
A Tighter Bound on the Information Bottleneck with Applications to Deep Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The Information Bottleneck, IB, provides a hypothetically optimal framework for
2 unsupervised data modeling, yet is often intractable. Recent efforts optimized
3 supervised DNNs with a variational upper bound to the IB objective, resulting in
4 improved robustness to adversarial attacks. However, when deriving the upper
5 bound, it is presumed that the supervisor distribution $p^*(y)$ is known, as it would
6 be in the unsupervised case. This work demonstrates that lifting this assumption
7 not only results in a tighter bound on the IB and improved empirical performance,
8 but also proposes a new motivation for regularization.

9

1 Introduction

10 Deep Neural Nets (DNNs) learn latent representations induced by their downstream task, objective
11 function, and other parameters. The quality of the learned representations impacts the DNN's
12 generalization ability and the coherence of the emerging latent space (Bengio 2009). A question
13 emerges regarding the extraction of an optimal latent representation for all data points from a restricted
14 set of training examples. Classic information theory provides rate-distortion (Shannon 1959) for
15 optimal compression of data. However, rate-distortion regards all information as equal, not taking
16 into account which information is more relevant to a specified downstream task, without constructing
17 tailored distortion functions. The Information Bottleneck, IB, (Tishby, Pereira, and Bialek 1999)
18 resolves this limitation by defining mutual information, MI, between the learned representation and a
19 designated downstream task as a universal distortion function. Yet, learning representations using
20 the IB method requires computations of mutual information, which is possible between discrete
21 distributions, and some continuous ones, but not in the general case (Chechik et al. 2003). Moreover,
22 MI is either difficult or impossible to optimize over when using deterministic models (Saxe et al.
23 2018; Amjad and Geiger 2020). Nonetheless, the promise of the IB remains alluring, and recent
24 works utilized VAE (Kingma and Welling 2014) inspired variational approximations to approximate
25 upper bounds to the IB objective, allowing its utilization as a loss function for DNNs, where the
26 underlying distributions are both continuous and unknown (Alemi et al. 2017; Fischer 2020; Cheng
27 et al. 2020). These approaches learn representations in supervised settings, without knowledge of the
28 underlying distribution $p^*(x, y)$, utilizing the learned variational conditional $p(y|x)$ to approximate
29 MI. In contrast, non variational IB methods learn representations in unsupervised settings, where the

30 stochastic process underlying the observed data is known (Tishby, Pereira, and Bialek 1999; Chechik
31 et al. 2003; Painsky and Tishby 2017). Nonetheless, when deriving the variational IB objectives,
32 previous research relax the problem by assuming that $p^*(y)$ is constant, while in practice it's an
33 optimized variational approximation. We derive a new upper bound for the IB objective, and a
34 subsequent variational approximation, by removing the relaxation. We show that our bound is tighter
35 than previous bounds, and that our proposed loss function is a tighter variational approximation,
36 when considering $p(y|x)$ as part of the optimization. We believe our new derivation is a better
37 adaptation of the IB for supervised tasks, and show empirical evidence of improved performance
38 across several challenging tasks over different modalities. We utilize previous studies on variational
39 representation learning and regularization (Alemi et al. 2018; Pereyra et al. 2017; Achille and Soatto
40 2018) to interpret our findings, and conclude that our proposed derivation applies regularization over
41 the variational classifier, preventing it from overfitting the learned representations and thus enabling
42 greater MI between learned representation and the real and unknown stochastic process.

43 The reader is encouraged to refer to the preliminaries provided in Appendix A before proceeding.

44 2 Related work

45 2.1 Deterministic Information Bottleneck

46 Classic information theory offers rate-distortion (Shannon 1959) to mitigate signal loss during
47 compression: A source X is compressed to an encoding Z , such that maximal compression is
48 achieved while keeping the encoding quality above a certain threshold. Encoding quality is measured
49 by a task specific distortion function: $d : X \times Z \mapsto \mathbb{R}^+$. Rate-distortion suggests a mapping that
50 minimizes the rate of bits to source sample, measured by $I(X; Z)$, that adheres to a chosen allowed
51 expected distortion $D \geq 0$. The Information Bottleneck, IB, (Tishby, Pereira, and Bialek 1999)
52 extends rate-distortion by replacing the tailored distortion functions with MI over a target distribution:
53 Let Y be the target signal for some specific downstream task, such that the joint distribution $P^*(x, y)$
54 is known, and define the distortion function as MI between Z and Y . The IB is the solution to
55 the optimization problem $Z : \min_{P(Z|X)} I(X; Z)$ subject to $I(Z; Y) \geq D$, that can be optimized by
56 minimizing the IB objective $\mathcal{L}_{IB} = I(X; Z) - \beta I(Z; Y)$ over $P(Z|X)$. The solution to the IB
57 objective is a function of the Lagrange multiplier β , and is a theoretical limit for representation
58 quality, given mutual information as an accepted metric, as elaborated in more detail in Appendix B.
59 The IB is in fact an unsupervised soft clustering problem, where each data point x is assigned a
60 probabilities to belong to the different clusters z (Slonim 2002). Chechik et al. (2003) showed that
61 computing the IB for continuous distributions is hard in the general case, and provided a method
62 to optimize the IB objective in the case where X, Y are known and jointly Gaussian. Painsky and
63 Tishby (2017) offered a limited linear approximation of the IB for any distribution by extracting
64 the jointly gaussian element of given distributions. Saxe et al. (2018) considered the application
65 of the IB functional as an objective for DNNs, and concluded that computing mutual information
66 in deterministic DNNs is problematic as the entropy term $H(Z|X)$ for a continuous Z is infinite.
67 Amjad and Geiger (2020) extended this observation and pointed out that for a discrete Z MI becomes
68 a piecewise constant function of its parameters, making gradient descent limiting and difficult.

69 2.2 Variational Information Bottleneck

70 Alemi et al. (2017) introduced the Variational Information Bottleneck (VIB) - a variational approxima-
71 tion for an upper bound to the IB objective for DNN optimization. Bounds for $I(X, Z)$ and $I(Z, Y)$

72 are derived from the non negativity of KL divergence, and are used to form an upper bound for the
 73 IB objective. A variational upper bound is derived by replacing intractable distributions with varia-
 74 tional approximations. Using the 'reparameterization trick' (Kingma and Welling 2014) a discrete
 75 empirical estimation of the variational upper bound is then used as a loss function for classifier DNN
 76 optimization. The subsequent loss function is equivalent to the β -autoencoder loss (Higgins et al.
 77 2017). VIB was evaluated over image classification tasks, and displayed substantial improvements in
 78 robustness to adversarial attacks, while causing a slight reduction in test set accuracy, comparing to
 79 equivalent deterministic models. (Achille and Soatto 2018) extended VIB with a total correlation
 80 term, designed to increase latent disentanglement. Fischer (2020) proposed an IB based loss func-
 81 tion named Conditional Entropy Bottleneck (CEB), in which the conditional mutual information of
 82 X and Z given Y is minimized, instead of the unconditional mutual information. The CEB loss,
 83 $L_{CEB} = \min_Z I(X; Z|Y) - \gamma I(Y; Z)$, is designed to minimize all information in Z that is not rele-
 84 vant to the downstream task Y , by conditioning over Y . Similarly to VIB, a variational approximation
 85 for CEB was proposed as $L_{VCEB} = \mathbb{E}_{x,y} \log(e(z|x)) - \mathbb{E}_{x,y} \log(b(z|y)) - \gamma \mathbb{E}_{x,y} \log(c(y|z))$ and
 86 tested over the FMNIST (Xiao, Rasul, and Vollgraf 2017) and CIFAR10 (Krizhevsky 2009) datasets.
 87 CEB was shown to be equivalent to IB for $\gamma = \beta - 1$ following the chain rule of mutual information
 88 and the IB Markov chain, as established in Appendix B. VIB was shown to be a special case of VCEB,
 89 where rate is approximated by the variational expression $\mathbb{E}_{x,y} \log(e(z|x)) - \mathbb{E}_{x,y} \log(b(z|y))$ instead
 90 of $\mathbb{E}_{x,y} \log(e(z|x)) - \mathbb{E}_{x,y} \log(r(z))$. Geiger and Fischer (2020) investigated whether VCEB is a
 91 tighter variational approximation to IB than VIB, and concluded that no ordering can be established
 92 in the general case, noting that any empirical improvement VCEB exhibits over VIB is not due to a
 93 tighter variational bound on the IB, but rather of VCEB being more amenable to optimization, or
 94 simply a successful loss function in its own regard. Cheng et al. (2020) proposed CLUB, an upper
 95 bound based MI estimator, that empirically outperformed the popular MINE estimator (Belghazi et al.
 96 2018). Since CLUB is an upper bound for MI, it was evaluated as a replacement to the upper bound
 97 for the IB rate term, $I(X; Z)$, proposed in VIB (Alemi et al. 2017). CLUB based VIB was tested it
 98 the over the MNIST dataset (Deng 2012), resulting in a slight improvement in accuracy compared to
 99 VIB, without reporting adversarial robustness. We note that CLUB does not prove a tighter bound
 100 on the VIB rate term , or on the IB objective. We also note, that our current work derives a tighter
 101 bound on IB through the IB distortion term, $I(Z; Y)$, and that combining our suggested method with
 102 a CLUB bound on rate is an interesting avenue for future work.

