REVISITING THE VARIATIONAL INFORMATION BOT TLENECK

Anonymous authors

Paper under double-blind review

Abstract

The Information Bottleneck (IB) framework offers a theoretically optimal approach to data modeling, though it is often intractable. Recent efforts have optimized supervised deep neural networks (DNNs) using a variational upper bound on the IB objective, leading to enhanced robustness to adversarial attacks. In these studies, supervision assumes a dual role: sometimes as a presumably constant and observed random variable, and at other times as its variational approximation. This work proposes an extension to the IB framework, and consequently to the derivation of its variational bound, that resolves this duality. Applying the resulting bound as an objective for supervised DNNs induces significant empirical improvements, and provides an information theoretic motivation for decoder regularization.

025 026 027

028

004

006

011

013

014

015

016

017 018

019

021

023

1 INTRODUCTION

029 The Variational Information Bottleneck, VIB, (Alemi et al. 2017) adapts the theoretically optimal¹, yet mostly intractable, Information Bottleneck, IB, (Tishby et al. 1999) to supervised DNNs. However, the IB is a method for unsupervised learning, and requires knowledge 032 of the underlying joint distribution p(x, y) (Slonim 2002). This requirement is relaxed in 033 the original VIB derivation, resulting in a duality in the usage of the downstream RV Y, 034 which is treated both as an observed RV when sampled from the training data, and as a variational approximation when optimized over. This work proposes a new adaptation of 037 the IB and VIB frameworks for supervised tasks, and consequently an information-theoretic 038 motivation for decoder regularization. 039

040 We begin by laying down what IB is, and how it can be adapted to DNNs. Classic 041 information theory provides rate-distortion (Shannon 1959) for optimal compression of data. 042 However, rate-distortion regards all information as equal, not taking into account which 043 information is more relevant to a specified downstream task, without constructing tailored 044 distortion functions. The Information Bottleneck (IB) (Tishby et al. 1999) resolves this 045 limitation by defining mutual information (MI) between the learned representation and 046 a designated downstream random variable (RV) as a universal distortion function. Yet, 047 048 learning representations using the IB method is possible given discrete distributions, and 049 some continuous ones, but not in the general case (Chechik et al. 2003). Moreover, MI is either difficult or impossible to optimize over when considering deterministic models, such as

¹Optimal data modeling with the IB method is established under the assumption that optimizing a precision-complexity trade-off will yield a model that is closer in nature to the real underlying process, and that mutual information is a sufficient metric for this purpose (Slonim 2002).

054 DNNs (Saxe et al. 2018; Amjad & Geiger 2020). Nonetheless, the promise of the IB remains alluring, and recent efforts utilized VAE (Kingma & Welling 2014) inspired variational methods to approximate upper bounds on the IB objective, allowing its utilization as a loss 057 function for DNNs, where the underlying distributions are both continuous and unknown 058 (Alemi et al. 2017; Fischer 2020; Cheng et al. 2020). These approaches learn representations 060 in supervised settings, without knowledge of the underlying distribution p(x, y), utilizing the learned variational conditional p(y|x) to approximate MI. In contrast, non variational 062 IB methods learn representations in unsupervised settings, where the stochastic process 063 underlying the observed data is known (Tishby et al. 1999; Chechik et al. 2003; Painsky & 064 Tishby 2017). Nonetheless, when deriving the variational IB objectives, previous research 065 considered the learned representation as the only optimized RV, when in practice a variational 066 classifier is also optimized. This work proposes a modification of the IB and variational 067 068 IB objectives, by setting the downstream RV as a parameterized model in the problem 069 definition. We believe our modification is a better adaptation of the IB for supervised tasks, and show empirical evidence of improved performance across several challenging tasks over 071 different modalities. Finally, we use our findings to propose a novel information theoretic 072 interpretation of overfitting in supervised DNNs. 073

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

077 078

079 080

081

2 Related work

2.1 Deterministic Information Bottleneck

082 Classic information theory offers rate-distortion (Shannon 1959) to mitigate signal loss during compression: A source X is compressed to an encoding Z, such that maximal compression is 084 achieved while keeping the encoding quality above a certain threshold. Encoding quality is 085 measured by a task specific distortion function: $d: X \times Z \mapsto \mathbb{R}^+$. Rate-distortion suggests a 086 mapping that minimizes the rate of bits to source sample, measured by I(X; Z), that adheres 087 to a chosen allowed expected distortion $D \ge 0$. The Information Bottleneck (IB) (Tishby 880 et al. 1999) extends rate-distortion by replacing the tailored distortion functions with MI 090 over a target distribution: Let Y be the target signal for some specific downstream task, such 091 that the joint distribution p(x, y) is known, and define the distortion function as MI between 092 Z and Y. The IB is the solution to the optimization problem $Z: \min_{\substack{p(z|x)}} I(X;Z)$ subject to 093 $I(Z;Y) \geq D$, that can be optimized by minimizing the IB objective $\mathcal{L}_{IB} = \beta I(X;Z) - I(Z;Y)$ 094 over p(z|x). The solution to this objective is a function of the Lagrange multiplier β , and 096 is a theoretical limit for representation quality, given mutual information as an accepted 097 metric, as elaborated in more detail in Appendix B. The IB is in fact an unsupervised soft 098 clustering problem, where each data point x is assigned a probability z to belong to different 099 clusters, given the joint distribution of the input and target tasks p(x, y) (Slonim 2002). 100 Chechik et al. (2003) showed that computing the IB for continuous distributions is hard 101 in the general case, and provided a method to optimize the IB objective in the case where 102 103 X, Y are jointly Gaussian and known. Painsky & Tishby (2017) offered a limited linear 104 approximation of the IB for any distribution by extracting the jointly Gaussian element of 105 given distributions. Saxe et al. (2018) considered the application of the IB objective as a 106 loss function for DNNs, and concluded that computing mutual information in deterministic 107 DNNs is problematic as the entropy term H(Z|X) for a continuous Z is infinite. Amjad &

Geiger (2020) extended this observation and pointed out that for a discrete Z MI becomes a piecewise constant function of its parameters, making gradient descent limited and difficult.

111 Considering the supervised problem, Geiger & Fischer (2020) suggested to consider the 112 classification output as an additional random variable, leading to an extended Markov chain 113 underlying the problem: $Y \leftrightarrow X \leftrightarrow Z \leftrightarrow \tilde{Y}$. A similar approach has also been suggested by 114 Piran et al. (2020) where a dual IB formulation was suggested that although still considers 115 the minimization of I(X;Z) replaces the constraints to one that takes into account \tilde{Y} . The 117 approach suggested here follows these ideas, but adds the additional objective of reducing 118 overfitting during the classification step.

119 120

121

2.2 VARIATIONAL INFORMATION BOTTLENECK

122 Alemi et al. (2017) introduced the Variational Information Bottleneck (VIB) - a variational 123 approximation for an upper bound to the IB objective for DNN optimization. Bounds for 124 I(X,Z) and I(Z,Y) are derived from the non negativity of KL divergence, and are used to 125 126 form an upper bound for the IB objective. Variational approximations are then used to replace 127 intractable distributions in the upper bound. Using the reparameterization trick (Kingma & 128 Welling 2014), a discrete empirical estimation of the variational upper bound is used as a 129 loss function for classifier DNN optimization, resulting in a loss function that is equivalent 130 to the β -autoencoder loss (Higgins et al. 2017). VIB was evaluated over image classification 131 tasks, and displayed substantial improvements in robustness to adversarial attacks, while 132 inflicting a slight reduction in test set accuracy, when compared to equivalent deterministic 133 134 models. The improved robustness is attributed to an improvement in representation quality, 135 and subsequently better generalization. Achille & Soatto (2018) extended VIB with a total 136 correlation term, designed to increase latent disentanglement. Fischer (2020) proposed an IB 137 based loss function named Conditional Entropy Bottleneck (CEB), in which the conditional 138 mutual information of X and Z given Y is minimized, instead of the unconditional mutual 139 information. The CEB loss, $L_{CEB} = \min_{Z} I(X; Z|Y) - \gamma I(Y; Z)$, is designed to minimize all 140 141 information in Z that is not relevant to the downstream task Y, by conditioning over Y. 142 CEB is equivalent to IB for $\gamma = \beta - 1$ following the chain rule of mutual information (Cover 143 1999) and the IB Markov chain, as established in Appendix B. However, its variational 144 approximation, VCEB, differs from VIB in the way the marginal is approximated. Geiger & 145 Fischer (2020) showed that VCEB is a tighter variational approximation for IB under certain 146 conditions, but not in the general case. Later work (Fischer & Alemi 2020) evaluated VCEB 147 on the ImageNet-A and ImageNet-C datasets, two flavors of ImageNet (Deng et al. 2009) that 148 149 assess model performance on challenging edge cases and robustness to common corruptions, 150 respectively. Results showed improved generalization, calibration, and robustness to the 151 targeted PGD (Madry et al. 2018) attack, in particular when model size was increased. 152

