
Track 1:

Achieving Domain-Independent Certified Robustness via *Knowledge Continuity*

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We present *knowledge continuity*, a novel definition inspired by Lipschitz continuity
2 which aims to certify the robustness of neural networks across input domains
3 (such as continuous and discrete domains in vision and language, respectively).
4 Most existing approaches that seek to certify robustness, especially Lipschitz
5 continuity, lie within the continuous domain with norm and distribution-dependent
6 guarantees. In contrast, our proposed definition yields certification guarantees that
7 depend only on the loss function and the intermediate learned metric spaces of the
8 neural network. These bounds are independent of domain modality, norms, and
9 distribution. We further demonstrate that the expressiveness of a model class is
10 not at odds with its knowledge continuity. This implies that achieving robustness
11 by maximizing knowledge continuity should not theoretically hinder inferential
12 performance. Finally, we present several applications of knowledge continuity such
13 as regularization and show that knowledge continuity can also localize vulnerable
14 components of a neural network.

15 1 Introduction

16 Deep neural networks (DNNs) have demonstrated remarkable generalization capabilities. Their
17 robustness, however, has been considerably more difficult to achieve. Robustness refers to the
18 preservation of model performance under natural or adversarial alterations of the input [14]. DNNs'
19 lack of robustness, highlighted by seminal works such as [19, 53] and recently [6, 4], poses signifi-
20 cant challenges to their adoption in critical applications, underscoring concerns for AI safety and
21 trustworthiness [15, 23, 7, 6].

22 Though issues of robustness emerged from computer vision applications, they have since spanned
23 multiple domains [1, 29, 59, 62, 6]. This research trajectory has not only prompted significant ad-
24 vancements in robustness improvements through architectural, procedural, and dataset augmentations,
25 but also unveiled the sophistication of adversarial attacks—the process through which counterex-
26 amples to robustness are generated [1, 29, 59, 62, 6]. In particular, a great deal of work has gone
27 into certified robustness which seeks to provide theoretical robustness guarantees. Certification is
28 desirable as it generally transcends any particular task, dataset, or model.

29 As a result, *Lipschitz continuity* has emerged, promising certified robustness by bounding the deriva-
30 tive of a neural network's output with respect to its input. In this way, Lipschitz continuity directly
31 captures the volatility of a model's performance, getting at the heart of robustness. Such an approach
32 has proven its merit in computer vision, facilitating robustness under norm and distributional assump-

33 tions [22, 50, 65, 63]. Its inherent ease and interpretability has lead to widespread adoption as a
34 means to measure and regulate robustness among practitioners as well [58, 10, 16, 55, 47].

35 Despite these successes in computer vision, there are fundamental obstacles when one tries to apply
36 Lipschitz continuity into discrete or non-metrizable domains such as natural language. Firstly,
37 characterizing distance in this input (and output) space is highly nontrivial, as language does not
38 have a naturally-endowed distance metric. Additionally, distance in this input (and output) space
39 cannot be task-invariant, as context could dramatically change the meaning of a sentence [41]. Lastly,
40 key architectures such as the Transformer [57] are provably *not* Lipschitz continuous [30]. Most
41 of these challenges are not unique to language and form the tip of the iceberg that represents the
42 strong divide of robustness between discrete/non-metrizable and continuous domains [17, 38]. For a
43 detailed summary of the related literature, see Appendix A.

44 To address these issues, we propose a new conceptual framework which we call *knowledge continuity*.
45 At its core, we adopt the following axiom:

46 *Robustness is the stability of a model’s performance with respect to its perceived*
47 *knowledge of input-output relations.*

48 Concretely, our framework is grounded on the premise that robustness is better achieved by focusing
49 on the probabilistic variation of a model’s loss with respect to its hidden representations, rather than
50 forcing arbitrary metrics on its inputs and outputs. Our approach results in certification guarantees
51 independent of domain modality, norms, and distribution. We demonstrate that the expressiveness of
52 a model class is not at odds with its knowledge continuity. In other words, achieving robustness by
53 improving knowledge continuity should not theoretically hinder inferential performance. We show
54 that in continuous settings (i.e. computer vision) knowledge continuity generalizes Lipschitz conti-
55 nuity and inherits its tight robustness bounds. Finally, we present an array of practical applications
56 using knowledge continuity both as an indicator to predict and characterize robustness as well as an
57 additional term in the loss function to train robust classifiers.

58 Although our results apply to all discrete/non-metrizable and continuous spaces, throughout the paper
59 we invoke examples from natural language as it culminates the aforementioned challenges. Further,
60 the ubiquity of large language models make their robustness a timely focus.

61 2 Knowledge Continuity

62 In this section, we provide a formulation of *knowledge continuity* and explore its theoretical properties.
63 Refer to Appendix B for all of the necessary background and notation.

64 We start by defining a model’s perceived knowledge through a rigorous treatment of its hidden
65 representation spaces. By considering the distance between inputs in some representation space in
66 conjunction with changes in loss, we result in a measure of *volatility* analogous to Lipschitz continuity.
67 Bounding this volatility in expectation then directly leads to our notion of knowledge continuity.
68 With these tools, we demonstrate a host of theoretical properties of knowledge continuity including
69 its certification of robustness, guarantees of expressiveness, and connections to Lipschitz continuity
70 in continuous settings. We summarize our theoretical contributions as follows:

- 71 • We *define* the perceived knowledge of a model as well as volatility and knowledge continuity
72 within a model’s representation space (see Def. 1, 2, 3, 4, respectively).
- 73 • We *prove* that knowledge continuity implies *probabilistic* certified robustness under perturbations in
74 the representation space and constraining knowledge continuity should not hinder the expressiveness
75 of the class of neural networks (see Thm. 2.1 and Prop. 2.2, 2.3, respectively).
- 76 • We *prove* that in some cases knowledge continuity is equivalent (in expectation) to Lipschitz
77 continuity. This shows that our axiomization of robustness aligns with existing results when
78 perturbation with respect to the input is well-defined (see Prop. 2.4, 2.6).

79 2.1 Defining Perceived Knowledge

80 Knowledge is generally accepted as a relational concept, as it arises from the connections we make
81 between ideas and experiences [21]. Herein, we capture the perceived knowledge of a model by
82 focusing on the relations it assigns to input-input pairs. Specifically, these relations are exposed by
83 decomposing a function $f : \mathcal{X} \rightarrow \mathcal{Y}$ into projections to intermediate metric spaces. Formally,

84 **Definition 1** (Metric Decomposition). We say that f admits a metric decomposition if there exists
 85 metric spaces $(\mathcal{Z}_1, d_1), \dots, (\mathcal{Z}_n, d_n)$ with metrics d_k for $k \in [n]$ such that

- 86 1. (\mathcal{Z}_k, d_k) is endowed with its Borel σ -algebra.
- 87 2. There exists measurable mappings h_0, h_1, \dots, h_n where $h_0 : \mathcal{X} \rightarrow \mathcal{Z}_1$, $h_k : \mathcal{Z}_k \rightarrow \mathcal{Z}_{k+1}$
 88 for $k \in [n-1]$, and $h_n : \mathcal{Z}_n \rightarrow \mathcal{Y}$.
- 89 3. $f = h_n \circ h_{n-1} \circ \dots \circ h_1 \circ h_0$.

90 To the best of our knowledge, all deep learning architectures admit metric decompositions, since
 91 their activations are generally real-valued. So, for all subsequent functions from \mathcal{X} to \mathcal{Y} , unless
 92 otherwise specified, we assume they are measurable and possess a metric decomposition. Further, we
 93 denote $f^k = h_k \circ h_{k-1} \circ \dots \circ h_1 \circ h_0$ and adopt the convention of calling h_k the k^{th} hidden layer. In
 94 Appendix C, we present several metric decompositions for a variety of architectures.

95 For any metric-decomposable function, an immediate consequence of our definition is that its metric
 96 decomposition may not be unique. However, in the context of neural networks, this is a desirable
 97 property. Seminal works from an array of deep learning subfields such as semi-supervised learn-
 98 ing [49], manifold learning [43], and interpretability [8] place great emphasis on the quality of learned
 99 representation spaces by examining the induced-topology of their metrics. This often does not affect
 100 the typical performance of the estimator, but has strong robustness implications [27]. Our results,
 101 which are dependent on particular metric decompositions, capture this trend. In Section 2.4, we
 102 discuss in detail the effects of various metric decompositions on our theoretical results.

103 2.2 Defining Knowledge Continuity

104 We first introduce what it means for a model’s performance to be volatile at a data point with respect
 105 to some learned representation of that model.

106 **Definition 2** (k -Volatility). Let $f : \mathcal{X} \rightarrow \mathcal{Y}$ and \mathcal{L} be any loss function. The k -volatility of a point
 107 $(x, y) \in \mathcal{X} \times \mathcal{Y}$ which we denote as $\sigma_f^k(x, y)$ is given by

$$\sigma_f^k(x, y) := \mathbb{E}_{\substack{(x', y') \sim \mathcal{D}_{\mathcal{X}, \mathcal{Y}} \\ f(x) \neq f(x')}} \left[\frac{\Delta \mathcal{L}_f^{(x, y)}(x', y')}{d_k(f^k(x), f^k(x'))} \right]. \quad (2.1)$$

108 By performing some algebra on the definition, we see that it decomposes nicely into two distinct
 109 terms: *sparsity* of the representation and *variation* in loss.

$$\begin{aligned} \sigma_f^k(x, y) &= \mathbb{E}_{(x', y') \sim \mathcal{D}_{\mathcal{X}, \mathcal{Y}}} \left[\frac{|\mathcal{L}(f(x), y) - \mathcal{L}(f(x'), y')|}{d_k(f^k(x), f^k(x'))} \right], \\ &= \mathcal{L}(f(x), y) \underbrace{\mathbb{E}_{(x', y')} \left[\frac{1}{d_k(f^k(x), f^k(x'))} \right]}_{\text{sparsity}} \cdot \underbrace{\left[1 - \frac{\mathcal{L}(f(x'), y')}{\mathcal{L}(f(x), y)} \right]}_{\text{variation in loss}}, \end{aligned} \quad (2.2)$$

110 Our notion of volatility essentially measures the derivative of performance with respect to perturba-
 111 tions to a model’s perceived knowledge. In particular, Eq. 2.2 reveals that there are two interactions
 112 in play which we illustrate in Fig. 1. Informally, we say that (x, y) is highly volatile if there is a
 113 large discrepancy in performance between it and points that are perceived to be conceptually similar.
 114 Therefore, highly volatile points capture inaccurate input-input knowledge relations. Additionally,
 115 (x, y) experiences low volatility if the space around it is sparse with respect to $\mathcal{D}_{\mathcal{X}, \mathcal{Y}}$. In other words,
 116 any set of perturbations applied in \mathcal{Z}_k would push (x, y) far away, with high probability. This makes
 117 (x, y) an isolated concept with little knowledge relationships associated with it.

118 Similar to Lipschitz continuity, the boundedness of the k -volatility of f across the data distribution is
 119 crucial and we denote this class of functions as *knowledge continuous*.

120 **Definition 3** (ε -Knowledge Continuity at a Point). We say that f is ε -knowledge continuous at
 121 $(x, y) \in \mathcal{X} \times \mathcal{Y}$ with respect to a function f , loss function \mathcal{L} , and hidden layer k if $\sigma_f^k(x, y) < \varepsilon$.

122 Conversely, we say that (x, y) is ε -knowledge discontinuous if the previous inequality does not hold.
 123 Further, (x, y) is simply knowledge discontinuous if $\sigma_f^k(x, y)$ is unbounded. Now, we extend this
 124 definition globally by considering the k -volatility between all pairs of points.