103 2.3 Information theoretic regularization

104 Label smoothing (Szegedy et al. 2016) and entropy regularization (Pereyra et al. 2017) both regularize
 105 classifier DNNs by increasing the entropy of their output. This is done either directly by inserting
 106 a scaled conditional entropy term to the loss function, $-\gamma \cdot H(p_\theta(y|x))$, or by smoothing the
 107 training data labels. Applying both methods was demonstrated to improve test accuracy and model
 108 calibration on various challenging classification tasks. Alemi et al. (2018) distinguishes between the
 109 IB theoretical limit for representation quality, and the limits of variational approximations in VAEs.
 110 The variational limit is looser, and depends on the chosen variational family. In the context of VAEs,
 111 the representation quality depends on choosing marginal and decoder distributions that are close to
 112 the true distributions. A common degeneration is the usage of overpowerful decoders that overfit
 113 encoder embeddings, inducing low quality latent representations. Alemi et al. (2018) claims that the
 114 VAE ELBO loss (Kingma and Welling 2014) is susceptible to this flaw, as its KL regularization term
 115 isn't directly affected by the quality of the reconstructed image, and no additional regularization is
 116 performed over the decoder directly. This work is further elaborated on in Appendix B. In the current
 117 study, a conditional entropy term (Pereyra et al. 2017) emerges from our proposed derivation of a

118 variational IB loss function, and we extend the theoretic framework proposed in (Alemi et al. 2018)
 119 to interpret why this new term facilitates a better variational approximation of the IB objective.

120 3 From VIB to VUB

121 As elaborated in Section 2.1, IB optimization is defined as $Z : \min_{P(Z|X)} I(X; Z)$ subject to $I(Z; Y) \geq$
 122 D . While this is true in the unsupervised case, adapting IB to supervised tasks admits the learned
 123 classifier as a new RV to the optimization problem. Geiger and Fischer (2020) suggested the Markov
 124 chain $Y \leftrightarrow X \rightarrow Z \rightarrow \hat{Y}$ for variational IB, distinguishing between the real unknown RV Y , and the
 125 variational approximation \hat{Y} . Following this logic, variational IB optimization should be defined as:
 126 $Z, \hat{Y} : \min_{P(Z|X), P(\hat{Y}|Z)} I(X; Z)$ subject to $I(Z; \hat{Y}) \geq D$. Previous research had relaxed the supervised
 127 optimization problem, in order to derive a tractable loss function. The VIB loss (Alemi et al. 2017)
 128 consists of a cross entropy (CE) term and a beta regulated KL regularization term, as in β -VAE loss
 129 (Higgins et al. 2017). The KL term is derived from a bound on the IB rate term $I(X; Z)$, while the
 130 CE term is derived from a bound on the IB distortion term $I(Z; Y) = H(Y) - H(Y|Z)$. When
 131 deriving the cross entropy, the term $H(Y)$ is ignored as it is assumed constant, and hence does not
 132 effect optimization, while the conditional entropy $H(Y|Z)$ is developed into the CE term. We derive
 133 a new upper bound for the IB objective by not omitting $H(Y)$ from the distortion term. Subsequently,
 134 the variational approximation for our proposed bound is a variationally tighter when taking into
 135 account that the optimization process is done over $P(\hat{Y}|Z)$ as well as $P(Z|X)$. This modification
 136 attains a tighter variational bound on the IB objective for any Y with positive entropy, and a tighter
 137 empirical bound for all Y .

138 3.1 IB upper bound

139 We begin by establishing a new upper bound for the IB functional by bounding the mutual information
 140 terms, using the same method as in VIB.

141 Consider $I(Z; X)$:

$$I(Z; X) = \int \int p^*(x, z) \log(p^*(z|x)) dx dz - \int p^*(z) \log(p^*(z)) dz \quad (\text{i})$$

142 For any probability distribution r we have that $D_{KL}(p^*(z) || r(z)) \geq 0$, it follows that:

$$\int p^*(z) \log(p^*(z)) dz \geq \int p^*(z) \log(r(z)) dz \quad (\text{ii})$$

143 And so, by Equation ii:

$$I(Z; X) \leq \int \int p^*(x) p^*(z|x) \log\left(\frac{p^*(z|x)}{r(z)}\right) dx dz \quad (\text{iii})$$

144 Consider $I(Z; Y)$:

145 For any probability distribution c we have that $D_{KL}(p^*(y|z) || c(y|z)) \geq 0$, it follows that:

$$\int p^*(y|z) \log(p^*(y|z)) dy \geq \int p^*(y|z) \log(c(y|z)) dy \quad (\text{iv})$$

¹⁴⁶ And so, by Equation iv:

$$\begin{aligned} I(Z;Y) &= \int \int p^*(y,z) \log \left(\frac{p^*(y,z)}{p^*(y)p^*(z)} \right) dydz \geq \int \int p^*(y|z)p^*(z) \log \left(\frac{c(y|z)}{p^*(y)} \right) dydz \\ &= \int \int p^*(y,z) \log(c(y|z)) dydz + H_{p^*}(Y) \end{aligned} \quad (\text{v})$$

¹⁴⁷ We now diverge from the original VIB derivation by replacing $H_{p^*}(Y)$ with $H_c(Y|Z)$ instead of
¹⁴⁸ omitting it. In addition, we limit the new term to make sure that the inequality holds.

$$I(Z;Y) \geq \int \int p^*(y,z) \log(c(y|z)) dydz + \min \{H_{p^*}(Y), H_c(Y|Z)\} \quad (\text{vi})$$

¹⁴⁹ We further develop this term using the IB Markov chain $Z \leftarrow X \leftarrow Y$ and total probability:

$$\begin{aligned} I(Z;Y) &\geq \int \int \int p^*(x)p^*(y|x)p^*(z|x) \log(c(y|z)) dx dy dz \\ &\quad + \min \left\{ H_{p^*}(Y), - \int \int c(y,z) \log(c(y|z)) dy dz \right\} \end{aligned} \quad (\text{vii})$$

¹⁵⁰ Finally, we join Equation iii with Equation vii to establish a new upper bound for the IB objective:

$$\begin{aligned} L_{IB} \leq L_{UB} &\equiv \int \int p^*(x)p^*(z|x) \log \left(\frac{p^*(z|x)}{r(z)} \right) dx dz \\ &\quad - \int \int \int p^*(x)p^*(y|x)p^*(z|x) \log(c(y|z)) dx dy dz - \min \left\{ H_{p^*}(Y), - \int \int c(y,z) \log(c(y|z)) dy dz \right\} \end{aligned} \quad (\text{viii})$$

¹⁵¹ 3.2 Variational approximation

¹⁵² Let $e(z|x)$ a variational encoder approximating $p^*(z|x)$ and $c(y|z)$ a variational classifier approxi-
¹⁵³ mating $p^*(y|z)$. We define the variational approximation L_{VUB} :

$$\begin{aligned} L_{UB} \approx L_{VUB} &\equiv \beta \int \int p^*(x)e(z|x) \log \left(\frac{e(z|x)}{r(z)} \right) dx dz \\ &\quad - \int \int \int p^*(x)p^*(y|x)e(z|x) \log(c(y|z)) dx dy dz \\ &\quad - \min \left\{ H_{p^*}(Y), - \int \int \int p^*(x)e(z|x)c(y|z) \log(c(y|z)) dx dy dz \right\} \end{aligned} \quad (\text{ix})$$

¹⁵⁴ 3.3 Empirical estimation

¹⁵⁵ The real and continuous distribution $p^*(x,y) = p^*(y|x)p^*(x)$ can be estimated by Monte Carlo
¹⁵⁶ sampling from a discrete dataset \mathcal{S} , and distributions featuring Z are sampled from a stochastic
¹⁵⁷ encoder. Let $e_\phi(z|x) \sim N(\mu, \Sigma)$ be a stochastic DNN encoder with parameters ϕ , and a final layer

158 of dimension $2K$, such that for each forward pass, the first K entries are used to encode μ , and the
 159 last K entries to encode a diagonal Σ , after a soft-plus transformation. For each $x_n \in \mathcal{S}$ we generate
 160 a sample \hat{z}_n from the encoder, using the reparameterization trick (Kingma and Welling 2014). Let
 161 C_λ be a discrete classifier neural net parameterized by λ , such that $C_\lambda(y|z) \sim \text{Categorical}$, let
 162 $H_S(Y)$ be the empirical entropy of the real RV Y as measured from the training data, and let \hat{Y} be
 163 the variational approximation for Y . We chose a standard Gaussian as a variational approximation
 164 for the marginal $r(z)$.

