see a divergence in the final test set loss as $\sigma_w$ decreases. For all of these networks, the training loss goes to zero.

In **??** we show the performance of these networks on the information plane. The $x$-axis shows a variational lower bound on the complexity of the learned representation: $I(Z; X|D)$. The $y$-axis shows a variational lower bound on learned relevant information: $I(Y; Z)$. For details on the calculation of the MI estimates see **??**. The curves show trajectories of the networks' representation as time varies from $\tau = 10^{-2}$ to $\tau = 10^{10}$ for different weight variances (the bias-variance in all networks was fixed to 0.01). The red line is the optimal theoretical IB bound.

(a) ReLU ............ (b) Erf

**Figure 1.12: Loss as function of time for different initial weights' variances on the Gaussian dataset.**

There are several features worth highlighting. First, we emphasize the somewhat surprising result that, as time goes to infinity, the MI between an infinite ensemble of infinitely-wide neural networks output and their input is finite and quite small. Even though every individual network provides a seemingly rich deterministic representation of the input, when we marginalize over the random initialization, the ensemble compresses the input quite strongly. The networks overfit at late times. For ERF networks, the more complex representations ($I(Z;X|D)$) overfit more. With optimal early stopping, over a wide range, these models achieve a near-optimal tradeoff in prediction versus compression. Varying the initial weight variance controls the amount of information the ensemble extracts.



(a) ReLU ............ (b) ERF

**Figure 1.13: Trajectories of the (bounds on) MI between the representation $Z$ and the input $X$ versus time. Curves differ only in their initial weight variance. The red line is the optimal IB as predicted by theory. Our estimate for $I(Z;X)$ is upper bounded by the log of the batch size ( $\log 1000 = 6.9$.)**

Next, we repeat the result of the previous section on the MNIST dataset (LeCun et al.; 2010). Unlike the typical setup, we turn MNIST into a binary regression task for the parity of the digit (even or odd). This time, the network is a standard two-layer convolutional

neural network with $5 \times 5$ filters and either RELU or ERF activation functions.

?? shows the results. Unlike in the jointly Gaussian dataset case, here both networks show some region of initial weight variances that do not overfit in the sense of demonstrating any advantage from early stopping. The ERF network at higher variances does show overfitting at low initial weight variances, but the RELU network does not. Notice that in the information plane, the ERF network shows overfitting at higher representational complexities ($I(Z; X)$ large), while the RELU network does not.



(a) ReLU                                    (b) ERF

(c) ReLU                                    (d) ERF

Figure 1.14: Loss as function of time and information plan trajectories for different initial weights' variances on MNIST.

### Conclusions

Infinite ensembles of infinitely-wide neural networks provide an interesting model family. Being linear in their parameters, they permit a high number of tractable calculations of information-theoretic quantities and their bounds. Despite their simplicity, they still can achieve good generalization performance (Arora et al.; 2019). This challenges existing claims for the purported connections between information theory and generalization in deep neural networks. In this preliminary work, we laid the groundwork for a larger-scale empirical and theoretical study of generalization in this simple model family. Given that real networks approach this Family in their infinite width limit, we believe a better

understanding of generalization in the NTK limit will shed light on generalization in deep neural networks.

### Bibliography

Achille, A. and Soatto, S. (2017). Emergence of Invariance and Disentangling in Deep Representations, *Proceedings of the ICML Workshop on Principled Approaches to Deep Learning* .

Achille, A. and Soatto, S. (2019). Where is the information in a deep neural network?

Alemi, A. A., Fischer, I., Dillon, J. V. and Murphy, K. (2016). Deep variational information bottleneck, *arXiv:1612.00410* .
**URL:** *http://arxiv.org/abs/1612.00410*

Amjad, R. A. and Geiger, B. C. (2018). How (not) to train your neural network using the information bottleneck principle, *arXiv preprint arXiv:1802.09766* .

Arora, S., Du, S. S., Li, Z., Salakhutdinov, R., Wang, R. and Yu, D. (2019). Harnessing the power of infinitely wide deep nets on small-data tasks.

Asadi, A., Abbe, E. and Verdu, S. (2018). Chaining mutual information and tightening generalization bounds, *in* S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi and R. Garnett (eds), *Advances in Neural Information Processing Systems 31*, Curran Associates, Inc., pp. 7234–7243.
**URL:** *http://papers.nips.cc/paper/7954-chaining-mutual-information-and-tightening-generalization-bounds.pdf*

Banerjee, A. (2006). On bayesian bounds, *Proceedings of the 23rd international conference on Machine learning*, ACM, pp. 81–88.

Bartlett, P. L. and Mendelson, S. (2002). Rademacher and gaussian complexities: Risk bounds and structural results, *Journal of Machine Learning Research* **3**(Nov): 463–482.

Bassily, R., Moran, S., Nachum, I., Shafer, J. and Yehudayoff, A. (2017). Learners that use little information.

Boucheron, S., Bousquet, O. and Lugosi, G. (2005). Theory of classification: A survey of some recent advances, *ESAIM: probability and statistics* **9**: 323–375.

Bousquet, O. and Elisseeff, A. (2002). Stability and generalization, *Journal of machine learning research* **2**(Mar): 499–526.

Chechik, G., Globerson, A., Tishby, N. and Weiss, Y. (2005). Information bottleneck for gaussian variables, *Journal of machine learning research* **6**(Jan): 165–188.

Chen, R. T. Q., Rubanova, Y., Bettencourt, J. and Duvenaud, D. (2018). Neural ordinary differential equations.

Cover, T. M. and Thomas, J. A. (2012). *Elements of information theory*, John Wiley & Sons.

Hardt, M., Recht, B. and Singer, Y. (2016). Train faster, generalize better: Stability of stochastic gradient descent, *International Conference on Machine Learning*, pp. 1225–1234.

Keskar, N. S., Mudigere, D., Nocedal, J., Smelyanskiy, M. and Tang, P. T. P. (2016). On large-batch training for deep learning: Generalization gap and sharp minima, *arXiv preprint arXiv:1609.04836* .

Kolchinsky, A., Tracey, B. D. and Van Kuyk, S. (2018). Caveats for information bottleneck in deterministic scenarios, *arXiv preprint arXiv:1808.07593* .

Kunstner, F., Balles, L. and Hennig, P. (2019). Limitations of the empirical fisher approximation.

LeCun, Y., Cortes, C. and Burges, C. (2010). Mnist handwritten digit database, *ATT Labs [Online]. Available: http://yann.lecun.com/exdb/mnist* **2**.

Lee, J., Xiao, L., Schoenholz, S. S., Bahri, Y., Novak, R., Sohl-Dickstein, J. and Pennington, J. (2019). Wide Neural Networks of Any Depth Evolve as Linear Models Under Gradient Descent, *arXiv e-prints* p. arXiv:1902.06720.

Liang, T., Poggio, T., Rakhlin, A. and Stokes, J. (2019). Fisher-rao metric, geometry, and complexity of neural networks, *The 22nd International Conference on Artificial Intelligence and Statistics*, pp. 888–896.

Negrea, J., Haghifam, M., Dziugaite, G. K., Khisti, A. and Roy, D. M. (2019). Information-theoretic generalization bounds for sgld via data-dependent estimates, *in* H. Wallach, H. Larochelle, A. Beygelzimer, F. d Alché-Buc, E. Fox and R. Garnett (eds), *Advances in Neural Information Processing Systems 32*, Curran Associates, Inc., pp. 11015–11025.
**URL:** *http://papers.nips.cc/paper/9282-information-theoretic-generalization-bounds-for-sgld-via-data-dependent-estimates.pdf*

Neyshabur, B., Tomioka, R. and Srebro, N. (2014). In search of the real inductive bias: On the role of implicit regularization in deep learning, *arXiv preprint arXiv:1412.6614* .

Neyshabur, B., Tomioka, R. and Srebro, N. (2015). Norm-based capacity control in neural networks, *Conference on Learning Theory*, pp. 1376–1401.

Novak, R., Xiao, L., Hron, J., Lee, J., Alemi, A. A., Sohl-Dickstein, J. and Schoenholz, S. S. (2019). Neural tangents: Fast and easy infinite neural networks in python, *arXiv preprint arXiv:1912.02803* .

Pensia, A., Jog, V. and Loh, P.-L. (2018). Generalization error bounds for noisy, iterative algorithms, *2018 IEEE International Symposium on Information Theory (ISIT)*, IEEE, pp. 546–550.

Poole, B., Ozair, S., van den Oord, A., Alemi, A. A. and Tucker, G. (2019). On variational bounds of mutual information, *CoRR* **abs/1905.06922**.
**URL:** *http://arxiv.org/abs/1905.06922*

Russo, D. and Zou, J. (2016). Controlling bias in adaptive data analysis using information theory, *Artificial Intelligence and Statistics*, pp. 1232–1240.

Russo, D. and Zou, J. (2019). How much does your data exploration overfit? controlling bias via information usage, *IEEE Transactions on Information Theory* **66**(1): 302–323.

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.

Shwartz-Ziv, R. and Tishby, N. (2017). Opening the black box of deep neural networks via information, *arXiv preprint arXiv:1703.00810* .

Slonim, N. (2002). *The information bottleneck: Theory and applications*, PhD thesis, Citeseer.

Steinke, T. and Zakynthinou, L. (2020). Reasoning about generalization via conditional mutual information, *arXiv preprint arXiv:2001.09122* .

Tishby, N., Pereira, F. C. and Bialek, W. (1999). The information bottleneck method, *In Proceedings of the 37-th Annual Allerton Conference on Communication, Control and Computing* .

Tishby, N. and Zaslavsky, N. (2015). Deep learning and the information bottleneck principle, *2015 IEEE Information Theory Workshop (ITW)*, IEEE, pp. 1–5.

Vapnik, V. N. (1998). Statistical learning theory, *Wiley* .

Watanabe, S. and Opper, M. (2010). Asymptotic equivalence of bayes cross validation and widely applicable information criterion in singular learning theory., *Journal of machine learning research* **11**(12).

Xu, A. and Raginsky, M. (2017). Information-theoretic analysis of generalization capability of learning algorithms, *Advances in Neural Information Processing Systems*, pp. 2524–2533.

Zhang, C., Bengio, S., Hardt, M., Recht, B. and Vinyals, O. (2016). Understanding deep learning requires rethinking generalization, *arXiv preprint arXiv:1611.03530* .

*Appendix*

## *Appendix A - Gaussian Dataset*

For our experiments we used a jointly Gaussian dataset, for which there is an analytic solution for the optimal representation (Chechik et al.; 2005).

Imagine a jointly Gaussian dataset, where we have $x_{ij} = L^x_{jk}\epsilon^x_{ik}$ with $\epsilon \sim \mathcal{N}(0,1)$. Make $y$ just an affine projection of $x$ with added noise.

$$y_{ij} = L^y_{jk}\epsilon_{ik} + A_{jk}x_{ik} = L^y_{jk}\epsilon^y_{ik} + A_{jk}L^x_{km}\epsilon^x_{im}. \tag{1.65}$$

Both $x$ and $y$ will be mean zero. We can compute their covariances.

$$\Sigma^x_{jk} = \langle x_{ij}x_{ik}\rangle = \langle L_{jm}\epsilon_{im}L_{kl}\epsilon_{il}\rangle = L_{jm}L_{kl}\delta_{ml} = L_{jm}L_{km} \tag{1.66}$$

Next look at the covariance of $y$.

$$\begin{aligned}
\Sigma^y_{jk} &= \langle y_{ij}y_{ik}\rangle \\
&= \left\langle \left(L^y_{jl}\epsilon^y_{il} + A_{jl}L^x_{lm}\epsilon^x_{im}\right)\left(L^y_{kn}\epsilon^y_{in} + A_{kn}L^x_{no}\epsilon^x_{io}\right)\right\rangle \\
&= L^y_{jl}L^y_{kn}\delta_{ln} + A_{jl}L^x_{lm}A_{kn}L_{no}\delta_{mo} \\
&= L^y_{jn}L^y_{kn} + A_{jl}\Sigma^x_{ln}A_{kn}
\end{aligned}$$

For the cross covariance:

$$\begin{aligned}
\Sigma^{xy}_{jk} &= \langle x_{ij}y_{ik}\rangle \\
&= \left\langle L^x_{jm}\epsilon^x_{im}\left(L^y_{kn}\epsilon^y_{in} + A_{kn}L^x_{no}\epsilon^x_{io}\right)\right\rangle \\
&= L^x_{jm}A_{kn}L^x_{no}\delta_{mo} \\
&= L^x_{jm}A_{kn}L^x_{nm} = \Sigma^x_{jn}A_{kn}
\end{aligned}$$

So we have for our entropy of $x$:

$$H(X) = \frac{n_x}{2}\log(2\pi e) + n_x \log\sigma_x$$

$$H(Y|X) = \frac{n_y}{2}\log(2\pi e) + n_y \log\sigma_y$$

as for the marginal entropy, we will assume the SVD decomposition $A = U\Sigma V^T$

$$H(Y) = \frac{n_y}{2}\log(2\pi e) + \frac{n_y}{2}\log\left|\sigma_y^2 I + \sigma_x^2 AA^T\right| = \frac{n_y}{2}\log(2\pi e) + \frac{1}{2}\sum_i \log\left(\sigma_y^2 + \sigma_x^2\Sigma_i^2\right)$$

So, solving for the mutual information between $x$ and $y$ we obtain:

$$I(X;Y) = H(Y) - H(Y|X) = \frac{1}{2} \sum_i \log \left( 1 + \frac{\sigma_x^2 \Sigma_i^2}{\sigma_y^2} \right)$$

## Appendix B - Additional Results

shows $I(X;Z|D)$, $I(\theta;D)$, $\frac{dI(\theta;D)}{dt}$ and the loss as function of time for a fixed initial weight's variance ($\sigma_w = 0.25$). (in log-log scale, notice that $y$-axes are different for each measure). For both the RELUand ERFnetworks, we see clear features in each plot near the optimal test loss.

### WAIC

We know that the WAIC is the diffrence betweem the Bayes and the Gibbs errors, namely -

$$WAIC = \sum_i \left( \log(\mathbb{E}\left[p(y_i|\theta)\right]) - \mathbb{E}[[\log(p(y_i|\theta))]] \right)$$

Training with squared loss is equivalent to assuming a Gaussian observation model $p(y|z) \sim \mathcal{N}(z, I)$. We can marginalize out our representation to obtain

$$\mathbb{E}\left[q(y|z)\right] = q(y|x) = \int dz\, q(y|z) p(z|x) \sim \mathcal{N}(\mu(x,\tau), I + \Sigma(x,\tau)). \qquad (1.67)$$

The expected log loss (Gibbs loss) has contributions both from the square loss of the mean prediction, as well as a term which couples to the trace of the covariance:

$$\mathbb{E}\left[\log q(y|z)\right] = \qquad\qquad (1.68)$$

$$-\frac{1}{2}\mathbb{E}\left[(y - z_t)^2\right] + \frac{1}{2} \log\left((2\pi)^k |I + \Sigma_t)|\right) = \qquad (1.69)$$

$$-\frac{1}{2}(y - \mu_t)^2 + \frac{1}{2} Tr(\Sigma) \qquad (1.70)$$

$$+\frac{k}{2} \log\left((2\pi)|I + \Sigma_t|\right) \qquad (1.71)$$

here $k$ is the dimensionality of $y$.