153 154

155

2.3 INFORMATION THEORETIC REGULARIZATION

Label smoothing (Szegedy et al. 2016) and entropy regularization (Pereyra et al. 2017)
regularize classifier DNNs by increasing classifier entropy, either by inserting a scaled
conditional entropy term to the objective, or by smoothing the training labels. Applying
either methods improved test accuracy and model calibration on various challenging tasks.
Alemi et al. (2018) extended the information plane (Tishby et al. 1999) to VAE (Kingma &
Welling 2014) settings, measuring distortion as MI between input and reconstructed images,

and rate as KL divergence between variational representation and marginal. The limits of 163 representation quality in VAEs are looser than the theoretical IB limits, and heavily depend 164 on the chosen variational families of the marginal and decoder distributions. The closer 165 the families are to the true distributions, the tighter the gap to the theoretical IB limit 166 for representation quality. Alemi et al. (2018) also showed that the ELBO loss is prone 167 168 to produce low quality representations: provided a strong enough decoder, the ELBO KL 169 regularization term might induce completely uninformative representations, that are then 170 overfitted by the powerful decoder, as elaborated in detail in Appendix B. In the current 171 study, a conditional entropy term (Pereyra et al. 2017) emerges during the derivation of our 172 proposed adaptation of the IB objective, providing a possible remedy to the discrepancies in 173 the ELBO loss, and subsequently VIB and VCEB loss, described in (Alemi et al. 2018). 174

175 176

177 178 179

180

181 182

183

FROM VIB TO SVIB 3

3.1PROBLEM DEFINITION

As elaborated in Section 2.1, the IB objective, $\mathcal{L}_{IB} = I(X;Z) - \beta I(Z;Y)$, is computed over the joint distribution p(x, y, z). When p(x, y) is given, this expression is optimized over the distribution p(z|x, y), as proposed by Tishby et al. (1999):

 $\min_{p(z|x,y)} I(Z;X)$

187

188 189

190

191

192

193

194

195

196

197

198 199

200

201

202

203

204

205 206

207

208

209

210

s.t. $I(Z;Y) \ge D_1$ (1)However, as mentioned, adapting IB to supervised tasks admits the learned classifier as a new RV to the optimization problem (Geiger & Fischer 2020; Piran et al. 2020). Thus, we consider the extended Markov chain $Y \leftrightarrow X \leftrightarrow Z \leftrightarrow \tilde{Y}$ for supervised IB, distinguishing between the true unknown RV Y, and the learned classifier \tilde{Y} . We follow this approach, and also assume that \tilde{Y} and Y share the same support. The IB framework connects the underlying joint distribution of the input and objective data, p(x, y), with a learned representation Z. We claim that when applying IB to supervised tasks, one must also consider the connection to the classifier defined by the output RV \tilde{Y} . Thus, we also want to consider the joint distribution over the pair Z, \tilde{Y} during optimization. Following the IB method logic, we seek a \tilde{Y} that will minimize mutual information with Z, whilst keeping below a defined distortion metric with the true Y. That is, we seek a second bottleneck that minimizes passage of information between Z and \tilde{Y} , so as to limit it to the minimum required to ensure that \tilde{Y} is similar enough to Y, given the transition through both X and Z. Since in this case we are optimizing over the joint conditional distribution $p(z, \tilde{y}|x, y)$, instead of the conditional $p(\tilde{y}|z, x, y)$, this problem is not simply an IB problem over the Markov chain $Y \leftrightarrow Z \leftrightarrow \tilde{Y}$. Moreover, contrary to the standard IB, X plays a significant role, controlling the distribution of Z, and the entire chain of four random variables must be taken into consideration. We thus define a second bottleneck for the true distribution p(x, y) and modeled distribution $c(\tilde{y}|z)p(z|x,y)$. We choose KL divergence as a distortion metric, as we assume Y and Y share the same support. For some positive scalar D_2 we have:

211 212

$$\min_{c(\tilde{y}|z)p(z|x,y)} I(Z; Y)$$

214
215 s.t.
$$D_{KL}\left(p(y=\tilde{y}|z,x)\middle|\middle|c(\tilde{y}|z)p(z|x)\right) \le D_2$$

(2)

Combining the two bottlenecks results in a new optimization problem, which we denote
 Supervised Information Bottleneck (SIB), which minimizes the following objective:

$$\mathcal{L}_{SIB} \equiv \beta I(X;Z) - I(Z;Y) + \lambda I(Z;\tilde{Y}) + D_{KL} \left(p(y=\tilde{y}|z,x) \middle\| c(\tilde{y}|z)p(z|x) \right)$$
(3)

3.2 Optimization Objective

We proceed to derive a tractable variational upper bound for \mathcal{L}_{SIB} , which we can use as an objective function for classifier DNNs. We begin by deriving the first bottleneck (1) as done in VIB (Alemi et al. 2017), and proceed to derive the second (2).

Consider I(Z; X):

$$I(Z;X) = \int \int p(x,z) \log \left(p(z|x) \right) \mathrm{d}x \,\mathrm{d}z - \int p(z) \log \left(p(z) \right) \mathrm{d}z \tag{4}$$

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

$$\int p(z) \log (p(z)) \, \mathrm{d}z \ge \int p(z) \log (r(z)) \, \mathrm{d}z \tag{5}$$

And so, by Equation 5:

$$I(Z;X) \le \int \int p(x)p(z|x) \log\left(\frac{p(z|x)}{r(z)}\right) dx dz$$
(6)

Consider I(Z;Y):

From the Barber-Agakov inequality (Barber & Agakov 2003), we have that for any probability distribution c:

$$I(Z;Y) \ge \int \int p(y,z) \log \left(c(y|z) \right) \mathrm{d}y \, \mathrm{d}z - \int p(y) \log \left(p(y) \right) \mathrm{d}y \tag{7}$$

Note that Equations 6 and 7 hold for any distribution r over the support of Z, and for any conditional distribution $c(\cdot|z)$ whose support equals the support of Y for every given value z in the support of Z. We link the two bottlenecks by choosing c to be $\tilde{Y}|Z = z \sim c(\cdot|z)$, meaning the variational classifier distribution. This connection is implicit in (Alemi et al. 2017), where \tilde{Y} is not formally defined. We now move on to the second bottleneck.

Consider $I(Z; \tilde{Y})$:

$$I(Z;\tilde{Y}) = H(\tilde{Y}) - H(\tilde{Y}|Z)$$
(8)

Choosing a discrete random variable for \tilde{Y} , as in labeled classification, we have $H(\tilde{Y}) \leq \log \| \tilde{\mathcal{Y}} \|$. Otherwise, choosing a continuous RV with finite support [a, b], we have that $H(\tilde{Y}) \leq \log(b-a)$. In both cases, $I(Z; \tilde{Y})$ is bounded from above by some constant

 $J = \log(b-a)$, or $J = \log || \tilde{Y} ||$, and the negative conditional entropy term $-H(\tilde{Y}|Z)$:

$$I(Z; \tilde{Y}) \le J - H(\tilde{Y}|Z) = J + \int \int p(\tilde{y}, z) \log\left(c(\tilde{y}|z)\right) d\tilde{y} dz \tag{9}$$

 Consider $D_{KL}\left(p(y=\tilde{y}|z,x) \middle| \middle| c(\tilde{y}|z)p(z|x)\right)$:

$$D_{KL}\left(p(y=\tilde{y}|z,x)\bigg|\bigg|c(\tilde{y}|z)p(z|x)\right) =$$
(10)

$$\int \int \int p(y,z,x) \log \left(p(y|z,x) \right) \, \mathrm{d}y \, \mathrm{d}x \, \mathrm{d}z - \int \int \int p(y,z,x) \log \left(c(y|z,x) \right) \, \mathrm{d}y \, \mathrm{d}x \, \mathrm{d}z \quad (11)$$

Applying the Markov chain $Y \leftrightarrow X \leftrightarrow Z \leftrightarrow \tilde{Y}$, and total probability, we get:

$$D_{KL}\left(p(y=\tilde{y}|z,x)\bigg|\bigg|c(\tilde{y}|z)p(z|x)\right) = \int \int p(y,x)\log\left(p(y|x)\right)\,\mathrm{d}y\,\mathrm{d}x - \int \int p(y,z)\log\left(c(y|z)\right)\,\mathrm{d}y\,\mathrm{d}z \quad (12)$$

Finally, we attain an upper bound for \mathcal{L}_{SIB} by combining Equations (6,7,9,12):

$$\mathcal{L}_{SIB} \leq \beta \int \int p(x)p(z|x) \log\left(\frac{p(z|x)}{r(z)}\right) dx dz - 2 \int \int p(y,z) \log\left(c(y|z)\right) dy dz +\lambda \int \int c(y|z)p(z) \log\left(c(y|z)\right) dy dz + \int \int p(y,x) \log\left(p(y|x)\right) dy dx + \int p(y) \log\left(p(y)\right) dy + \lambda J$$
(13)

Note that p(x, y) and J are constants, and so the last three terms in Equation (13) can be ignored in the course of optimization.