$$\hat{L}_{VUB} \equiv \frac{1}{N} \sum_{n=1}^N \left[\beta D_{KL} \left(e_\phi(z|x_n) \middle\| r(z) \right) - \log(C_\lambda(y_n|\hat{z}_n)) - \min \left\{ H_S(Y), H_{C_\lambda}(\hat{Y}|Z) \right\} \right]_{(x)}$$

165 3.4 Interpretation

166 The conditional entropy term that follows from our proposed derivation provides strong classifier
 167 regularization, as shown in (Pereyra et al. 2017). It was also shown that the quality of representations
 168 learned by ELBO loss functions degenerates as a consequence of overfitting decoders (Alemi et al.
 169 2018). As VIB is in fact a beta regulated ELBO loss, it follows that regulating its classifier can
 170 lead to better representations learned. We propose an additional theoretical motivation for the
 171 effects of entropy regularization on supervised learning: Following the variational IB Markov chain
 172 $Y \leftrightarrow X \rightarrow Z \rightarrow \hat{Y}$ (Geiger and Fischer 2020) we have that \hat{Y} is a deterministic function of Z , and so
 173 $H(Z) \geq H(\hat{Y})$ by the data processing inequality (Cover 1999). The following inequality holds:

$$2H(\hat{Y}) - H(\hat{Y}, Z) \stackrel{\text{DPI}}{\leq} H(\hat{Y}) + H(Z) - H(\hat{Y}, Z) = I(Z; \hat{Y}) \stackrel{(v)}{\leq} I(Z; Y) \leq \min(H(Z), H(Y))$$

174 For the left most expression, we have that any increase in $H(\hat{Y})$ will increase the joint entropy
 175 $H(\hat{Y}, Z)$ by at most the same magnitude, leading to the proportionality $H(\hat{Y}) \propto 2H(\hat{Y}) - H(\hat{Y}, Z)$.
 176 It follows that the true mutual information $I(Z; Y)$ is squeezed between $\min(H(Z), H(Y))$ and
 177 $H(\hat{Y}) \geq H(\hat{Y}|Z)$. We propose that the increase in variational entropy, that follows from VUB, can
 178 skew the true mutual information $I(Z; Y)$ up. Since the real Y is constant, this increase can only be
 179 caused by learning a better representation Z . Figure 1 illustrates the possible effect of an increase
 180 in variational entropy: The left hand diagram suggest a model with low variational entropy, and
 181 a low variational mutual information, and the right hand diagram suggest a model with increased
 182 variational entropy. The cross entropy term pushes $I(Z; \hat{Y})$ up, and the real mutual information
 183 $I(Z; Y)$ increases as a result of the Barber-Agakov inequality (Barber and Agakov 2003).

184 4 Experiments

185 We follow the experimental setup proposed by Alemi et al. (2017), extending it to NLP tasks as
 186 well. We trained image classification models on the ImageNet 2012 dataset (Deng et al. 2009), and
 187 text classification on the IMDB sentiment analysis dataset (Maas et al. 2011). For each dataset,
 188 we compared a competitive Vanilla model with a VIB and a VUB model trained with beta values
 189 of $\beta = 10^{-i}$ for $i \in \{1, 2, 3\}$. Each model was trained and evaluated 5 times per β value, with
 190 consistent performance and statistical significance shown by a Wilcoxon rank sum test. Each model
 191 was evaluated using test set accuracy, and robustness to various adversarial attacks. For image

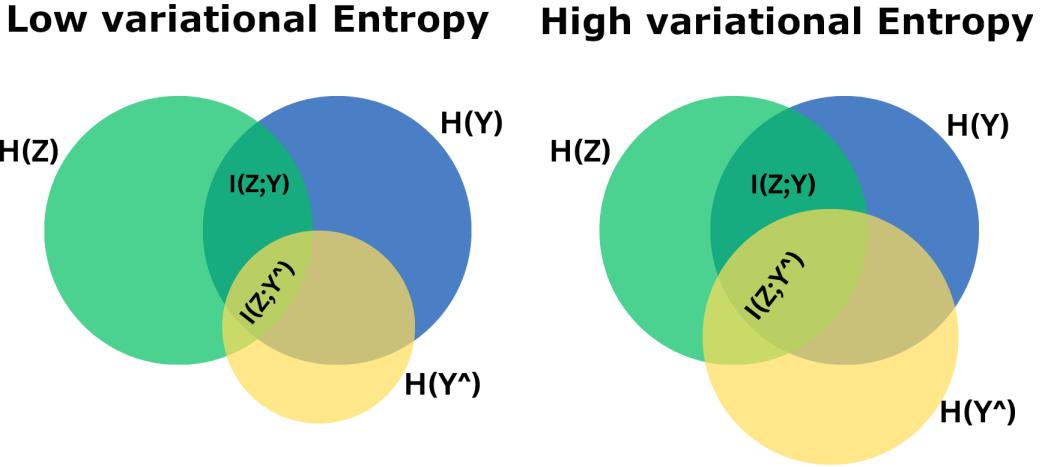


Figure 1: Venn diagrams illustrating the possible increase in real mutual information as a result of increased variational entropy. The right diagram features higher variational entropy, and higher variational mutual information that induce higher real mutual information as a result of the Barber-Agakov inequality (Barber and Agakov 2003).

classification, we employed the untargeted Fast Gradient Sign (FGS) attack (Goodfellow, Shlens, and Szegedy 2015) as well as the targeted CW L_2 attack (Carlini and Wagner 2017), (Kaiwen 2018). For text classification, we used the untargeted Deep Word Bug attack (Gao et al. 2018), (Morris et al. 2020) as well as the untargeted PWWS attack (Ren et al. 2019). Elaboration on the experimental setup, results and further insights from the experiments are available in Appendix C. Code to reconstruct the experiments is provided to the supplementary materials of this paper.

198 4.1 Image classification

199 A pre-trained inceptionV3 (Szegedy et al. 2016) base model was used and achieved a 77.21% accuracy
200 on the ImageNet 2012 validation set (Test set for ImageNet is unavailable). Image classification
201 evaluation results are shown in Table 1, examples of successful attacks are shown in Figures 5, 6 in
202 Appendix C. The empirical results presented in Table 1 confirm that while VIB reduces performance
203 on the validation set, it substantially improves robustness to adversarial attacks. Moreover, these
204 results demonstrate that VUB significantly outperforms VIB in terms of validation accuracy, while
205 providing competitive robustness to attacks similarly to VIB. A comparison of the best VIB and VUB
206 models further substantiates these findings, with statistical significance confirmed by a p-value of
207 less than 0.05 in a Wilcoxon rank sum test.

208 4.2 Text classification

209 A fine tuned BERT uncased (Devlin et al. 2019) base model was used and achieved a 93.0%
210 accuracy on the IMDB sentiment analysis test set. Text classification evaluation results are shown in
211 Table 2, examples of successful attacks are shown in Figure 4 in Appendix C. In this modality, VUB
212 significantly outperforms VIB in both test set accuracy and robustness to the two attacks. Moreover,
213 VUB also outperformed the original model in terms of test set accuracy. A comparison of the best
214 VIB and VUB models further substantiates these findings, with statistical significance confirmed by a
215 p-value of less than 0.05 in a Wilcoxon rank sum test.

β	Val \uparrow	FGS $\downarrow_{\epsilon=0.1}$	FGS $\downarrow_{\epsilon=0.5}$	CW \uparrow
Vanilla model				
-	77.2%	68.9%	67.7%	788
VIB models				
10^{-3}	73.7% $\pm .1\%$	59.5% $\pm .2\%$	63.9% $\pm .2\%$	3917 ± 291
10^{-2}	72.8% $\pm .1\%$	53.5% $\pm .2\%$	62.0% $\pm .1\%$	3318 ± 293
10^{-1}	72.1% $\pm .01\%$	58.4% $\pm .1\%$	62.0% $\pm .1\%$	3318 ± 293
VUB models				
10^{-3}	75.5% $\pm .03\%$	62.8% $\pm .1\%$	66.4% $\pm .1\%$	2666 ± 140
10^{-2}	75.0% $\pm .05\%$	57.6% $\pm .2\%$	64.3% $\pm .1\%$	1564 ± 218
10^{-1}	74.8% $\pm .09\%$	57.9% $\pm .5\%$	64.8% $\pm .5\%$	3575 ± 456

Table 1: ImageNet evaluation scores for vanilla, VIB and VUB models, average over 5 runs with standard deviation. First column is performance on the ImageNet validation set (higher is better \uparrow), second and third columns are the % of successful FGS attacks at $\epsilon = 0.1, 0.5$ (lower is better \downarrow), and the fourth column is the average L_2 distance for a successful Carlini Wagner L_2 targeted attack (higher is better \uparrow). VUB attains significantly higher accuracy over unseen data in all settings, while preserving competitive robustness to adversarial attacks.