For the Bayes loss

$$\log \mathbb{E}\left[q(y|z)\right] = -\frac{1}{2}(y - \mu_t)^T (I + \Sigma_t))^{-1} (y - \mu_t)$$

$$+\frac{k}{2} \log\left((2\pi)|I + \Sigma_t|\right)$$

and by the Woodbury matrix identity -

$$(I + V)^{-1} = I - (I + V)^{-1} V$$

and we get

$$log\mathbb{E}\left[q(y|z)\right] = -\frac{1}{2}\left(y - \mu_t\right)^2 + \frac{1}{2}\left(y - \mu_t\right)\left(I + \Sigma_t\right)^{-1}\Sigma\left(y - \mu_t\right)$$
$$+ \frac{k}{2}\log\left((2\pi)|I + \Sigma_t|\right)$$

Combining all together

$$WAIC = \frac{1}{2}\left(y - \mu_t\right)\left(I + \Sigma_t\right)^{-1}\Sigma\left(y - \mu_t\right) + \frac{1}{2}Tr(\Sigma)$$

$I(\theta; D$

$$I(\theta; D) =$$
$$\mathbb{E}_{p(\theta_t)}\left[\log p(\theta_t|D) - \log p(\theta_t)\right]$$
$$= \mathbb{E}_{p(\theta_t)}\left[\log p(\theta_t|D) - \log\int dD' p(D')p(\theta_t|D')\right]$$

# The Dual Information Bottleneck

# The Dual Information Bottleneck

Zoe Piran [1] Ravid Shwartz-Ziv [2] Naftali Tishby[2,3]

[1] Racah Institute of Physics
The Hebrew University of Jerusalem
Jerusalem, Israel
[2] The Edmond and Lilly Safra Center for Brain Sciences, The Hebrew University,
Jerusalem, Israel.
[3] School of Computer Science and Engineering,
The Hebrew University,
Jerusalem, Israel.

## Abstract

The Information Bottleneck (IB) framework is a general characterization of optimal representations obtained using a principled approach for balancing accuracy and complexity. Here we present a new framework, the Dual Information Bottleneck (dualIB), which resolves some of the known drawbacks of the IB. We provide a theoretical analysis of the dualIB framework; (i) solving for the structure of its solutions (ii) unraveling its superiority in optimizing the *mean prediction error exponent* and (iii) demonstrating its ability to preserve exponential forms of the original distribution. To approach large scale problems, we present a novel variational formulation of the dualIB for Deep Neural Networks. In experiments on several data-sets, we compare it to a variational form of the IB. This exposes superior Information Plane properties of the dualIB and its potential in improvement of the error.

## *Introduction*

The Information Bottleneck (IB) method Tishby et al. (1999), is an information-theoretic framework for describing efficient representations of a "input" random variable $X$ that preserve the information on an "output" variable $Y$. In this setting the joint distribution of $X$ and $Y$, $p(x, y)$, defines the problem, or rule, and the training data are a finite sample from this distribution. The stochastic nature of the label is essential for the analytic regularity of the IB problem. In the case of deterministic labels, we assume a *noise model* which induces a distribution. The representation variable $\hat{X}$ is in general a stochastic

93

function of $X$ which forms a Markov chain $Y \to X \to \hat{X}$, and it depends on $Y$ only through the input $X$. We call the map $p(\hat{x} \mid x)$ the *encoder* of the representation and denote by $p(y \mid \hat{x})$ the *Bayes optimal decoder* for this representation; i.e., the best possible prediction of the *desired label* $Y$ from the representation $\hat{X}$.

The IB has direct successful applications for representation learning in various domains, from vision and speech processing Ma et al. (2019), to neuroscience Schneidman et al. (2001), and Natural Language Processing Li and Eisner (2019). Due to the notorious difficulty in estimating mutual information in high dimension, variational approximations to the IB have been suggested and applied also to Deep Neural Networks (DNNs) (e.g., Alemi et al.; 2016; Achille and Soatto; 2018a; Parbhoo et al.; 2018; Poole et al.; 2019). Additionally, following (Shwartz-Ziv and Tishby; 2017), several recent works tackled the problem of understanding DNNs using the IB principle (Nash et al.; 2018; Goldfeld et al.; 2018)

Still, there are several drawbacks to the IB framework which motivated this work. While the standard approach in representation learning is to use the topology or a specific parametric model over the input, the IB principle and equations are completely non-parametric, and operate directly on the encoder and decoder probability distributions. Moreover, the original IB formulation, as common in Information Theory, assumes full knowledge of the pattern-label distribution and does not relate directly to the task of prediction the label for unseen input patterns, when trained from finite samples. These issues were addressed before for general learning with the IB (Slonim et al.; 2006; Shamir et al.; 2010) and by extending it in the context of large DNNs by Achille et. al. (Achille and Soatto; 2018a; Achille et al.; 2019).

Here, we address the above drawbacks by introducing a novel theoretical framework, the Dual Information Bottleneck (dualIB). The dualIB can account for known features of the data and use them to make better predictions over unseen examples, from small samples for large scale problems. Further, it emphasizes the prediction problem, inferring $\hat{Y}$, which wasn't present in the original IB formulation due to the complete distributional knowledge assumption.

### 1.7.4 Contributions of this work

We present here the Dual Information Bottleneck (dualIB) aiming to obtain optimal representations, which resolves the IB drawbacks:

- We provide a theoretical analysis which obtains an analytical solution to the framework and compare its behaviour to the IB.

- For data which can be approximated by exponential families we provide closed

form solutions, dualExpIB, which preserves the sufficient statistics of the original distribution.

- We show that by accounting for the prediction variable, the dualIB formulation optimizes a bound over the error exponent.

- We present a novel variational form of the dualIB for Deep Neural Networks (DNNs) allowing its application to real world problems. Using it we empirically investigate the dynamics of the dualIB and validate the theoretical analysis.

### Background

The Information Bottleneck (IB) framework is defined as the trade off between the encoder and decoder mutual information values. It is defined by the minimization of the Lagrangian:

$$\mathcal{F}[p_\beta(\hat{x} \mid x); p_\beta(y \mid \hat{x})] = I(X; \hat{X}) - \beta I(Y; \hat{X}) \,, \tag{1.72}$$

independently over the convex sets of the normalized distributions, $\{p_\beta(\hat{x} \mid x)\}$, $\{p_\beta(\hat{x})\}$ and $\{p_\beta(y \mid \hat{x})\}$, given a positive Lagrange multiplier $\beta$ constraining the information on $Y$, while preserving the Markov Chain $Y \to X \to \hat{X}$. Three self-consistent equations for the optimal encoder-decoder pairs, known as the IB *equations*, define the solutions to the problem. An important characteristic of the equations is the existence of critical points along the optimal line of solutions in the *information plane* (presenting $I(Y; \hat{X})$ vs. $I(X; \hat{X})$) (Wu and Fischer; 2020; Parker et al.; 2003). The IB optimization trade off can be considered as a generalized rate-distortion problem Cover and Thomas (2006) with the distortion function, $d_{\text{IB}}(x, \hat{x}) = D[p(y \mid x)||p_\beta(y|\hat{x})]$. For more background on the IB framework see §1.7.7.3.

### The Dual Information Bottleneck

Supervised learning is generally separated into two phases: the training phase, in which the internal representations are formed from the training data, and the prediction phase, in which these representations are used to predict labels of new input patterns (Shalev-Shwartz and Ben-David; 2014). To explicitly address these different phases we add to the IB Markov chain another variable, $\hat{Y}$, the *predicted label* from the trained representation, which obtains the same values as $Y$ but is distributed differently:

$$\overbrace{Y \to \underbrace{X \to \hat{X}_\beta}_{\text{}} \to \hat{Y}}^{\text{training}}_{\text{prediction}} . \tag{1.73}$$

The left-hand part of this chain describes the representation training, while the right-hand part is the Maximum Likelihood prediction using these representations Slonim et al. (2006). So far the prediction variable $\hat{Y}$ has not been a part of the IB optimization problem. It has been implicitly assumed that the *Bayes optimal decoder*, $p_\beta(y \mid \hat{x})$, which maximizes the full representation-label information, $I(Y; \hat{X})$, for a given $\beta$, is also the best choice for making predictions. Namely, the prediction of the label, $\hat{Y}$, from the representation $\hat{X}_\beta$ through the right-hand Markov chain by the mixture using the internal representations, $p_\beta(\hat{y} \mid x) \equiv \sum_{\hat{x}} p_\beta(y = \hat{y} \mid \hat{x})p_\beta(\hat{x} \mid x)$, is optimal when $p_\beta(y \mid \hat{x})$ is the *Bayes optimal decoder*. However, this is not necessarily the case, for example when we train from finite samples Shamir et al. (2010).

Focusing on the prediction problem, we define the dualIB distortion by switching the order of the arguments in the KL-divergence of the original IB distortion, namely:

$$d_{\text{dualIB}}(x, \hat{x}) = D[p_\beta(y \mid \hat{x}) \| p(y \mid x)] = \sum_y p_\beta(y \mid \hat{x}) \log \frac{p_\beta(y \mid \hat{x})}{p(y \mid x)} \ . \qquad (1.74)$$

In geometric terms this is known as the *dual* distortion problem Felice and Ay (2019). The dualIB optimization can then be written as the following rate-distortion problem:

$$\mathcal{F}^*[p_\beta(\hat{x} \mid x); p_\beta(y \mid \hat{x})] = I(X; \hat{X}) + \beta \mathbb{E}_{p_\beta(x, \hat{x})}[d_{\text{dualIB}}(x, \hat{x})] \ . \qquad (1.75)$$

As the decoder defines the prediction stage ($p_\beta(y = \hat{y} \mid \hat{x})$) we can write (see proof in §1.7.10) the average distortion in terms of mutual information on $\hat{Y}$, $I(\hat{X}; \hat{Y})$ and $I(X; \hat{Y})$:

$$\mathbb{E}_{p_\beta(x, \hat{x})}[d_{\text{dualIB}}(x, \hat{x})] = \underbrace{I(\hat{X}; \hat{Y}) - I(X; \hat{Y})}_{(a)} + \underbrace{\mathbb{E}_{p(x)}[D[p_\beta(\hat{y} \mid x) \| p(y = \hat{y} \mid x)]]}_{(b)}. \qquad (1.76)$$

This is similar to the known IB relation: $\mathbb{E}_{p_\beta(x, \hat{x})}[d_{\text{IB}}(x, \hat{x})] = I(Y; X) - I(Y; \hat{X})$ with an extra positive term $(b)$. Both terms, $(a)$ and $(b)$, vanish precisely when $\hat{X}$ is a sufficient statistic for $X$ with respect to $\hat{Y}$. In such a case we can reverse the order of $X$ and $\hat{X}$ in the Markov chain equation **??**. This replaces the roles of $Y$ and $\hat{Y}$ as the variable for which the representations, $\hat{X}_\beta$, are approximately minimally sufficient statistics. In that sense the dualIB shifts the emphasis from the training phase to the prediction phase. This implies that minimizing the dualIB functional maximizes a lower bound on the mutual information $I(X; \hat{Y})$.

### 1.7.5 *The* dualIB *equations*

Solving the dualIB minimization problem equation **??**, we obtain a set of self consistent equations. Generalized Blahut-Arimoto (BA) iterations between them converges to an optimal solution. The equations are similar to the original IB equations with the following

modifications: (i) Replacing the distortion by its dual in the encoder update; (ii) Updating the decoder by the encoder's geometric mean of the data distributions $p(y \mid x)$.

**Theorem 1.7.2.** *The* dualIB *equations are given by:*

$$
\begin{cases}
(i) & p_\beta(\hat{x} \mid x) = \frac{p_\beta(\hat{x})}{Z_{\hat{\mathbf{x}}\mid\mathbf{x}}(x;\beta)} e^{-\beta D[p_\beta(y\mid\hat{x})\|p(y\mid x)]} \\
(ii) & p_\beta(\hat{x}) = \sum_x p_\beta(\hat{x} \mid x) p(x) \\
(iii) & p_\beta(y \mid \hat{x}) = \frac{1}{Z_{\mathbf{y}\mid\hat{\mathbf{x}}}(\hat{x};\beta)} \prod_x p(y \mid x)^{p_\beta(x\mid\hat{x})}
\end{cases}
, \qquad (1.77)
$$

*where $Z_{\hat{\mathbf{x}}\mid\mathbf{x}}(x;\beta), Z_{\mathbf{y}\mid\hat{\mathbf{x}}}(\hat{x};\beta)$ are normalization terms.*

The proof is given in §1.7.10. It is evident that the basic structure of the equations of the dualIB and IB is similar and they approach each other for large values of $\beta$. In the following sections we explore the implication of the differences on the properties of the new framework.

### 1.7.6  The critical points of the dualIB

As mentioned in §1.7.4 and Wu and Fischer (2020), the "skeleton" of the IB optimal bound (the information curve) is constituted by the critical points in which the topology (cardinality) of the representation changes. Using perturbation analysis over the dualIB optimal representations we find that small changes in the encoder and decoder that satisfy equation ?? for a given $\beta$ are approximately determined through a nonlinear eigenvalues problem. [2]

**Theorem 1.7.3.** *The* dualIB *critical points are given by non-trivial solutions of the nonlinear eigenvalue problem:*

$$
\left[I - \beta C_{xx'}^{\text{dualIB}}(\hat{x},\beta)\right]\delta\log p_\beta(x' \mid \hat{x}) = 0 \ , \qquad \left[I - \beta C_{yy'}^{\text{dualIB}}(\hat{x},\beta)\right]\delta\log p_\beta(y' \mid \hat{x}) = 0. \tag{1.78}
$$

*The matrices $C_{xx'}^{\text{dualIB}}(\hat{x};\beta), C_{yy'}^{\text{dualIB}}(\hat{x};\beta)$ have the same eigenvalues $\{\lambda_i\}$, with $\lambda_1(\hat{x}) = 0$. With binary $y$, the critical points are obtained at $\lambda_2(\hat{x}) = \beta^{-1}$.*

The proof to Theorem 1.7.3 along with the structure of the matrices $C_{xx'}^{\text{dualIB}}(\hat{x};\beta), C_{yy'}^{\text{dualIB}}(\hat{x};\beta)$ is given in §1.7.10. We found that this solution is similar to the nonlinear eigenvalues problem for the IB, given in §1.7.7.3. As in the IB, at the critical points we observe cusps with an undefined second derivative in the mutual information values as functions of $\beta$ along the optimal line. That is the general structure of the solutions is preserved between the frameworks, as can be seen in Figure 1.15c.

---

[2] For simplicity we ignore here the possible interactions between the different representations.

The *Information Plane*, $I_y = I(\hat{X}; Y)$ vs. $I_x = I(\hat{X}; X)$, is the standard depiction of the compression-prediction trade-off of the IB and has known analytic properties(Giladbachrach et al.; 2003). First, we note that both curves obey similar constraints, as given in *lemma* 1.7.4 below.