3.3 VARIATIONAL APPROXIMATION AND EMPIRICAL ESTIMATION

We further develop the upper bound in Equation (13) using the IB Markov chain $Y \leftrightarrow X \leftrightarrow Z \leftrightarrow \tilde{Y}$ and total probability, and define tractable variational distributions to replace intractable ones. Let e(z|x) a variational encoder approximating the conditional p(z|x), let r(z) be a variational approximation for the marginal, and let let c(y|z) a variational classifier approximating p(y|z). We define the variational approximation L_{SVIB} :

$$\mathcal{L}_{SVIB} \equiv \beta \int \int p(x)e(z|x) \log\left(\frac{e(z|x)}{r(z)}\right) dx dz$$

$$-2 \int \int \int p(x)p(y|x)e(z|x) \log\left(c(y|z)\right) dx dy dz$$

$$+\lambda \int \int \int p(x)e(z|x)c(y|z) \log\left(c(y|z)\right) dx dy dz$$
(14)

324 As is common in VIB and VAE literature, we chose a standard Gaussian for the variational 325 marginal r(z), a spherical Gaussian for the variational encoder e(z|x), and a categorical 326 distribution for the variational classifier c(y|z). We use DNNs to model these distributions 327 as follows: Let $e_{\phi}(z|x) \sim N(\mu, \Sigma)$ be a stochastic DNN encoder with parameters ϕ , and a 328 final layer of dimension 2K, such that for each forward pass, the first K entries are used to 329 330 encode μ , and the last K entries to encode a diagonal Σ , after a soft-plus transformation. Let C_{γ} be a discrete classifier neural net parameterized by γ , such that $C_{\gamma}(y|z) \sim Categorical$. 332 r(z) is constant and unparameterized. We use Monte Carlo sampling over some discrete 333 dataset \mathcal{S} to empirically estimate \mathcal{L}_{SVIB} . The true and possibly continuous distribution 334 p(x,y) = p(y|x)p(x) can be sampled from S. Distributions featuring Z are samples from the 335 stochastic encoder using the reparameterization trick (Kingma & Welling 2014), such that 336 for each $x_n \in \mathcal{S}$ we generate a sample \hat{z}_n . Finally, we use the variational classifier to attain 338 instances \tilde{y}_n , given an instance \hat{z}_n . 339

$$\widehat{\mathcal{L}}_{SVIB} \equiv \frac{1}{N} \sum_{n=1}^{N} \left[\beta D_{KL} \left(e_{\phi}(z|x_n) \middle\| r(z) \right) - \log \left(C_{\gamma} \left(y_n \middle| \hat{z}_n \right) \right) + \lambda \log \left(C_{\gamma} \left(\tilde{y}_n \middle| \hat{z}_n \right) \right) \right]$$
(15)

3.4 Motivation

347 Tishby et al. (1999) proposed that representations are optimal if they contain just enough information for a required downstream task, and proposed the information bottleneck as 349 a method to obtain such representations. However, in the supervised case an additional 350 information processing stage is added, where representations are decoded by a learned 351 $decoder^2$, in a joint training process. As mentioned in Section 2.3, Alemi et al. (2018) 352 observed that the ELBO loss function (Kingma & Welling 2014) may learn uninformative 353 representations even when strong KL regularization is imposed, since an overpowerful decoder 354 355 can overfit the learned embeddings. This observation holds for all VIB loss functions (Alemi 356 et al. 2017; Fischer 2020; Cheng et al. 2020), as VIB is equivalent to the ELBO loss, as 357 shown in (Alemi et al. 2017). Our proposed extension to the IB and VIB frameworks asks 358 to resolve this conflict. By appending an additional bottleneck between the representation 359 Z, and learned classifier Y, we learn a classifier that holds the minimal information about 360 the representation that is required to meet a designated distortion target over the true 361 downstream RV. Extending the work in (Alemi et al. 2018). We propose to define a decoder 362 363 \tilde{Y} as overfitting, if a substantial amount of its information about Z lacks relevance about Y. 364 The conditional MI $I(Z; \tilde{Y}|Y)$ measures the amount of information Z and \tilde{Y} share, that is 365 uninformative about about Y. Hence, we have that \tilde{Y} overfits Z if: 366

367 368

369 370

371

377

340 341

343 344 345

346

 $I(Z; \tilde{Y}) \gg I(Z; \tilde{Y}) - I(Z; \tilde{Y}|Y)$ $H(\tilde{Y}|Y) \gg H(\tilde{Y}|Z)$ (16)

Where the last line follows from the SIB Markov chain.

By deriving the second bottleneck, \mathcal{L}_{SVIB} introduces a modulated conditional entropy term to the loss function: $-\lambda H(\tilde{Y}|Z)$, inducing an increase in the right hand side of Equation 16. At the same time, we expect that the left hand side conditional entropy will be reduced, by

²Here decoder in the general sense, including classifiers and other decoders

the power of the cross entropy term. Applying these two forces together prevents decoders
 from overfitting embeddings, as is illustrated in Figure 1.



Figure 1: Venn diagrams illustrating decoder overfitting. The left diagram depicts an overfitted decoder where \tilde{Y} holds no information about Y, and $H(\tilde{Y}|Y) \gg H(\tilde{Y}|Z)$. The right diagram depicts a regularized decoder where $H(\tilde{Y}|Y)$ is not much greater than $H(\tilde{Y}|Z)$.

4 Experiments

We follow the experimental setup proposed by Alemi et al. (2017), extending it to NLP tasks as well. We trained image classification models on the ImageNet 2012 dataset (Deng et al. 2009), and text classification models on the IMDB sentiment analysis dataset (Maas et al. 2011). For each dataset, we compared a competitive Vanilla model with VIB models trained over 8 different β values ranging from 10^{-4} to 0.5, a VCEB model trained with ρ values ranging from 1 to 7, and an SVIB model trained with different combinations of β and λ values. All models were trained over a frozen encoder of the vanilla model, to allow faster experimentation. Each model was trained and evaluated 5 times per setting. Models were evaluated over test set accuracy and robustness to various adversarial attacks, showing consistent performance. For image classification, we employed the untargeted Fast Gradient Sign (FGS) attack (Goodfellow et al. 2015), as well as the targeted CW L_2 attack (Carlini & 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). The empirical results presented in Figure 2 confirms that while VIB, VCEB and SVIB models mostly decrease test set accuracy compared to the vanilla model, they significantly improve robustness to the applied adversarial attacks. SVIB attains significantly higher test set accuracy over VIB and VCEB, notably outperforming the vanilla model for IMDB, while scoring the highest robustness in all attacks, apart from the CW attack. A comparison of the best VIB, VCEB and SVIB models further substantiates these findings, with statistical significance confirmed by a p-value of less than 0.05 on a Wilcoxon rank sum test. We note that our experiments compare identical models, varying only in objective functions and scaling parameters. This design highlights performance differences solely due to these factors. Methods like training from scratch to boost overall performance were omitted to ensure a robust comparison. Elaboration on the experimental setup, detailed results, and further insights from the experiments are available in Appendix C. Code to reconstruct the experiments is provided in the supplementary materials of this submission.

432 4.1 IMAGE CLASSIFICATION

 A pre-trained inceptionV3 (Szegedy et al. 2016) base model was used and achieved a 77.21% accuracy on the ImageNet 2012 validation set (Test set for ImageNet is unavailable). Image classification evaluation results are shown in Figure 2, examples of successful attacks are shown in Figures 6, 7 in Appendix C.

4.2 Text classification

A fine tuned BERT uncased (Devlin et al. 2019) base model was used, and achieved a 93.0% accuracy on the IMDB sentiment analysis test set. Text classification evaluation results are shown in Figure 2, examples of successful attacks are shown in Figures 1,2 in Appendix C.



Figure 2: Performance comparison across models and metrics for IMDB and ImageNet. **Higher is better** \uparrow **in all plots**. Analyzing accuracy and robustness against adversarial attacks for vanilla, SVIB, VIB, and VCEB models under Varying β and ρ values, average over 5 runs with standard deviation. Left column features IMDB tasks, right column features ImageNet tasks. Upper row shows accuracy over test set, and bottom rows depict robustness under various adversarial attacks, presented as the rate of deflected attacks, or as the average L_2 distance required for a successful CW attack. Results show that SVIB attains significantly higher test set accuracy, outperforming the vanilla model for IMDB, while attaining better robustness in all attacks apart from the CW attack. ρ values apply to CEB models, while β values apply for SVIB and VIB models. SVIB results are presented for $\lambda = 1$ in IMDB and $\lambda = 2$ in ImageNet. For all experimental results please see the results Section in Appendix C.