β	Test \uparrow	DWB \downarrow	PWWS \downarrow
Vanilla model			
-	93.0%	54.3%	100%
VIB models			
10^{-3}	91.0% $\pm 1.0\%$	35.1% $\pm 4.4\%$	41.6% $\pm 6.6\%$
10^{-2}	90.8% $\pm 0.5\%$	41.0% $\pm 4.8\%$	62.9% $\pm 14.3\%$
10^{-1}	89.4% $\pm .9\%$	90.0% $\pm 8.0\%$	99.1% $\pm 0.9\%$
VUB models			
10^{-3}	93.2% $\pm .5\%$	27.5% $\pm 2.0\%$	28.4% $\pm 1.3\%$
10^{-2}	92.6% $\pm .8\%$	30.8% $\pm 2.0\%$	50.0% $\pm 4.8\%$
10^{-1}	89.2% $\pm 2.0\%$	99.2% $\pm 0.5\%$	100% $\pm 0\%$

Table 2: IMDB evaluation scores for vanilla, VIB and VUB models, average over 5 runs with standard deviation. First column is performance over the test set (higher is better \uparrow), second is % of successful Deep Word Bug attacks (lower is better \downarrow) and third is % of successful PWWS attacks (lower is better \downarrow). In almost all cases VUB attains significantly higher accuracy over unseen data, as well as significantly higher robustness to adversarial attacks. For this modality, VUB also outperforms the vanilla model in terms of test set accuracy for $\beta = 10^{-3}$.

216 **5 Discussion**

217 The IB is a private case of rate distortion, and is designed to learn representations in unsupervised
218 settings. Adapting the IB to supervised settings requires maximization of mutual information between
219 compressed representation and learned discriminator, rather than the real and unknown downstream
220 RV. This distinction requires lifting the assumption that the downstream distribution is known during
221 the derivation of a tractable objective. Lifting this assumption allows us to derive a tighter variational
222 bound on the IB, which we show to produce better classification accuracy, with equivalent or superior
223 robustness to adversarial attacks, over high dimensional tasks of different modalities. We extend
224 previous theoretical studies on representation learning (Alemi et al. 2018) to provide a possible
225 insight to the effects of our proposed derivation on the quality of the learned representation. We
226 show that, apart from the regularization described by (?), increasing classifier entropy enables the
227 variational mutual information between representation and classifier to better approximate the real
228 mutual information between representation and the real underlying distribution, bringing the attained
229 rate distortion ratio closer to its theoretical IB limit. While other advancements have been done
230 in recent years, (Fischer 2020; Cheng et al. 2020), none show a tighter bound than VIB, and all
231 modify the derivation of the rate term, while our work both derives an upper bound by modifying the
232 distortion term.

233 While providing a complete framework for optimal data modeling, the IB, and its variational ap-
234 proximations, rely on three assumptions: (1) It suffices to optimize the mutual information metric
235 to optimize a model’s performance; (2) Forgetting more information about the input, while keeping
236 relevant information about the output, induces better generalization; (3) Mutual information between
237 the input, output and latent representation can be either computed or approximated to a desired level
238 of accuracy. Our improved empirical results, induced by a tighter bound, suggest better data modeling,
239 and strengthen the cause for the IB, and its variational approximations, as a learning framework for
240 classifier DNNs. Conversely, one might argue that the improved adversarial robustness of both VIB
241 and VUB is an artifact of latent geometry: Both methods promote a disentangled latent space by
242 using a stochastic factorized prior, as suggested by Chen et al. (2018). In addition, both utilize KL
243 regularization, enforcing clustering around a 0 mean, which might increase latent smoothness. These
244 geometric traits of the latent space can make it difficult for minor perturbations to significantly alter
245 latent semantics, possibly making the models more robust to attacks.

246 In conclusion, while the IB and its variational approximations do not provide a complete theoretical
247 framework for DNN data modeling and regularization, they offer a strong, measurable, and theo-
248 retically grounded approach. VUB is presented as a tractable and tighter upper bound of the IB
249 functional, that can be easily adapted to any classifier DNN to significantly increase robustness to
250 various adversarial attacks, while inflicting minimal decrease in test set performance, and in some
251 cases even increasing it.

252 This study opens many opportunities for further research: Further improvements to the upper bound,
253 including combining VUB with proposed new bounds for the IB rate term such as CLUB (Cheng
254 et al. 2020); Applying VUB in self-supervised learning, and in particular to measure whether
255 representations learned with VUB capture better semantics than representations learned with non
256 IB inspired loss functions; Finally, comparing the effects of latent geometry vs rate distortion on
257 adversarial robustness is left to future work.

258 **References**

- 259 Achille, A.; and Soatto, S. 2018. Information dropout: Learning optimal representations through noisy
260 computation. *IEEE transactions on pattern analysis and machine intelligence*, 40(12): 2897–2905.
- 261 Alemi, A. A.; Fischer, I.; Dillon, J. V.; and Murphy, K. 2017. Deep Variational Information
262 Bottleneck. In *Proceedings of the International Conference on Learning Representations (ICLR)*.
263 Google Research.
- 264 Alemi, A. A.; Poole, B.; Fischer, I.; Dillon, J. V.; Saurous, R. A.; and Murphy, K. 2018. Fixing a
265 Broken ELBO. In *Proceedings of Machine Learning Research*, volume 80, 159–168. PMLR.
- 266 Amjad, R. A.; and Geiger, B. C. 2020. Learning Representations for Neural Network-Based
267 Classification Using the Information Bottleneck Principle. *IEEE Trans. Pattern Anal. Mach. Intell.*,
268 42(9): 2225–2239.
- 269 Barber, D.; and Agakov, F. V. 2003. The IM algorithm: a variational approach to Information
270 Maximization. In *Neural Information Processing Systems*.
- 271 Belghazi, M. I.; Baratin, A.; Rajeswar, S.; Ozair, S.; Bengio, Y.; Hjelm, R. D.; and Courville, A. C.
272 2018. Mutual Information Neural Estimation. In *International Conference on Machine Learning*.
- 273 Bengio, Y. 2009. Learning Deep Architectures for AI. *Foundations and Trends in Machine Learning*,
274 2(1): 1–127.
- 275 Carlini, N.; and Wagner, D. A. 2017. Towards Evaluating the Robustness of Neural Networks. In
276 *IEEE Symposium on Security and Privacy*, 39–57. IEEE Computer Society.
- 277 Chechik, G.; et al. 2003. Gaussian information bottleneck. In *Advances in Neural Information
278 Processing Systems*.
- 279 Chen, T. Q.; Li, X.; Grosse, R. B.; and Duvenaud, D. 2018. Isolating Sources of Disentanglement
280 in Variational Autoencoders. In *Proceedings of the 32nd International Conference on Neural
Information Processing Systems*, 2615–2625.
- 282 Cheng, P.; Hao, W.; Dai, S.; Liu, J.; Gan, Z.; and Carin, L. 2020. CLUB: A Contrastive Log-ratio
283 Upper Bound of Mutual Information. In III, H. D.; and Singh, A., eds., *Proceedings of the 37th
284 International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning
285 Research*, 1779–1788. PMLR.
- 286 Cover, T. M. 1999. *Elements of information theory*. John Wiley & Sons.
- 287 Deng, J.; Dong, W.; Socher, R.; Li, L.-J.; Li, K.; and Fei-Fei, L. 2009. Imagenet: A large-scale
288 hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*,
289 248–255. Ieee.
- 290 Deng, L. 2012. The mnist database of handwritten digit images for machine learning research. *IEEE
291 Signal Processing Magazine*, 29(6): 141–142.
- 292 Devlin, J.; Chang, M.-W.; Lee, K.; and Toutanova, K. 2019. BERT: Pre-training of Deep Bidirectional
293 Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North
294 American Chapter of the Association for Computational Linguistics: Human Language Technologies,
295 Volume 1 (Long and Short Papers)*, 4171–4186. Minneapolis, Minnesota.
- 296 Fischer, I. 2020. The Conditional Entropy Bottleneck. *Entropy*, 22(9): 999.