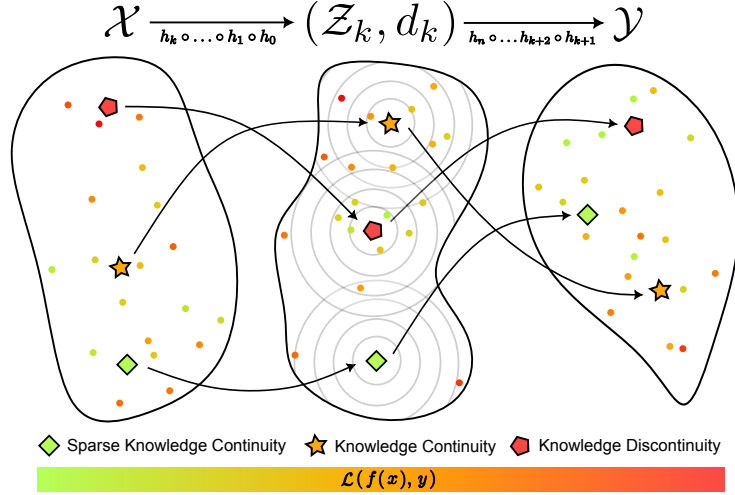


Figure 1: Various types of knowledge (dis)continuities. $f : \mathcal{X} \rightarrow \mathcal{Y}$ is a measurable map, and (\mathcal{Z}_k, d_k) is one of its hidden representation. \blacklozenge denotes knowledge continuity from sparsity: an isolated concept with no knowledge relations close to it. So, any perturbation moves \blacklozenge far away with high probability. Smooth changes in loss around \star implies knowledge continuity. Finally, \blacklozenge lacks continuity due to drastic changes in loss nearby.

125 **Definition 4** (ε -Knowledge Continuity). We say that f is ε -knowledge continuous with respect to a
 126 loss function \mathcal{L} and hidden layer k if

$$\mathbb{E}_{(x,y) \sim \mathcal{D}}[\sigma_f^k(x, y)] < \varepsilon. \quad (2.3)$$

127 Though the functional forms of Lipschitz continuity and knowledge continuity are similar, there are
 128 important differences that allow us to prove more general results. Firstly, unlike Lipschitz continuity
 129 which is an analytical property of the model f , knowledge continuity is a statistical one. In this way,
 130 non-typical data points, even if they are volatile, are ignored, whereas Lipschitz continuity treats all
 131 points equally. This is necessary in many discrete applications, as projecting a countable input space
 132 onto a non-countable metric space inevitably results in a lack of correspondence thereof. Moreover,
 133 the ground-truth function from $\mathcal{X} \rightarrow \mathcal{Y}$ may not be well-defined on *all* of \mathcal{X} : consider sentiment
 134 classification of an alpha-numeric UUID string or dog-cat classification of Gaussian noise. Secondly,
 135 knowledge continuity of an estimator is measured with respect to the loss function rather than its
 136 output. This property allows us to achieve the expressiveness guarantees in Section 2.4, since it
 137 places no restrictions on the function class of estimators. Lastly, knowledge continuity measures the
 138 distance between inputs with the endowed metric in its hidden layers. This flexibility allows us to
 139 define knowledge continuity even when the input domain is not a metric space.

140 2.3 Certification of Robustness

141 Our first main result demonstrates that ε -knowledge continuity implies probabilistic certified robust-
 142 ness in the hidden representation space. In Theorem 2.1, given some reference set $A \subset \mathcal{X} \times \mathcal{Y}$, we
 143 bound the probability that a δ -sized perturbation in the representation space away from A will result
 144 in an η change in loss. In other words, knowledge continuity is able to characterize the robustness of
 145 any subset of data points with positive measure.

146 **Theorem 2.1.** Let $A \subset \mathcal{X} \times \mathcal{Y}$ such that $\mathbb{P}_{\mathcal{D}_{\mathcal{X}, \mathcal{Y}}}[A] > 0$ and $\delta, \eta > 0$. Let $A' = \{(x', y') \in \mathcal{X} \times \mathcal{Y} :$
 147 $\exists (x, y) \in A, \Delta \mathcal{L}_f^{(x, y)}(x', y') > \eta\}$. If $f : \mathcal{X} \rightarrow \mathcal{Y}$ is ε -knowledge continuous with respect to the
 148 hidden layer indexed by k and (\mathcal{Z}_k, d_k) is bounded by $B > 0$, then

$$\mathbb{P}_{(x,y) \sim \mathcal{D}_{\mathcal{X}, \mathcal{Y}}}[A' \mid d_k(f^k(x), f^k(A)) < \delta] \leq \frac{\varepsilon \delta}{\eta \left(1 - \exp \left[-\Omega \left(\frac{\delta}{B} - \sqrt{\log \frac{1}{\mathbb{P}[A]}} \right)^2 \right] \right)}. \quad (2.4)$$

149 See Appendix D for the proof. We can lose the assumptions of boundedness and knowledge of $\mathbb{P}[A]$
 150 by taking limits of Eq. D.11 with respect to B and $\mathbb{P}[A]$. This result is shown in Appendix D.

151 **2.4 Expressiveness**

152 Our second main result demonstrates that ε -knowledge continuity can be achieved without theoret-
 153 ically compromising the accuracy of the model. In other words, universal function approximation is
 154 an invariant property with respect to ε -knowledge continuity. Universal approximation puts limits
 155 on what neural networks can learn [12, 24, 37]. A major limitation of Lipschitz functions is that
 156 they are not universal function approximators of arbitrary functions (see Appendix A for a detailed
 157 discussion). However, we show that this is achievable with knowledge continuity.

158 First, we formally define a universal function approximator.

159 **Definition 5** (Universal Function Approximator). *Suppose that \mathcal{L} is Lebesgue-integrable in both
 160 coordinates. Let $\mathcal{F} \subset \mathcal{Y}^{\mathcal{X}}$ be a set of measurable functions from $\mathcal{X} \rightarrow \mathcal{Y}$ such that for any $f \in \mathcal{F}$,
 161 there exists $\mu_f \ll \mathcal{D}_{\mathcal{X}, \mathcal{Y}}$ such that $\mu_f(\text{graph}(f)) = 1$. Then, $\mathcal{U} \subset \mathcal{F}$ is a universal function
 162 approximator of \mathcal{F} if for every $f \in \mathcal{F}$ and every $\varepsilon > 0$, there exists $\hat{f} \in \mathcal{U}$ such that*

$$\int \mathcal{L}(\hat{f}(x), y) d\mu_f < \varepsilon. \quad (2.5)$$

163 We now show that it is always possible to learn some hidden representation that is perfectly robust.

164 **Proposition 2.2.** *Let $\mathcal{U} \subset \mathcal{Y}^{\mathcal{X}}$ be a universal function approximator of $\mathcal{Y}^{\mathcal{X}}$ with respect to some loss
 165 function \mathcal{L} . Then, for any $f \in \mathcal{Y}^{\mathcal{X}}$ and sequence $\varepsilon_1, \varepsilon_2, \dots$ such that $\varepsilon_n \rightarrow 0$ there are a sequence of
 166 ε_n -knowledge continuous functions in \mathcal{U} such that $\int \mathcal{L}(f_n(x), y) d\mu_f < \varepsilon_n$, for $n \in \mathbb{N}$.*

167 *Proof.* Choose $f_n \in \mathcal{U}$ such that $\int \mathcal{L}(f_n(x), y) d\mu_f < \frac{1}{2}\varepsilon_n$. Consider the 1-layer metric decomposi-
 168 tion of f , $h_1 : \mathcal{X} \rightarrow \mathcal{Z}_1$ where $\mathcal{Z}_1 = \mathcal{X}$ equipped with the trivial metric ($d_1(x, y) = 1$ if $x \neq y$ and 0
 169 otherwise). Then, $f_n = f_n \circ h_1$. So, it follows that

$$\mathbb{E} \sigma_{f_n}^1(x, y) = \int \frac{\Delta \mathcal{L}_{f_n}^{(x, y)}(x', y')}{d_1(h_1(x), h_1(x'))} d\mu_f \leq \int \Delta \mathcal{L}_{f_n}^{(x, y)}(x', y') d\mu_f \leq \varepsilon_n. \quad (2.6)$$

170 and by the construction of f_n , the proof is completed. ■

171 In other words, if our estimator was given “infinite representational capacity,” robustness can be
 172 trivially achieved by isolating every point as its own concept (as discussed in Section 2.2). We can,
 173 however, construct a tighter model by relaxing the assumptions on the input-output metric spaces.
 174 These added constraints make it so that trivial metric decompositions are no longer possible unless
 175 the metric in \mathcal{X} is also trivial. We state this formally below, note the highlighted differences between
 176 this and Prop. 2.2.

177 **Proposition 2.3.** *Suppose $(\mathcal{X}, d_{\mathcal{X}}), (\mathcal{Y}, d_{\mathcal{Y}}) := (\mathcal{X}, d_{\mathcal{X}})$ are **compact** metric spaces, $\mathcal{F} \subset \mathcal{Y}^{\mathcal{X}}$ is the
 178 **set of all continuous functions** from \mathcal{X} to \mathcal{Y} such that $\int d_{\mathcal{X}}(x, x')^{-1} d\mu_f < \infty$ and \mathcal{L} be Lipschitz
 179 continuous in both coordinates. Then, there exists a universal function approximator \mathcal{U} of \mathcal{F} that is
 180 knowledge continuous (i.e. $\mathbb{E} \sigma_f^k(x, y) < \infty$ for some k).*

181 See Appendix E for the proof.

182 **2.5 Connections to Lipschitz Continuity**

183 We now demonstrate that our axiomization of robustness presented in Section 1 aligns with the notion
 184 of robustness¹ commonly prescribed in vision [14]. This unifies the certified robustness bounds with
 185 respect to the representation space derived in Thm. 2.1 with existing work certifying robustness with
 186 respect to the input space in continuous applications such as vision.

187 Our first result identifies conditions under which knowledge continuity, implies Lipschitz continuity.

188 **Proposition 2.4.** *Suppose that $(\mathcal{X}, d_{\mathcal{X}}), (\mathcal{Y}, d_{\mathcal{Y}})$ are metric spaces. Let the first n metric decomposi-
 189 tions of $f : \mathcal{X} \rightarrow \mathcal{Y}$ be K_i -Lipschitz continuous, for $i \in [n]$. If f is ε -knowledge continuous with
 190 respect to the n^{th} hidden layer and $d_{\mathcal{Y}}(f(x), f(x')) \leq \eta \Delta \mathcal{L}_f^{(x, y)}(x', y)$ for all $x, x' \in \mathcal{X}, y \in \mathcal{Y}$,
 191 and some $\eta > 0$, then f is Lipschitz continuous in expectation. That is,*

$$\mathbb{E}_{(x, y), (x', y') \sim \mathcal{D}_{\mathcal{X}, \mathcal{Y}}} \frac{d_{\mathcal{Y}}(f(x), f(x'))}{d_{\mathcal{X}}(x, x')} \leq \varepsilon \eta \prod_{j=1}^n K_j. \quad (2.7)$$

¹Small perturbations on the input result in small changes in performance which implies small changes in output when the loss function is Lipschitz continuous.

192 The proof is presented in Appendix F and follows easily through some algebraic manipulation. Next,
 193 combining this proposition with an auxiliary result from [74], we directly yield a certification on the
 194 input space.

195 **Corollary 2.5.** *Suppose that assumptions of Prop. 2.4 are true. And also assume that $(\mathcal{X}, d_{\mathcal{X}}) =$
 196 (\mathbb{R}^n, ℓ_p) , $(\mathcal{Y}, d_{\mathcal{Y}}) = (\mathbb{R}^m, \ell_p)$, for $1 \leq p \leq \infty$. Define a classifier from $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$, g ,
 197 where $g(x) := \arg \max_{k \in [m]} f_k(x)$ for any $x \in \mathbb{R}^n$. Then, with probability $1 - \frac{\varepsilon \eta}{t} \prod_{j=1}^n K_j$,
 198 $g(x) = g(x + \delta)$ for all $\|\delta\|_p < \frac{\sqrt[3]{2}}{2t} \text{margin}(f(x))$ and $t > 0$. $f_k(x)$ is the k^{th} coordinate of $f(x)$
 199 and $\text{margin}(f(x))$ denotes the difference between the largest and second-largest output logits.*

200 See Appendix F for the proof. Our second result identifies conditions under which Lipschitz
 201 continuity, implies knowledge continuity.