486 5 DISCUSSION

488 The IB is a special case of rate-distortion, and was initially designed to optimize compressed 489 representations. Applying the IB objective for supervised tasks results in optimization of a 490 classifier distribution as well, and requires a reformulation of the initial problem to include 491 492 both representation and classification. We propose Supervised IB (SIB), an extension to the 493 original IB that considers the classifier distribution as well, and adds an additional bottleneck 494 to mitigate information flow between representations and classifier. We derive a tractable 495 variational approximation for SIB, SVIB, and show that it outperforms VIB and VCEB in 496 terms of classification accuracy and robustness to adversarial attacks, over high dimensional 497 tasks of different modalities, with high statistical significance. We use previous information 498 theoretic frameworks for deep learning (Alemi et al. 2018; Perevra et al. 2017; Szegedy et al. 499 500 2016) to interpret our findings, and propose a definition for decoder overfitting, and a new 501 motivation for conditional entropy regularization. While other advancements have been 502 achieved in recent years, (Fischer 2020; Cheng et al. 2020; Achille & Soatto 2018), none 503 propose a reformulation for IB, as is required in our opinion. 504

This study opens many opportunities for further research: Applying SVIB in self-supervised learning, and in particular measuring whether representations learned with SVIB capture better semantics than representations learned with non IB inspired loss functions, empirical studies with a full covariance matrix SVIB, a GMM model SVIB, adding β and λ annealing to SVIB, and combining SVIB with CEB are left for future work.

511

512 REFERENCES

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

Alexander A. Alemi, Ian Fischer, Joshua V. Dillon, and Kevin Murphy. Deep variational
information bottleneck. In *Proceedings of the International Conference on Learning Representations (ICLR)*, Google Research, 2017.

Alexander A. Alemi, Ben Poole, Ian Fischer, Joshua V. Dillon, Rif A. Saurous, and Kevin
Murphy. Fixing a broken elbo. In *Proceedings of Machine Learning Research*, volume 80, pp.
159–168, PMLR, 2018. URL http://dblp.uni-trier.de/db/conf/icml/icml2018.html#
AlemiPFDS018.

527

Rana Ali Amjad and Bernhard C. Geiger. Learning representations for neural networkbased classification using the information bottleneck principle. *IEEE Trans. Pattern Anal. Mach. Intell.*, 42(9):2225-2239, 2020. URL http://dblp.uni-trier.de/db/journals/
pami/pami42.html#AmjadG20.

David Barber and Felix V. Agakov. The im algorithm: a variational approach to in formation maximization. In Neural Information Processing Systems, 2003. URL https:
 //api.semanticscholar.org/CorpusID:14633080.

Nicholas Carlini and David A. Wagner. Towards evaluating the robustness of neural networks.
In *IEEE Symposium on Security and Privacy*, pp. 39–57. IEEE Computer Society, 2017.
URL http://dblp.uni-trier.de/db/conf/sp/sp2017.html#Carlini017.

Gal Chechik et al. Gaussian information bottleneck. In Advances in Neural Information
 Processing Systems, 2003. URL https://proceedings.neurips.cc/paper/2003/hash/
 7e05d6f828574fbc975a896b25bb011e-Abstract.html.

Pengyu Cheng, Weituo Hao, Shuyang Dai, Jiachang Liu, Zhe Gan, and Lawrence Carin.
CLUB: A contrastive log-ratio upper bound of mutual information. In Hal Daumé III and
Aarti Singh (eds.), Proceedings of the 37th International Conference on Machine Learning,
volume 119 of Proceedings of Machine Learning Research, pp. 1779–1788. PMLR, 13–18 Jul
2020. URL https://proceedings.mlr.press/v119/cheng20b.html.

- Thomas M Cover. *Elements of information theory*. John Wiley & Sons, 1999.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A largescale hierarchical image database. In 2009 IEEE conference on computer vision and pattern
 recognition, pp. 248–255. Ieee, 2009.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguis-tics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 4171–4186, Minneapolis, Minnesota, 2019. URL https://www.aclweb.org/anthology/N19-1423.
- Ian Fischer. The conditional entropy bottleneck. *Entropy*, 22(9):999, 2020. URL http:
 //dblp.uni-trier.de/db/journals/entropy/entropy22.html#Fischer20.
- Ian Fischer and Alexander A. Alemi. Ceb improves model robustness. *Entropy*, 22(10), 2020.
 ISSN 1099-4300. doi: 10.3390/e22101081. URL https://www.mdpi.com/1099-4300/22/
 10/1081.
- Ji Gao, Jack Lanchantin, Mary Lou Soffa, and Yanjun Qi. Black-box generation of adversarial
 text sequences to evade deep learning classifiers. In 2018 IEEE Security and Privacy
 Workshops (SPW), pp. 50–56. IEEE, 2018.
- Bernhard Geiger and Ian Fischer. A comparison of variational bounds for the information
 bottleneck functional. *Entropy*, 2020. URL https://www.mdpi.com/1099-4300/22/11/
 1229.
- Ziv Goldfeld and Yury Polyanskiy. The information bottleneck problem and its applications
 in machine learning. *IEEE Journal on Selected Areas in Information Theory*, 1(1):19–38,
 2020. doi: 10.1109/JSAIT.2020.2991561.
- Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing
 adversarial examples. In *ICLR (Poster)*, 2015. URL http://dblp.uni-trier.de/db/conf/
 iclr/iclr2015.html#GoodfellowSS14.
- Irina Higgins, Loïc Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew
 Botvinick, Shakir Mohamed, and Alexander Lerchner. beta-vae: Learning basic visual
 concepts with a constrained variational framework. In *ICLR (Poster)*, 2017.
- 592

563

566

579

583

593 Kaiwen. pytorch-cw2, 2018. URL https://github.com/kkew3/pytorch-cw2. GitHub repository.

607

⁵⁹⁴ Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. In 2nd International
⁵⁹⁵ Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014,
⁵⁹⁷ Conference Track Proceedings, 2014.

Diederik P. Kingma and Max Welling. An introduction to variational autoencoders. Foundations and Trends in Machine Learning, 12(4):307-392, 2019. URL http://dblp.uni-trier.
de/db/journals/ftml/ftml12.html#KingmaW19.

Andrew Maas, Raymond E Daly, Peter T Pham, Dan Huang, Andrew Y Ng, and Christopher
Potts. Learning word vectors for sentiment analysis. In *Proceedings of the 49th annual meeting of the association for computational linguistics: Human language technologies*, pp.
142–150, 2011.

Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu.
 Towards deep learning models resistant to adversarial attacks. In International Conference
 on Learning Representations, 2018. URL https://openreview.net/forum?id=rJzIBfZAb.

John Morris, Eli Lifland, Jin Yong Yoo, Jake Grigsby, Di Jin, and Yanjun Qi. Textattack: A
framework for adversarial attacks, data augmentation, and adversarial training in nlp. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing:*System Demonstrations, pp. 119–126, 2020.

Amichai Painsky and Naftali Tishby. Gaussian lower bound for the information bottleneck
 limit. J. Mach. Learn. Res., 18:213:1-213:29, 2017. URL http://dblp.uni-trier.de/db/
 journals/jmlr/jmlr18.html#PainskyT17.

Gabriel Pereyra, George Tucker, Jan Chorowski, Lukasz Kaiser, and Geoffrey E. Hinton.
Regularizing neural networks by penalizing confident output distributions. In *Proceedings*of the International Conference on Learning Representations, OpenReview.net, 2017. URL
http://dblp.uni-trier.de/db/conf/iclr/iclr2017w.html#PereyraTCKH17.

Zoe Piran, Ravid Shwartz-Ziv, and Naftali Tishby. The dual information bottleneck, 2020.
 URL https://arxiv.org/abs/2006.04641.

Shuhuai Ren, Yihe Deng, Kun He, and Wanxiang Che. Generating natural language
adversarial examples through probability weighted word saliency. In *Proceedings of the 57th annual meeting of the association for computational linguistics*, pp. 1085–1097, 2019.

Andrew M Saxe, Yamini Bansal, Joel Dapello, and Madhu Advani. On the information
bottleneck theory of deep learning. 2018.

Claude E. Shannon. Coding theorems for a discrete source with a fidelity criterion. In *IRE National Convention*, 1959.

Ravid Shwartz-Ziv and Naftali Tishby. Opening the black box of deep neural networks via
 information, 2017. URL http://arxiv.org/abs/1703.00810. 19 pages, 8 figures.

Noam Slonim. The information bottleneck: Theory and applications. PhD thesis, Hebrew University of Jerusalem Jerusalem, Israel, 2002.