- 297 Gao, J.; Lanchantin, J.; Soffa, M. L.; and Qi, Y. 2018. Black-box generation of adversarial text
298 sequences to evade deep learning classifiers. In *2018 IEEE Security and Privacy Workshops (SPW)*,
299 50–56. IEEE.
- 300 Geiger, B.; and Fischer, I. 2020. A Comparison of Variational Bounds for the Information Bottleneck
301 Functional. *Entropy*.
- 302 Goldfeld, Z.; and Polyanskiy, Y. 2020. The Information Bottleneck Problem and its Applications in
303 Machine Learning. *IEEE Journal on Selected Areas in Information Theory*, 1(1): 19–38.
- 304 Goodfellow, I. J.; Shlens, J.; and Szegedy, C. 2015. Explaining and Harnessing Adversarial Examples.
305 In *ICLR (Poster)*.
- 306 Higgins, I.; Matthey, L.; Pal, A.; Burgess, C.; Glorot, X.; Botvinick, M.; Mohamed, S.; and Lerchner,
307 A. 2017. beta-VAE: Learning Basic Visual Concepts with a Constrained Variational Framework. In
308 *ICLR (Poster)*.
- 309 Kaiwen. 2018. pytorch-cw2. GitHub repository.
- 310 Kingma, D. P.; and Welling, M. 2014. Auto-Encoding Variational Bayes. In *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*.
- 313 Kingma, D. P.; and Welling, M. 2019. An Introduction to Variational Autoencoders. *Foundations and Trends in Machine Learning*, 12(4): 307–392.
- 315 Krizhevsky, A. 2009. Learning Multiple Layers of Features from Tiny Images. 32–33.
- 316 Maas, A.; Daly, R. E.; Pham, P. T.; Huang, D.; Ng, A. Y.; and Potts, C. 2011. Learning word
317 vectors for sentiment analysis. In *Proceedings of the 49th annual meeting of the association for
318 computational linguistics: Human language technologies*, 142–150.
- 319 Morris, J.; Lifland, E.; Yoo, J. Y.; Grigsby, J.; Jin, D.; and Qi, Y. 2020. TextAttack: A Framework for
320 Adversarial Attacks, Data Augmentation, and Adversarial Training in NLP. In *Proceedings of the
321 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*,
322 119–126.
- 323 Painsky, A.; and Tishby, N. 2017. Gaussian Lower Bound for the Information Bottleneck Limit. *J.
324 Mach. Learn. Res.*, 18: 213:1–213:29.
- 325 Pereyra, G.; Tucker, G.; Chorowski, J.; Kaiser, L.; and Hinton, G. E. 2017. Regularizing Neural Net-
326 works by Penalizing Confident Output Distributions. In *Proceedings of the International Conference
327 on Learning Representations*. OpenReview.net.
- 328 Ren, S.; Deng, Y.; He, K.; and Che, W. 2019. Generating natural language adversarial examples
329 through probability weighted word saliency. In *Proceedings of the 57th annual meeting of the
330 association for computational linguistics*, 1085–1097.
- 331 Saxe, A. M.; Bansal, Y.; Dapello, J.; and Advani, M. 2018. On the information bottleneck theory of
332 deep learning.
- 333 Shannon, C. E. 1959. Coding Theorems for a Discrete Source With a Fidelity Criterion. In *IRE
334 National Convention*.

- 335 Shwartz-Ziv, R.; and Tishby, N. 2017. Opening the Black Box of Deep Neural Networks via
336 Information. 19 pages, 8 figures, arXiv:arXiv:1703.00810.
- 337 Slonim, N. 2002. *The information bottleneck: Theory and applications*. Ph.D. thesis, Hebrew
338 University of Jerusalem Jerusalem, Israel.
- 339 Szegedy, C.; Vanhoucke, V.; Ioffe, S.; Shlens, J.; and Wojna, Z. 2016. Rethinking the inception
340 architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and*
341 *pattern recognition*, 2818–2826.
- 342 Tishby, N.; Pereira, F. C.; and Bialek, W. 1999. The Information Bottleneck Method. In *The*
343 *37th annual Allerton Conference on Communication, Control, and Computing*. Hebrew University,
344 Jerusalem 91904, Israel.
- 345 Tishby, N.; and Zaslavsky, N. 2015. Deep Learning and the Information Bottleneck Principle.
346 arXiv:1503.02406.
- 347 Vapnik, V. N. 1995. *The nature of statistical learning theory*. Springer-Verlag New York, Inc. ISBN
348 0-387-94559-8.
- 349 Xiao, H.; Rasul, K.; and Vollgraf, R. 2017. Fashion-MNIST: a Novel Image Dataset for Benchmarking
350 Machine Learning Algorithms. Cite arxiv:1708.07747Comment: Dataset is freely available at
351 <https://github.com/zalandoresearch/fashion-mnist> Benchmark is available at <http://fashion-mnist.s3-website.eu-central-1.amazonaws.com/>.

353 **Appendix A - Preliminaries**

354 We denote random variables (RVs) with upper cased letters X, Y , and their realizations in lower case
 355 x, y . Denote discrete Probability Mass Functions (PMFs) with an upper case $P(x)$ and continuous
 356 Probability Density Functions (PDFs) with a lower case $p(x)$. Hat notation denotes empirical
 357 measurements.

358 Let X, Y be two observed random variables with an unknown joint distribution $p^*(x, y)$ and marginals
 359 $p^*(x), p^*(y)$. We can attempt to approximate these distributions using a model p_θ with parameters
 360 θ , such that for generative tasks $p_\theta(x) \approx p^*(x)$, and for discriminative tasks $p_\theta(y|x) \approx p^*(y|x)$,
 361 using a dataset $\mathcal{S} = \{(x_1, y_1), \dots, (x_N, y_N)\}$ to fit our model. One can also assume the existence
 362 of an additional unobserved RV $Z \sim p^*(z)$ that influences or generates the observed RVs X, Y .
 363 Since Z is unobserved, it is absent from the dataset \mathcal{S} , and so cannot be modeled directly. Denote
 364 $p_\theta(x) = \int p_\theta(x|z)p_\theta(z)dz = \int p_\theta(x, z)dz$ the marginal, $p_\theta(z)$ the prior as it is not conditioned over
 365 any other RV, and $p_\theta(z|x)$ the posterior following Bayes' rule.

366 When modeling an unobserved variable of an unknown distribution, we encounter a problem as the
 367 marginal $p_\theta(x) = \int p_\theta(x, z)dz$ doesn't have an analytic solution. This intractability can be overcome
 368 by choosing some tractable parametric variational distribution $q_\phi(z|x)$ to approximate the posterior
 369 $p_\theta(z|x)$, such that $q_\phi(z|x) \approx p_\theta(z|x)$, and estimate $p_\theta(x, z)$ or $p_\theta(x, z|y)$ by fitting the dataset \mathcal{S}
 370 (Kingma and Welling 2019).

371 Vapnik (1995) defines *supervised* learning as follows:

372 • A generator of random vectors $x \in \mathbb{R}^d$, drawn independently from an unknown probability
 373 distribution $p^*(x)$.

374 • A supervisor who returns a scalar output value $y \in \mathbb{R}$, according to an unknown conditional
 375 probability distribution $p^*(y|x)$. We note that these probabilities can indeed be soft labels,
 376 where y is a continuous probability vector, rather than the more commonly used hard labels.

377 • A learning machine capable of implementing a predefined set of functions, $f(x, \theta) : \mathbb{R}^d \times$
 378 $\Theta \mapsto \mathbb{R}$, where Θ is a set of parameters.

379 The problem of learning is that of choosing from the given set of functions, the one that best approximates
 380 the supervisor's response. The choice is typically based on a training set of n independent and
 381 identically distributed pairs of observations drawn according to $p(x, y) = p(x)p(y|x) : \mathcal{S}$.

382 Given a set of unlabeled data points $\{x_1, \dots, x_n\}, x_i \in \mathbb{R}^d$, Slonim (2002) defines *unsupervised*
 383 learning as the task of constructing a compact representation of these points, which in some sense
 384 reveals their hidden structure. This representation can be used further to achieve a variety of goals,
 385 including reasoning, prediction, communication etc. In particular, unsupervised clustering partitions
 386 the data points into exhaustive and mutually exclusive clusters, where each cluster can typically be
 387 represented by a centroid, typically a weighted average of the cluster's members. Soft clustering
 388 assigns cluster probabilities for each data point, and fits an assignment by minimizing the expected
 389 loss for these probabilities, usually a distance metric such as MSE.