**Lemma 1.7.4.** *along the optimal lines of $I_x(\beta)$ and $I_y(\beta)$ the curves are non-decreasing piece-wise concave functions of $\beta$. When their second derivative (with respect to $\beta$) is defined, it is strictly negative.*

Next, comparing the dualIB's and IB's information planes we find several interesting properties which are summarized in the following Theorem (see §1.7.11 for the proof).
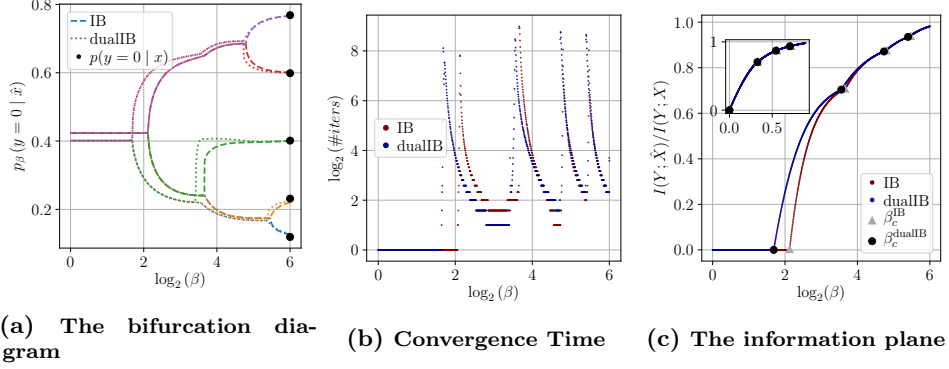
**Theorem 1.7.5.** *(i) The critical points of the two algorithms alternate, $\forall i, i+1$, $\beta_{c,i}^{\text{dualIB}} \leq \beta_{c,i}^{\text{IB}} \leq \beta_{c,i+1}^{\text{dualIB}} \leq \beta_{c,i}^{\text{IB}}$. (ii) The distance between the two information curves is minimized at $\beta_c^{\text{dualIB}}$. (iii) The two curves approach each other as $\beta \to \infty$.*

From Theorem 1.7.5 we deduce that as the dimensionality of the problem increases (implying the number of critical points grows) the dualIB's approximation of the IB's information plane becomes tighter. We illustrate the behavior of the dualIB's solutions in comparison to the IB's on a low-dimensional problem that is easy to analyze and visualize, with a binary $Y$ and only 5 possible inputs $X$ (the complete definition is given in §1.7.11). For any given $\beta$, the encoder-decoder iterations converge to stationary solutions of the dualIB or IB *equations*. The evolution of the optimal decoder, $p_\beta(y = 0 \mid \hat{x})$, $\forall \hat{x} \in \hat{X}$, as a function of $\beta$, forms a *bifurcation diagram* (Figure 1.15a), in which the critical points define the location of the bifurcation. At the critical points the number of iterations diverges (Figure 1.15b). While the overall structure of the solutions is similar, we see a "shift" in the appearance of the representation splits between the two frameworks. Specifically, as predicted by Theorem 1.7.5 the dualIB bifurcations occur at lower $\beta$ values than those of the IB. The inset of Figure 1.15c depicts this comparison between the two information curves. While we know that $I_y^{\text{IB}}(\beta)$ is always larger, we see that for this setting the two curves are almost indistinguishable. Looking at $I_y$ as a function of $\beta$ (Figure 1.15c) the importance of the critical points is revealed as the corresponding cusps along these curves correspond to "jumps" in the accessible information. Furthermore, the distance between the curves is minimized precisely at the dual critical points, as predicted by the theory.

### The Exponential Family dualIB

One of the major drawbacks of the IB is that it fails to capture an existing parameterization of the data, that act as minimal sufficient statistics for it. Exponential families are the class of parametric distributions for which minimal sufficient statistics exist, forming an elegant theoretical core of parametric statistics and often emerge as maximum entropy Jaynes (1957) or stochastic equilibrium distributions, subject to observed expectation constraints.

**(a) The bifurcation diagram**

**(b) Convergence Time**

**(c) The information plane**

**Figure 1.15:** (*a*) The *bifurcation diagram*; each color corresponds to one component of the representation $\hat{x} \in \hat{X}$ and depicts the decoder $p_\beta(y = 0 \mid \hat{x})$. Dashed lines represent the IB's solution and dotted present the dualIB's. The black dots denote the input distribution $p(y = 0 \mid x)$. (*b*) Convergence time the BA algorithms as a function of $\beta$. (*c*) The desired label Information $I_y^{\mathrm{IB}}(\beta)$ and $I_y^{\mathrm{dualIB}}(\beta)$ as functions of $\beta$. The inset shows the information plane, $I_y$ vs. $I_x$. The black dots are the dualIB's critical points, $\beta_c^{\mathrm{dualIB}}$, and the grey triangles are the IB's, $\beta_c^{\mathrm{IB}}$

As the IB ignores the structure of the distribution, given data from an exponential family it won't consider these known features. Contrarily, the dualIB accounts for this structure and its solution are given in terms of these features, defining the dualExpIB equations.

We consider the case in which the rule distribution is of the form, $p(y \mid x) = e^{-\sum_{r=0}^{d} \lambda^r(y) A_r(x)}$, where $A_r(x)$ are $d$ functions of the input $x$ and $\lambda^r(y)$ are functions of the label $y$, or the parameters of this exponential family [3]. For exponential forms the mutual information, $I(X;Y)$, is fully captured by the $d$ conditional expectations. This implies that all the relevant information (in the training sample) is captured by $d$-dimensional empirical expectations which can lead to a reduction in computational complexity.

Next we show that for the dualIB, for all values of $\beta$, this dimension reduction is preserved or improved along the dual information curve. The complete derivations are given in §1.7.13.

**Theorem 1.7.6.** (*dualExpIB*) *For data from an exponential family the optimal encoder-decoder of the* dualIB *are given by:*

$$p_\beta(\hat{x} \mid x) = \frac{p_\beta(\hat{x}) e^{\beta \lambda_\beta^0(\hat{x})}}{Z_{\hat{\mathbf{x}} \mid \mathbf{x}}(x; \beta)} e^{-\beta \sum_{r=1}^{d} \lambda_\beta^r(\hat{x}) [A_r(x) - A_{r,\beta}(\hat{x})]}$$

$$p_\beta(y \mid \hat{x}) = e^{-\sum_{r=1}^{d} \lambda^r(y) A_{r,\beta}(\hat{x}) - \lambda_\beta^0(\hat{x})} \ , \quad \lambda_\beta^0(\hat{x}) = \log\left(\sum_y e^{-\sum_{r=1}^{d} \lambda^r(y) A_{r,\beta}(\hat{x})}\right), \quad (1.79)$$

---

[3] The normalization factors, $Z_{\mathbf{y} \mid \mathbf{x}}(x)$, are written, for brevity, as $\lambda_{\mathbf{x}}^0 \equiv \log(\sum_y \prod_{r=1}^{d} e^{-\lambda^r(y) A_r(x)})$ with $A_0(x) \equiv 1$. We do not constrain the marginal $p(x)$.

*with the constraints and multipliers expectations,*

$$A_{r,\beta}(\hat{x}) \equiv \sum_x p_\beta(x \mid \hat{x}) A_r(x) \; , \; \lambda_\beta^r(\hat{x}) \equiv \sum_y p_\beta(y \mid \hat{x}) \lambda^r(y) \; , \; 1 \le r \le d \; . \qquad (1.80)$$

This defines a simplified iterative algorithm for solving the dualExpIB problem. Given the mapping of $x \in X$ to $\{A_r(x)\}_{r=1}^d$ the problem is completely independent of $x$ and we can work in the lower dimensional embedding of the features, $A_r(x)$. Namely, the update procedure is reduced to the dimensions of the sufficient statistics. Moreover, the representation is given in terms of the original features, a desirable feature for any model based problem.

### Optimizing the Error Exponent

The dualIB optimizes an upper bound on the error exponent of the representation multi class testing problem. The error exponent accounts for the decay of the prediction error as a function of data size $n$. This implies the dualIB tends to minimize the prediction error. For the classical binary hypothesis testing, the classification Bayes error, $P_e^{(n)}$, is the weighted sum of type 1 and type 2 errors. For large $n$, both errors decay exponentially with the test size $n$, and the best error exponent, $D^*$, is given by the Chernoff information. The Chernoff information is also a measure of distance defined as, $C(p_0, p_1) = \min_{0 < \lambda < 1} \{\log \sum_x p_0^\lambda(x) p_1^{1-\lambda}(x)\}$, and we can understand it as an optimization on the log-partition function of $p_\lambda$ to obtain $\lambda$ (for further information see Cover and Thomas (2006) and §1.7.13).

The intuition behind the optimization of $D^*$ by the dualIB is in its distortion, the order of the prediction and the observation which implies the use of geometrical mean. The best achievable exponent (see Cover and Thomas (2006)) in Bayesian probability of error is given by the KL-distortion between $p_{\lambda^*}$ ($\propto p_0^{\lambda^*}(x) p_1^{1-\lambda^*}(x)$) to $p_0$ or $p_1$, such that $p_{\lambda^*}$ is the mid point between $p_0$ and $p_1$ on the geodesic of their geometric means. By mapping the dualIB decoder to $\lambda$, it follows that the above minimization is proportional to the log-partition function of $p_\beta(x \mid \hat{x})$, namely we obtain the mapping $p_\beta(x \mid \hat{x}) = p_\lambda$.

In the generalization to multi class classification the error exponent is given by the pair of hypotheses with the minimal Chernoff information Westover (2008). However, finding this value is generally difficult as it requires solving for each pair in the given classes. Thus, we consider an upper bound to it, the mean of the Chernoff information terms over classes. The representation variable adds a new dimension on which we average on and we obtain a bound on the optimal (in expectation over $\hat{x}$) achievable exponent, $\hat{D}_\beta = \min_{p_\beta(y|\hat{x}), p_\beta(\hat{x}|x)} \mathbb{E}_{p_\beta(y, \hat{x})}[D[p_\beta(x \mid \hat{x}) \mid p(x \mid y)]]$. This expression is bounded from above by the dualIB minimization problem. Thus, the dualIB (on expectation) minimizes the prediction error for every $n$. A formal derivation of the above along with an analytical

example of a multi class classification problem is given in §1.7.13. In §1.7.7.3 we experimentally demonstrate that this also holds for a variational dualIB framework using a DNN.

### Variational Dual Information Bottleneck

The Variational Information Bottleneck (VIB) approach introduced by Achille and Soatto (Achille and Soatto; 2018b) and Alemi et al. (Alemi et al.; 2016) allows to parameterize the IB model using Deep Neural Networks (DNNs). The variational bound of the IB is obtained using DNNs for both the encoder and decoder. Since then, various extensions have been made (Strouse and Schwab; 2017; Elad et al.; 2019) demonstrating promising attributes. Recently, along this line, the Conditional Entropy Bottleneck (CEB) (Fischer; 2018) was proposed. The CEB provides variational optimizing bounds on $I(Y; \hat{X})$, $I(X; \hat{X})$ using a variational decoder $q(y \mid \hat{x})$, variational conditional marginal, $q(\hat{x} \mid y)$, and a variational encoder, $p(\hat{x} \mid x)$, all implemented by DNNs.

Here, we present the variational dualIB (VdualIB), which optimizes the variational dualIB objective for using in DNNs. Following the CEB formalism, we bound the dualIB objective. We develop a variational form of the dualIB distortion and combine it with the bound for $I(X; \hat{X})$ (as in the CEB). This gives us the following Theorem (for the proof see §1.7.14.).

**Theorem 1.7.7.** *The* VdualIB *objective is given by:*

$$\min_{q(\hat{x}|y), p(\hat{x}|x)} \left\{ \mathbb{E}_{\tilde{p}(y|x)p(\hat{x}|x)p(x)} \left[ \log \frac{p(\hat{x} \mid x)}{q(\hat{x} \mid y)} \right] + \beta \mathbb{E}_{p(y|\hat{x})p(\hat{x}|x)} \left[ \log \frac{p(y \mid \hat{x})}{\tilde{p}(y \mid x)} \right] \right\}, \qquad (1.81)$$

*where $\tilde{p}(y \mid x)$ is a distribution based on the given labels of the data-set, which we relate to as the noise model. Under the assumption that the noise model captures the distribution of the data the above provides a variational upper bound to the* dualIB *functional equation* **??**.

Due to the nature of its objective the dualIB requires a noise model. We must account for the contribution to the objective arising from $\tilde{p}(y \mid x)$ which could be ignored in the VIB case. The noise model can be specified by its assumptions over the data-set. In §1.7.7.2 we elaborate on possible noise models choices and their implications on the performance. Notice that the introduction of $\tilde{p}(y \mid x)$ implies that, unlike most machine learning models, the VdualIB does not optimize directly the error between the predicted and desired labels in the training data. Instead, it does so implicitly with respect to the noisy training examples. This is not unique to the VdualIB, as it is equivalent to training with noisy labels, often done to prevent over-fitting. For example, in (Müller et al.; 2019) the authors show that label noise can improve generalization that results in a reduction in the mutual information between the input and the output.

In practice, similarly to the CEB, for the stochastic encoder, $p(\hat{x} \mid x)$, we use the original architecture, replacing the final softmax layer with a linear dense layer with $d$ outputs. These outputs are taken as the means of a multivariate Gaussian distribution with unit diagonal covariance. For the variational decoder, $q(y \mid \hat{x})$, any classifier network can be used. We take a linear softmax classifier which takes the encoder as its input. The reverse decoder $q(\hat{x} \mid y)$ is implemented by a network which maps a one-hot representation of the labels to the $d$-dimensional output interpreted as the mean of the corresponding Gaussian marginal.
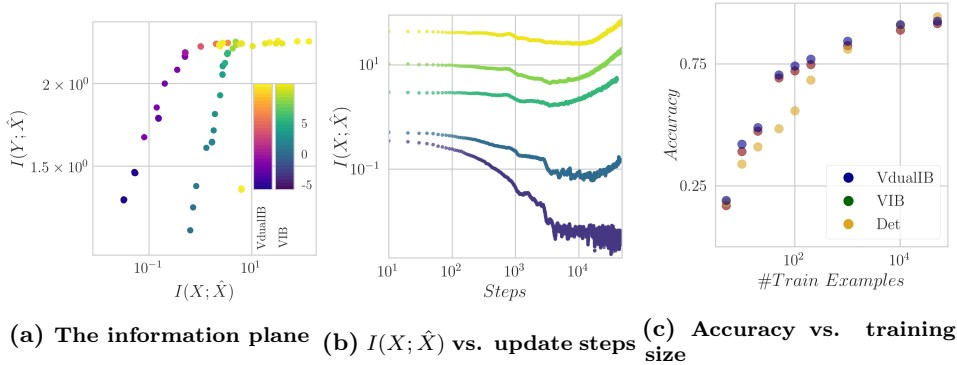
### 1.7.7   Experiments

To investigate the VdualIB on real-world data we compare it to the CEB model using a DNN over two data-sets, FasionMNIST and CIFAR10. For both, we use a Wide ResNet $28 - 10$ (Zagoruyko and Komodakis; 2016) as the encoder, a one layer Gaussian decoder and a single layer linear network for the reverse decoder (similarly to the setup in (Fischer and Alemi; 2020)). We use the same architecture to train networks with VdualIB and VIB objectives. (See §1.7.14 for details on the experimental setup). We note that in our attempts to train over the CIFAR100 data-set the results did not fully agree with the results on the above data-sets (for more information see §1.7.14.1). An open source implementation is available here.

#### 1.7.7.1   The variational information plane

As mentioned, the information plane describes the compression-prediction trade-off. It enables us to compare different models and evaluate their "best prediction level" in terms of the desired label information, for each compression level. In (Fischer and Alemi; 2020) the authors provide empirical evidence that information bottlenecking techniques can improve both generalization and robustness. Other works (Fischer; 2018; Achille and Soatto; 2018a,b) provide both theoretical and conceptual insights into why these improvements occur.