202 **Proposition 2.6.** *Let $(\mathcal{X}, d_{\mathcal{X}}), (\mathcal{Y}, d_{\mathcal{Y}})$ be a metric spaces. Let $f : \mathcal{X} \rightarrow \mathcal{Y}$ be ε -Lipschitz continuous
 203 and $\mathcal{L}(f(x), y)$ be η -Lipschitz continuous with respect to both coordinates. If the first n metric
 204 decompositions of f are K_i -Lipschitz continuous, then f is knowledge continuous with respect to the
 205 n^{th} hidden layer. That is,*

$$\mathbb{E}_{(x,y) \sim \mathcal{D}_{\mathcal{X},\mathcal{Y}}} \sigma_f^n(x, y) \leq \varepsilon \eta \prod_{j=1}^n \frac{1}{K_j}. \quad (2.8)$$

206 See Appendix F for the proof. In continuous applications such as computer vision, the assumptions of
 207 both propositions are generally met (i.e. our input-output spaces are metric spaces, all hidden layers
 208 are Lipschitz, and loss functions are locally Lipschitz). Furthermore, common architectures such as
 209 fully connected networks, CNNs, RNNs, and even vision transformers are Lipschitz continuous [58,
 210 48]. This implies that our notion of robustness is indeed an appropriate generalization that transcends
 211 domain modality since in continuous settings we can recover the strong bounds of Lipschitz continuity
 212 while expanding into new discrete and non-metrizable territory.

213 3 Practical Applications

214 In addition to the theoretical guarantees yielded by knowledge continuity in Section 2, we now
 215 demonstrate that knowledge continuity can be easily applied in practice.

216 **Using knowledge continuity to predict adversarial robustness.** For a given model, f , and hidden
 217 representation, k , we first determine the smallest ε_k such that f is ε_k -knowledge continuous. Then,
 218 we collate all ε_k through a simple average. When we regress these scores from a series of model
 219 families and sizes against their empirical adversarial robustness strong correlation is observed. In
 220 particular, knowledge continuity alone is able to explain 35% of the variance in adversarial attack
 221 success rate. We present a detailed discussion of these experiments in Appendix G.

222 **Knowledge continuity can localize vulnerable hidden representations.** Since knowledge continuity
 223 is layer-specific, we repeat the previous experiment, but holding the index of the hidden representation
 224 constant. We plot the relationship between explained variance of adversarial robustness and layer
 225 index. We find that models belonging to different families result in dramatically different curves.
 226 We tune our regularization hyperparameters according to these curves and find they yield superior
 227 performance. These results are represented in Appendix G, H, and I.

228 **Regulating knowledge continuity.** Motivated by the theoretical results in Section 2, we devise
 229 algorithms to estimate the k -volatility of a given model during training. These algorithms are
 230 described in Appendix I along with guarantees on their convergence rate and unbiasedness. Then,
 231 we directly append this estimate of volatility to our loss function as a regularization term. By
 232 minimizing this regularized loss, we find that the adversarial robustness of the resulting model
 233 significantly improves. Moreover, our method outperforms existing works both in terms of robustness
 234 and training speed (up to $2\times$ for TextFooler [29] and $3\times$ for ALUM [36]). These results are presented
 235 in Appendix I, Table 1.

236 4 Conclusion

237 In this paper, we propose a novel definition, *knowledge continuity*, which addresses key limitations
 238 associated with Lipschitz robustness. We demonstrate that our definition certifies robustness across
 239 domain modality, distribution, and norms. We also show that knowledge continuity, in contrast to
 240 Lipschitz continuity, does not affect the universal approximation property of neural networks. We
 241 further establish conditions under which knowledge continuity and Lipschitz continuity are equivalent.
 242 Lastly, we present several practical applications that directly benefit the practitioner.

243 **References**

- 244 [1] M. Alzantot, Y. Sharma, A. Elgohary, B.-J. Ho, M. Srivastava, and K.-W. Chang. Generating
245 natural language adversarial examples. *arXiv preprint arXiv:1804.07998*, 2018.
- 246 [2] C. Anil, J. Lucas, and R. Grosse. Sorting out lipschitz function approximation. In *International
247 Conference on Machine Learning*, pages 291–301. PMLR, 2019.
- 248 [3] P. L. Bartlett, D. J. Foster, and M. J. Telgarsky. Spectrally-normalized margin bounds for neural
249 networks. *Advances in neural information processing systems*, 30, 2017.
- 250 [4] S. Biderman, U. PRASHANTH, L. Sutawika, H. Schoelkopf, Q. Anthony, S. Purohit, and
251 E. Raff. Emergent and predictable memorization in large language models. *Advances in Neural
252 Information Processing Systems*, 36, 2023.
- 253 [5] S. Boucheron, G. Lugosi, and P. Massart. Concentration inequalities: A nonasymptotic theory
254 of independence, 2013.
- 255 [6] N. Carlini, M. Nasr, C. A. Choquette-Choo, M. Jagielski, I. Gao, A. Awadalla, P. W. Koh,
256 D. Ippolito, K. Lee, F. Tramer, et al. Are aligned neural networks adversarially aligned? *arXiv
257 preprint arXiv:2306.15447*, 2023.
- 258 [7] P. Chao, A. Robey, E. Dobriban, H. Hassani, G. J. Pappas, and E. Wong. Jailbreaking black box
259 large language models in twenty queries. *arXiv preprint arXiv:2310.08419*, 2023.
- 260 [8] C. Chen, O. Li, D. Tao, A. Barnett, C. Rudin, and J. K. Su. This looks like that: deep learning
261 for interpretable image recognition. *Advances in neural information processing systems*, 32,
262 2019.
- 263 [9] M. Cisse, P. Bojanowski, E. Grave, Y. Dauphin, and N. Usunier. Parseval networks: Improving
264 robustness to adversarial examples. In *International conference on machine learning*, pages
265 854–863. PMLR, 2017.
- 266 [10] J. M. Cohen, E. Rosenfeld, and J. Z. Kolter. Certified adversarial robustness via randomized
267 smoothing, 2019.
- 268 [11] Z. Cranko, Z. Shi, X. Zhang, R. Nock, and S. Kornblith. Generalised lipschitz regularisation
269 equals distributional robustness. In *International Conference on Machine Learning*, pages
270 2178–2188. PMLR, 2021.
- 271 [12] G. Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics of control,
272 signals and systems*, 2(4):303–314, 1989.
- 273 [13] L. Devroye, M. Lerasle, G. Lugosi, and R. I. Oliveira. Sub-gaussian mean estimators. 2016.
- 274 [14] N. Drenkow, N. Sani, I. Shpitser, and M. Unberath. A systematic review of robustness in deep
275 learning for computer vision: Mind the gap?, 2022.
- 276 [15] K. Eykholt, I. Evtimov, E. Fernandes, B. Li, A. Rahmati, C. Xiao, A. Prakash, T. Kohno, and
277 D. Song. Robust physical-world attacks on deep learning visual classification. In *Proceedings
278 of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- 279 [16] M. Fazlyab, A. Robey, H. Hassani, M. Morari, and G. Pappas. Efficient and accurate estimation
280 of lipschitz constants for deep neural networks. *Advances in Neural Information Processing
281 Systems*, 32, 2019.
- 282 [17] F. Gama, J. Bruna, and A. Ribeiro. Stability properties of graph neural networks. *IEEE
283 Transactions on Signal Processing*, 68:5680–5695, 2020.
- 284 [18] S. Garg and G. Ramakrishnan. BAE: BERT-based adversarial examples for text classification.
285 *arXiv preprint arXiv:2004.01970*, 2020.
- 286 [19] I. J. Goodfellow, J. Shlens, and C. Szegedy. Explaining and harnessing adversarial examples.
287 *Proceedings of 3rd International Conference on Learning Representations*, 2014.

- 288 [20] H. Gouk, E. Frank, B. Pfahringer, and M. J. Cree. Regularisation of neural networks by
289 enforcing lipschitz continuity. *Machine Learning*, 110:393–416, 2021.
- 290 [21] G. S. Halford, W. H. Wilson, and S. Phillips. Relational knowledge: the foundation of higher
291 cognition. *Trends in Cognitive Sciences*, 14(11):497–505, 2010.
- 292 [22] M. Hein and M. Andriushchenko. Formal guarantees on the robustness of a classifier against
293 adversarial manipulation, 2017.
- 294 [23] D. Hendrycks, K. Zhao, S. Basart, J. Steinhardt, and D. Song. Natural adversarial examples. In
295 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*,
296 pages 15262–15271, June 2021.
- 297 [24] K. Hornik, M. Stinchcombe, and H. White. Multilayer feedforward networks are universal
298 approximators. *Neural Networks*, 2(5):359–366, 1989.
- 299 [25] H.-Y. Huang, R. Kueng, and J. Preskill. Predicting many properties of a quantum system from
300 very few measurements. *Nature Physics*, 16(10):1050–1057, 2020.
- 301 [26] Y. Huang, H. Zhang, Y. Shi, J. Z. Kolter, and A. Anandkumar. Training certifiably robust neural
302 networks with efficient local lipschitz bounds. *Advances in Neural Information Processing*
303 *Systems*, 34:22745–22757, 2021.
- 304 [27] A. Ilyas, S. Santurkar, D. Tsipras, L. Engstrom, B. Tran, and A. Madry. Adversarial examples
305 are not bugs, they are features. *Advances in neural information processing systems*, 32, 2019.
- 306 [28] R. Jia, A. Raghunathan, K. Göksel, and P. Liang. Certified robustness to adversarial word
307 substitutions. In K. Inui, J. Jiang, V. Ng, and X. Wan, editors, *Proceedings of the 2019*
308 *Conference on Empirical Methods in Natural Language Processing and the 9th International*
309 *Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4129–4142, Hong
310 Kong, China, Nov. 2019. Association for Computational Linguistics.
- 311 [29] D. Jin, Z. Jin, J. T. Zhou, and P. Szolovits. Is BERT really robust? a strong baseline for natural
312 language attack on text classification and entailment. In *Proceedings of the AAAI conference on*
313 *artificial intelligence*, volume 34, pages 8018–8025, 2020.
- 314 [30] H. Kim, G. Papamakarios, and A. Mnih. The lipschitz constant of self-attention. In M. Meila
315 and T. Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*,
316 volume 139 of *Proceedings of Machine Learning Research*, pages 5562–5571. PMLR, 18–24
317 Jul 2021.
- 318 [31] M. Lecuyer, V. Atlidakis, R. Geambasu, D. Hsu, and S. Jana. Certified robustness to adversarial
319 examples with differential privacy. In *2019 IEEE symposium on security and privacy (SP)*,
320 pages 656–672. IEEE, 2019.
- 321 [32] K. Leino, Z. Wang, and M. Fredrikson. Globally-robust neural networks. In *International*
322 *Conference on Machine Learning*, pages 6212–6222. PMLR, 2021.
- 323 [33] B. Li, C. Chen, W. Wang, and L. Carin. Certified adversarial robustness with additive noise.
324 *Advances in neural information processing systems*, 32, 2019.
- 325 [34] L. Li, R. Ma, Q. Guo, X. Xue, and X. Qiu. Bert-attack: Adversarial attack against bert using
326 bert. *arXiv preprint arXiv:2004.09984*, 2020.
- 327 [35] J. Lin, C. Song, K. He, L. Wang, and J. E. Hopcroft. Nesterov accelerated gradient and scale
328 invariance for adversarial attacks. In *International Conference on Learning Representations*,
329 2020.
- 330 [36] X. Liu, H. Cheng, P. He, W. Chen, Y. Wang, H. Poon, and J. Gao. Adversarial training for large
331 neural language models. *arXiv preprint arXiv:2004.08994*, 2020.
- 332 [37] Z. Lu, H. Pu, F. Wang, Z. Hu, and L. Wang. The expressive power of neural networks: A view
333 from the width. *Advances in neural information processing systems*, 30, 2017.