390 In this work, information theoretic functions share the same notation for discrete and continuous
 391 settings and are denoted as follows:

392

	Notation	Differential	Discrete
Entropy	$H_p(X)$	$-\int p(x)\log(p(x))dx$	$-\sum_{x \in X} P(x)\log(P(x))$
Conditional entropy	$H_p(X Y)$	$-\int \int p(x,y)\log(p(x y))dxdy$	$-\sum_{x \in X} \sum_{y \in Y} P(x,y)\log(P(x y))$
Cross entropy	$CE(p, q)$	$-\int p(x)\log(q(x))dx$	$-\sum_{x \in X} P(x)\log(Q(x))$
Joint entropy	$H_p(X, Y)$	$-\int \int p(x,y)\log(p(x,y))dxdy$	$-\sum_{x \in X} \sum_{y \in Y} P(x,y)\log(P(x,y))$
KL divergence	$D_{KL}(p q)$	$\int p(x)\log\left(\frac{p(x)}{q(x)}\right)dx$	$\sum_{x \in X} P(x)\log\left(\frac{P(x)}{Q(x)}\right)$
Mutual information (MI)	$I(X; Y)$	$\int \int p(x,y)\log\left(\frac{p(x,y)}{p(x)p(y)}\right)dxdy$	$\sum_{x \in X} \sum_{y \in Y} P(x,y)\log\left(\frac{P(x,y)}{P(x)P(y)}\right)$

393

394 Appendix B - Related work elaboration

395 This appendix expands the related work Section 2, by providing a deeper review of the IB, the IB
 396 theory of deep learning, and variational approximations for the IB.

397 The Information Plane

398 As mentioned in Section 2.1, the solution to the IB objective, $\mathcal{L}_{IB} = I(X; Z) - \beta I(Z; Y)$, depends
 399 on the Lagrange multiplier β . Hence, the IB objective has no one unique solution, and can thus
 400 be plotted as a function of β and of Z 's cardinality, over a Cartesian system composed of the axes
 401 $I(X; Z)$ (rate) and $I(Z; Y)$ (distortion). When plotted as a function of β the IB functional is called
 402 the 'information curve' and its Cartesian system the 'information plane' (Tishby, Pereira, and Bialek
 403 1999), as illustrated in Figure 2. When β approaches 0 the distortion term is nullified and we learn a
 404 representation that has no information over the down stream task, and maximal compression (such a
 405 representation may be a null vector). When β approaches ∞ we learn a representation that has the
 406 maximal possible information over the downstream task, but that also holds the maximal information
 407 over the training data, hence overfitted. The region above the information curve is unreachable by
 408 any possible representation. The different bifurcation of the information curve, illustrated in Figure 2,
 409 correspond to the different possible cardinalities of the compressed representation.

410 Fixing a Broken Elbo

411 As mentioned in 2.3, Alemi et al. (2018) extends the information plane with an additional theoretical
 412 bound for rate distortion ratio, imposed by the usage of finite parametric families of variational approx-
 413 imations. Alemi et al. (2018) observed that the KL regularization term in the ELBO loss function isn't
 414 directly affected by the quality of the reconstructed image, and vice versa. Following this logic, given
 415 a powerful enough decoder, optimizing the ELBO loss can result in poor latent representation, com-

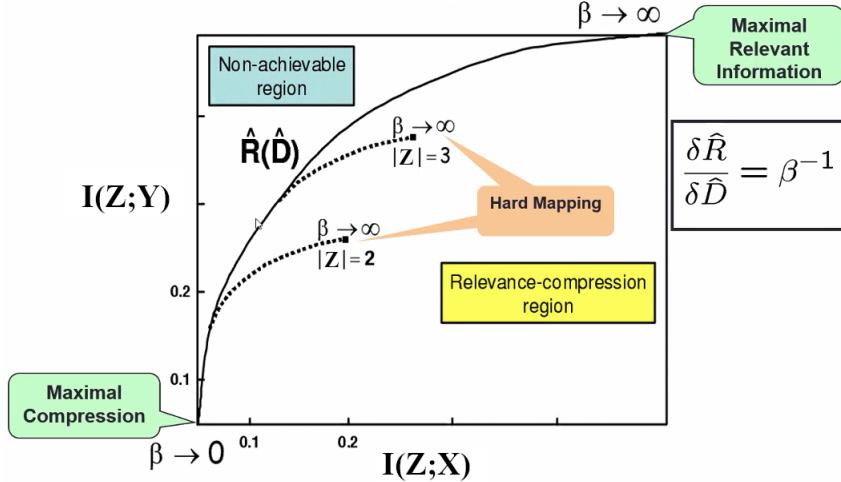


Figure 2: The information plane and curve: rate-distortion ratio over β . At $\beta = 0$ the representation is compressed but uninformative (maximal compression), at $\beta \rightarrow \infty$ the representation is informative but potentially overfitted (maximal information). Taken from (Slonim 2002).

416 compensated for by an overfitted decoder. Mutual information is used to measure representation quality:
 417 $I_e(X; Z) \equiv \int \int p_e(x, z) \log \left(\frac{p_e(x, z)}{p^*(x)p_e(z)} \right) dx dz$, for a stochastic encoder e . I_e can be bounded from
 418 both sides as follows: $H_{p^*}(X) - I_d(X; \hat{X}) \leq I_e(X; Z) \leq D_{KL}(q_\phi(z|x) || p_\theta(z))$, where $H_{p^*}(X)$
 419 is the true data entropy, I_d is the variational decoder's distortion. It is observed, that the KL term does
 420 not depend on the variational decoder distribution, and can be minimized regardless of reconstruction
 421 quality. Equivalently, good reconstruction does not directly depend on good representation. Figure 3
 422 is suggested as an extension to the information plane (Tishby, Pereira, and Bialek 1999), by replacing
 423 real rate and distortion with their variational approximations. The new plane is divided into three
 424 sub planes: (1) Infeasible: This is the IB theoretical limit (As per Figure 2); (2) Feasible: Attainable
 425 given an infinite model family, and complete variety of: $e(z|x)$, $d(x|z)$ and $p_\theta(z)$; (3) Realizable:
 426 Attainable given a finite parametric and tractable variational family.

 427 The black diagonal line at the lower left of Figure 3 satisfies $H_{p^*}(X) - I_d(X; \hat{X}) = D_{KL}(q_\phi(z|x) || p_\theta(z))$, resulting in tight variational bounds on the mutual information: $R \leq I_e \leq R$.
 428 (Alemi et al. 2018) notes that a common degeneration of representations in VAEs is a decrease in rate
 429 and distortion over the training data, together with an increase in both over unseen data. This is caused
 430 by overpowered decoders that overfit an uninformative learned representation. Low ELBO loss can be
 431 attained as both the distortion and rate terms are minimized during training. $D_{KL}(q_\phi(z|x) || p_\theta(z))$
 432 approaches 0 iff $e(z|x) \rightarrow p_\theta(z)$. In this case, $e(z|x)$ is close to independence of x , and the latent
 433 representation fails to encode information about the input. However, a suitably powerful decoder
 434 could non the less learn to overfit encoded traces of the training examples, and reach a low distortion
 435 score during optimization. In the current study, we extend this theoretical framework to explain the
 436 advancements of our proposed loss function.

438 IB theory of deep learning

439 The following is a summary of work leveraging the IB framework for deterministic DNN optimization
 440 and interpretation. For a more comprehensive review of this opinion-splitting topic, the reader is
 441 advised to consult the work of Goldfeld and Polyanskiy (2020).

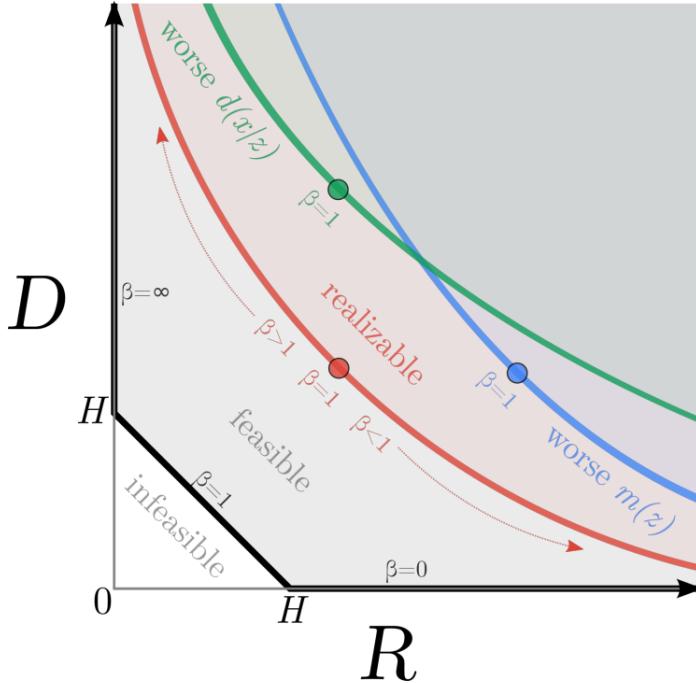


Figure 3: Phase diagram, a proposed information plane interpretation of VAEs. Axes are variational rate and distortion. The IB theoretical limit is extended by an additional limit induced by the constraint of a finite parametric variational family. Once a family is chosen, we seek to learn an optimal marginal $m(z)$ and decoder $d(x|z)$ in order to approach the new limit. β modulation controls the tradeoff between rate and distortion, regardless of the variational family. Note that this figure is inverted in orientation to Figure 2, i.e. low distortion corresponds to better performance, and not to lower MI. Taken from (Alemi et al. 2018).