In Figure 1.16 we present the information plane of the VdualIB where the distribution $\tilde{p}(y \mid x)$ (the noise model) is a learnt confusion matrix, ConfVdualIB (similarly to (Wu and Fischer; 2020)). We compare it to the VIB over a range of $\beta$ values ($-5 \leq \log \beta \leq 5$). Figure 1.16a validates that, as expected, the information growth is approximately monotonic with $\beta$. Comparing the VdualIB to the VIB model, we can see significant differences between their representations. The VdualIB successfully obtains better compressed representations in comparison to the VIB performance, where only for large values of $I(X; \hat{X})$ their performances match. As predicted by the theory, in the limit $\beta \to \infty$ the models behaviour match. Furthermore, the VdualIB values are smoother and they are spread over the information plane, making it easier to optimize for a specific value in it. In Figure 1.16b
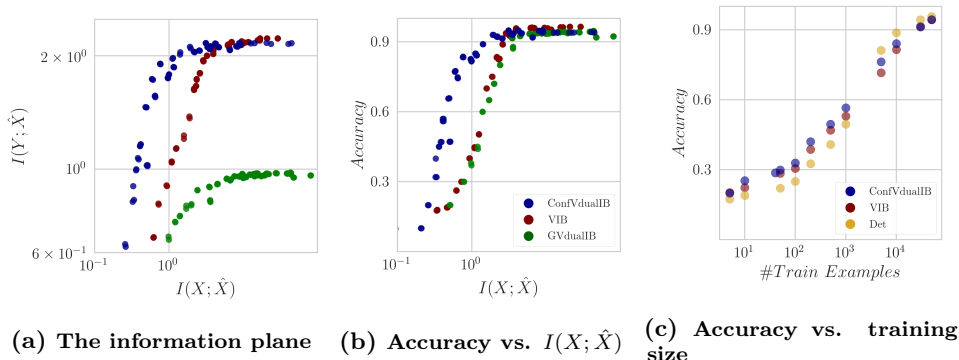
**(a) The information plane**    **(b)** $I(X; \hat{X})$ **vs. update steps**    **(c) Accuracy vs. training size**

**Figure 1.16: Experiments over FashionMNIST.** (*a*) **The information plane of the** ConfVdualIB **and VIB for a range of** $\beta$ **values at the final training step.** (*b*) **The evolution of the the** ConfVdualIB**'s** $I(X; \hat{X})$ **along the optimization update steps.** (*c*) **The models accuracy as a function of the training set size.**

we consider the dynamics of $I(X; \hat{X})$ for several values of $\beta$. Interestingly, at the initial training stage the representation information for all values of $\beta$ decreases. However, as the training continues, the information increases only for high $\beta$s.

### 1.7.7.2 *The* VdualIB *noise model*

As mentioned above, learning with the VdualIB objective requires a choice of a noise model for the distribution $\tilde{p}(y \mid \hat{x})$. To explore the influence of different models on the learning we evaluate four types, with different assumptions on the access to the data. (i) Adding Gaussian noise to the one-hot vector of the true label (GVdualIB); (ii) An analytic Gaussian integration of the log-loss around the one-hot labels; (iii) A pre-computed confusion matrix for the labels (ConfVdualIB) as in (Wu and Fischer; 2020); (iv) Using predictions of another trained model as the induced distribution. Where for (i) and (ii) the variance acts as a free parameter determining the noise level. The complexity of the noise models can be characterized by the additional prior knowledge on our data-set they require. While adding Gaussian noise does not require prior knowledge, using a trained model requires access to the prediction for every data sample. The use of a confusion matrix is an intermediate level of prior knowledge requiring access only to the $|\mathcal{Y}| \times |\mathcal{Y}|$ pre-computed matrix. Here we present cases (i) and (iii) (see §1.7.14 for (ii) and (iv)). Note that although using the VIB does not require the use of a noise model we incorporate it by replacing the labels with $\tilde{p}(y \mid x)$. In the analysis below, the results are presented with the VIB trained with the same noise model as the VdualIB (see §1.7.14 for a comparison between training VdualIB with noise and VIB without it).

Figure 1.17 depicts the information plane of the CIFAR10 data-set. Figure 1.17a shows the information obtained from a range of $\beta$. The colors depict the different models

(a) **The information plane**    (b) **Accuracy vs.** $I(X;\hat{X})$    (c) **Accuracy vs. training size**

**Figure 1.17: Experiments over CIFAR10.** ($a$) **The information plane of the** VIB**,** ConfVdualIB**,** GVdualIB **and** VIB **for a range of** $\beta$ **values.** ($b$) **The accuracy of the models as a function of the mutual information,** $I(X;\hat{X})$**.** ($c$) **The accuracy of the models as a function of the training set size.**

ConfVdualIB, GVdualIB and the VIB. As we can see, training a VdualIB with Gaussian noise achieves much less information with the labels at any given $I(X;\hat{X})$. We note that we verified that this behaviour is consistent over a wide range of variances. The ConfVdualIB model performance is similar to the VIB's with the former showcasing more compressed representations. When we present the prediction accuracy (Figure 1.17b), here all models attain roughly the same accuracy values. The discrepancy between the accuracy and information, $I(Y;\hat{X})$, is similar to the one discussed in (Dusenberry et al.; 2020).

#### 1.7.7.3  *Performance with different training set sizes*

Our theoretical analysis (§1.7.6) shows that under given assumptions the dualIB bounds the optimal achievable error exponent on expectation hence it optimizes the error for a given data size $n$. We turn to test this in the VdualIB setting. We train the models on a subset of the training set and evaluate them on the test set. We compare the VIB and the VdualIB to a deterministic network (Det; Wide Res Net 28-10). Both the VIB and VdualIB are trained over a wide range of $\beta$ values ($-5 \leq \log \beta \leq 6$). Presented is the best accuracy value for each model at a given $n$. Figure 1.16c and Figure 1.17c show the accuracy of the models as a function of the training set size over FashionMNIST and CIFAR10 respectively. The VdualIB performance is slightly better in comparison to the VIB, while the accuracy of the deterministic network is lower for small training sets. The superiority of the variational models over the deterministic network is not surprising as minimizing $I(X;\hat{X})$ acts as regularization.

104

## Conclusions

We present here the Dual Information Bottleneck (dualIB), a framework resolving some of the known drawbacks of the IB obtained by a mere switch between the terms in the distortion function. We provide the dualIB self-consistent equations allowing us to obtain analytical solutions. A local stability analysis revealed the underlying structure of the critical points of the solutions, resulting in a full bifurcation diagram of the optimal pattern representations. The study of the dualIB objective reveals several interesting properties. First, when the data can be modeled in a parametric form the dualIB preserves this structure and it obtains the representation in terms of the original parameters, as given by the dualExpIB equations. Second, it optimizes the mean prediction error exponent thus improving the accuracy of the predictions as a function of the data size. In addition to the dualIB analytic solutions, we provide a variational dualIB (VdualIB) framework, which optimizes the functional using DNNs. This framework enables practical implementation of the dualIB to real world data-sets. While a broader analysis is required, the VdualIB experiments shown validate the theoretical predictions. Our results demonstrate the potential advantages and unique properties of the framework.

## Bibliography

Achille, A., Paolini, G. and Soatto, S. (2019). Where is the information in a deep neural network?, *arXiv preprint arXiv:1905.12213* .

Achille, A. and Soatto, S. (2018a). Emergence of invariance and disentanglement in deep representations, *The Journal of Machine Learning Research* **19**(1): 1947–1980.

Achille, A. and Soatto, S. (2018b). Information dropout: Learning optimal representations through noisy computation, *IEEE transactions on pattern analysis and machine intelligence* **40**(12): 2897–2905.

Alemi, A. A., Fischer, I., Dillon, J. V. and Murphy, K. (2016). Deep variational information bottleneck, *ArXiv* **abs/1612.00410**.

Brown, L. D. (1986). Fundamentals of statistical exponential families with applications in statistical decision theory, *Lecture Notes-Monograph Series* **9**: i–279.
**URL:** *http://www.jstor.org/stable/4355554*

Cover, T. M. and Thomas, J. A. (2006). *Elements of Information Theory (Wiley Series in Telecommunications and Signal Processing)*, Wiley-Interscience, New York, NY, USA.

Dusenberry, M. W., Jerfel, G., Wen, Y., Ma, Y.-a., Snoek, J., Heller, K., Lakshminarayanan, B. and Tran, D. (2020). Efficient and scalable bayesian neural nets with rank-1 factors, *arXiv preprint arXiv:2005.07186* .

Elad, A., Haviv, D., Blau, Y. and Michaeli, T. (2019). Direct validation of the information bottleneck principle for deep nets, *Proceedings of the IEEE International Conference on Computer Vision Workshops.*

Felice, D. and Ay, N. (2019). Divergence functions in information geometry, *in* F. Nielsen and F. Barbaresco (eds), *Geometric Science of Information - 4th International Conference, GSI 2019, Toulouse, France, August 27-29, 2019, Proceedings*, Vol. 11712 of *Lecture Notes in Computer Science*, Springer, pp. 433–442.
**URL:** *https://doi.org/10.1007/978-3-030-26980-7_45*

Fischer, I. (2018). The conditional entropy bottleneck, *URL openreview. net/forum.*

Fischer, I. and Alemi, A. A. (2020). Ceb improves model robustness, *arXiv preprint arXiv:2002.05380 .*

Gilad-bachrach, R., Navot, A. and Tishby, N. (2003). An information theoretic tradeoff between complexity and accuracy, *In Proceedings of the COLT*, Springer, pp. 595–609.

Goldfeld, Z., Berg, E. v. d., Greenewald, K., Melnyk, I., Nguyen, N., Kingsbury, B. and Polyanskiy, Y. (2018). Estimating information flow in deep neural networks, *arXiv preprint arXiv:1810.05728 .*

Jaynes, E. T. (1957). Information theory and statistical mechanics, *Phys. Rev.* **106**: 620–630.
**URL:** *https://link.aps.org/doi/10.1103/PhysRev.106.620*

Kingma, D. P. and Ba, J. (2014). Adam: A method for stochastic optimization, *arXiv preprint arXiv:1412.6980 .*

Li, X. L. and Eisner, J. (2019). Specializing word embeddings (for parsing) by information bottleneck, *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 2744–2754.

Ma, S., McDuff, D. and Song, Y. (2019). Unpaired image-to-speech synthesis with multimodal information bottleneck, *Proceedings of the IEEE International Conference on Computer Vision*, pp. 7598–7607.

Müller, R., Kornblith, S. and Hinton, G. E. (2019). When does label smoothing help?, *Advances in Neural Information Processing Systems*, pp. 4696–4705.

Nash, C., Kushman, N. and Williams, C. K. (2018). Inverting supervised representations with autoregressive neural density models, *arXiv preprint arXiv:1806.00400 .*

Painsky, A. and Wornell, G. W. (2018). Bregman Divergence Bounds and the Universality of the Logarithmic Loss, *arXiv e-prints* p. arXiv:1810.07014.

Parbhoo, S., Wieser, M. and Roth, V. (2018). Causal deep information bottleneck, *ArXiv* **abs/1807.02326**.

Parker, A. E., Gedeon, T. and Dimitrov, A. G. (2003). Annealing and the rate distortion problem, *in* S. Becker, S. Thrun and K. Obermayer (eds), *Advances in Neural Information Processing Systems 15*, MIT Press, pp. 993–976.
**URL:** *http://papers.nips.cc/paper/2264-annealing-and-the-rate-distortion-problem.pdf*

Poole, B., Ozair, S., Oord, A. v. d., Alemi, A. A. and Tucker, G. (2019). On variational bounds of mutual information, *arXiv preprint arXiv:1905.06922* .

Schneidman, E., Slonim, N., Tishby, N., van Steveninck, R. d. and Bialek, W. (2001). Analyzing neural codes using the information bottleneck method, *Advances in Neural Information Processing Systems, NIPS* .

Shalev-Shwartz, S. and Ben-David, S. (2014). *Understanding machine learning: From theory to algorithms*, Cambridge university press.

Shamir, O., Sabato, S. and Tishby, N. (2010). Learning and generalization with the information bottleneck, *Theor. Comput. Sci.* **411**: 2696–2711.

Shwartz-Ziv, R. and Tishby, N. (2017). Opening the Black Box of Deep Neural Networks via Information, *arXiv e-prints* p. arXiv:1703.00810.

Slonim, N., Friedman, N. and Tishby, N. (2006). Multivariate information bottleneck, *Neural Computation* **18**(8): 1739–1789. PMID: 16771652.
**URL:** *https://doi.org/10.1162/neco.2006.18.8.1739*

Strouse, D. and Schwab, D. J. (2017). The deterministic information bottleneck, *Neural computation* **29**(6): 1611–1630.

Tishby, N., Pereira, F. C. and Bialek, W. (1999). The information bottleneck method, *Proc. of the 37-th Annual Allerton Conference on Communication, Control and Computing*, pp. 368–377.

Tredicce, J. R., Lippi, G. L., Mandel, P., Charasse, B., Chevalier, A. and Picqué, B. (2004). Critical slowing down at a bifurcation, *American Journal of Physics* **72**(6): 799–809.

Tusnady, G. and Csiszar, I. (1984). Information geometry and alternating minimization procedures, *Statistics & Decisions: Supplement Issues* **1**: 205–237.

Westover, M. B. (2008). Asymptotic geometry of multiple hypothesis testing, *IEEE transactions on information theory* **54**(7): 3327–3329.

Wu, T. and Fischer, I. (2020). Phase transitions for the information bottleneck in representation learning, *arXiv preprint arXiv:2001.01878* .

Zagoruyko, S. and Komodakis, N. (2016). Wide residual networks, *arXiv preprint arXiv:1605.07146* .

Zaslavsky, N. and Tishby, N. (2019). Deterministic annealing and the evolution of optimal information bottleneck representations, *Preprint* .

*Appendix*

## Appendix A - The Information Bottleneck method

The Information Bottleneck (IB) trade off between the encoder and decoder mutual information values is defined by the minimization of the Lagrangian:

$$\mathcal{F}[p_\beta(\hat{x} \mid x); p_\beta(y \mid \hat{x})] = I(X; \hat{X}) - \beta I(Y; \hat{X}) \ , \tag{1.82}$$

independently over the convex sets of the normalized distributions, $\{p_\beta(\hat{x} \mid x)\}$, $\{p_\beta(\hat{x})\}$ and $\{p_\beta(y \mid \hat{x})\}$, given a positive Lagrange multiplier $\beta$. As shown in Tishby et al. (1999); Shamir et al. (2010), this is a natural generalization of the classical concept of *Minimal Sufficient Statistics* Cover and Thomas (2006), where the estimated parameter is replaced by the output variable $Y$ and *exact* statistical sufficiency is characterized by the mutual information equality: $I(\hat{X}; Y) = I(X; Y)$. The minimality of the statistics is captured by the minimization of $I(X; \hat{X})$, due to the Data Processing Inequality (DPI). However, non-trivial minimal sufficient statistics only exist for very special parametric distributions known as exponential families Brown (1986). Thus in general, the IB relaxes the minimal sufficiency problem to a continuous family of representations $\hat{X}$ which are characterized by the trade off between compression, $I(X; \hat{X}) \equiv I_X$, and accuracy, $I(Y; \hat{X}) \equiv I_Y$, along a convex line in the *Information-Plane* ($I_Y$ vs. $I_X$). When the rule $p(x, y)$ is strictly stochastic, the convex optimal line is smooth and each point along the line is uniquely characterized by the value of $\beta$. We can then consider the optimal representations $\hat{x} = \hat{x}(\beta)$ as encoder-decoder pairs: $(p_\beta(x \mid \hat{x}), p_\beta(y \mid \hat{x}))$[4] - a point in the continuous manifold defined by the Cartesian product of these distribution simplexes. We also consider a small variation of these representations, $\delta\hat{x}$, as an infinitesimal change in this (encoder-decoder) continuous manifold (not necessarily on the optimal line(s)).

### 1.7.8  IB and Rate-Distortion Theory

The IB optimization trade off can be considered as a generalized rate-distortion problem Cover and Thomas (2006) with the distortion function between a data point, $x$ and a representation point $\hat{x}$ taken as the KL-divergence between their predictions of the desired label $y$:

$$
\begin{aligned}
d_{\mathrm{IB}}(x, \hat{x}) &= D[p(y \mid x) || p_\beta(y|\hat{x})] \\
&= \sum_y p(y \mid x) \log \frac{p(y \mid x)}{p_\beta(y \mid \hat{x})}.
\end{aligned}
\tag{1.83}
$$

---

[4] Here we use the *inverse encoder*, which is in the fixed dimension simplex of distributions over $X$.