- 334 [38] B. Lütjens, M. Everett, and J. P. How. Certified adversarial robustness for deep reinforcement
335 learning. In L. P. Kaelbling, D. Kragic, and K. Sugiura, editors, *Proceedings of the Conference*
336 *on Robot Learning*, volume 100 of *Proceedings of Machine Learning Research*, pages 1328–
337 1337. PMLR, 30 Oct–01 Nov 2020.
- 338 [39] A. L. Maas, R. E. Daly, P. T. Pham, D. Huang, A. Y. Ng, and C. Potts. Learning word vectors
339 for sentiment analysis. In *Proceedings of the 49th Annual Meeting of the Association for*
340 *Computational Linguistics: Human Language Technologies*, pages 142–150, Portland, Oregon,
341 USA, June 2011. Association for Computational Linguistics.
- 342 [40] C. McDiarmid et al. On the method of bounded differences. *Surveys in combinatorics*,
343 141(1):148–188, 1989.
- 344 [41] T. Mikolov, K. Chen, G. Corrado, and J. Dean. Efficient estimation of word representations in
345 vector space, 2013.
- 346 [42] T. Miyato, S.-i. Maeda, M. Koyama, and S. Ishii. Virtual adversarial training: a regularization
347 method for supervised and semi-supervised learning. *IEEE transactions on pattern analysis*
348 *and machine intelligence*, 41(8):1979–1993, 2018.
- 349 [43] M. Moor, M. Horn, B. Rieck, and K. Borgwardt. Topological autoencoders. In *International*
350 *conference on machine learning*, pages 7045–7054. PMLR, 2020.
- 351 [44] A. S. Nemirovskij and D. B. Yudin. Problem complexity and method efficiency in optimization.
352 1983.
- 353 [45] Y. Nie, A. Williams, E. Dinan, M. Bansal, J. Weston, and D. Kiela. Adversarial nli: A new
354 benchmark for natural language understanding. *arXiv preprint arXiv:1910.14599*, 2019.
- 355 [46] Y. Oren, S. Sagawa, T. B. Hashimoto, and P. Liang. Distributionally robust language modeling.
356 In K. Inui, J. Jiang, V. Ng, and X. Wan, editors, *Proceedings of the 2019 Conference on*
357 *Empirical Methods in Natural Language Processing and the 9th International Joint Conference*
358 *on Natural Language Processing (EMNLP-IJCNLP)*, pages 4227–4237, Hong Kong, China,
359 Nov. 2019. Association for Computational Linguistics.
- 360 [47] P. Pauli, A. Koch, J. Berberich, P. Kohler, and F. Allgöwer. Training robust neural networks
361 using lipschitz bounds. *IEEE Control Systems Letters*, 6:121–126, 2021.
- 362 [48] X. Qi, J. Wang, Y. Chen, Y. Shi, and L. Zhang. Lipsformer: Introducing lipschitz continuity to
363 vision transformers. In *The Eleventh International Conference on Learning Representations*,
364 2022.
- 365 [49] A. Radford, K. Narasimhan, T. Salimans, I. Sutskever, et al. Improving language understanding
366 by generative pre-training, 2018.
- 367 [50] W. Ruan, X. Huang, and M. Kwiatkowska. Reachability analysis of deep neural networks with
368 provable guarantees. *arXiv preprint arXiv:1805.02242*, 2018.
- 369 [51] A. Shafahi, M. Najibi, M. A. Ghiasi, Z. Xu, J. Dickerson, C. Studer, L. S. Davis, G. Taylor,
370 and T. Goldstein. Adversarial training for free! *Advances in Neural Information Processing*
371 *Systems*, 32, 2019.
- 372 [52] M. H. Stone. The generalized weierstrass approximation theorem. *Mathematics Magazine*,
373 21(5):237–254, 1948.
- 374 [53] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus.
375 Intriguing properties of neural networks, 2014.
- 376 [54] N. Tishby, F. C. Pereira, and W. Bialek. The information bottleneck method. *arXiv preprint*
377 *physics/0004057*, 2000.
- 378 [55] Y. Tsuzuku, I. Sato, and M. Sugiyama. Lipschitz-margin training: Scalable certification of
379 perturbation invariance for deep neural networks. In S. Bengio, H. Wallach, H. Larochelle,
380 K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Pro-*
381 *cessing Systems*, volume 31. Curran Associates, Inc., 2018.

- 382 [56] M. Usama and D. E. Chang. Towards robust neural networks with lipschitz continuity. In
383 *Digital Forensics and Watermarking: 17th International Workshop, IWDW 2018, Jeju Island,*
384 *Korea, October 22-24, 2018, Proceedings 17*, pages 373–389. Springer, 2019.
- 385 [57] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and
386 I. Polosukhin. Attention is all you need. *Advances in neural information processing systems*,
387 30, 2017.
- 388 [58] A. Virmaux and K. Scaman. Lipschitz regularity of deep neural networks: analysis and efficient
389 estimation. *Advances in Neural Information Processing Systems*, 31, 2018.
- 390 [59] E. Wallace, S. Feng, N. Kandpal, M. Gardner, and S. Singh. Universal adversarial triggers for
391 attacking and analyzing nlp, 2021.
- 392 [60] B. Wang, S. Wang, Y. Cheng, Z. Gan, R. Jia, B. Li, and J. Liu. Infobert: Improving robustness of
393 language models from an information theoretic perspective. *arXiv preprint arXiv:2010.02329*,
394 2020.
- 395 [61] W. Wang, P. Tang, J. Lou, and L. Xiong. Certified robustness to word substitution attack
396 with differential privacy. In K. Toutanova, A. Rumshisky, L. Zettlemoyer, D. Hakkani-Tur,
397 I. Beltagy, S. Bethard, R. Cotterell, T. Chakraborty, and Y. Zhou, editors, *Proceedings of*
398 *the 2021 Conference of the North American Chapter of the Association for Computational*
399 *Linguistics: Human Language Technologies*, pages 1102–1112, Online, June 2021. Association
400 for Computational Linguistics.
- 401 [62] A. Wei, N. Haghtalab, and J. Steinhardt. Jailbroken: How does llm safety training fail?, 2023.
- 402 [63] L. Weng, H. Zhang, H. Chen, Z. Song, C.-J. Hsieh, L. Daniel, D. Boning, and I. Dhillon.
403 Towards fast computation of certified robustness for ReLU networks. In J. Dy and A. Krause,
404 editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of
405 *Proceedings of Machine Learning Research*, pages 5276–5285. PMLR, 10–15 Jul 2018.
- 406 [64] L. Weng, H. Zhang, H. Chen, Z. Song, C.-J. Hsieh, L. Daniel, D. Boning, and I. Dhillon.
407 Towards fast computation of certified robustness for relu networks. In *International Conference*
408 *on Machine Learning*, pages 5276–5285. PMLR, 2018.
- 409 [65] T.-W. Weng, H. Zhang, P.-Y. Chen, J. Yi, D. Su, Y. Gao, C.-J. Hsieh, and L. Daniel. Evalu-
410 ating the robustness of neural networks: An extreme value theory approach. In *International*
411 *Conference on Learning Representations*, 2018.
- 412 [66] E. Wong and Z. Kolter. Provable defenses against adversarial examples via the convex outer
413 polytope. In *International conference on machine learning*, pages 5286–5295.
414 PMLR, 2018.
- 415 [67] E. Wong, L. Rice, and J. Z. Kolter. Fast is better than free: Revisiting adversarial training, 2020.
- 416 [68] E. Wong, F. Schmidt, J. H. Metzen, and J. Z. Kolter. Scaling provable adversarial defenses.
417 *Advances in Neural Information Processing Systems*, 31, 2018.
- 418 [69] X. Xu, L. Li, Y. Cheng, S. Mukherjee, A. H. Awadallah, and B. Li. Certifiably robust transform-
419 ers with 1-lipschitz self-attention, 2023.
- 420 [70] J. Y. Yoo and Y. Qi. Towards improving adversarial training of nlp models, 2021.
- 421 [71] Y. Yoshida and T. Miyato. Spectral norm regularization for improving the generalizability of
422 deep learning. *arXiv preprint arXiv:1705.10941*, 2017.
- 423 [72] A. Zhang, A. Chan, Y. Tay, J. Fu, S. Wang, S. Zhang, H. Shao, S. Yao, and R. K.-W. Lee. On
424 orthogonality constraints for transformers. In C. Zong, F. Xia, W. Li, and R. Navigli, editors,
425 *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and*
426 *the 11th International Joint Conference on Natural Language Processing (Volume 2: Short*
427 *Papers)*, pages 375–382, Online, Aug. 2021. Association for Computational Linguistics.
- 428 [73] B. Zhang, D. Jiang, D. He, and L. Wang. Boosting the certified robustness of l-infinity distance
429 nets. *arXiv preprint arXiv:2110.06850*, 2021.

- 430 [74] B. Zhang, D. Jiang, D. He, and L. Wang. Rethinking lipschitz neural networks and certified
431 robustness: A boolean function perspective. *Advances in Neural Information Processing*
432 *Systems*, 35:19398–19413, 2022.
- 433 [75] H. Zhang, H. Chen, C. Xiao, S. Gowal, R. Stanforth, B. Li, D. Boning, and C.-J. Hsieh.
434 Towards stable and efficient training of verifiably robust neural networks. *arXiv preprint*
435 *arXiv:1906.06316*, 2019.
- 436 [76] A. Zou, Z. Wang, N. Carlini, M. Nasr, J. Z. Kolter, and M. Fredrikson. Universal and transferable
437 adversarial attacks on aligned language models, 2023.

438	Appendix Table of Contents	
439	A Related Works	13
440	B Notations and Background	13
441	C Examples of Metric Decompositions	14
442	D Proof of Robustness	14
443	E Proof of Expressiveness	17
444	F Proof of Equivalence Between Lipschitz Continuity and Knowledge Continuity	18
445	G Predicting Adversarial Robustness with Volatility	20
446	H Localizing Volatile Hidden Representations	22
447	H.1 Per-Layer Volatility	22
448	H.2 Per-Model Volatility	23
449	I Regularizing Knowledge Continuity	23
450	I.1 Estimating Knowledge Continuity Algorithmically	23
451	I.2 Theoretical Guarantees of Knowledge Continuity Regulation	25
452	I.3 Regulating Knowledge Continuity “In the Wild”	27
453	I.4 Ablation Studies	27
454	I.5 Training Details	30
455	J Limitations	31
456	K Broader Impacts	31
457	L Reproducibility	31

458 A Related Works

459 There have been extensive studies on developing robust neural networks with theoretical guarantees.
460 These approaches with respect to our contribution can be organized into the following categories.

461 **Certified robustness with Lipschitz continuity.** The exploration of Lipschitz continuity as a
462 cornerstone for improving model robustness has yielded significant insights, particularly in the
463 domain of computer vision. This principle, which ensures bounded derivatives of the model’s output
464 with respect to its input, facilitates a smoother model behavior and inherently encourages robustness
465 against adversarial perturbations. This methodology, initially suggested by [19], has since been
466 rigorously analyzed and expanded upon. Most theoretical results in this area focus on certifying
467 robustness with respect to the ℓ_2 -norm [9, 71, 20, 2, 32, 22, 3]. A recent push, fueled by new
468 architectural developments, has also expanded these results into ℓ_∞ -norm perturbations [74, 73, 75].
469 Further, Lipschitz continuity also serves practitioners as a computationally effective way to train
470 more robust models [55, 65, 56, 11]. This stands in contrast to (virtual) adversarial training methods
471 which brute-force the set of adversarial examples, then iteratively re-trains on them [42, 51, 67].
472 Though Lipschitz continuity has seen much success in continuous domains, it does not apply to
473 non-metrizable domains such as language. Further, architectural limitations of prevalent models such
474 as the Transformer [57, 30] exacerbate this problem. These challenges highlight a critical need for a
475 new approach that can accommodate the specificities of discrete and non-metrizable domains while
476 providing robustness guarantees.

477 **Achieving robustness in discrete/non-metrizable spaces.** Non-metrizable spaces, where it is non-
478 trivial to construct a distance metric on the input/output domains, pose a unique challenge to certified
479 robustness. Instead of focusing on point-wise perturbations, many studies have opted to examine
480 how the output probability distribution of a model changes with respect to input distribution shifts
481 by leveraging information bottleneck methods [54, 60, 46]. Most of these bounds lack granularity
482 and cannot often be expressed in closed-form. In contrast to these theoretical approaches, recent
483 efforts have refocused on directly adapting the principles underlying Lipschitz continuity to language.
484 Virtual adversarial training methods such as [36, 70] mimic the measurement of Lipschitz continuity
485 by comparing the textual embeddings with the KL-divergence of the output logits. Along these lines,
486 techniques akin to those used in adversarial training in vision have also been translated to language,
487 reflecting a shift towards robustness centered around the learned representation space [34, 18, 29].
488 Though these approaches have seen empirical success, they lack theoretical guarantees. As a result,
489 their implementations and success rate is heavily task-dependent [36, 70]. There have also been
490 attempts to mitigate the non-Lipschitzness of Transformers [72, 69] by modifying its architecture.
491 These changes, however, add significant computational overhead.