645
 646
 647
 648
 647
 649
 647
 648
 648
 649
 649
 649
 649
 649
 641
 641
 642
 642
 643
 644
 644
 644
 645
 645
 646
 647
 648
 648
 649
 649
 649
 649
 641
 642
 642
 642
 643
 644
 644
 645
 645
 646
 647
 648
 649
 649
 649
 649
 649
 649
 649
 641
 642
 642
 642
 644
 645
 645
 646
 647
 648
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649
 649

⁶⁴⁸ Naftali Tishby and Noga Zaslavsky. Deep learning and the information bottleneck principle, 2015.

⁶⁵¹ Naftali Tishby, Fernando C. Pereira, and William Bialek. The information bottleneck method.
⁶⁵² In *The 37th annual Allerton Conference on Communication, Control, and Computing.*,
⁶⁵⁴ Hebrew University, Jerusalem 91904, Israel, 1999.

Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-SNE. Journal of
Machine Learning Research, 9:2579-2605, 2008. URL http://www.jmlr.org/papers/v9/
vandermaaten08a.html.

Vladimir N. Vapnik. The nature of statistical learning theory. Springer-Verlag New York,
Inc., 1995. ISBN 0-387-94559-8.

702 APPENDIX A - PRELIMINARIES

704 705 Notation

We denote random variables (RVs) with upper cased letters X, Y, and their realizations in lower case x, y. Denote discrete Probability Mass Functions (PMFs) with an upper case P(x)and continuous Probability Density Functions (PDFs) with a lower case p(x). Subscripts are written where the RVs identities are not clear from the context, and hat notation denotes empirical measurements.

712 Let X, Y be two observed random variables with a true and unknown joint distribution 713 p(x, y), and true marginals p(x), p(y). We can attempt to approximate these distributions 714 using a model p_{θ} with parameters θ , such that for generative tasks $p_{\theta}(x) \approx p(x)$, and 715 716 for discriminative tasks $p_{\theta}(y|x) \approx p(y|x)$, using a dataset of N i.i.d observation pairs 717 $\mathcal{S} = \{(x_1, y_1), ..., (x_N, y_N)\}$ to fit our model. One can also assume the existence of an 718 additional unobserved RV $Z \sim p(z)$ that influences or generates the observed RVs X, Y. 719 Since Z is unobserved, it is absent from the dataset \mathcal{S} , and so cannot be modeled directly. 720 Denote $p_{\theta}(x) = \int p_{\theta}(x|z)p_{\theta}(z) dz = \int p_{\theta}(x,z) dz$ the marginal, $p_{\theta}(z)$ the prior as it is not 721 conditioned over any other RV, and $p_{\theta}(z|x)$ the posterior following Bayes' rule. 722

723 724

725

734

735 736

737

738 739

740

741

742

743 744

745

746

VARIATIONAL APPROXIMATIONS

When modeling an unobserved variable of an unknown distribution, we encounter a problem as the marginal $p_{\theta}(x) = \int p_{\theta}(x, z) dz$ doesn't have an analytic solution. This intractability can be overcome by choosing some tractable parametric variational distribution $q_{\phi}(z|x)$ to approximate the posterior $p_{\theta}(z|x)$, 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 S (Kingma & Welling 2019).

733 Learning tasks

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

- A generator of random vectors $x \in \mathbb{R}^d$, drawn independently from an unknown probability distribution p(x).
- A supervisor who returns a scalar output value $y \in \mathbb{R}$, according to an unknown conditional probability distribution p(y|x). We note that these probabilities can indeed be soft labels, where y is a continuous probability vector, rather the more commonly used hard labels.
- A learning machine capable of implementing a predefined set of functions, $f(x, \theta)$: $\mathbb{R}^d \times \Theta \mapsto \mathbb{R}$, where Θ is a set of parameters.

The problem of supervised learning is that of choosing from the given set of functions, the one that best approximates the supervisor's response, based on observation pairs from the training set S, drawn according to p(x, y) = p(x)p(y|x).

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

INFORMATION THEORETIC FUNCTIONS

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

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

Appendix B - Related work elaboration

This appendix supplements the related work presented in Section 2, by providing a deeper review of the IB, the IB theory of deep learning, and variational approximations for the IB.

816 THE INFORMATION PLANE

As mentioned in Section 2.1, the solution to the IB objective, $\mathcal{L}_{IB} = I(X; Z) - \beta I(Z; Y)$, depends on the Lagrange multiplier β . Hence, the IB objective has no one unique solution, and can thus be plotted as a function of β and of Z's cardinality, over a Cartesian system composed of the axes I(X; Z) (rate) and I(Z; Y) (distortion). We denote the resulting curve the *information curve*, and its Cartesian system the *information plane* (Tishby et al. 1999), as illustrated in Figure 3. When β approaches 0 the distortion term is nullified and we learn a representation that has maximal compression but no information over the down stream task (such a representation may be a null vector), and when β approaches ∞ we learn a representation that has the maximal possible information over the downstream task, but minimal compression. The region above the information curve is unreachable by any possible representation. The different bifurcation of the information curve, illustrated in Figure 3, correspond to the different possible cardinalities of the compressed representation.



Figure 3: The information plane and curve: rate-distortion ratio over β . At $\beta = 0$ the representation is compressed but uninformative (maximal compression), at $\beta \to \infty$ the representation is informative but potentially overfitted (maximal information). Taken from (Slonim 2002).

856 FIXING A BROKEN ELBO

Kingma & Welling (2014) introduced variational auto encoders (VAEs) as a latent model based generative DNN architecture. In VAEs, an unobserved RV Z is assumed to generate evidence X, a variational DNN encoder e(z|x) is used to approximate the intractable posterior p(z|x), and a variational DNN decoder $d(\hat{x}|z)$ is used to reconstruct X. The log probability log(p(x)) is developed in to the tractable Evidence Lower Bound (ELBO) loss: $log(p(x)) \leq \mathcal{L}_{\text{ELBO}}(x) \equiv -\mathbb{E}_{e(z|x)} [log(d(x|z))] + D_{KL} (e(z|x)||m(z))$, consisting of a reconstruction error term (cross entropy), and a KL regularization term between encoder and variational marginal m(z).

867 Alemi et al. (2018) adapt the information plane (Tishby et al. 1999) to VAEs by defining an 868 additional theoretical bound for the ratio between rate and distortion, imposed by the limits 869 of finite parametric families of variational approximations. Instead of true rate and distortion, 870 the proposed information plane features variational rate as $R \equiv D_{KL} \left(e(z|x) || m(z) \right)$, and 871 variational distortion as $D \equiv -\int \int p(x)e(z|x)\log(d(x|z)) dxdz$. Figure 4 illustrates the 872 suggested information plane, which is divided into three sub planes: (1) Infeasible: This is 873 874 the IB theoretical limit (As per Figure 3); (2) Feasible: Attainable given an infinite model 875 family, and complete variety of e(z|x), d(x|z) and m(z); (3) Realizable: Attainable given a 876 finite parametric and tractable variational family. The black diagonal line at the lower left 877 satisfies $H_p(X) - D = R$, resulting in tight variational bounds on the mutual information. 878

879 Alemi et al. (2018) observe that the variational rate R does not depend on the variational 880 decoder distribution d(x|z). As R is used as the ELBO KL regularizer, high variational 881 compression rates can be attained regardless of MI between decoder and learned representa-882 tion. Equivalently, good reconstruction does not directly depend on good representation. 883 Empirical evidence suggest that VAEs are prone to learn uninformative representations 884 while still achieving low ELBO loss, a degeneration made possible by overpowerful decoders 885 that are able to overfit the little information captured by the encoder. $D_{KL}(e(z|x)||m(z))$ 886 887 approaches 0 iff $e(z|x) \to m(z)$, making e(z|x) close to independence from x, resulting in a latent representation that fails to encode information about the input. However, a suitably 889 powerful decoder could possibly learn to overfit encoded traces of the training examples, and 890 reach a low distortion score during optimization. 891

In the current study, we extend this theoretical framework to explain the advancements of our proposed loss function.

892

893

906 907 908



912 913

Figure 4: 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. Taken from (Alemi et al. 2018).

 \overline{H}

R

D

Η

918 IB THEORY OF DEEP LEARNING 919

924

925

926

927

928

929 930

931

932

933

934

935

936

937

953

954

955

956 957

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

Tishby & Zaslavsky (2015) proposed a representation-learning interpretation of DNNs using the IB framework, regarding DNNs as Markov cascades of intermediate representations between hidden layers. Under this notion, comparing the optimal and the achieved ratedistortion ratios between DNN layers will indicate if a model is too complex or too simple for a given task and training set. Shwartz-Ziv & Tishby (2017) visualized and analyzed the information plane behavior of DNNs over a toy problem with a known joint distribution. Mutual information of the different layers was estimated and used to analyze the training process. The learning process over Stochastic Gradient Descent (SGD) exhibited two separate and sequential behaviors: A short Empirical Error Minimization phase (ERM) characterized by a rapid decrease in distortion, followed by a long compression phase with an increase in rate until convergence to an optimal IB limit, as demonstrated in Figure 5. Similar, yet repetitive behavior was observed in the current study, as elaborated in Section 5.