The expected distortion $\mathbb{E}_{p_\beta(x,\hat{x})}[d_{\mathrm{IB}}(x,\hat{x})]$ for the optimal decoder is simply the label-information loss: $I(X;Y) - I(\hat{X};Y)$, using the Markov chain condition. Thus minimizing the expected IB distortion is equivalent to maximizing $I(\hat{X};Y)$, or minimizing equation **??**. Minimizing this distortion is equivalent to minimizing the cross-entropy loss, and it provides an upper-bound to other loss functions such as the $\mathcal{L}_1$-loss (due to the Pinsker inequality, see also Painsky and Wornell (2018)). Pinsker implies that both orders of the cross-entropy act as an upper bound to the $\mathcal{L}_1$-loss, $\min\{D[q||p], D[p||q]\} \geq \frac{1}{2\log 2}\|p - q\|_1^2$ .

### 1.7.9   The IB Equations

For discrete $X$ and $Y$, a necessary condition for the IB (local) minimization is given by the three self-consistent equations for the optimal encoder-decoder pairs, known as the IB *equations*:

$$\begin{cases} (i) & p_\beta(\hat{x} \mid x) = \frac{p_\beta(\hat{x})}{Z(x;\beta)}e^{-\beta D[p(y|x)\|p_\beta(y|\hat{x})]} \\ (ii) & p_\beta(\hat{x}) = \sum_x p_\beta(\hat{x} \mid x)p(x) \\ (iii) & p_\beta(y \mid \hat{x}) = \sum_x p(y \mid x)p_\beta(x \mid \hat{x}) \end{cases} , \qquad (1.84)$$

where $Z(x;\beta)$ is the normalization function. Iterating these equations is a generalized, Blahut-Arimoto, alternating projection algorithm Tusnady and Csiszar (1984); Cover and Thomas (2006) and it converges to a stationary point of the Lagrangian, equation **??** Tishby et al. (1999). Notice that the minimizing decoder, (equation **??**-$(iii)$), is precisely the *Bayes optimal decoder* for the representation $\hat{x}(\beta)$, given the Markov chain conditions.

### 1.7.10   Critical points and critical slowing down

One of the most interesting aspects of the IB equations is the existence of critical points along the optimal line of solutions in the information plane (i.e. the information curve). At these points the representations change topology and cardinality (number of clusters) Zaslavsky and Tishby (2019); Parker et al. (2003) and they form the skeleton of the information curve and representation space. To identify such points we perform a perturbation analysis of the IB equations:[5]:

$$\delta \log p_\beta(x \mid \hat{x}) = \beta \sum_y p(y \mid x)\delta \log p_\beta(y \mid \hat{x}), \qquad (1.85)$$

$$\delta \log p_\beta(y \mid \hat{x}) = \frac{1}{p_\beta(y \mid \hat{x})} \sum_x p(y \mid x)p_\beta(x \mid \hat{x})\delta \log p_\beta(x \mid \hat{x}). \qquad (1.86)$$

---

[5] We ignore here the possible interaction between the different representations, for simplicity.

Substituting equation **??** into equation **??** and vice versa one obtains:

$$\delta \log p_\beta(x \mid \hat{x}) = \beta \sum_{y,x'} p(y \mid x') \frac{p(y \mid x)}{p_\beta(y \mid \hat{x})} p_\beta(x' \mid \hat{x}) \delta \log p_\beta(x' \mid \hat{x})$$

$$\delta \log p_\beta(y \mid \hat{x}) = \beta \sum_{x,y'} p(y \mid x) \frac{p_\beta(x \mid \hat{x})}{p_\beta(y \mid \hat{x})} p(y' \mid x) \delta \log p_\beta(y' \mid \hat{x})$$

Thus by defining the matrices:

$$C_{xx'}^{\text{IB}}(\hat{x}, \beta) = \sum_y p(y \mid x) \frac{p_\beta(x' \mid \hat{x})}{p_\beta(y \mid \hat{x})} p(y \mid x') \ , \ C_{yy'}^{\text{IB}}(\hat{x}, \beta) = \sum_x p(y \mid x) \frac{p_\beta(x \mid \hat{x})}{p_\beta(y \mid \hat{x})} p(y' \mid x).$$
$$(1.87)$$

We obtain the following nonlinear eigenvalues problem:

$$\left[ I - \beta C_{xx'}^{\text{IB}}(\hat{x}, \beta) \right] \delta \log p_\beta(x' \mid \hat{x}) = 0 \ , \quad \left[ I - \beta C_{yy'}^{\text{IB}}(\hat{x}, \beta) \right] \delta \log p_\beta(y' \mid \hat{x}) = 0, \quad (1.88)$$

These two matrices have the same eigenvalues and have non-trivial eigenvectors (i.e., different co-existing optimal representations) at the critical values of $\beta$, the bifurcation points of the IB solution. At these points the cardinality of the representation $\hat{X}$ (the number of "IB-clusters") changes due to splits of clusters, resulting in topological phase transitions in the encoder. These critical points form the "skeleton" of the topology of the optimal representations. Between critical points the optimal representations change continuously (with $\beta$). The important computational consequence of critical points is known as *critical slowing down* Tredicce et al. (2004). For binary $y$, near a critical point the convergence time, $\tau_\beta$, of the iterations of equation **??** scales like: $\tau_\beta \sim 1/(1 - \beta\lambda_2)$, where $\lambda_2$ is the second eigenvalue of either $C_{yy'}^{\text{IB}}$ or $C_{xx'}^{\text{IB}}$. At criticality, $\lambda_2(\hat{x}) = \beta^{-1}$ and the number of iterations diverges. This phenomenon dominates any local minimization of equation **??** which is based on alternate encoder-decoder optimization.

The appearance of the critical points and the critical slowing-down is visualized in Figure 1 in the main text.

### Appendix B - The dualIB mathematical formulation

The dualIB is solved with respect to the full Markov chain $(Y \to X \to \hat{X}_\beta \to \hat{Y})$ in which we introduce the new variable, $\hat{y}$, the *predicted label*. Thus, in analogy to the IB we want to write the optimization problem in term of $\hat{Y}$.

Developing the expected distortion we find:

$$\mathbb{E}_{p_\beta(x,\hat{x})}[d_{\mathrm{dualIB}}(x,\hat{x})] = \sum_{x,\hat{x}} p_\beta(x,\hat{x}) \sum_{\hat{y}} p_\beta(y=\hat{y}\mid\hat{x}) \log \frac{p_\beta(y=\hat{y}\mid\hat{x})}{p(y=\hat{y}\mid x)}$$

$$= \sum_{\hat{x},\hat{y}} p_\beta(\hat{x})p_\beta(\hat{y}\mid\hat{x}) \log \frac{p_\beta(\hat{y}\mid\hat{x})}{p_\beta(\hat{y})} - \sum_{x,\hat{y}} p(x)p_\beta(\hat{y}\mid x) \log \frac{p_\beta(\hat{y}\mid x)}{p_\beta(\hat{y})}$$

$$+ \sum_{x,\hat{y}} p(x)p_\beta(\hat{y}\mid x) \log \frac{p_\beta(\hat{y}\mid x)}{p(y=\hat{y}\mid x)}$$

$$= I(\hat{X};\hat{Y}) - I(X;\hat{Y}) + \mathbb{E}_{p(x)}[D[p_\beta(\hat{y}\mid x)\|p(y=\hat{y}\mid x)]].$$

Allowing the dual optimization problem to be written as:

$$\mathcal{F}^*[p(\hat{x}\mid x); p(y\mid\hat{x})] = I(X;\hat{X}) - \beta\Big\{I(X;\hat{Y}) - I(\hat{X};\hat{Y}) - \mathbb{E}_{p(x)}[D[p_\beta(\hat{y}\mid x)\|p(y=\hat{y}\mid x)]]\Big\}.$$

### *Appendix C - The DualIB solutions*

To prove *theorem* 2 we want to obtain the normalized distributions minimizing the dualIB rate-distortion problem.

*Proof.* (*i*) Given that the problem is formulated as a rate-distortion problem the encoder's update rule must be the known minimizer of the distortion function. Cover and Thomas (2006). Thus the IB encoder with the dual distortion is plugged in. (*ii*) For the decoder, by considering a small perturbation in the distortion $d_{\mathrm{dualIB}}(x,\hat{x})$, with $\alpha(\hat{x})$ the normalization Lagrange multiplier, we obtain:

$$\delta d_{\mathrm{dualIB}}(x,\hat{x}) = \delta\left(\sum_y p_\beta(y\mid\hat{x}) \log \frac{p_\beta(y\mid\hat{x})}{p(y\mid x)} + \alpha(\hat{x})\left(\sum_y p_\beta(y\mid\hat{x}) - 1\right)\right)$$

$$\frac{\delta d_{\mathrm{dualIB}}(x,\hat{x})}{\delta p_\beta(y\mid\hat{x})} = \log \frac{p_\beta(y\mid\hat{x})}{p(y\mid x)} + 1 + \alpha(\hat{x}).$$

Hence, minimizing the expected distortion becomes:

$$0 = \sum_x p_\beta(x\mid\hat{x})\left[\log \frac{p_\beta(y\mid\hat{x})}{p(y\mid x)} + 1\right] + \alpha(\hat{x})$$

$$= \log p_\beta(y\mid\hat{x}) - \sum_x p_\beta(x\mid\hat{x}) \log p(y\mid x) + 1 + \alpha(\hat{x}),$$

which yields Algorithm 1, row 6. $\qquad\square$

Considering the dualIB encoder-decoder, Algorithm 1, we find that $\mathbb{E}_{p_\beta(x,\hat{x})}[d_{\mathrm{dualIB}}(x,\hat{x})]$

reduces to the expectation of the decoder's log partition function:

$$\mathbb{E}_{p_\beta(x,\hat{x})}[d_{\text{dualIB}}(x,\hat{x})] = \sum_{x,\hat{x}} p_\beta(x,\hat{x}) \sum_y p_\beta(y \mid \hat{x}) \log \frac{p_\beta(y \mid \hat{x})}{p(y \mid x)}$$

$$= -\mathbb{E}_{p_\beta(\hat{x})}\big[\log Z_{\mathbf{y}|\hat{\mathbf{x}}}(\hat{x};\beta)\big] + \sum_{\hat{x}} p_\beta(\hat{x})\Bigg[\sum_{\hat{x},y} p_\beta(x' \mid \hat{x}) \log p(y \mid x') - \sum_x p_\beta(x \mid \hat{x}) \log p(y \mid x)\Bigg]$$

$$= -\mathbb{E}_{p_\beta(\hat{x})}\big[\log Z_{\mathbf{y}|\hat{\mathbf{x}}}(\hat{x};\beta)\big].$$

### Appendix D - Stability analysis

Here we provide the detailed stability analysis allowing the definition of the matrices $C_{xx'}^{\text{dualIB}}, C_{yy'}^{\text{dualIB}}$ (*theorem* 4) which allows us to claim that they obey the same rules as the $C$ matrices of the IB. Similarly to the IB in this calculation we ignore second order contributions which arise form the normalization terms. Considering a variation in $\hat{x}$ we get:

$$\delta \log p_\beta(x \mid \hat{x}) = \beta \sum_y p_\beta(y \mid \hat{x})\left(\log \frac{p(y \mid x)}{p_\beta(y \mid \hat{x})} - 1\right)\delta \log p_\beta(y \mid \hat{x})$$

$$= \beta \sum_y p_\beta(y \mid \hat{x})\left[\log p(y \mid x) - \sum_{\tilde{x}} p_\beta(\tilde{x} \mid \hat{x}) \log p(y \mid \tilde{x})\right]\delta \log p_\beta(y \mid \hat{x})$$

$$+ \beta \sum_y \log Z_{\mathbf{y}|\hat{\mathbf{x}}}(\hat{x};\beta)\frac{\partial p_\beta(y \mid \hat{x})}{\partial \hat{x}}$$

$$= \beta \sum_{y,\tilde{x}} p_\beta(y \mid \hat{x})p_\beta(\tilde{x} \mid \hat{x}) \log \frac{p(y \mid x)}{p(y \mid \tilde{x})}\delta \log p_\beta(y \mid \hat{x}), \tag{1.89}$$

$$\delta \log p_\beta(y \mid \hat{x}) = -\frac{1}{Z_{\mathbf{y}|\hat{\mathbf{x}}}(\hat{x};\beta)}\frac{\partial Z_{\mathbf{y}|\hat{\mathbf{x}}}(\hat{x};\beta)}{\partial \hat{x}} + \sum_x p_\beta(x \mid \hat{x}) \log p(y \mid x)\delta \log p_\beta(x \mid \hat{x})$$

$$= -\sum_{\tilde{y}} p_\beta(\tilde{y} \mid \hat{x}) \sum_x p_\beta(x \mid \hat{x}) \log p(\tilde{y} \mid x)\delta \log p_\beta(x \mid \hat{x})$$

$$+ \sum_x p_\beta(x \mid \hat{x}) \log p(y \mid x)\delta \log p_\beta(x \mid \hat{x})$$

$$= \sum_{x,\tilde{y}} p_\beta(x \mid \hat{x})p_\beta(\tilde{y} \mid \hat{x}) \log \frac{p(y \mid x)}{p(\tilde{y} \mid x)}\delta \log p_\beta(x \mid \hat{x}). \tag{1.90}$$

Substituting equation **??** into equation **??** and vice versa one obtains:

$$\delta \log p_\beta(x \mid \hat{x}) = \beta \sum_{x',y,\tilde{y},\tilde{x}} p_\beta(y \mid \hat{x})p_\beta(\tilde{x} \mid \hat{x}) \log \frac{p(y \mid x)}{p(y \mid \tilde{x})}$$

$$\cdot p_\beta(x' \mid \hat{x})p_\beta(\tilde{y} \mid \hat{x}) \log \frac{p(y \mid x')}{p(\tilde{y} \mid x')} \delta \log p_\beta(x' \mid \hat{x})$$

$$\delta \log p_\beta(y \mid \hat{x}) = \beta \sum_{x,y',\tilde{x},\tilde{y}} p_\beta(x \mid \hat{x})p_\beta(\tilde{y} \mid \hat{x}) \log \frac{p(y \mid x)}{p(\tilde{y} \mid x)}$$

$$\cdot p_\beta(y' \mid \hat{x})p_\beta(\tilde{x} \mid \hat{x}) \log \frac{p(y' \mid x)}{p(y' \mid \tilde{x})} \delta \log p_\beta(y' \mid \hat{x}).$$

We now define the $C^{\text{dualIB}}$ matrices as follows:

$$C_{xx'}^{\text{dualIB}}(\hat{x};\beta) = \sum_{y,\tilde{y},\tilde{x}} p_\beta(y \mid \hat{x})p_\beta(\tilde{x} \mid \hat{x}) \log \frac{p(y \mid x)}{p(y \mid \tilde{x})} \cdot p_\beta(x' \mid \hat{x})p_\beta(\tilde{y} \mid \hat{x}) \log \frac{p(y \mid x')}{p(\tilde{y} \mid x')}$$

$$C_{yy'}^{\text{dualIB}}(\hat{x};\beta) = \sum_{x,\tilde{x},\tilde{y}} p_\beta(x \mid \hat{x})p_\beta(\tilde{y} \mid \hat{x}) \log \frac{p(y \mid x)}{p(\tilde{y} \mid x)} \cdot p_\beta(y' \mid \hat{x})p_\beta(\tilde{x} \mid \hat{x}) \log \frac{p(y' \mid x)}{p(y' \mid \tilde{x})}.$$

Using the above definition we have an equivalence to the IB stability analysis in the form of:

$$\left[I - \beta C_{xx'}^{\text{dualIB}}(\hat{x},\beta)\right]\delta \log p_\beta(x' \mid \hat{x}) = 0 \ , \quad \left[I - \beta C_{yy'}^{\text{dualIB}}(\hat{x},\beta)\right]\delta \log p_\beta(y' \mid \hat{x}) = 0.$$

Note that for the binary case, the matrices may be simplified to:

$$C_{xx'}^{\text{dualIB}}(\hat{x};\beta) = \sum_{y,\tilde{x}} p_\beta(y \mid \hat{x})p_\beta(\tilde{x} \mid \hat{x}) \log \frac{p(y \mid x)}{p(y \mid \tilde{x})} \cdot p_\beta(x' \mid \hat{x})(1 - p_\beta(y \mid \hat{x})) \log \frac{p(y \mid x')}{1 - p(y \mid x')}$$

$$C_{yy'}^{\text{dualIB}}(\hat{x};\beta) = \sum_{x,\tilde{x}} p_\beta(x \mid \hat{x})(1 - p_\beta(y \mid \hat{x})) \log \frac{p(y \mid x)}{1 - p(y \mid x)} \cdot p_\beta(y' \mid \hat{x})p_\beta(\tilde{x} \mid \hat{x}) \log \frac{p(y' \mid x)}{p(y' \mid \tilde{x})}.$$

We turn to show that the $C^{\text{dualIB}}$ matrices share the same eigenvalues with $\lambda_1(\hat{x}) = 0$.