492 **Other robustness approaches.** In parallel, other certified robustness approaches such as randomized
493 smoothing [10, 33, 31] give state-of-the-art certification for ℓ_2 -based perturbations. Notable works
494 such as [28, 61] have sought to generalize these techniques into language, but their guarantees
495 strongly depend on the type of perturbation being performed. On the other hand, analytic approaches
496 through convex relaxation inductively bound the output of neurons in a ReLU network across
497 layers [66, 68, 64]. These works, however, are difficult to scale and also do not transfer easily to
498 discrete/non-metrizable domains.

499 Our approach, inspired by Lipschitz continuity, distills the empirical intuition from the works
500 of [36, 70] and provides theoretical certification guarantees independent of perturbation-type [28, 61]
501 and domain modality. We demonstrate that knowledge continuity yields many practical applications
502 analogous to Lipschitz continuity which are easy to implement and are computationally competitive.

503 B Notations and Background

504 **Notations.** Let $\mathbb{R}^{\geq 0} := [0, \infty)$. For any function $f : \mathcal{X} \rightarrow \mathcal{Y}$, we denote $\text{graph}(f) := \{(x, y) \in$
505 $\mathcal{X} \times \mathcal{Y} : f(x) = y\}$. Let $[n]$ denote the set $\{1, 2, \dots, n\}$ for $n \in \mathbb{N}$. $(\mathcal{X}, \mathcal{F}_{\mathcal{X}}, \mathbb{P}_{\mathcal{X}})$, $(\mathcal{Y}, \mathcal{F}_{\mathcal{Y}}, \mathbb{P}_{\mathcal{Y}})$
506 are probability spaces and $(\mathcal{X} \times \mathcal{Y}, \mathcal{F}_{\mathcal{X}} \otimes \mathcal{F}_{\mathcal{Y}}, \mathbb{P}_{\mathcal{X}} \times \mathbb{P}_{\mathcal{Y}})$ denotes the product measurable space of
507 the probability spaces \mathcal{X}, \mathcal{Y} . Since our contribution focuses on the supervised learning regime, we
508 colloquially refer to \mathcal{X}, \mathcal{Y} as the input and labels, respectively. We call any probability measure $\mathbb{P}_{\mathcal{X} \times \mathcal{Y}}$
509 absolutely continuous to $\mathbb{P}_{\mathcal{X}} \times \mathbb{P}_{\mathcal{Y}}$ (i.e. $\mathcal{D}_{\mathcal{X}, \mathcal{Y}}(E) = 0$ for every $E \in \mathcal{X} \times \mathcal{Y}$ with $(\mathbb{P}_{\mathcal{X}} \times \mathbb{P}_{\mathcal{Y}})(E) = 0$)
510 a data-distribution and denote it as $\mathcal{D}_{\mathcal{X}, \mathcal{Y}}$. If $(\mathcal{Z}, d_{\mathcal{Z}})$ is a metric space with metric d and $A \subset \mathcal{Z}$,

511 then for any $z \in \mathcal{Z}$, $d_{\mathcal{Z}}(z, A) = \inf_{a \in A} d_{\mathcal{Z}}(a, z)$. We say that a metric space is bounded by some
512 $B \in \mathbb{R}^{\geq 0}$, if $\sup_{x', x \in \mathcal{X}} d(x, x') < B$. Denote by $\text{Id}_{\mathcal{Z}} : \mathcal{Z} \rightarrow \mathcal{Z}$ the identity function of \mathcal{Z} . Let
513 $\mathcal{L} : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}^{\geq 0}$ be a loss function such that $\mathcal{L}(y, y') = 0$ if and only if $y = y'$. For any $f : \mathcal{X} \rightarrow \mathcal{Y}$
514 and $(x, y), (x', y') \in \mathcal{X} \times \mathcal{Y}$, we denote $\Delta \mathcal{L}_f^{(x, y)}(x', y') := |\mathcal{L}(f(x), y) - \mathcal{L}(f(x'), y')|$. Unless
515 otherwise specified, it will be assumed that f is a measurable function from \mathcal{X} to \mathcal{Y} with a metric
516 decomposition (see Def. 1).

517 **Lipschitz continuity.** Given two metric spaces $(\mathcal{X}, d_{\mathcal{X}}), (\mathcal{Y}, d_{\mathcal{Y}})$ a function $f : \mathcal{X} \rightarrow \mathcal{Y}$ is K -
518 Lipschitz continuous if there exists $K \in \mathbb{R}^{\geq 0}$ such that for all $x, x' \in \mathcal{X}$, $d_{\mathcal{Y}}(f(x), f(x')) \leq$
519 $K d_{\mathcal{X}}(x, x')$.

520 C Examples of Metric Decompositions

521 We present several common neural network architectures and some possible metric decompositions
522 for them.

523 *Example 1.* Suppose that $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is a fully-connected neural network with n hidden layers,
524 weight matrices W_i , biases b_i for $i \in [n]$, and activation functions σ_i for $i \in [n - 1]$. Let $r(W_i)$
525 denote the number of rows of W_i . Then, for any $1 < p \leq \infty$ consider the set of metric spaces
526 $(\mathbb{R}^n, \ell_p), (\mathbb{R}^{r(W_1)}, \ell_p), \dots, (\mathbb{R}^{r(W_n)}, \ell_p), (\mathbb{R}^m, \ell_p)$, we define the metric decomposition of h_i such
527 that $h_0 : \mathbb{R}^n \rightarrow \mathbb{R}^{r(W_1)}$, $h_i : \mathbb{R}^{r(W_i)} \rightarrow \mathbb{R}^{r(W_{i+1})}$, and $h_n : \mathbb{R}^{r(W_n)} \rightarrow \mathbb{R}^m$. Each of these functions
528 are simply the hidden layers in the fully-connected network. That is,

$$h_i(x) = \sigma_i(W_i x + b_i). \quad (\text{C.1})$$

529 *Example 2.* If f is a convolutional neural network, we can decompose it in the same way as before.
530 Except, h_i is now the convolution operation.

531 *Example 3.* We present two distinct metric decompositions of a residual network. Consider two
532 fully-connected layers A, B such that $x \xrightarrow{A} A(x) \xrightarrow{B} B(A(x)) \xrightarrow{x} B(A(x)) + x$. Here, the input x
533 feeds back into the output layer B creating a residual block (the set of layers between the input and
534 the residual connection).

535 We can aggregate each residual block as one metric decomposition. That is, let $h = B(A(x)) + x$.
536 Then, $x \xrightarrow{h} h(x)$. Clearly, this is the same function as before; moreover, we yield the metric decom-
537 position of $h(x)$. This is the approach we use in practice when dealing with residual connections.
538 Moreover, this is also the standard way to count layers in computer vision and natural language
539 processing.

540 We can also represent each layer within the residual block as a part of a metric decomposition. Define
541 $A' : x \mapsto (A(x), x)$, $B' : (A(x), x) \mapsto (B(A(x)), x)$ and $x' : (B(A(x)), x) \mapsto (B(A(x)) + x, x)$.
542 Then, it follows that $x \rightarrow A' \rightarrow B' \rightarrow x'$ forms a metric decomposition. Here, the metric in each
543 layer is with respect to the quotient space where $(a, a') \sim (b, b')$ if and only if $a = b$. Therefore, we
544 also recover the same vector space structure.

545 We again emphasize that any particular metric decomposition does not affect our theoretical results.
546 Our propositions and theorems only rely on the fact that a metric decompositions exist.

547 D Proof of Robustness

548 **Lemma D.1.** Let (X, d) be a metric space. Suppose that $x \in X$ and $f_x(z) = d(x, z)$. Then, f_x is
549 1-Lipschitz with respect to the metric d . Moreover, if $A \subset X$ and $f_A(z) = \inf_{a \in A} d(x, a)$. Then, f_A
550 is also 1-Lipschitz.

551 *Proof.* Fix some $x \in X$. By the definition of 1-Lipschitzness, it suffices to show that for all $z, y \in X$,
 552 $|f_x(z) - f_x(y)| \leq d(z, y)$. Thus,

$$\begin{aligned} |f_x(z) - f_x(y)| &= |d(x, z) - d(x, y)|, \\ &= |d(x, z) + d(z, y) - d(x, y) - d(x, y)|, \\ &\leq |d(x, y) - d(z, y) - d(x, y)|, \\ &\leq d(z, y). \end{aligned}$$

553 The latter statement follows from the same argument above with obvious modification. \blacksquare

554 Next, we state the McDiarmid's Inequality [40] and Lévy's Inequalities [5] without proof.

555 **Definition 6.** A function $f : X_1 \times X_2 \times \dots \times X_n \rightarrow \mathbb{R}$ satisfies the bounded differences property if
 556 there are constants c_1, c_2, \dots, c_n such that for all $1 \leq i \leq n$ and $x_1 \in X_1, x_2 \in X_2, \dots, x_n \in X_n$,

$$\sup_{x'_i \in X_i} |f(x_1, \dots, x_i, \dots, x_n) - f(x_1, \dots, x'_i, \dots, x_n)| \leq c_i. \quad (\text{D.1})$$

557 **Theorem D.2** (McDiarmid's Inequality). Assume that the function $f : X_1 \times X_2 \times \dots \times X_n \rightarrow \mathbb{R}$
 558 satisfy the bounded differences property with bounds c_1, \dots, c_n . Consider the independent random
 559 variables Y_1, \dots, Y_n where $Y_i \in X_i$ for all $1 \leq i \leq n$. Then, for any $\varepsilon > 0$,

$$\mathbb{P}[f(Y_1, \dots, Y_n) - \mathbb{E}[f(Y_1, \dots, Y_n)] \geq \varepsilon] \leq \exp\left(-\frac{2\varepsilon^2}{\sum_{i=1}^n c_i^2}\right), \quad (\text{D.2})$$

$$\mathbb{P}[f(Y_1, \dots, Y_n) - \mathbb{E}[f(Y_1, \dots, Y_n)] \leq -\varepsilon] \leq \exp\left(-\frac{2\varepsilon^2}{\sum_{i=1}^n c_i^2}\right). \quad (\text{D.3})$$

560 Let (X, d) be a metric space and $f : X \rightarrow \mathbb{R}$ be a Lipschitz function f with Lipschitz constant C .
 561 Consider the measure space formed by X , the σ -algebra of all Borel sets of X and a probability
 562 measure \mathbb{P} . Let Y be a random variable taking values in X and distributed according to \mathbb{P} .

563 **Definition 7** (Concentration Functions). For all $t > 0$, the concentration functions of X is defined by

$$\alpha(t) = \sup_{A \subset X: \mathbb{P}(A) \geq 1/2} \mathbb{P}[d(Y, A) \geq t], \quad (\text{D.4})$$

564 where $d(Y, A) = \inf_{x \in A} d(x, Y)$.

565 Informally, the concentration function $\alpha(t)$ represents the scatter of the random variable Y in the
 566 metric space. Specifically, for a fixed $t > 0$, if $\alpha(t) \approx 0$, then Y is dispersed throughout the metric
 567 space since for any subset of X that has significant probability mass, with high probability Y is t
 568 away from this subset.

569 **Theorem D.3** (Lévy's Inequalities). For any Lipschitz function f with Lipschitz constant $C > 0$,

$$\mathbb{P}[f(Y) \geq \mathbf{M}f(Y) + t] \leq \alpha\left(\frac{t}{C}\right) \quad \text{and} \quad \mathbb{P}[f(Y) \leq \mathbf{M}f(Y) - t] \leq \alpha\left(\frac{t}{C}\right), \quad (\text{D.5})$$

570 where $\mathbf{M}f(Y)$ is the median of $f(Y)$. That is, $\mathbb{P}[f(Y) \geq \mathbf{M}f(Y)] \geq 1/2$ and $\mathbb{P}[f(Y) \leq$
 571 $\mathbf{M}f(Y)] \geq 1/2$.

572 *Proof.* Here, we directly quote a proof from [5]. Consider the set $A = \{x \in X : f(x) \leq \mathbf{M}f(Y)\}$.
 573 By the definition of a median, $\mathbb{P}[A] \geq 1/2$. On the other hand, by the Lipschitz property of f ,

$$A_t = \left\{x \in X : d(x, A) \leq \frac{t}{C}\right\} \subset \left\{x \in X : f(x) < \mathbf{M}f(Y) + \frac{t}{C}\right\}.$$

574 The inequalities now follow from the definition of the concentration function (the second follows
 575 from the first by considering $-f$). \blacksquare

576 Leveraging Thm. D.2 and Thm. D.3, we can bound the distance between a fixed non-measure zero
 577 set and a point. Other than the boundedness of the metric space, we assume all the same notation as
 578 before.