Figure 5: Information plane scatters of different DNN layers (colors) in 50 randomized networks. Left are initial weights, center are at 400 epochs, and right at 9000 epochs. Taken from Shwartz-Ziv & Tishby (2017).

Saxe et al. (2018) reproduced the experiments described in (Shwartz-Ziv & Tishby 2017), 958 959 expanding them to different activation functions, different datasets and different methods to 960 estimate mutual information. It was found that double-sided saturated nonlinear activations, 961 such as the tanh, produced a distinct compressions stage when mutual information was 962 measured by binning, as performed in (Shwartz-Ziv & Tishby 2017), while other activations 963 did not. It was also shown that DNN generalization did not depend on a distinct compression 964 stage, and that DNNs do forget task irrelevant information, but this happens concurrently 965 966 to the learning of task relevant information, and not necessarily separately. Amjad & 967 Geiger (2020) argued against the use of the IB as an objective for deterministic DNNs, 968 as mutual information in deterministic DNNs is either infinite or step like, because of 969 mutual information's invariance to invertible transformations, and because of the absence 970 of a decision function in the objective. Using IB as an objective in stochastic DNNs, such 971 as of the variational IB family, is suggested as a possible solution. When examining the

information plane behavior in the current study, we notice recurring patterns of distortion
reduction followed by rate increase, resembling the ERM and representation compression
stages described by Shwartz-Ziv & Tishby (2017), as elaborated in Appendix 5.

977 CONDITIONAL ENTROPY BOTTLENECK

As mentioned in Section 2.2, Fischer (2020) showed that the conditional entropy bottleneck is equivalent to IB for $\gamma = \beta - 1$ following the chain rule of mutual information (Cover 1999), and the IB Markov chain. We develop this equivalence in detail:

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

¹⁰²⁶ Appendix C - Experiments elaboration

1027 1028

Image classification models were trained on the first 500,000 samples of the ImageNet 1029 2012 dataset (Deng et al. 2009), and text classification over the entire IMDB sentiment 1030 analysis dataset (Maas et al. 2011). For each dataset, a competitive pre-trained model 1032 (Vanilla model) was evaluated and then used to encode embeddings. These embeddings were then used as a dataset for a new stochastic classifier net with either a VIB or a SVIB loss 1034 function. Stochastic classifiers consisted of two ReLU activated linear layers of the same 1035 dimensions as the pre-trained model's logits (2048 for image and 768 for text classification), 1036 followed by reparameterization and a final softmax activated FC layer. Learning rate was 1037 10^{-4} and decaying exponentially with a factor of 0.97 every two epochs. Batch sizes were 32 1038 for ImageNet and 16 for IMDB. All models were trained using an Nvidia RTX3080 GPU 1039 1040 with approximately 1-2 days per a single experiment run. Beta values of $\beta = 10^{-i}$ for 1041 $i \in \{1, 2, 3\}$ were tested, and we used a single forward pass per sample for inference, since 1042 previous studies indicated that these are the best range and sample rate for VIB (Alemi et al. 1043 2017; 2018). Each model was trained and evaluated 5 times per β value, with consistent 1044 performance. Statistical significance was demonstrated in all comparisons using the Wilcoxon 1045 rank sum test with all metrics compared attaining a p-value of less than 0.05. Rank sum 1046 was computed as follows: A sorted vector of results was prepared for each compared 1047 1048 metric, where each entry featured the attained result in each of the 5 i.i.d. experiments 1049 per algorithm, and a boolean indicator value for the algorithm type. For example, let r :=1050 ((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))1051 be a sorted vector of (test accuracy, algorithm) tuples, 1 being SVIB, 0 VIB. We compute 1052 the rank-sum as follows: 1053

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

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

In practice, these were computed with the Python Scipy library as follows:

1055 1056 1057

1054

1058 1059

1060 1061

1062 1063 1064

1066

1068

```
import scipy.stats as stats
vib_scores = [0.915, 0.915, 0.91, 0.92, 0.89]
svib_scores = [0.93, 0.935, 0.925, 0.93, 0.94]
pvalue = stats.ranksums(svib_scores, vib_scores, 'greater').pvalue
assert pvalue < 0.05</pre>
```

```
1069
1070 IMAGE CLASSIFICATION
```

1071 The ImageNet 2012 validation set was used for evaluation as the test set for ImageNet 1072 is unavailable. InceptionV3 yields a slightly worse single shot accuracy than inceptionV2 1073 (80.4%) when run in a single model and single crop setting, however we've used InceptionV3 1074 1075 over V2 for simplicity. Each model was trained for 100 epochs. The entire validation set 1076 was used to measure accuracy and robustness to FGS attacks, while only 1% of it was used 1077 for CW attacks, as they are computationally expensive. Complete results are available in 1078 Section 5. Examples of successful attacks are shown in Figures 6.7. t-SNE (van der Maaten 1079 & Hinton 2008) visualization of the latent space of each model is presented in Figure 8.



Figure 6: 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 β =0.01. Target: Soccer ball



Figure 7: 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 6.



t-SNE visualization of ImageNet embedding across different models

Figure 8: ImageNet embeddings of the different models casted to 2D using the t-SNE algorithm (van der Maaten & Hinton 2008). 5000 datapoints of the first 500 ImageNet labels. The VIB and CEB castings share similar traits of well separated clusters, while the vanilla casting shows some clustering that that seems less formed and unseparated. The SVIB casting shows very little clustering and features the most dispersed distribution. The visualization suggests that the conditional entropy term in SVIB has negated the clustering effect of the ELBO loss, and induced a more uniform representation.

1188 TEXT CLASSIFICATION 1189

Each model was trained for 150 epochs. The entire test set was used to measure accuracy,
while only the first 200 entries in the test set were used for adversarial attacks, as they
are computationally expensive. Complete results are available in Section 5. Examples of
successful attacks are shown in Tables 1,2.