*Proof.* The matrices, $C_{xx'}^{\text{dualIB}}(\hat{x};\beta)$, $C_{yy'}^{\text{dualIB}}(\hat{x};\beta)$, are given by:

$$C_{xx'}^{\text{dualIB}}(\hat{x};\beta) = A_{xy}(\hat{x};\beta)B_{yx'}(\hat{x};\beta) \ , \quad C_{yy'}^{\text{dualIB}}(\hat{x};\beta) = B_{yx}(\hat{x};\beta)A_{xy'}(\hat{x};\beta),$$

with:

$$A_{xy}(\hat{x};\beta) = p_\beta(y \mid \hat{x}) \sum_{\tilde{x}} p_\beta(\tilde{x} \mid \hat{x}) \log \frac{p(y \mid x)}{p(y \mid \tilde{x})} \ , \quad B_{yx}(\hat{x};\beta) = p_\beta(x \mid \hat{x}) \sum_{\tilde{y}} p_\beta(\tilde{y} \mid \hat{x}) \log \frac{p(y \mid x)}{p(\tilde{y} \mid x)}.$$

Given that the matrices are obtained by the multiplication of the same matrices, it follows that they have the same eigenvalues $\{\lambda_i(\hat{x};\beta)\}$.

To prove that $\lambda_1(\hat{x}; \beta) = 0$ we show that $\det(C_{yy'}^{\text{dualIB}}) = 0$. We present the exact calculation for a binary label $y \in \{y_0, y_1\}$ (the argument for general $y$ follows by encoding the label as a sequence of bits and discussing the first bit only, as a binary case):

$$
\begin{aligned}
\det(C_{yy'}^{\text{dualIB}}(\hat{x}; \beta)) &= \sum_{x,\tilde{x}} p_\beta(x \mid \hat{x}) p_\beta(y_1 \mid \hat{x}) \log \frac{p(y_0 \mid x)}{p(y_1 \mid x)} \cdot p_\beta(y_0 \mid \hat{x}) p_\beta(\tilde{x} \mid \hat{x}) \log \frac{p(y_0 \mid x)}{p(y_0 \mid \tilde{x})} \\
&\quad \cdot \sum_{x',\tilde{x}',} p_\beta(x' \mid \hat{x}) p_\beta(y_0 \mid \hat{x}) \log \frac{p(y_1 \mid x')}{p(y_0 \mid x')} \cdot p_\beta(y_1 \mid \hat{x}) p_\beta(\tilde{x}' \mid \hat{x}) \log \frac{p(y_1 \mid x')}{p(y_1 \mid \tilde{x}')} \\
&\quad - \sum_{x,\tilde{x}} p_\beta(x \mid \hat{x}) p_\beta(y_0 \mid \hat{x}) \log \frac{p(y_1 \mid x)}{p(y_0 \mid x)} \cdot p_\beta(y_0 \mid \hat{x}) p_\beta(\tilde{x} \mid \hat{x}) \log \frac{p(y_0 \mid x)}{p(y_0 \mid \tilde{x})} \\
&\quad \cdot \sum_{x',\tilde{x}'} p_\beta(x' \mid \hat{x}) p_\beta(y_1 \mid \hat{x}) \log \frac{p(y_0 \mid x')}{p(y_1 \mid x')} \cdot p_\beta(y_1 \mid \hat{x}) p_\beta(\tilde{x}' \mid \hat{x}) \log \frac{p(y_1 \mid x')}{p(y_1 \mid \tilde{x}')} \\
&= \sum_{x,x',\tilde{x},\tilde{x}'} p_\beta(x \mid \hat{x}) p_\beta(x' \mid \hat{x}) p_\beta^2(y_0 \mid \hat{x}) p_\beta^2(y_1 \mid \hat{x}) p_\beta(\tilde{x} \mid \hat{x}) \log \frac{p(y_0 \mid x)}{p(y_0 \mid \tilde{x})} p_\beta(\tilde{x}' \mid \hat{x}) \log \frac{p(y_1 \mid x')}{p(y_1 \mid \tilde{x}')} \\
&\quad \cdot \left[ \log \frac{p(y_0 \mid x)}{p(y_1 \mid x)} \log \frac{p(y_1 \mid x')}{p(y_0 \mid x')} - \log \frac{p(y_0 \mid x)}{p(y_1 \mid x)} \log \frac{p(y_1 \mid x')}{p(y_0 \mid x')} \right] = 0.
\end{aligned}
$$

Given that the determinant is 0 implies that $\lambda_1(\hat{x}) = 0$. $\qquad\square$

For a binary problem we can describe the non-zero eigenvalue using $\lambda_2(\hat{x}) = \text{Tr}(C_{yy'}^{\text{dualIB}}(\hat{x}; \beta))$. That is:

$$
\begin{aligned}
\lambda_2(\hat{x}) &= \sum_{x,\tilde{x}} p_\beta(x \mid \hat{x}) p_\beta(y_1 \mid \hat{x}) \log \frac{p(y_0 \mid x)}{p(y_1 \mid x)} \cdot p_\beta(y_0 \mid \hat{x}) p_\beta(\tilde{x} \mid \hat{x}) \log \frac{p(y_0 \mid x)}{p(y_0 \mid \tilde{x})} \\
&\quad + \sum_{x,\tilde{x}} p_\beta(x \mid \hat{x}) p_\beta(y_0 \mid \hat{x}) \log \frac{p(y_1 \mid x)}{p(y_0 \mid x)} \cdot p_\beta(y_1 \mid \hat{x}) p_\beta(\tilde{x} \mid \hat{x}) \log \frac{p(y_1 \mid x)}{p(y_1 \mid \tilde{x})} \\
&= p_\beta(y_1 \mid \hat{x}) p_\beta(y_0 \mid \hat{x}) \sum_{x,\tilde{x}} p_\beta(x \mid \hat{x}) p_\beta(\tilde{x} \mid \hat{x}) \log \frac{p(y_0 \mid x)}{p(y_1 \mid x)} \left[ \log \frac{p(y_0 \mid x)}{p(y_0 \mid \tilde{x})} - \log \frac{p(y_1 \mid x)}{p(y_1 \mid \tilde{x})} \right].
\end{aligned}
$$

*1.7.11  Definition of the sample problem*

We consider a problem for a binary label $Y$ and 5 possible inputs $X$ uniformly distributed, i.e. $\forall x \in \mathcal{X}, p(x) = 1/5$ and the conditional distribution, $p(y \mid x)$, given by:

|         | $x = 0$ | $x = 1$ | $x = 2$ | $x = 3$ | $x = 4$ |
|---------|---------|---------|---------|---------|---------|
| $y = 0$ | 0.12    | 0.23    | 0.4     | 0.6     | 0.76    |
| $y = 1$ | 0.88    | 0.77    | 0.6     | 0.4     | 0.24    |

**Appendix E - Information plane analysis**

We rely on known results for the rate-distortion problem and the information plane:

**Lemma 1.7.8.** *$I(x; \hat{X})$ is a non-increasing convex function of the distortion $\mathbb{E}_{p_\beta(x,\hat{x})}[d(x,\hat{x})]$ with a slope of $-\beta$.*

We emphasis that this is a general result of rate-distortion thus holds for the dualIB as well.

**Lemma 1.7.9.** *For a fixed encoder $p_\beta(\hat{x} \mid x)$ and the Bayes optimal decoder $p_\beta(y \mid \hat{x})$:*

$$\mathbb{E}_{p_\beta(x,\hat{x})}[d_{\mathrm{IB}}(x,\hat{x})] = I(X;Y) - I(\hat{X};Y).$$

*Thus, the information curve, $I_y$ vs. $I_x$, is a non-decreasing concave function with a positive slope, $\beta^{-1}$. The concavity implies that $\beta$ increases along the curve.*

Cover and Thomas (2006); Gilad-bachrach et al. (2003).

*1.7.12   Proof of Lemma 3*

In the following section we provide a proof to *lemma* 3, for the IB and dualIB problems.

*Proof.* We want to analyze the behavior of $I_x(\beta)$, $I_y(\beta)$, that is the change in each term as a function of the corresponding $\beta$. From *lemma* 1.7.9, the concavity of the information curve, we can deduce that both are non-decreasing functions of $\beta$. As the two $\beta$ derivatives are proportional it's enough to discuss the first one.

Next, we focus on their behavior between two critical points. That is, where the cardinality of $\hat{X}$ is fixed (clusters are "static"). For "static" clusters, the $\beta$ derivative of $I_x$, along the optimal line is given by:

$$\begin{aligned}
\frac{\partial I(X;\hat{X})}{\partial \beta} &= -\frac{\partial}{\partial \beta}\left[\sum_{x,\hat{x}} p_\beta(x,\hat{x})\big(\log Z_{\hat{\mathbf{x}}|\mathbf{x}}(x;\beta) + \beta d(x,\hat{x})\big)\right] \\
&= -\beta\left\langle d(x,\hat{x})\frac{\partial \log p_\beta(\hat{x} \mid x)}{\partial \beta}\right\rangle_{p_\beta(x,\hat{x})} \\
&\approx \beta\left\langle d(x,\hat{x})\left[\frac{\partial \log Z_{\hat{\mathbf{x}}|\mathbf{x}}(x;\beta)}{\partial \beta} + d(x,\hat{x})\right]\right\rangle_{p_\beta(x,\hat{x})} \\
&\approx \beta\left\langle \underbrace{\left\langle d^2(x,\hat{x})\right\rangle_{p_\beta(\hat{x}|x)} - \left\langle d(x,\hat{x})\right\rangle^2_{p_\beta(\hat{x}|x)}}_{\mathrm{Var}(d(x))}\right\rangle_{p(x)} .
\end{aligned}$$

This first of all reassures that the function is non-decreasing as $\mathrm{Var}(d(x)) \geq 0$.

The piece-wise concavity follows from the fact that when the number of clusters is fixed (between the critical points) - increasing $\beta$ decreases the clusters conditional entropy $H(\hat{X} \mid x)$, as the encoder becomes more deterministic. The mutual information is bounded

116

by $H(\hat{X})$ and it's $\beta$ derivative decreases. Further, between the critical points there are no sign changes in the second $\beta$ derivative. □

### 1.7.13 Proof of Theorem 4

*Proof.* The proof follows from *lemma* 3 together with the critical points analysis above, and is only sketched here. As the encoder and decoder at the critical points, $\beta_c^{\mathrm{IB}}$ and $\beta_c^{\mathrm{dualIB}}$, have different left and right derivatives, they form cusps in the curves of the mutual information ($I_x$ and $I_y$) as functions of $\beta$. These cusps can only be consistent with the optimality of the IB curves ( implying that sub-optimal curves lie below it; i.e, the IB slope is steeper) if $\beta_c^{\mathrm{dualIB}} < \beta_c^{\mathrm{IB}}$ (this is true for any sub-optimal distortion), otherwise the curves intersect.

Moreover, at the dualIB critical points, the distance between the curves is minimized due to the strict concavity of the functions segments between the critical points. As the critical points imply discontinuity in the derivative, this results in a "jump" in the information values. Therefore, at any $\beta_c^{\mathrm{dualIB}}$ the distance between the curves has a (local) minimum. This is depicted in Figure 4 (in the main text), comparing $I_x(\beta)$ and $I_y(\beta)$ and their differences for the two algorithms.

The two curves approach each other for large $\beta$ since the two distortion functions become close in the low distortion limit (as long as $p(y \mid x)$ is bounded away from 0). □

### Appendix F - Derivation of the dualExpIB

We provide elaborate derivations to *theorem* 9; that is, we obtain the dualIB optimal encoder-decoder under the exponential assumption over the data. We use the notations defined in §*The Exponential Family dualIB*.

- The *decoder*. Substituting the exponential assumption into the dualIB log-decoder yields:

$$\log p_\beta(y \mid \hat{x}) = \sum_x p_\beta(x \mid \hat{x}) \log p(y \mid x) - \log Z_{\mathbf{y}|\hat{\mathbf{x}}}(\hat{x}; \beta)$$

$$= -\sum_x \sum_{r=0}^{d} p_\beta(x \mid \hat{x}) \lambda^r(y) A_r(x) - \log Z_{\mathbf{y}|\hat{\mathbf{x}}}(\hat{x}; \beta)$$

$$= -\sum_{r=1}^{d} \lambda^r(y) A_{r,\beta}(\hat{x}) - \mathbb{E}_{p_\beta(x|\hat{x})}\left[\lambda_{\mathbf{x}}^0\right] - \log Z_{\mathbf{y}|\hat{\mathbf{x}}}(\hat{x}; \beta).$$

Taking a closer look at the normalization term:

$$Z_{\mathbf{y}|\hat{\mathbf{x}}}(\hat{x}; \beta) = \sum_y e^{\sum_x p_\beta(x|\hat{x}) \log p(y|x)} = e^{-\mathbb{E}_{p_\beta(x|\hat{x})}\left[\lambda_{\mathbf{x}}^0\right]} \sum_y e^{-\sum_{r=1}^d \lambda^r(y) A_{r,\beta}(\hat{x})}$$

$$\log Z_{\mathbf{y}|\hat{\mathbf{x}}}(\hat{x}; \beta) = -\mathbb{E}_{p_\beta(x|\hat{x})}\left[\lambda_{\mathbf{x}}^0\right] + \log\left(\sum_y e^{-\sum_{r=1}^d \lambda^r(y) A_{r,\beta}(\hat{x})}\right).$$

From which it follows that $\lambda_\beta^0(\hat{x})$ is given by:

$$\lambda_\beta^0(\hat{x}) = \log\left(\sum_y e^{-\sum_{r=1}^d \lambda^r(y) A_{r,\beta}(\hat{x})}\right),$$

and we can conclude that the dualExpIB decoder takes the form:

$$\log p_\beta(y \mid \hat{x}) = -\sum_{r=1}^d \lambda^r(y) A_{r,\beta}(\hat{x}) - \lambda_\beta^0(\hat{x}).$$

- The *encoder*.