579 **Lemma D.4.** Suppose that (X, d) is a bounded metric space such that $\sup_{x, x' \in X} d(x, x') < B$ for
 580 some $B > 0$. Let $A \subset X$ such that $\mathbb{P}[A] > 0$ and $\delta > 0$. Then,

$$\mathbb{P}[d(Y, A) \geq \delta] \leq \exp\left(-\frac{2}{B^2} \left(\delta - B \sqrt{\frac{1}{2} \log \frac{1}{\mathbb{P}[A]}}\right)^2\right).$$

581 *Proof.* Let $f_A(x) = d(x, A)$. Then, by Lem. D.1, $f_A(\cdot)$ is 1-Lipschitz with respect to d . Since the
 582 metric space is bounded by constant $B > 0$, f satisfies the bounded differences property with constant
 583 B . By Theorem D.2,

$$\mathbb{P}[\mathbb{E}[f_A(Y)] - f_A(Y) \geq \delta] \leq e^{-2\delta^2/B^2}. \quad (\text{D.6})$$

584 If $\delta = \mathbb{E}[f_A(Y)]$, then the left-hand side becomes $\mathbb{P}[f_A(Y) \leq 0] \geq \mathbb{P}[A]$. Therefore, by the
 585 previous inequality,

$$\mathbb{P}[A] \leq \mathbb{P}[f_A(Y) \leq 0], \quad (\text{D.7})$$

$$\mathbb{P}[A] \leq \exp(-2 \mathbb{E}[f_A(Y)]^2 / B^2), \quad (\text{D.8})$$

$$\mathbb{E}[f_A(Y)] \geq B \sqrt{\frac{1}{2} \log \frac{1}{\mathbb{P}[A]}}. \quad (\text{D.9})$$

586 Therefore,

$$\mathbb{P}\left[d(Y, A) \geq \underbrace{\delta + B \sqrt{\frac{1}{2} \log \frac{1}{\mathbb{P}[A]}}}_{(a)}\right] \leq e^{-2\delta^2/B^2}. \quad (\text{D.10})$$

587 The statement of the theorem then follows directly from a substitution of the term labeled (a)
 588 above. ■

589 Now, we are ready to prove the main regarding robustness (Thm. 2.1). For coherence, we restate the
 590 statement of the theorem before detailing the proof.

591 **Theorem.** Let $A \subset \mathcal{X} \times \mathcal{Y}$ such that $\mathbb{P}_{\mathcal{X}}[A] > 0$ and $\delta, \eta > 0$. Let $A' = \{(x', y') \in \mathcal{X} \times \mathcal{Y} :$
 592 $\exists(x, y) \in A, \Delta \mathcal{L}_f^{(x, y)}(x', y') > \eta\}$. If $f : \mathcal{X} \rightarrow \mathcal{Y}$ is ε -knowledge continuous with respect to the
 593 hidden layer indexed by k and (Z_k, d_k) is bounded by $B > 0$, then

$$\mathbb{P}_{(x, y) \sim \mathcal{D}_{\mathcal{X}, \mathcal{Y}}}[A' \mid d_k(f^k(x), f^k(A)) < \delta] \leq \frac{\varepsilon \delta}{\eta \left(1 - \exp\left[-\Omega\left(\frac{\delta}{B} - \sqrt{\log \frac{1}{\mathbb{P}[A]}}\right)^2\right]\right)}. \quad (\text{D.11})$$

594 *Proof.* By the definition of conditional probability, we have that

$$\mathbb{P}_{(x', y') \sim \mathcal{D}}[A' \mid d_k(f^k(x), f^k(x')) < \delta] = \frac{\mathbb{P}_{(x', y') \sim \mathcal{D}}[A' \text{ and } d_k(f^k(x), f^k(x')) < \delta]}{\mathbb{P}_{(x', y') \sim \mathcal{D}}[d_k(f^k(x), f^k(x')) < \delta]}. \quad (\text{D.12})$$

595 We start by bounding the numerator of Eq. D.12. By the definition of ε -knowledge continuity,

$$\mathbb{E} \sigma_f^k(x, y) = \iint \frac{\Delta \mathcal{L}_f^{(x, y)}(x', y')}{d_k(f^k(x), f^k(x'))} d(\mathbb{P} \times \mathbb{P}), \quad (\text{D.13})$$

$$\geq \iint_{d_k(f^k(x), f^k(x')) < \delta} \frac{\Delta \mathcal{L}_f^{(x, y)}(x', y')}{d_k(f^k(x), f^k(x'))} d(\mathbb{P} \times \mathbb{P}), \quad (\text{D.14})$$

$$\geq \frac{1}{\delta} \iint_{d_k(f^k(x), f^k(x')) < \delta} \Delta \mathcal{L}_f^{(x, y)}(x', y') d(\mathbb{P} \times \mathbb{P}), \quad (\text{D.15})$$

$$\delta \mathbb{E} \sigma_f^k(x, y) \geq \iint_{\substack{d_k(f^k(x), f^k(x')) < \delta \\ (x, y) \in A}} \Delta \mathcal{L}_f^{(x, y)}(x', y') d(\mathbb{P} \times \mathbb{P}). \quad (\text{D.16})$$

596 This gives us an upper-bound of expectation of $\Delta\mathcal{L}_f^{(x,y)}(x', y')$ over the set of all points that are
 597 within δ -radius from A . Next, by Markov's inequality,

$$\mathbb{P}[A' \text{ and } d_k(f^k(x), f^k(x')) < \delta] \leq \frac{\delta \mathbb{E} \sigma_f^k(x, y)}{\eta}, \quad (\text{D.17})$$

$$\leq \frac{\delta \varepsilon}{\eta}. \quad (\text{D.18})$$

598 The last inequality follows from the fact that f is ε -knowledge continuous. Now, by applying the
 599 complement of Lem. D.4, we lower-bound the denominator and yield the following

$$\mathbb{P}_{(x', y') \sim \mathcal{D}} [A' \mid d_k(f^k(x), f^k(x')) < \delta] \leq \frac{\varepsilon \delta}{\eta \left(1 - \exp \left(-\frac{2}{B^2} \left(\delta - B \sqrt{\frac{1}{2} \log \frac{1}{\mathbb{P}[A]}} \right)^2 \right) \right)}. \quad (\text{D.19})$$

600 The proof is concluded by applying big-Omega notation to the exponentiated. ■

601 **Corollary D.5.** *If (Z_k, d_k) is unbounded, then*

$$\mathbb{P}_{(x, y) \sim \mathcal{D}_{\mathcal{X}, \mathcal{Y}}} [A' \mid d_k(f^k(x), f^k(A)) < \delta] \leq \frac{\varepsilon \delta}{\eta(1 - \mathbb{P}[A])}. \quad (\text{D.20})$$

602 *If $\mathbb{P}[A] = 0$, then*

$$\mathbb{P}_{(x, y) \sim \mathcal{D}_{\mathcal{X}, \mathcal{Y}}} [A' \mid d_k(f^k(x), f^k(A)) < \delta] \leq \frac{\varepsilon \delta}{\eta}. \quad (\text{D.21})$$

603 *Proof.* These results follow from directly taking the limit as $B \rightarrow \infty$ and applying some of the
 604 bounds acquired in the proof of Thm. 2.1. This yields Eq. D.20. Next, setting $\mathbb{P}[A] = 0$ easily results
 605 in Eq. D.21. ■

606 E Proof of Expressiveness

607 Here, we show the main result regarding the expressiveness of ε -knowledge continuous estimators
 608 (Prop. 2.3). For completeness, we restate the statement of the proposition before proceeding with the
 609 proof.

610 **Proposition.** *Suppose $(\mathcal{X}, d_{\mathcal{X}}), (\mathcal{Y}, d_{\mathcal{Y}}) := (\mathcal{X}, d_{\mathcal{X}})$ are **compact** metric spaces, $\mathcal{F} \subset \mathcal{Y}^{\mathcal{X}}$ is the
 611 **set of all continuous functions** from \mathcal{X} to \mathcal{Y} such that $\int d_{\mathcal{X}}(x, x')^{-1} d\mu_f < \infty$ and \mathcal{L} be Lipschitz
 612 continuous in both coordinates. Then, there exists a universal function approximator \mathcal{U} of \mathcal{F} that is
 613 knowledge continuous (i.e. $\mathbb{E} \sigma_f^k(x, y) < \infty$ for some k).*

614 *Proof.* By assumption, $\mathcal{X} = \mathcal{Y}$ and $d_{\mathcal{X}} = d_{\mathcal{Y}}$. First, we consider the set of all Lipschitz continuous
 615 functions from $\mathcal{X} \rightarrow \mathcal{X}$. Clearly, the set of all Lipschitz continuous functions separate points in \mathcal{X}
 616 by the fact that the $d_{\mathcal{X}}$ is Lipschitz continuous (see Lem. D.1). Thus, since \mathcal{X} is compact, by the
 617 Stone-Weierstrass Theorem [52] the set of Lipschitz continuous functions must be dense in the set of
 618 all continuous functions from \mathcal{X} to \mathcal{X} . This implies that for any sequence $\varepsilon_1, \varepsilon_2, \dots$ we can choose
 619 Lipschitz continuous functions f_1, f_2, \dots such that $\int \mathcal{L}(f_n(x), y) d\mu_f < \varepsilon_n$. It remains to show that
 620 each of these functions are in fact knowledge continuous. Since \mathcal{X} is a metric space, we consider the
 621 trivial metric decomposition of our sequence of functions (see Remark ??). Specifically, we denote

622 $h_1 = \text{Id}_{\mathcal{X}}$ and proceed to bound $\mathbb{E} \sigma_f^1(x, y)$.

$$\mathbb{E} \sigma_{f_n}^1(x, y) = \iint \frac{\Delta \mathcal{L}_{f_n}^{(x,y)}(x', y')}{d_{\mathcal{X}}(x, x')} (d\mu_f \times d\mu_f), \quad (\text{E.1})$$

$$\leq \iint \frac{|\mathcal{L}(f_n(x), y) - \mathcal{L}(f_n(x'), y) + \mathcal{L}(f_n(x'), y) - \mathcal{L}(f_n(x'), y')|}{d_{\mathcal{X}}(x, x')} (d\mu_f \times d\mu_f), \quad (\text{E.2})$$

$$\leq \iint \frac{|\mathcal{L}(f_n(x), y) - \mathcal{L}(f_n(x'), y)|}{d_{\mathcal{X}}(x, x')} d(\mu_f \times \mu_f) \quad (\text{E.3})$$

$$+ \iint \frac{|\mathcal{L}(f_n(x'), y) - \mathcal{L}(f_n(x'), y')|}{d(x, x')} (d\mu_f \times d\mu_f), \quad (\text{E.4})$$

$$\leq \iint \frac{L d_{\mathcal{X}}(f(x), f(x'))}{d_{\mathcal{X}}(x, x')} d(\mu_f \times \mu_f) + \iint \frac{L d_{\mathcal{X}}(y, y')}{d_{\mathcal{X}}(x, x')} d(\mu_f \times \mu_f), \quad (\text{E.5})$$

$$\leq \iint LK d(\mu_f \times \mu_f) + LB \int \frac{1}{d_{\mathcal{X}}(x, x')} d\mu_f, \quad (\text{E.6})$$

$$= LK + LB \int d_{\mathcal{X}}(x, x')^{-1} d\mu_f, \quad (\text{E.7})$$

623 where L is the Lipschitz constant of \mathcal{L} , K is the Lipschitz constant of the f_n , and B bounds the metric
 624 space \mathcal{X} (since any compact metric space is bounded). The remaining assumption in the proposition
 625 concludes the proof of the proposition. \blacksquare

626 F Proof of Equivalence Between Lipschitz Continuity and Knowledge 627 Continuity

628 We present the proofs of the results that establish conditions when knowledge continuity implies
 629 Lipschitz continuity and vice versa. As before, we restate all of the statements before providing their
 630 proof. First, we identify conditions under which knowledge continuity implies Lipschitz continuity
 631 (Prop. 2.4).