1196		Original text	
1197	:	the acting costumes music cinematography and	
1198		sound are all <i>astounding</i> given the production's	
1199		austere locales.	
1200		Perturbed text	
1201			
1202		the acting, costumes, music, cinematography and	
1203		sound are all <i>dumbjounding</i> given the production's	
1204		austere locales.	
1205	Table 1. Example	of a successful PWWS attack on a vanilla Bert model fine tuned over	the
1206	IMDB dataset. Th	he original label is 'Positive sentiment'. The substituted word, marke	d in
1207	italic font, changed	d the classification to 'Negative sentiment'. SVIB and VIB classifiers	are
1200	far less susceptible	e to these perturbations as shown in Figure 2.	
1205			
1210			
1212		Original text	
1213	:	areat historical movie, will not allow a viewer to	
1214		leave once you begin to watch. View is presented	
1215		differently than displayed by most school books on	
1216		this <i>subject</i> . My only fault for this movie is it was	
1217		photographed in black and white; wished it had been	
1218		in color wow !	
1219		Perturbed text	
1220		anreat historical movie will not allow a viewer to	
-			
1221		leave once you begin to watch. View is presented	
1221 1222		leave once you begin to watch. View is presented differently than displayed by most school books on	
1221 1222 1223		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	
1221 1222 1223 1224		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	
1221 1222 1223 1224 1225 1226		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 !	
1221 1222 1223 1224 1225 1226 1227	Table 2. Example	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 !	ined
1221 1222 1223 1224 1225 1226 1227 1228	Table 2: Example over the IMDB da	of a successful Deep Word Bug attack on a vanilla Bert model, fine tu taset. The original label is 'Positive sentiment'. Perturbations, marke	uned d in
1221 1222 1223 1224 1225 1226 1227 1228 1229	Table 2: Example over the IMDB da italic font, change	of a successful Deep Word Bug attack on a vanilla Bert model, fine tw taset. The original label is 'Positive sentiment'. SVIB and VIB classifiers are	uned d in e far
1221 1222 1223 1224 1225 1226 1227 1228 1229 1230	Table 2: Example over the IMDB da italic font, change less susceptible to	<i>givent</i> information in order, with 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 ! of a successful Deep Word Bug attack on a vanilla Bert model, fine tw taset. The original label is 'Positive sentiment'. Perturbations, marke the classification to 'Negative sentiment'. SVIB and VIB classifiers are these perturbations, as shown in Figure 2.	uned d in e far
1221 1222 1223 1224 1225 1226 1227 1228 1229 1230 1231	Table 2: Example over the IMDB da italic font, change less susceptible to	of a successful Deep Word Bug attack on a vanilla Bert model, fine tu taset. The original label is 'Positive sentiment'. Perturbations, marke the classification to 'Negative sentiment'. SVIB and VIB classifiers are these perturbations, as shown in Figure 2.	uned d in e far
1221 1222 1223 1224 1225 1226 1227 1228 1229 1230 1231 1232	Table 2: Example over the IMDB da italic font, change less susceptible to	of a successful Deep Word Bug attack on a vanilla Bert model, fine tu taset. The original label is 'Positive sentiment'. SVIB and VIB classifiers are these perturbations, as shown in Figure 2.	uned d in e far
1221 1222 1223 1224 1225 1226 1227 1228 1229 1230 1231 1232 1233	Table 2: Example over the IMDB da italic font, change less susceptible to	I ave 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 ! of a successful Deep Word Bug attack on a vanilla Bert model, fine tu taset. The original label is 'Positive sentiment'. Perturbations, marke the classification to 'Negative sentiment'. SVIB and VIB classifiers are these perturbations, as shown in Figure 2.	uned d in e far
1221 1222 1223 1224 1225 1226 1227 1228 1229 1230 1231 1232 1233 1234	Table 2: Example over the IMDB da italic font, change less susceptible to	of a successful Deep Word Bug attack on a vanilla Bert model, fine tu taset. The original label is 'Positive sentiment'. Perturbations, marke the classification to 'Negative sentiment'. SVIB and VIB classifiers are these perturbations, as shown in Figure 2.	uned d in e far
1221 1222 1223 1224 1225 1226 1227 1228 1229 1230 1231 1232 1233 1234 1235	Table 2: Example over the IMDB da italic font, change less susceptible to	of a successful Deep Word Bug attack on a vanilla Bert model, fine tu taset. The original label is 'Positive sentiment'. Perturbations, marke the classification to 'Negative sentiment'. SVIB and VIB classifiers are these perturbations, as shown in Figure 2.	uned d in e far
1221 1222 1223 1224 1225 1226 1227 1228 1229 1230 1231 1232 1233 1234 1235 1236	Table 2: Example over the IMDB da italic font, change less susceptible 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 ! of a successful Deep Word Bug attack on a vanilla Bert model, fine tu taset. The original label is 'Positive sentiment'. Perturbations, marke the classification to 'Negative sentiment'. SVIB and VIB classifiers are these perturbations, as shown in Figure 2.	uned d in e far
1221 1222 1223 1224 1225 1226 1227 1228 1229 1230 1231 1232 1233 1234 1235 1236 1237	Table 2: Example over the IMDB da italic font, change less susceptible 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 ! of a successful Deep Word Bug attack on a vanilla Bert model, fine tu taset. The original label is 'Positive sentiment'. Perturbations, marke the classification to 'Negative sentiment'. SVIB and VIB classifiers are these perturbations, as shown in Figure 2.	uned d in e far
1221 1222 1223 1224 1225 1226 1227 1228 1229 1230 1231 1232 1233 1234 1235 1236 1237 1238	Table 2: Example over the IMDB da italic font, change less susceptible 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 ! 	uned d in e far
1221 1222 1223 1224 1225 1226 1227 1228 1229 1230 1231 1232 1233 1234 1235 1236 1237 1238 1238 1239	Table 2: Example over the IMDB da italic font, change less susceptible to	I ave 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 ! of a successful Deep Word Bug attack on a vanilla Bert model, fine tu taset. The original label is 'Positive sentiment'. Perturbations, marke the classification to 'Negative sentiment'. SVIB and VIB classifiers are these perturbations, as shown in Figure 2.	uned d in e far
1221 1222 1223 1224 1225 1226 1227 1228 1229 1230 1231 1232 1233 1234 1235 1236 1237 1238 1239 1240	Table 2: Example over the IMDB da italic font, change less susceptible 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 ! of a successful Deep Word Bug attack on a vanilla Bert model, fine tu taset. The original label is 'Positive sentiment'. Perturbations, marke the classification to 'Negative sentiment'. SVIB and VIB classifiers are these perturbations, as shown in Figure 2.	nned d in e far

1242 Complete empirical results 1243

1245 The following tables contain the results of all experiments run in this study.

1240
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1000

1244

1246 1247

1248								
1249 1250	eta	λ	$\mathbf{Val} \uparrow$	$\mathop{\mathbf{FGS}}_{\epsilon=0.1}\uparrow$	$\mathop{\mathbf{FGS}}\limits_{\epsilon=0.5}\uparrow$	$\mathbf{CW}\uparrow$		
1251	Vanilla model							
1252			77.2%	31.1%	30.3%	788		
1253		_	011.270	51.170	52.570	100		
1254								
1255	10^{-4}	2	$75.4\% \pm .01\%$	$40.1\% \pm .08\%$	$33.7\% \pm 2.1\%$	$3401 \\ \pm 267$		
1257 1258	10^{-3}	0.5	$74.9\% \pm .06\%$	$38.4\% \pm .06\%$	$33.8\% \pm .1\%$	$3293 \\ \pm 140$		
1259 1260	10^{-3}	1	$75.5\% \pm .03\%$	$37.2\% \pm .1\%$	$33.6\% \pm .1\%$	2666 ± 140		
1262	10^{-3}	2.0	75.4% +.07\%	38.1% +.1%	33.7% +.1%	2981 + 260		
1264 1265	10^{-3}	2.5	75.3% $\pm.01\%$	38.3% $\pm .2\%$	33.8% $\pm .15\%$	3095 ± 407		
1266 1267	10^{-3}	3.0	$75.3\% \pm .03\%$	$38.5\% \pm .2\%$	$33.9\% \pm .16\%$	$\begin{array}{r} 3078 \\ \pm 443 \end{array}$		
1268 1269	10^{-2}	0.5	$74.2\% \pm .11\%$	$42.0\% \pm .13\%$	$35.2\% \pm .06\%$	$2354 \\ \pm 394$		
1270 1271 1272	10^{-2}	1	$75.0\% \pm .05\%$	$42.4\% \pm .2\%$	$35.7\% \pm .1\%$	$1564 \\ \pm 218$		
1273 1274	10^{-2}	2.0	$75.3\% \pm .07\%$	$43.1\% \pm .1\%$	$36.3\% \pm .1\%$	$1748 \\ \pm 160$		
1275 1276	10^{-2}	2.5	$75.4\% \pm .06\%$	$43.0\% \pm .13\%$	$36.0\% \pm .1\%$	$1814 \\ \pm 144$		
1277 1278	10^{-2}	3.0	$75.4\% \pm .07\%$	$42.9\% \pm .18\%$	$36.2\% \pm .12\%$	$1749 \\ \pm 138$		
1280 1281	10^{-1}	0.5	$73.1\% \pm .04\%$	$39.1\% \pm .2\%$	$32.6\% \pm .19\%$	$3738 \\ \pm 138$		
1282 1283	10^{-1}	1	$74.8\% \pm .09\%$	$42.1\% \pm .5\%$	$35.2\% \pm .5\%$	$\begin{array}{c} 3575 \\ \pm 456 \end{array}$		
1284 1285	10^{-1}	2.0	$75.4\% \pm .03\%$	$46.6\% \pm 1.8\%$	$39.8\% \pm 2.1\%$	$\begin{array}{c} 3332 \\ \pm 443 \end{array}$		
1286 1287	10^{-1}	2.5	$75.4\% \pm .03\%$	$45.6\% \pm 1.2\%$	$38.7\% \pm 1.3\%$	$\begin{array}{c} 3581 \\ \pm 243 \end{array}$		
1200 1289 1290	10^{-1}	3.0	75.1% $\pm .09\%$	46.0% +.8%	39.3% +1.0%	3536 + 315		
1291								

Table 3: Complete ImageNet evaluation scores for vanilla and SVIB models, average over 5 1292 runs with standard deviation. First column is performance on the ImageNet validation set, 1293 second and third columns are the percent of unsuccessful FGS attacks at $\epsilon=0.1, 0.5,$ and 1294 the fourth column is the average \hat{L}_2 distance for a successful Carlini Wagner L_2 targeted 1295 attack. For all columns higher is better $\uparrow.$