  The core of the encoder is the dual distortion function which may now be written as:

  $$\begin{aligned} d_{\text{dualIB}}(x, \hat{x}) &= \sum_y p_\beta(y \mid \hat{x}) \log \frac{p_\beta(y \mid \hat{x})}{p(y \mid x)} \\ &= \sum_y p_\beta(y \mid \hat{x})\left[\left(\lambda_{\mathbf{x}}^0 - \lambda_\beta^0(\hat{x})\right) + \sum_{r=1}^d \lambda^r(y)(A_r(x) - A_{r,\beta}(\hat{x}))\right] \\ &= \lambda_{\mathbf{x}}^0 - \lambda_\beta^0(\hat{x}) + \sum_{r=1}^d \lambda_\beta^r(\hat{x})(A_r(x) - A_{r,\beta}(\hat{x})), \end{aligned}$$

  substituting this into the encoder's definition we obtain:

  $$\begin{aligned} p_\beta(\hat{x} \mid x) &= \frac{p_\beta(\hat{x})}{Z_{\hat{\mathbf{x}}|\mathbf{x}}(x; \beta)} e^{-\beta\left[\lambda_{\mathbf{x}}^0 - \lambda_\beta^0(\hat{x}) + \sum_{r=1}^d \lambda_\beta^r(\hat{x})[A_r(x) - A_{r,\beta}(\hat{x})]\right]} \\ &= \frac{p_\beta(\hat{x}) e^{\beta \lambda_\beta^0(\hat{x})}}{Z_{\hat{\mathbf{x}}|\mathbf{x}}(x; \beta)} e^{-\beta \sum_{r=1}^d \lambda_\beta^r(\hat{x})[A_r(x) - A_{r,\beta}(\hat{x})]}. \end{aligned}$$

We can further write down the information quantities under these assumptions:

$$I(X; \hat{X}) = \sum_{x,\hat{x}} p_\beta(x, \hat{x}) \log \frac{p_\beta(x \mid \hat{x})}{p(x)}$$

$$= H(X) - \beta \sum_{r=1}^{d} \sum_{\hat{x}} p_\beta(\hat{x}) \lambda_\beta^r(\hat{x}) \left[ \sum_x p_\beta(x \mid \hat{x}) A_r(x) - A_{r,\beta}(\hat{x}) \right] +$$

$$\beta \mathbb{E}_{p_\beta(\hat{x})} \left[ \lambda_\beta^0(\hat{x}) \right] - \mathbb{E}_{p(x)} \left[ \log Z_{\hat{\mathbf{x}} \mid \mathbf{x}}(x; \beta) \right]$$

$$= H(X) + \beta \mathbb{E}_{p_\beta(\hat{x})} \left[ \lambda_\beta^0(\hat{x}) \right] - \mathbb{E}_{p(x)} \left[ \log Z_{\hat{\mathbf{x}} \mid \mathbf{x}}(x; \beta) \right]$$

$$I(Y; \hat{X}) = \sum_{y,\hat{x}} p_\beta(y, \hat{x}) \log \frac{p_\beta(y \mid \hat{x})}{p(y)}$$

$$= H(Y) - \sum_{r=1}^{d} \sum_{\hat{x}} p_\beta(\hat{x}) \sum_y p_\beta(y \mid \hat{x}) \lambda^r(y) A_{r,\beta}(\hat{x}) - \mathbb{E}_{p_\beta(\hat{x})} \left[ \lambda_\beta^0(\hat{x}) \right]$$

$$= H(Y) - \mathbb{E}_{p_\beta(\hat{x})} \left[ \sum_{r=1}^{d} \lambda_\beta^r(\hat{x}) A_{r,\beta}(\hat{x}) + \lambda_\beta^0(\hat{x}) \right]$$

### Appendix G - Optimizing the error exponent

We start by to expressing the Chernoff information for the binary hypothesis testing problem using $p(y \mid x)$:

$$C(p_0, p_1) = \min_{\lambda \in [0,1]} \log \left( \sum_x p(x \mid y_0)^{q_\lambda(y_0)} p(x \mid y_1)^{q_\lambda(y_1)} \right)$$

$$= \min_{\lambda \in [0,1]} \log \left( \sum_x p(y = 0 \mid x)^\lambda p(x)^\lambda p(y = 0)^{-\lambda} p(y = 1 \mid x)^{1-\lambda} p(x)^{1-\lambda} p(y = 1)^{\lambda - 1} \right)$$

$$= \min_{\lambda \in [0,1]} \log \left( \sum_x p(x) p(y = 0 \mid x)^\lambda p(y = 1 \mid x)^{1-\lambda} \right) - \log \left( p(y = 0)^\lambda p(y = 1)^{1-\lambda} \right)$$

$$= \min_{q_\lambda(y)} \log \left( \sum_x e^{q_\lambda(y_0) \log p(x \mid y_0) + q_\lambda(y_1) \log p(x \mid y_1)} \right)$$

$$= \min_{q_\lambda(y)} \log \left( \sum_x e^{-D[q_\lambda(y) \| p(y \mid x)] + D[q_\lambda(y) \| py] + \log p(x)} \right)$$

$$= \min_{q_\lambda(y)} \log \left( e^{D[q_\lambda(y) \| p(y)]} \sum_x e^{-D[q_\lambda(y) \| p(y \mid x)] + \log p(x)} \right)$$

$$= \min_{q_\lambda(y)} \left\{ \log \left( \sum_x p(x) e^{-D[q_\lambda(y) \| p(y \mid x)]} \right) + D[q_\lambda(y) \| p(y)] \right\},$$

where $q_\lambda(y_0) = \lambda, q_\lambda(y_1) = 1 - \lambda$. Now, if we consider the mapping, $q_\lambda(y) = p_\beta(y \mid \hat{x})$ we can write the above as:

$$C(p_0, p_1) = \min_{p_\beta(y|\hat{x})} \left\{ \log \left( \sum_x p(x) e^{-D[p_\beta(y|\hat{x})\|p(y|x)]} \right) + D[p_\beta(y \mid \hat{x})\|p(y)] \right\}.$$

The above term in minimization is proportional to log-partition function of $p_\beta(x \mid \hat{x})$, namely we get the mapping $p_\beta(x \mid \hat{x}) = p_\lambda$. Next we shall generalize the setting to the $M$-hypothesis testing problem. Having that solving for the Chernoff information is notoriously difficult we consider an upper bound to it, taking the expectation over the classes. Instead of choosing $p_{\lambda^*}$ as the maximal value of the minimimum $\{D[p_{\lambda^*}\|p_0], D[p_{\lambda^*}\|p_1]\}$ we consider it w.r.t the full set $\{D[p_{\lambda^*}\|p_i]\}_{i=1}^M$. Using the above mapping we must take the expectation also over the representation variable $\hat{x}$. Thus we get the expression:

$$D^*(\beta) = \min_{p_\beta(y|\hat{x}), p_\beta(\hat{x}|x)} \mathbb{E}_{p_\beta(y,\hat{x})}[D[p_\beta(x \mid \hat{x}) \mid p(x \mid y)]].$$

From the definition of $D^*(\beta)$ we obtain the desired bound of the dualIB:

$$D^*(\beta) = \min_{p_\beta(y|\hat{x}), p_\beta(\hat{x}|x)} \mathbb{E}_{p_\beta(y,\hat{x})}[D[p_\beta(x \mid \hat{x})\|p(x \mid y)]]$$

$$= \min_{p_\beta(y|\hat{x}), p_\beta(\hat{x}|x)} \sum_{x,y,\hat{x}} p_\beta(y \mid \hat{x}) p_\beta(\hat{x}) [D[p_\beta(x \mid \hat{x}) \mid p(x \mid y)]]$$

$$= \min_{p_\beta(y|\hat{x}), p_\beta(\hat{x}|x)} \sum_{x,y,\hat{x}} p_\beta(y \mid \hat{x}) p_\beta(\hat{x}) p_\beta(x \mid \hat{x}) \left\{ \log \frac{p_\beta(y \mid \hat{x})}{p(y \mid x)} + \log \frac{p_\beta(x \mid \hat{x})}{p_\beta(y \mid \hat{x})} + \log \frac{p(y)}{p(x)} \right\}$$

$$= \min_{p_\beta(y|\hat{x}), p_\beta(\hat{x}|x)} \left\{ I(X; \hat{X}) + \mathbb{E}_{p_\beta(x,\hat{x})}[D[p_\beta(y \mid \hat{x})\|p(y \mid x)]] + H(Y \mid \hat{X}) + \mathbb{E}_{p_\beta(y)}[\log p(y)] \right\}$$

$$\leq \min_{p_\beta(y|\hat{x}), p_\beta(\hat{x}|x)} \left\{ I(X; \hat{X}) + \mathbb{E}_{p_\beta(x,\hat{x})}[D[p_\beta(y \mid \hat{x})\|p(y \mid x)]] \right\}$$

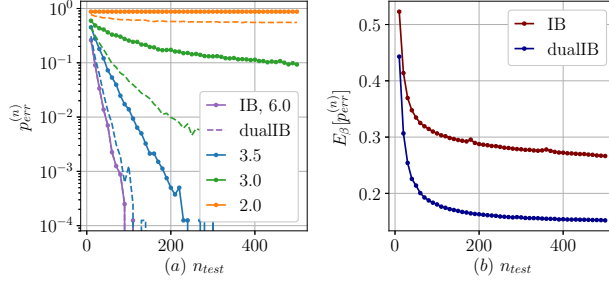$$\leq \mathcal{F}^*[p(\hat{x} \mid x); p(y \mid \hat{x})].$$

*1.7.14   Error exponent optimization example*

To demonstrate the above properties we consider a classification problem with $M = 8$ classes, each class characterized by $p_i = p(x \mid y_i)$. The *training* is performed according to the above algorithms to obtain the IB (dualIB) encoder and decoder. For the prediction, given a new sample $x^{(n)} \overset{i.i.d}{\sim} p(x \mid y)$ defining an empirical distribution $\hat{p}(x)$ the prediction is done by first evaluating $\hat{p}_\beta(\hat{x}) = \sum_x p_\beta(\hat{x} \mid x) \hat{p}(x)$. Next, using the (representation) optimal decision rule, we obtain the prediction:

$$\hat{H}_\beta = \arg \min_i D[\hat{p}_\beta(\hat{x})\|p_\beta(\hat{x} \mid y_i)],$$

and we report $p_{err}^{(n)}$, the probability of miss-classification. This represents the most general classification task; the distributions $p_i$ represent the empirical distributions over a training

data-set and then testing is performed relative to a test set. Looking at the results, Figure 1.18, it is evident that indeed the dualIB improves the prediction error (at $log_2(\beta) = 6$ the algorithms performance is identical due to the similarity of the algorithms behavior as $\beta$ increases).



**Figure 1.18: The probability of error, $p_{err}^{(n)}$, as a function of test sample size, $n_{test}$. (a) The exponential decay of error for representative $\beta$ values ($\log_2(\beta)$ reported in the legend). For a given $\beta$ the IB performance is plotted in solid line and the dualIB in dashed (for $\log_2(\beta) = 6$ the lines overlap). (b) The expectation of the error over all $\beta$'s ($\log_2 \beta \in [1, 6]$).**

### *Appendix H - The variational* dualIB

*Derivation of the* VdualIB *objective*

Just as (Fischer and Alemi; 2020) did, we can variationally upper bound the information of the input with the representation variable using:

$$I(\hat{X}; X \mid Y) = \mathbb{E}_{p(x,y)p(\hat{x}|x)}\left[\log \frac{p(\hat{x} \mid x, y)}{p(\hat{x} \mid y)}\right] \leq \mathbb{E}_{\tilde{p}(y|x)p(x)p(\hat{x}|x)}\left[\log \frac{p(\hat{x} \mid x)}{q(\hat{x} \mid y)}\right]$$

where $q(\hat{x} \mid y)$ is a variational class conditional marginal. In contradiction to the CEB, in order to bound the dualIB distortion, we replace the bound on $I(\hat{X}; Y)$ with a bound over the expected dualIB distortion. Here, given the assumption of a noise model $\tilde{p}(y \mid x)$ which we evaluate the expected distortion with respect to it:

$$\mathbb{E}_{p(x,\hat{x})}[d_{\text{dualIB}}(x, \hat{x})] = \mathbb{E}_{p(y|\hat{x})p(\hat{x}|x)p(x)}\left[\log \frac{p(y \mid \hat{x})}{\tilde{p}(y \mid x)}\right]$$

Combining the above together gives the variational upper bound to the dualIB as the following objective:

$$I(X; \hat{X}) + \beta \mathbb{E}_{p(x,\hat{x})}[d_{\text{dualIB}}(x, \hat{x})] \leq \mathbb{E}_{\tilde{p}(y|x)p(\hat{x}|x)p(x)}\left[\log \frac{p(\hat{x} \mid x)}{q(\hat{x} \mid y)}\right] + \beta \mathbb{E}_{p(y|\hat{x})p(\hat{x}|x)p(x)}\left[\log \frac{p(y \mid \hat{x})}{\tilde{p}(y \mid x)}\right]$$

*Experimental setup*

For both CIFAR10 and FasionMNIST We trained a set of $30$ $28 - 10$ Wide ResNet models in a range of values of $\beta$ ($-5 \le \log \beta \le 5$). The training was doneusing Adam (Kingma and Ba; 2014) at a base learning rate of $10^{-4}$. We lowered the learning rate two times by a factor of 0.3 each time. Additionally, following Fischer and Alemi (2020), we use a jump-start method for $\beta < 100$. We start the training with $\beta = 100$, anneal down to the target $\beta$ over 1000 steps. The training includes data augmentation with horizontal flip and width height shifts. Note, that we exclude from the analysis runs that didn't succeed to learn at all (for which the results look as random points).