632 **Proposition.** *Suppose that $(\mathcal{X}, d_{\mathcal{X}})$, $(\mathcal{Y}, d_{\mathcal{Y}})$ are metric spaces. Let the first n metric decompositions
 633 of $f : \mathcal{X} \rightarrow \mathcal{Y}$ be K_i -Lipschitz continuous, for $i \in [n]$. If f is ε -knowledge continuous with respect
 634 to the n^{th} hidden layer and $d_{\mathcal{Y}}(f(x), f(x')) \leq \eta \Delta \mathcal{L}_f^{(x,y)}(x', y)$ for all $x, x' \in \mathcal{X}$, $y \in \mathcal{Y}$, and some
 635 $\eta > 0$, then f is Lipschitz continuous in expectation. That is,*

$$\mathbb{E}_{(x,y),(x',y') \sim \mathcal{D}_{\mathcal{X},\mathcal{Y}}} \frac{d_{\mathcal{Y}}(f(x), f(x'))}{d_{\mathcal{X}}(x, x')} \leq \varepsilon \eta \prod_{j=1}^n K_j. \quad (\text{F.1})$$

636 *Proof.* We proceed to bound the knowledge continuity of f from below.

$$\mathbb{E} \sigma_f^k(x, y) \geq \mathbb{E}_{(x,y) \sim \mathcal{D}_{\mathcal{X},\mathcal{Y}}} \mathbb{E}_{\substack{(x',y') \sim \mathcal{D}_{\mathcal{X},\mathcal{Y}} \\ y'=y}} \frac{\Delta \mathcal{L}_f^{(x,y)}(x', y)}{d_k(f^k(x), f^k(x'))}, \quad (\text{F.2})$$

$$\geq \mathbb{E}_{(x,y) \sim \mathcal{D}_{\mathcal{X},\mathcal{Y}}} \mathbb{E}_{\substack{(x',y') \sim \mathcal{D} \\ y'=y}} \frac{\Delta \mathcal{L}_f^{(x,y)}(x', y)}{\prod_{j=1}^n K_j d_{\mathcal{X}}(x', x)}, \quad (\text{F.3})$$

$$\geq \mathbb{E}_{(x,y) \sim \mathcal{D}_{\mathcal{X},\mathcal{Y}}} \mathbb{E}_{\substack{(x',y') \sim \mathcal{D} \\ y'=y}} \frac{\frac{1}{\eta} d_{\mathcal{Y}}(f(x), f(x'))}{\prod_{j=1}^n K_j d_{\mathcal{X}}(x, x')}, \quad (\text{F.4})$$

$$= \mathbb{E}_{(x,y),(x',y') \sim \mathcal{D}_{\mathcal{X},\mathcal{Y}}} \frac{\frac{1}{\eta} d_{\mathcal{Y}}(f(x), f(x'))}{\prod_{j=1}^n K_j d_{\mathcal{X}}(x, x')}. \quad (\text{F.5})$$

637 Eq. F.2 comes from the fact that we take the expectation only over pairs of points $(x, y), (x', y')$
 638 where $y = y'$ and also because the summand is always nonnegative. Then, we inductively apply the

639 definition of K_i -Lipschitz continuity to yield Eq. F.3. Eq. F.4 follows directly from the assumption
 640 in the statement of the proposition. Since the expression in Eq. F.4 now has no dependence on the
 641 label distribution, we may expand the expectation which results in Eq. F.5. Lastly, by the definition
 642 of ε -knowledge continuity,

$$\varepsilon \geq \mathbb{E}_{(x,y),(x',y') \sim \mathcal{D}_{\mathcal{X},\mathcal{Y}}} \frac{\frac{1}{\eta} d_{\mathcal{Y}}(f(x), f(x'))}{\prod_{j=1}^n K_j d_{\mathcal{X}}(x, x')},$$

$$\varepsilon \eta \prod_{j=1}^n K_j \geq \mathbb{E}_{(x,y),(x',y') \sim \mathcal{D}_{\mathcal{X},\mathcal{Y}}} \frac{d_{\mathcal{Y}}(f(x), f(x'))}{d_{\mathcal{X}}(x, x')},$$

643 and this concludes the proof of the proposition. \blacksquare

644 To prove Cor. 2.5, we need the following auxiliary result from [74].

645 **Proposition F.1.** *For a neural network $f : \mathbb{R}^n \rightarrow \mathbb{R}^K$ with Lipschitz constant L under ℓ_p -norm,
 646 define the resulting classifier as $g(x) := \arg \max_{k \in [K]} f_k(x)$ for an input x . Then, g is provably
 647 robust under perturbations $\|\delta\|_p < \frac{\sqrt[p]{2}}{2L} \text{margin}(f(x))$, i.e.*

$$g(x + \delta) = g(x) \quad \text{for all } \|\delta\|_p < \frac{\sqrt[p]{2}}{2L} \text{margin}(f(x)). \quad (\text{F.6})$$

648 Here, $\text{margin}(f(x))$ is the difference between the largest and second largest output logit.

649 The following proof is from [74].

650 *Proof.* Let $f_j(x)$ denote the j^{th} coordinate of $f(x)$. We proceed by way of contraposition. Suppose
 651 that $g(x) \neq g(x + \delta)$ for some $\delta \in \mathbb{R}^n$. We show that $\|\delta\|_p \geq \frac{\sqrt[p]{2}}{2L} \text{margin}(f(x))$. Let $g(x) = \alpha$ and
 652 $g(x + \delta) = \beta$. Then,

$$\|f(x + \delta) - f(x)\|_p = \left(\sum_{k=1}^K |f_k(x + \delta) - f(x)_k|^p \right)^{1/p}, \quad (\text{F.7})$$

$$\geq (|f_{\alpha}(x + \delta) - f_{\alpha}(x)|^p + |f_{\beta}(x + \delta) - f_{\beta}(x)|^p)^{1/p}. \quad (\text{F.8})$$

653 The minimum of Eq. F.8 is achieved when $f_{\alpha}(x + \delta) = f_{\beta}(x + \delta) = (f_{\alpha}(x) + f_{\beta}(x))/2$. Then,
 654 through a direct substitution we have that

$$\|f(x + \delta) - f(x)\|_p \geq \frac{\sqrt[p]{2}}{2} (f_{\alpha}(x) - f_{\beta}(x)), \quad (\text{F.9})$$

655 by the definition of $\text{margin}(f(x))$, $f_{\alpha}(x) - f_{\beta}(x) \geq \text{margin}(f(x))$. Lastly, by the definition of
 656 L -Lipschitz continuity, we have that

$$L \|\delta\|_p \geq \|f(x + \delta) - f(x)\|_p \geq \frac{\sqrt[p]{2}}{2} \text{margin}(f(x)). \quad (\text{F.10})$$

657 Rearranging this expression results in the proposition. \blacksquare

658 We are now ready for the proof of Cor. 2.5. We simply Prop. F.1 in conjunction with Markov's
 659 inequality to bound the Lipschitz constant.

660 **Corollary.** *Suppose that assumptions of Prop. 2.4 are true. And also assume that $(\mathcal{X}, d_{\mathcal{X}}) = (\mathbb{R}^n, \ell_p)$,
 661 $(\mathcal{Y}, d_{\mathcal{Y}}) = (\mathbb{R}^m, \ell_p)$, for $1 \leq p \leq \infty$. Define a classifier from $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$, g , where $g(x) :=$
 662 $\arg \max_{k \in [m]} f_k(x)$ for any $x \in \mathbb{R}^n$. Then, with probability $1 - \frac{\varepsilon \eta}{t} \prod_{j=1}^n K_j$, $g(x) = g(x + \delta)$
 663 for all $\|\delta\|_p < \frac{\sqrt[p]{2}}{2t} \text{margin}(f(x))$ and $t > 0$. $f_k(x)$ is the k^{th} coordinate of $f(x)$ and $\text{margin}(f(x))$
 664 denotes the difference between the largest and second-largest output logits.*

665 *Proof.* By Prop. 2.4, we have that

$$\mathbb{E}_{(x,y),(x',y') \sim \mathcal{D}_{\mathcal{X},\mathcal{Y}}} \frac{d_{\mathcal{Y}}(f(x), f(x'))}{d_{\mathcal{X}}(x, x')} \leq \varepsilon \eta \prod_{j=1}^n K_j. \quad (\text{F.11})$$

666 By Markov's inequality,

$$\mathbb{P}_{(x,y),(x',y') \sim \mathcal{D}_{\mathcal{X},\mathcal{Y}}} \left[\frac{d_{\mathcal{Y}}(f(x), f(x'))}{d_{\mathcal{X}}(x, x')} \geq t \right] \leq \frac{\varepsilon \eta}{t} \prod_{j=1}^n K_j. \quad (\text{F.12})$$

667 We yield the corollary by directly applying Prop. F.1 assuming that f is t -Lipschitz continuous. ■

668 Next, we establish conditions under which Lipschitz continuity implies knowledge continuity
669 (Prop. 2.6).

670 **Proposition.** Let $(\mathcal{X}, d_{\mathcal{X}}), (\mathcal{Y}, d_{\mathcal{Y}})$ be a metric spaces. Let $f : \mathcal{X} \rightarrow \mathcal{Y}$ be ε -Lipschitz continuous
671 and $\mathcal{L}(f(x), y)$ be η -Lipschitz continuous with respect to both coordinates. If the first n metric
672 decompositions of f are K_i -Lipschitz continuous, then f is knowledge continuous with respect to the
673 n^{th} hidden layer. That is,

$$\mathbb{E}_{(x,y) \sim \mathcal{D}_{\mathcal{X},\mathcal{Y}}} \sigma_f^n(x, y) \leq \varepsilon \eta \prod_{j=1}^n \frac{1}{K_j}. \quad (\text{F.13})$$

674 *Proof.* Let us start with the definition of ε -Lipschitz continuity and lower-bound it. For any
675 $(x, y), (x', y') \in \mathcal{X} \times \mathcal{Y}$,

$$\frac{d_{\mathcal{Y}}(f(x), f(x'))}{d_{\mathcal{X}}(x, x')} \leq \varepsilon, \quad (\text{F.14})$$

$$\frac{d_{\mathcal{Y}}(f(x), f(x'))}{\prod_{j=1}^n \frac{1}{K_j} d_k(f^k(x), f^k(x'))} \leq \varepsilon, \quad (\text{F.15})$$

$$\frac{\frac{1}{\eta} |\mathcal{L}(x, y) - \mathcal{L}(x', y')|}{\prod_{j=1}^n \frac{1}{K_j} d_k(f^k(x), f^k(x'))} \leq \varepsilon, \quad (\text{F.16})$$

$$\frac{|\mathcal{L}(x, y) - \mathcal{L}(x', y')|}{d_k(f^k(x), f^k(x'))} \leq \varepsilon \eta \prod_{j=1}^n \frac{1}{K_j}. \quad (\text{F.17})$$

676 Eq. F.15 follows from inductively applying the definition of Lipschitz continuity on the metric
677 decompositions of f . Specifically, $d_{i+1}(f^{i+1}(x), f^{i+1}(x')) \leq K_i d_i(f^i(x), f^i(x'))$. Then, by the
678 Lipschitz continuity of \mathcal{L} in both coordinates we yield Eq. F.16. Since the Lebesgue integral preserves
679 order, Eq. F.17 directly implies the statement of the proposition and this concludes the proof. ■

680 G Predicting Adversarial Robustness with Volatility

681 As discussed in Section 3, we regress k -volatility scores for a variety of models across all layers
682 against their empirical adversarial robustness. Herein, we describe this experimental procedure and
683 detail the results. Throughout this section, we adopt the shorthand $\text{KVS} := \mathbb{E} \sigma_f^k(x, y)$ and refer to
684 this as the knowledge volatility score.

685 We run all our experiments against the IMDB dataset [39] with TextFooler [29] as the benchmark
686 adversarial attack. We run linear regression to predict the number of successful adversarial attacks,
687 using model type and model size. We then incorporate our vulnerability score, calculated over all
688 layers, and notice how our R^2 changes.

689 For our linear regression, we use the LinearRegression class from sklearn (version 1.3.2), and default
690 hyperparameters ($\alpha = 1.0$, $\text{max_iter} = 1000$). To calculate the number of adversarial attacks, we
691 use TextFooler algorithm [29] on a holdout test set with respect to a pretrained model. We say that
692 an adversarial attack is successful if the model previously characterized it correctly, but under the
693 perturbation of TextFooler, the model now classifies it incorrectly.