442 Tishby and Zaslavsky (2015) proposed a representation-learning interpretation of DNNs using the IB
 443 framework, regarding DNNs as Markov cascades of intermediate representations between hidden
 444 layers. Under this notion, comparing the optimal and the achieved rate-distortion ratios between DNN
 445 layers will indicate if a model is too complex or too simple for a given task and training set. Shwartz-
 446 Ziv and Tishby (2017) visualized and analyzed the information plane behavior of DNNs over a toy
 447 problem with a known joint distribution. Mutual information of the different layers was estimated
 448 and used to analyze the training process. The learning process over Stochastic Gradient Descent
 449 (SGD) exhibited two separate and sequential behaviors: A short Empirical Error Minimization phase
 450 (ERM) characterized by a rapid decrease in distortion, followed by a long compression phase with an
 451 increase in rate until convergence to an optimal IB limit as demonstrated in Figure 4. Similar yet
 452 repetitive behavior was observed in the current study, as elaborated in Section ??.

453 Saxe et al. (2018) reproduced the experiments described in (Shwartz-Ziv and Tishby 2017), expanding
 454 them to different activation functions, different datasets and different methods to estimate mutual
 455 information. It was found that double-sided saturated nonlinear activations, such as the tanh, produced
 456 a distinct compression stage when mutual information was measured by binning, as performed
 457 in (Shwartz-Ziv and Tishby 2017), while other activations did not. It was also shown that DNN
 458 generalization did not depend on a distinct compression stage, and that DNNs do forget task irrelevant
 459 information, but this happens concurrently to the learning of task relevant information, and not
 460 necessarily in a distinct phase. Amjad and Geiger (2020) also showed that equivalent representations
 461 might yield the same IB loss while one achieved better classification rate than the other. These

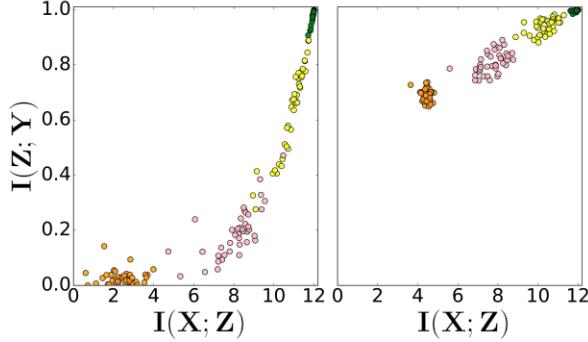


Figure 4: Information plane scatters of different DNN layers (colors) in 50 randomized networks. From Shwartz-Ziv and Tishby (2017). Left are initial weights, Right are at 400 epochs. Our study reproduced similar yet repetitive behavior on complicated high dimensional tasks, as elaborated in Section ?? and in Figure 7.

462 discrepancies are caused because mutual information in deterministic DNNs is either infinite or step
 463 like, because of mutual information's invariance to invertible transformations and because of the
 464 absence of a decision function in the objective.
 465 When examining the information plane behavior in the current study, we notice recurring patterns of
 466 distortion reduction followed by rate increase, resembling the ERM and representation compression
 467 stages described by Shwartz-Ziv and Tishby (2017), as elaborated in Appendix 5.

468 **Conditional Entropy Bottleneck**

469 As mentioned in Section 2.2, Fischer (2020) showed that the conditional entropy bottleneck is
 470 equivalent to IB for $\gamma = \beta - 1$ following the chain rule of mutual information and the IB Markov
 471 chain. We develop this equivalence in detail:

$$\begin{aligned}
 CEB &= I(X; Z|Y) - \gamma I(Z; Y) \\
 &\stackrel{\text{MI chain rule}}{=} H(Z|Y) - H(Z|X, Y) - \gamma I(Z; Y) \\
 &\stackrel{Z \leftarrow X \leftrightarrow Y}{=} H(Z|Y) - H(Z|X) - \gamma I(Z; Y) \\
 &\stackrel{\gamma := \beta - 1}{\implies} H(Z|Y) - H(Z|X) - (\beta - 1)I(Z; Y) \\
 &= H(Z|Y) - H(Z|X) - \beta I(Z; Y) + I(Z; Y) \\
 &= H(Z|Y) - H(Z|X) - \beta I(Z; Y) + H(Z) - H(Z|Y) \\
 &= H(Z) - H(Z|X) + H(Z|Y) - H(Z|Y) - \beta I(Z; Y) \\
 &= I(X; Z) - \beta I(Z; Y)
 \end{aligned}$$

472 **Appendix C - Experiments elaboration**

473 Image classification models were trained on the first 500,000 samples of the ImageNet 2012 dataset
 474 (Deng et al. 2009), and text classification over the entire IMDB sentiment analysis dataset (Maas
 475 et al. 2011). For each dataset, a competitive pre-trained model (Vanilla model) was evaluated and
 476 then used to encode embeddings. These embeddings were then used as a dataset for a new stochastic
 477 classifier net with either a VIB or a VUB loss function. Stochastic classifiers consisted of two ReLU

478 activated linear layers of the same dimensions as the pre-trained model's logits (2048 for image
 479 and 768 for text classification), followed by reparameterization and a final softmax activated FC
 480 layer. Learning rate was 10^{-4} and decaying exponentially with a factor of 0.97 every two epochs.
 481 Batch sizes were 32 for ImageNet and 16 for IMDB. All models were trained using an Nvidia
 482 RTX3080 GPU with approximately 1-2 days per a single experiment run. Beta values of $\beta = 10^{-i}$
 483 for $i \in \{1, 2, 3\}$ were tested, and we used a single forward pass per sample for inference, since
 484 previous studies indicated that these are the best range and sample rate for VIB (Alemi et al. 2017,
 485 2018). Each model was trained and evaluated 5 times per β value, with consistent performance.
 486 Statistical significance was demonstrated in all comparisons using the Wilcoxon rank sum test as
 487 follows: For each compared metric a sorted vector of results was prepared, where each entry featured
 488 the attained result in each of the 5 i.i.d. experiments per algorithm, and a boolean indicator value
 489 for the algorithm type. All metrics compared attained a p-value of less than 0.05. For example:
 490 $r = ((0.94, 1) (0.935, 1) (0.93, 1) (0.93, 1) (0.925, 1) (0.92, 0) (0.915, 0) (0.915, 0) (0.91, 0) (0.89, 0))^{-1}$
 491 be a sorted vector of (test accuracy, algorithm) tuples, 1 being VUB, 0 VIB. We computed the
 492 rank-sum as follows:
 493

$$\mu_T = \frac{5 \cdot 11}{2} = 27.5, \sigma_T = \sqrt{\frac{5 \cdot 5 \cdot 11}{12}} \approx 4.78, Z(T) = \frac{15 - 27.5}{4.78} \approx -2.61$$

$$\Phi^{-1}(pval) = -2.61, pval = 0.0045 \leq 0.05$$

493 In practice, these were computed with the Python Scipy library:
 494

```

494     import scipy.stats as stats
495     vib_scores = [0.915, 0.915, 0.91, 0.92, 0.89]
496     vub_scores = [0.93, 0.935, 0.925, 0.93, 0.94]
497     pvalue = stats.ranksums(vub_scores, vib_scores, 'greater').pvalue
498     assert pvalue < 0.05
  
```

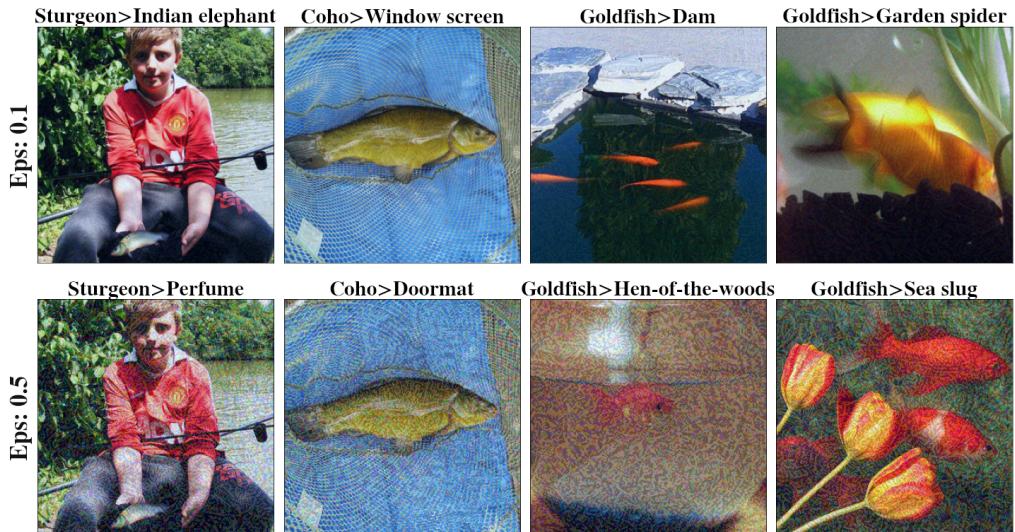
499 **Image classification**

500 The ImageNet 2012 validation set was used for evaluation as the test set for ImageNet is unavailable.
 501 InceptionV3 yields a slightly worse single shot accuracy than inceptionV2 (80.4%) when run in
 502 a single model and single crop setting, however we've used InceptionV3 over V2 for simplicity.
 503 Each model was trained for 100 epochs. The entire validation set was used to measure accuracy and
 504 robustness to FGS attacks, while only 1% of it was used for CW attacks, as they are computationally
 505 expensive.