1296 1297	eta	ρ	$\mathbf{Val} \uparrow$	$\mathop{\mathbf{FGS}}_{\epsilon=0.1}\uparrow$	$\mathop{\mathbf{FGS}}_{\epsilon=0.5}\uparrow$	$\mathbf{CW}\uparrow$	
1298	VIB models						
1300 1301	10^{-4}	-	$74.8\% \pm .01\%$	$28.3\% \pm .2\%$	$29.3\% \pm .2\%$	$1554 \\ \pm 280$	
1302 1303 1304	$\frac{5}{10^{-4}}$	-	$74.1\% \pm .01\%$	$37.7\% \pm .01\%$	$34.8\% \pm .01\%$	$3104 \\ \pm 529$	
1305 1306	10^{-3}	-	$73.7\% \pm .1\%$	$40.5\% \pm .2\%$	$36.1\% \pm .2\%$	$3917 \\ \pm 291$	
1307 1308	$\frac{5\cdot}{10^{-3}}$	-	$73.0\% \pm .04\%$	$44.9\% \pm .13\%$	$37.8\% \pm .21\%$	$3358 \\ \pm 245$	
1309 1310	10^{-2}	-	$72.8\% \pm .1\%$	$46.5\% \pm .2\%$	$38.0\% \pm .1\%$	$3318 \\ \pm 293$	
1311 1312	$\frac{5\cdot}{10^{-2}}$	-	$72.3\% \pm .07\%$	$44.7\% \pm .3\%$	$34.9\% \pm .32\%$	$3654 \\ \pm 333$	
1314 1315	10^{-1}	-	$72.1\% \pm .01\%$	$41.6\% \pm .1\%$	$38.0\% \pm .1\%$	$\begin{array}{c} 3318 \\ \pm 293 \end{array}$	
1316 1317	$\frac{5\cdot}{10^{-1}}$	-	$0.1\% \pm 0\%$	$0\% \pm 0\%$	$0\% \pm 0\%$	$\begin{array}{c} 0 \\ \pm 0 \end{array}$	
1318			CEB n	nodels			
1320 1321	-	1	$73.0\% \pm .07\%$	$26.5\% \pm .22\%$	$28.7\% \pm .15\%$	$4527 \\ \pm 64$	
1322 1323	-	2	$73.2\% \pm 0\%$	$26.4\% \pm .21\%$	$29.0\% \pm .03\%$	$4342 \\ \pm 173$	
1324 1325	-	3	$73.4\% \pm 0\%$	$26.7\% \pm .12\%$	$29.3\% \pm .18\%$	4556 ± 177	
1326 1327	-	4	$73.8\% \pm .08\%$	$27.0\% \pm 0\%$	$29.9\% \\ \pm .07\%$	$3689 \\ \pm 347$	
1328 1329 1330	-	5	$74.3\% \pm .05\%$	$27.6\% \pm .13\%$	$30.1\% \pm .22\%$	$\begin{array}{c} 1776 \\ \pm 146 \end{array}$	
1331 1332	-	6	$74.6\% \pm .03\%$	$27.7\% \pm .35\%$	$30.0\% \pm .13\%$	$1103 \\ \pm 154$	
1333 1334	-	7	$74.6\% \pm .04\%$	$28.0\% \pm .02\%$	$30.1\% \pm .02\%$	$\begin{array}{c} 847 \\ \pm 16 \end{array}$	

Table 4: Complete ImageNet evaluation scores for VIB and CEB models, average over 5 runs with standard deviation. First column is performance on the ImageNet validation set, second and third columns are the percent of unsuccessful FGS attacks at $\epsilon = 0.1, 0.5$, and the fourth column is the average L_2 distance for a successful Carlini Wagner L_2 targeted attack. For all columns higher is better \uparrow .

1350	β	λ	$\mathbf{Test}\uparrow$	DWB↑	PWWS ↑			
1351			Vanilla model					
1353			02.007	AE 707	0.007			
1354	-	-	95.0%	43.770	0.0%			
1355	SVIB models							
1356	10^{-4}	1	92.4%	68.4%	63.9%			
1357			$\pm.01\%$	$\pm 1.7\%$	$\pm 3.3\%$			
1358	10^{-3}	0.5	92.3%	70.7%	68.3%			
1359		0.0	$\pm .07\%$	$\pm 2.3\%$	$\pm 3.3\%$			
1360	10-3	1	03.2%	72 5%	71.6%			
1361	10	T	+.5%	+2.0%	+1.3%			
1362	10-3	2.0	2:070	±1.070				
1363	10 5	2.0	92.3% $\pm 0.7\%$	74.7% $\pm 2.5\%$	73.1% $\pm 2.4\%$			
1364			土.0770	±3.370	±3.470			
1365	10^{-3}	2.5	92.4%	75.9%	72.4%			
1300			±.07%	$\pm 1.9\%$	±1.8%			
1307	10^{-3}	3.0	92.3%	74.5%	74.4%			
1300			$\pm.04\%$	$\pm 1.7\%$	$\pm.9\%$			
1309	10^{-2}	0.5	92.4%	66.1%	68.3%			
1371			$\pm .06\%$	$\pm 4.2\%$	$\pm 3.3\%$			
1372	10^{-2}	1	92.6%	69.2%	50.0%			
1373	10	I	$\pm .8\%$	$\pm 2.0\%$	$\pm 4.8\%$			
1374	10-2	2.0	00.407	C1 007	40.207			
1375	10 -	2.0	92.4% + 1%	04.8% +4.7%	40.3% +7.4%			
1376			1.170	±4.170	±1.470			
1377	10^{-2}	2.5	92.3%	58.1%	28.9%			
1378			$\pm .1\%$	$\pm 2.5\%$	$\pm 2.45\%$			
1379	10^{-2}	3.0	92.3%	54.0%	22.5%			
1380			$\pm 0.1\%$	$\pm 3.3\%$	$\pm 2.6\%$			
1381	10^{-1}	0.5	92.2%	1.1%	0.0%			
1382			$\pm 0.02\%$	$\pm 1.1\%$	$\pm 0\%$			
1383	10^{-1}	1	89.2%	0.8%	0.0%			
1384			$\pm 2.0\%$	$\pm 0.5\%$	$\pm 0\%$			
1385	10-1	2.0	02.3%	0.0%	0.0%			
1386	10	2.0	+ 2%	+0%	+0%			
1387	10-1	25	±:2%	2.070	<u> </u>			
1380	10 1	2.5	92.4% + 1%	0.0%	0.0%			
1300	1		1.1/0	±070	<u></u>			
1391	10^{-1}	3.0	92.4%	0.0%	0.0%			
1392			±.1%	土0%	±0%			

Table 5: Complete IMDB evaluation scores for vanilla and SVIB models, average over 5 runs
with standard deviation. First column is performance over the test set, second is percent of
unsuccessful Deep Word Bug attacks, and third column is percent of unsuccessful PWWS
attacks. For all columns higher is better ↑.

04 05	β	ρ	$\mathbf{Test}\uparrow$	$\mathbf{DWB}\uparrow$	PWWS ↑
06			VIB models		
07	10-4		92.1%	67.0%	60.8%
08	10		$\pm 1.1\%$	$\pm 3.2\%$	$\pm 1.4\%$
10	$\overline{5 \cdot 10^{-4}}$	_	92.2%	68.2%	64.3%
11			$\pm.07\%$	$\pm 3.0\%$	$\pm 1.3\%$
12	10^{-3}	-	91.0%	64.9%	58.4%
3			$\pm 1.0\%$	$\pm 4.4\%$	$\pm 6.6\%$
	$\overline{5 \cdot 10^{-3}}$	-	92.2%	62.9%	48.3%
			$\pm .07\%$	$\pm 3.9\%$	$\pm 7.5\%$
,	10^{-2}	-	90.8%	59.0%	37.1%
			$\pm 0.5\%$	$\pm 4.8\%$	$\pm 14.3\%$
	$5 \cdot 10^{-2}$		92.4%	14.4%	1.0%
			$\pm .1\%$	$\pm 5.5\%$	$\pm 0.3\%$
	10^{-1}	-	89.4%	10.0%	0.9%
			$\pm.9\%$	$\pm 8.0\%$	$\pm 0.9\%$
			CEB models		
	_	0.1	92.7%	46.7%	1.65%
			$\pm.04\%$	$\pm 0.68\%$	$\pm 0.27\%$
	-	1	92.7%	43.2%	1.53%
			$\pm 0\%$	$\pm 1.45\%$	$\pm 0.8\%$
	-	2	92.5%	40.8%	0%
			$\pm 0\%$	$\pm.72\%$	$\pm 0\%$
	-	3	92.3%	36.2%	0%
			$\pm 0\%$	$\pm 0\%$	±0%
	-	4	92.1%	38.8%	0%
			$\pm 0\%$	$\pm 0\%$	±0%
	-	5	92.2%	39.6%	1.0%
			$\pm 0\%$	$\pm 0\%$	±0%
	-	6	92.1%	41.9%	0%
			$\pm 0\%$	$\pm 0\%$	±0%
	-	7	92.2%	41.9%	2.15%
			±0%	$\pm 0\%$	±0%
	-	8	92.2%	45.9%	0%
			$\pm 0\%$	$\pm 0\%$	±0%

Table 6: Complete IMDB evaluation scores for VIB and CEB models, average over 5 runs with standard deviation. First column is performance over the test set, second is percent of unsuccessful Deep Word Bug attacks, and third column is percent of unsuccessful PWWS attacks. For all columns higher is better \uparrow .