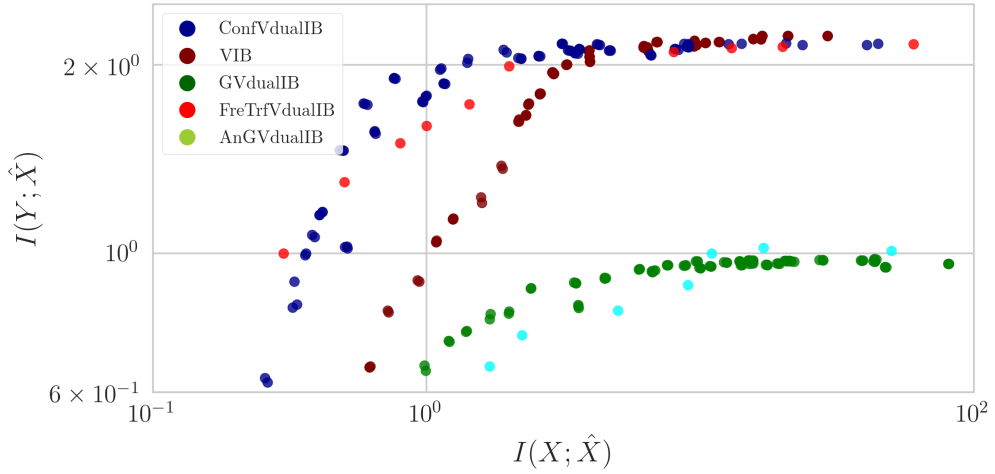*The variational information plane*

Note that, for the information plane analysis, there were several runs that failed to achieved more than random accuracy. In such cases, we remove them. The confusion matrix used for the FashionMNIST data-set is:

$$
\begin{pmatrix}
0.828 & 0.013 & 0.012 & 0.011 & 0.018 & 0. & 0.002 & 0.004 & 0.085 & 0.027 \\
0.01 & 0.91 & 0. & 0.005 & 0.001 & 0.001 & 0. & 0.001 & 0.011 & 0.061 \\
0.047 & 0.001 & 0.708 & 0.064 & 0.088 & 0.014 & 0.063 & 0.004 & 0.008 & 0.003 \\
0.003 & 0.004 & 0.016 & 0.768 & 0.033 & 0.093 & 0.05 & 0.019 & 0.004 & 0.01 \\
0.01 & 0. & 0.039 & 0.043 & 0.788 & 0.012 & 0.057 & 0.043 & 0.006 & 0.002 \\
0.002 & 0. & 0.01 & 0.137 & 0.029 & 0.777 & 0.008 & 0.033 & 0. & 0.004 \\
0.007 & 0.002 & 0.01 & 0.054 & 0.029 & 0.007 & 0.888 & 0.001 & 0.001 & 0.001 \\
0.024 & 0.002 & 0.014 & 0.039 & 0.076 & 0.017 & 0.004 & 0.818 & 0.002 & 0.004 \\
0.027 & 0.013 & 0. & 0.007 & 0.003 & 0. & 0.003 & 0. & 0.933 & 0.014. \\
0.019 & 0.064 & 0.001 & 0.007 & 0.002 & 0.001 & 0.001 & 0. & 0.018 & 0.887
\end{pmatrix}
$$

*The* VdualIB *noise models*

As described in the main text we consider two additional noise models; (i) An analytic Gaussian integration of the log-loss around the one-hot labels (AnGVdualIB) (ii) Using predictions of another trained model as the induced distribution (PrdTrVdualIB). In this case, we use a deterministic wide ResNet $28 - 10$ network that achieved 95.8% accuracy on CIFAR10. In Figure 1.19 we can see all the different models, 4 noise models for the VdualIB and the VIB). As expected, we can see that analytic Gaussian integration noise model obtains similar results to adding Gaussian noise to the one-hot vector of the true label, while the performance of the noise models that are based on a trained network are similar to the ConfVdualIB.
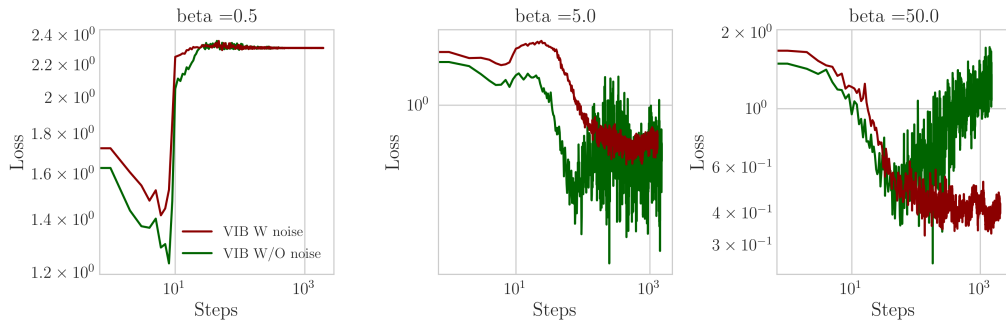
**Figure 1.19: The information plane for the different noise models.**

### 1.7.14.1  *Training* VIB *model with noise*

In our analysis, we train a VIB model with the same noise model as the VdualIB. Namely, instead of training with a deterministic label (one-hot vector of zeros and ones), we use our noise model also for the VIB. As mentioned in the text, this training procedure is closely related to label smoothing. In Figure 1.20, we present the loss function of the VIB on CIFAR10 with and without the noise models along the training process for 3 different values of $\beta$. For a small $\beta$ (left) both regimes under-fit the data as expected. However, when we enlarge $\beta$, we can see that the labels' noise makes the training more stable and for a high value of $\beta$ (right) training without noise over-fits the data and the loss increases.



**Figure 1.20: The influence of a noise model on the VIB performance. Loss as function of the update steps for different values of $\beta$, $\beta = 0.5, 5.0. 50.0$ from left to right.**
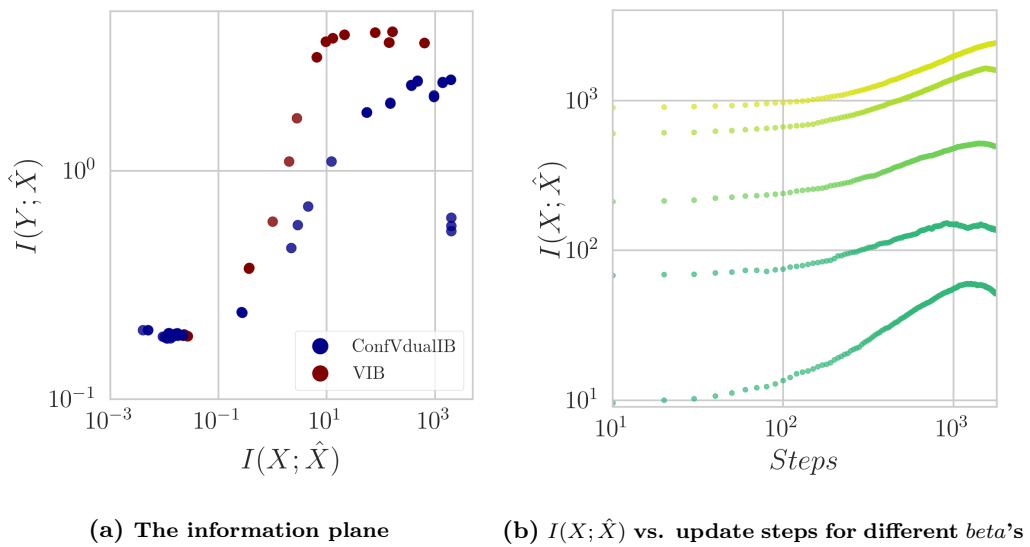
123

The confusion matrix for the CIFAR10 data set is:

$$
\begin{pmatrix}
0.878 & 0. & 0.017 & 0.013 & 0.002 & 0.001 & 0.082 & 0. & 0.007 & 0. \\
0. & 0.984 & 0.002 & 0.009 & 0.001 & 0. & 0.003 & 0. & 0. & 0. \\
0.013 & 0.001 & 0.896 & 0.009 & 0.038 & 0. & 0.043 & 0. & 0. & 0. \\
0.022 & 0.004 & 0.011 & 0.913 & 0.023 & 0. & 0.027 & 0. & 0.001 & 0. \\
0. & 0. & 0.072 & 0.022 & 0.85 & 0. & 0.058 & 0. & 0. & 0. \\
0. & 0. & 0. & 0. & 0. & 0.982 & 0. & 0.011 & 0. & 0.007 \\
0.099 & 0.001 & 0.049 & 0.021 & 0.055 & 0. & 0.768 & 0. & 0.005 & 0. \\
0. & 0. & 0. & 0. & 0. & 0.006 & 0. & 0.976 & 0. & 0.019 \\
0.004 & 0.001 & 0.001 & 0.001 & 0.004 & 0.002 & 0.003 & 0.001 & 0.98 & 0.001 \\
0. & 0. & 0. & 0. & 0. & 0.004 & 0. & 0.02 & 0.001 & 0.974
\end{pmatrix}
$$

*CIFAR100 results*

As mentioned in the text, we trained VdualIB networks also on CIFAR100. For this, we used the same $28 - 10$ Wide ResNet with a confusion matrix as our noise model. The confusion matrix was calculated based on the predictions of a deterministic network. The deterministic network achieved $80.2\%$ accuracy on CIFAR100. In 1.21a, we can see the information plane for both VdualIB and the VIB models. As we can see, both models are monotonic with $I(X; \hat{X})$, however the VIB's performance is better. The VIB achieves higher values of information with the labels along with more compressed representation at any given level of predication. Although a broader analysis is required, and possible further parameter tuning of the architecture, we hypothesize that the caveat is in the noise model used for the VdualIB. Using a noise model which is based on a network that achieves almost $20\%$ error might be insufficient in this case. It might be that "errors" in the noise model becomes similar to random errors, similar to the Gaussian case, and hence depicting similar learning performance to the GVdualIB case.

When we look at the information with the input as a function of time ,Figure 1.21b, we see that similar to the FasionMNIST and CIFAR10 results, the information saturates for small values of $\beta$, but over-fits for higher values of it.

**(a)** The information plane

**(b)** $I(X;\hat{X})$ vs. update steps for different *beta*'s

**Figure 1.21: Experiments over CIFAR100.** (*a*) **The information plane of the** ConfVdualIB **and** VIB **for a range of** $\beta$ **values at the final training step.** (*b*) **The evolution of the the** ConfVdualIB**'s** $I(X;\hat{X})$ **along the optimization update steps.**

# General Discussion

This thesis explored DNNs via an information theory perspective. Using the Information Bottleneck (IB) principle, we explored the underlying behavior of DNNs. We formulated deep learning as an information-theoretic trade-off between compression and prediction, which gives an optimal representation for each layer. Our theory is rooted in the idea that DNN learning with SGD aims to learn optimal representations in the IB sense. At its core, our theory aims to summarize each hidden layer using mutual information with the input and output. In the first section of this thesis, we combined the novel perspective with an empirical case study to make claims about phase transitions in SGD optimization dynamics, the computational benefits of deep networks, and the relation between generalization and compressed representation. These observations were collected to a new information-theoretic paradigm to explain deep learning, which inspired multiple follow-up works. As several researchers have noticed, measuring information in deterministic DNNs is hard (Goldfeld et al.; 2018). To overcome this problem, in the second section of this thesis, we utilized the NTK framework to derive many information-theoretic quantities in infinitely wide neural networks. These quantities allow us to explore where the information is in DNNs and the relationship between generalization, compression, and information. By ensemble over different initial conditions, we tried to find the hyperparameters' effect on the network's information. Our analysis revealed several interesting connections between the different information-theoretic quantities in this framework, the optimally of DNNs, their generalization ability, and capacity.

In the last section of this thesis, we presented the dualIB framework, which enabled an optimal representation that resolves some of the original IB's drawbacks. We provided the dualIB self-consistent equations, which allowed us to obtain analytical solutions; we characterized the structure of the critical points of the solutions, resulting in a full bifurcation diagram of its representation; we derived several interesting properties of the dualIB: First, when the data can be modeled in a parametric form, it preserves the original distribution's statistics. Second, it optimizes the mean prediction error exponent, therefore improving the predictions as a function of the data size. Additionally, we provide a variational dualIB framework. By optimizing its functional using DNNs, we can apply the dualIB for real-world datasets. These results demonstrate the potential advantages of

the framework in the context of information in DNNs.

Although many works have been done based on ideas from this thesis, the framework, and the experiments laid out in it open several more important avenues for future research. A few examples are outlined below.

- **Generalization and compression** – One of the most important questions in DNNs is their generalization ability. The followed-up works that have presented empirically investigations of the connection between compression of the information and generalization have attained mixed results. While the authors of Cheng et al. (2019) saw a clear connection, the authors of Gabrié et al. (2018) concluded that compression and generalization may not be linked. This contradicts many PAC-Bayes bounds on the information, which bound the generalization gap between train and test error by these information-theoretic quantities. The natural question that arises is the reason for this discrepancy. Is it due to poor information estimators in DNNs? Are these bounds not tight enough? This research line can shed light on the important information quantities and how useful they can explain DNNs. Our work involving infinitely-wide neural networks is the first step in this direction. However, a broader analysis is needed to understand the connections between the different factors.

- **Multi-domain, semi-supervised representation learning** – Semi-supervised learning takes advantage of a large amount of unlabeled data that are available for many uses in addition to typically smaller sets of labeled data (Van Engelen and Hoos; 2020). However, semi-supervised learning's advantages are unclear from an information and probabilistic perspective. Taking advantage of this approach to multi-domain data allows us to better control information in our network. Using information-theoretic principles and self-supervised approaches, we can create compressed representation learning from multi-domain and multi-task data. Several questions remain to be examined: How should data from different modalities be stored? For each domain, what are the main factors leading to better generalizations? By compressing irrelevant information within each domain, how can we make our models more robust by using semi-supervised learning? The generalization and robustness should be improved with the integration of multimodality models using semi-supervised learning.

- **The benefits of the hidden layers in DNNs** – This thesis suggested that one of the benefits of hidden layers is computational; by adding layers, the amount of compression each layer needs to do is reduced, resulting in a dramatic reduction in training time. Nevertheless, training very neural networks is challenging, and the network's performance does not monotonically increase with layers (Zagoruyko and Komodakis; 2017). Understanding this trade-off between learning from finite samples and computational benefit will allow us to develop better principle designs for DNNs.

- **Finite-sample information plane** – In sections 1 and 2, we illustrated the network's behavior in the information plane for different dataset sizes. It remains a challenge, however, to pinpoint the exact relationship between the two. We need to investigate further how the finite-sample problem affects the optimization dynamics and whether a corresponding IB problem exists for each dataset size.

- **Variational dualIB** – Our thesis presented a variational dualIB formulation. A study comparing the differences between the IB and the dualIB models would be an interesting future direction. According to our initial analysis, the variational dualIB model compresses better than the regular VIB model. It will be interesting to explore it on larger datasets, more noise models, and different architectures. In addition, the question whether the differences between the models affect network properties, such as robustness, remains open.

- **Biologically plausible models** – The question of how the brain processes sensory input and elevates it was at the heart of much of the early interest in neural networks. In spite of being inspired by brains, deep artificial neural networks do not exhibit brain-like characteristics. Adapting biologically plausible deep learning algorithms to information-theoretic learning principles would be an interesting future direction. Gradient back-propagation, for example, relies on mechanisms that seem biologically implausible. It is not possible to alternate between a bottom-up forward pass and a top-down backward pass, and the labels are not available for each example. These problems may be solveable by the IB principle. For example, by combining a bottleneck objective with self-supervised learning, neural networks can be trained layer-wise without labels presented and through alternative bottom-up forward passes and top-down backward passes (Pogodin and Latham; 2020).

More generally, in this thesis, we outlined the foundation for a novel, comprehensive theory of large-scale learning via deep neural networks that builds upon the correspondence between deep learning and the information bottleneck framework. One of the most important and challenging directions in this field is to use theoretical tools from other fields to analyze deep networks. Even though most deep learning studies use practical applications, they can only be effective if they are backed up by good theoretical knowledge. Our theory offers a number of benefits, such as providing a deeper understanding of the information that resides within DNNs, identifying different explanations for their behavior, and opening up theoretical and practical research opportunities in the field. Hence, many open questions remain for further research, and by combining these directions, we may arrive at a stronger theory of the field that can also bring practical benefits and specific design principles.

### Bibliography

Cheng, H., Lian, D. and Gao, S.and Geng, Y. (2019). Utilizing information bottleneck to evaluate the capability of deep neural networks for image classification, *Entropy* .

Gabrié, M., Manoel, A., Luneau, C., Barbier, J., Macris, N., Krzakala, F. and Zdeborová, L. (2018). Entropy and mutual information in models of deep neural networks, *arXiv preprint arXiv:1805.09785* .

Goldfeld, Z., Berg, E. v. d., Greenewald, K., Melnyk, I., Nguyen, N., Kingsbury, B. and Polyanskiy, Y. (2018). Estimating information flow in deep neural networks, *arXiv preprint arXiv:1810.05728* .

Pogodin, R. and Latham, P. E. (2020). Kernelized information bottleneck leads to biologically plausible 3-factor hebbian learning in deep networks, *arXiv preprint arXiv:2006.07123* .

Van Engelen, J. E. and Hoos, H. H. (2020). A survey on semi-supervised learning, *Machine Learning* **109**(2): 373–440.

Zagoruyko, S. and Komodakis, N. (2017). Diracnets: Training very deep neural networks without skip-connections, *arXiv preprint arXiv:1706.00388* .