506 **Text classification**

507 Each model was trained for 150 epochs. The entire test set was used to measure accuracy, while
 508 only the first 200 entries in the test set were used for adversarial attacks as they are computationally
 509 expensive.

Untargeted FGS attacks for VIB $\beta=0.01$



Untargeted FGS attacks for VUB $\beta=0.01$

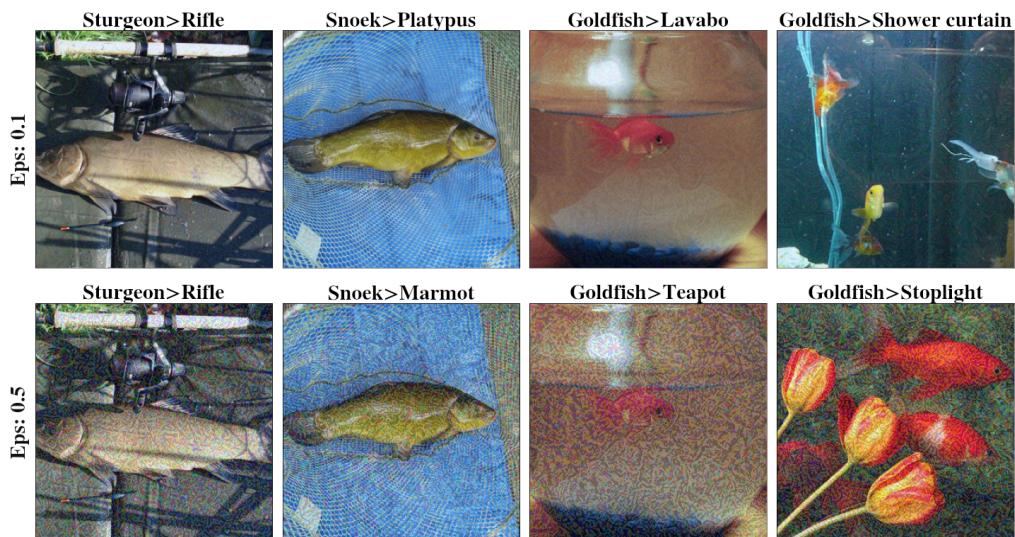


Figure 5: Successful untargeted FGS attack examples. Images are perturbations of previously successfully classified instances from the ImageNet validation set. Perturbation magnitude is determined by the parameter ϵ shown on the left, the higher, the more perturbed. Original and wrongly assigned labels are listed at the top of each image. Notice the deterioration of image quality as ϵ increases.

Targeted CW attacks for VIB $\beta=0.01$. Target: Soccer ball



Targeted CW attacks for VUB $\beta=0.01$. Target: Soccer ball



Figure 6: Successful targeted CW attack examples. Images are perturbations of previously successfully classified instances from the ImageNet validation set. The target label is 'Soccer ball'. Average L_2 distance required for a successful attack is shown on the left. The higher the required L_2 distance, the greater the visible change required to fool the model. Original and wrongly assigned labels are listed at the top of each image. Mind the difference in noticeable change as compared to the FGS perturbations presented in Figure 5.

Original text
the acting , costumes , music , cinematography and sound are all <i>astounding</i> given the production 's austere locales .
Perturbed text
the acting , costumes , music , cinematography and sound are all <i>dumbfounding</i> given the production 's austere locales .

Table 3: Example of a successful PWWS attack on a vanilla Bert model, fine tuned over the IMDB dataset. The original label is 'Positive sentiment'. The substituted word, marked in italic font, changed the classification to 'Negative sentiment'. VUB and VIB classifiers are far less susceptible to these perturbations as shown in Table 2.

510 In addition to the above evaluation metrics, we also measured approximated rate and distortion
 511 throughout training, and plotted them on the information plane as shown in Figure 7 in Appendix 5.
 512 Examining the resulting curve, we notice recurring patterns of distortion reduction followed by rate
 513 increase, resembling the ERM and representation compression stages described by Shwartz-Ziv and
 514 Tishby (2017). While previous research has documented the occurrence of error minimization and
 515 representation compression phases (Shwartz-Ziv and Tishby 2017), our work revealed that these
 516 phases can occur in cycles throughout training. This finding is particularly noteworthy, because
 517 previous studies observed this phenomenon in simple toy problems, whereas our research demon-
 518 strated it in complex tasks of high dimensionality, with unknown distributions. This suggests that this

Original text
<i>great</i> historical movie, will not allow a viewer to leave once you begin to watch. View is presented differently than displayed by most school books on this <i>subject</i> . My only fault for this movie is it was photographed in black and white; wished it had been in color ... wow !
Perturbed text
<i>gnreat</i> historical movie, will not allow a viewer to leave once you begin to watch. View is presented differently than displayed by most school books on this <i>sSbject</i> . My only fault for this movie is it was photographed in black and white; wished it had been in color ... wow !

Table 4: Example of a successful Deep Word Bug attack on a vanilla Bert model, fine tuned over the IMDB dataset. The original label is 'Positive sentiment'. Perturbations, marked in italic font, change the classification to 'Negative sentiment'. VUB and VIB classifiers are far less susceptible to these perturbations, as shown in Table 2.

Estimated information plane

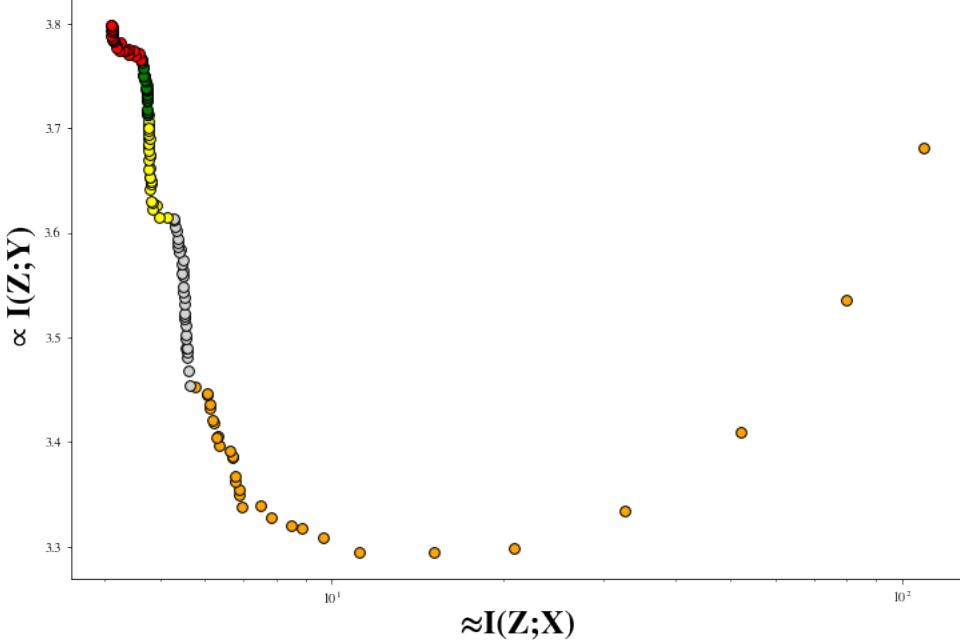


Figure 7: Estimated information plane metrics per epoch for VUB trained on IMDB with $\beta = 0.001$. $I(X; Z)$ is approximated by $H(R) - H(Z|X)$ and $\frac{1}{CE(Y; \hat{Y})}$ is used as an analog for $I(Z; Y)$. The epochs have been grouped and color-coded in intervals of 30 epochs in the order: Orange (0-30), gray (30-60), yellow (60-90), green (90-120) and red (120-150). We notice recurring patterns of distortion reduction followed by rate increase, resembling the ERM and representation compression stages described by Shwartz-Ziv and Tishby (2017).

519 information plane behavior is not limited to simplified scenarios, but is a possible characteristic of
520 the learning process in more challenging tasks as well.