694 For each model, we lay out our features as follows:

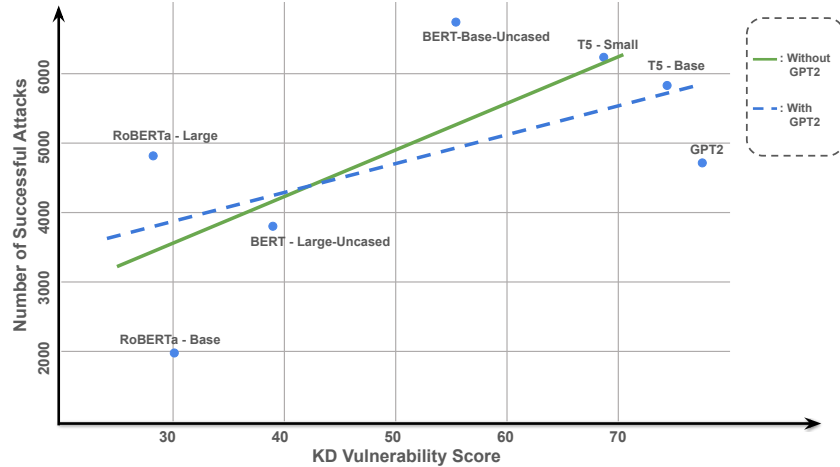


Figure 2: Regression analysis of Knowledge discontinuity vs. number of successful attacks under TextFooler. Knowledge discontinuities alone can explain 35% of the variance of successful adversarial attacks against a model ($R^2 = 0.35$). The line of best fit is given by $SuccessfulAttacks = 49(KVS) + 2254$.

- 695 1. A 0 or 1 representing whether this model is encoder-only
- 696 2. A 0 or 1 representing whether this model is decoder-only
- 697 3. A 0 or 1 representing whether this model is encoder-decoder
- 698 4. A floating point representing the natural log of the number of parameters in this model
- 699 5. The vulnerability score associated with this model

700 For example, the following vector represents bert-large:

$$[1, 0, 0, 19.630, 54.044]$$

701 We choose to use the logarithm of the model size. Intuitively, we expect that past a certain size,
 702 a well-trained model will perform so well that it essentially masters the task, and there is little
 703 adversarial robustness to be gained by adding more parameters.

704 After running our linear regression, we proceed to obtain the coefficients, and then calculate the per-
 705 mutation importance of each of our features. We get the following results below for our coefficients:

	Without vulnerability score	With vulnerability score
encoder	-548.43	1484.91
decoder	-556.89	-2816.49
encoder-decoder	1105.32	1331.57
ln(num_params)	-362.59	65.50
vulnerability score	N/A	95.74

707 We calculate importance values using 100 random permutations. We ultimately get the following
 708 table:

	Without vulnerability score	With vulnerability score
encoder	0.0652	0.403
decoder	0.0195	0.712
encoder-decoder	0.177	0.291
ln(num_params)	0.0442	-6.08e-05
vulnerability score	—	2.57
R^2	0.282	0.479

710 Notice the importance of the vulnerability score, especially in proportion to the other features. Clearly,
 711 this illustrates both the predictive power and importance of our vulnerability score.

712 H Localizing Volatile Hidden Representations

713 We seek to localize volatile hidden representations, both in the sense of which layers are more volatile,
714 and which areas of the representation space for a given layer are more volatile. We consider the same
715 selection of models in Appendix G, the same dataset (IMDB), and the same attack (TextFooler).

716 Throughout this section, we adopt the shorthand $KVS := \mathbb{E} \sigma_f^k(x, y)$ and refer to this as the knowledge
717 volatility score.

718 H.1 Per-Layer Volatility

719 We start by plotting the KVS for each of our models, against the actual number of successful
720 adversarial attacks. We use this as a proxy for analyzing volatility, since the more volatile, the higher
721 the correlation between these two variables.

722 Then, to analyze this on a per-layer basis, we notice that KVS can be calculated independently for
723 any given layer, since each layer emits its own distance metric.

724 Thus, we ultimately plot R^2 vs relative depth for our given models. We notice that the foremost and
725 final hidden layers are most explanatory (see Fig. 4). However, we see that GPT2 admits a surprising
726 behavior, in that its middle hidden layers are most participatory in adversarial vulnerability. We
727 now specifically look at this as a case study. To do this, we repeat the experiments in Appendix G
728 across 10 relative depths and plot the R^2 with and without GPT2 (see Fig. 3). Indeed, without GPT2
729 we see that the trend of R^2 seems to be more linear. These results directly inform the choice of
730 hyperparameters in Appendix I since we want to minimize KVS only over the highly salient layer,
731 rather than all of them.

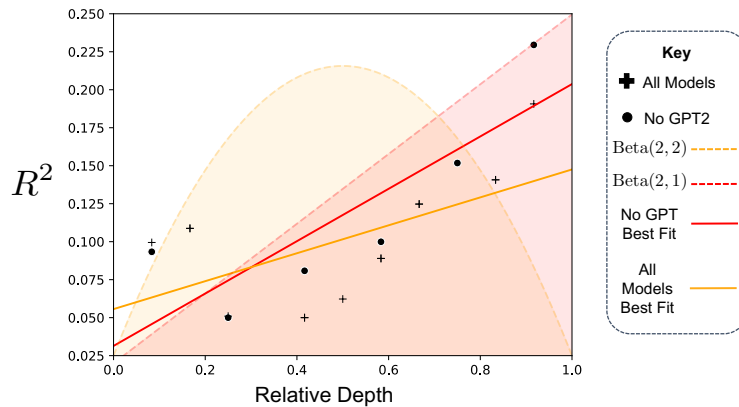


Figure 3: The explained variance of knowledge continuities for each relative depth across all models and without GPT2. The distribution of points warrant the use of various parameterizations of the Beta distribution in Alg. 2 for different models.

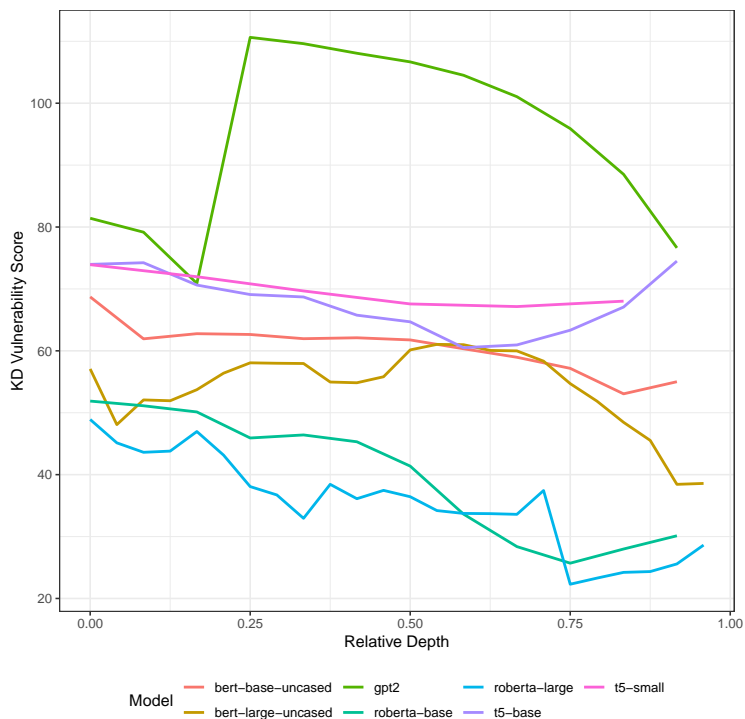


Figure 4: The KVS versus relative depth for BERTBase-Uncased, BERT-Large-Uncased, T5-Base, T5-Small, RoBERTa-Base, RoBERTa-Large, GPT2. Notice that we can track which layers are responsible for what portion of the vulnerability score of each model. Notice that GPT2 has a spike toward the middle, and teeters out toward the end— perhaps this is because the deeper layers are responsible for decoding, and have less of an effect on classification performance. Such a plot could be a useful for both practical applications and future research, as a computationally efficient method to roughly assess how different layers may contribute to adversarial vulnerability.

732 H.2 Per-Model Volatility

733 We start by exploring the KVS of each of our test models. We notice that KVS cannot be predicted
 734 by surface-level features such as size or model type alone. This is shown clearly in Fig. 5. Yet, as
 735 discussed in Appendix G, it is still able to predict actual adversarial vulnerability with moderate
 736 power. Thus, we conjecture that KVS captures a complex aspect of the model’s vulnerability which
 737 cannot be solely attributed to its size or type.

738 I Regularizing Knowledge Continuity

739 In this section, we provide a comprehensive overview of regulating knowledge continuity to achieve
 740 robustness. We first show a simple algorithm that estimates k -volatility. Then, we demonstrate how
 741 this can be used to augment any loss function to achieve regularization. We present some theoretical
 742 guarantees that revolve around the unbiasedness of our estimation algorithm and some guarantees
 743 of its rate of convergence. Lastly, we present detailed discussion of the results shown in Table 1
 744 including training details and ablation studies over the hyperparameters.

745 I.1 Estimating Knowledge Continuity Algorithmically

746 We first present a method for estimating the knowledge continuity of a hidden representation space.
 747 This is shown in Alg. 1. In the following subsection, we provide some guidance to choosing the
 748 subsampling hyperparameter M . In theory, one should choose $M = N$. However, if $N \gg 1$, this can
 749 become quickly intractable. Therefore, we multiplicatively bound the error of the unbiased estimator
 750 with respect to M and the variance of k -volatility. As discussed in the main text, the choice of metric

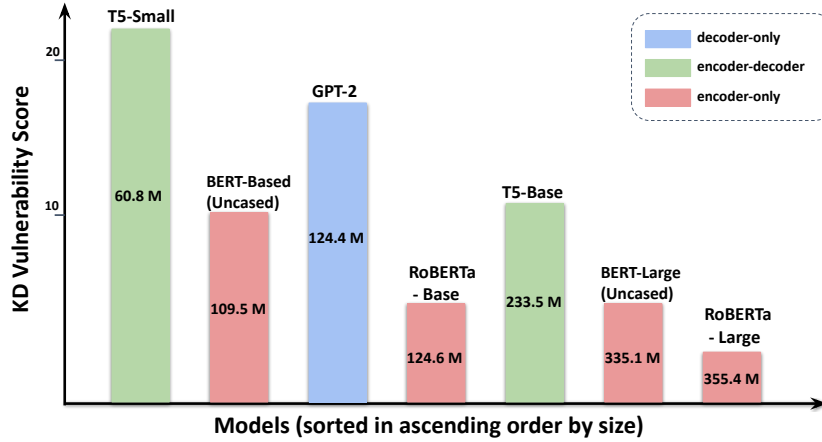


Figure 5: The KVS of each model, in the ascending order of model size. As shown, a model’s KVS cannot be solely attributed to its size or type.

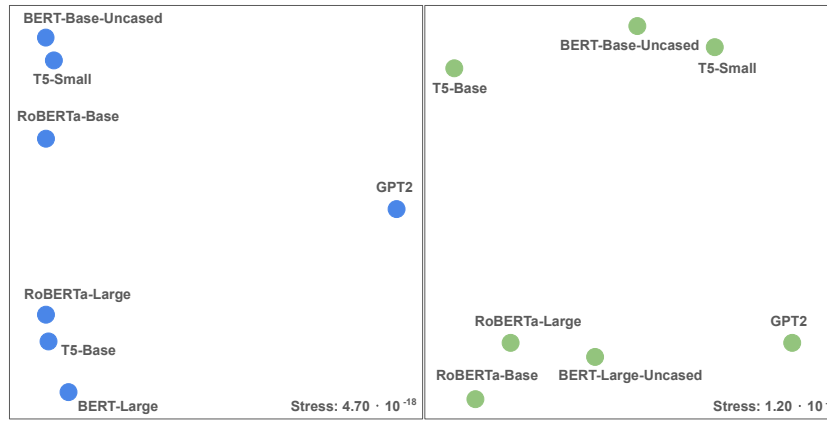


Figure 6: Left: Actual Adversarial Attacks, Right: Predicted Vulnerabilities using KVS

751 (or representation space) which we enforce knowledge continuity against is crucial as it determines
 752 the type of robustness we will achieve. Therefore, in Alg. 2, we incorporate this detail by sampling
 753 the index of the hidden layer using some Beta distribution specified by hyperparameters α, β . Note
 754 that we choose the Beta distribution for simplicity, however, it can be replaced by any distribution
 755 like a mixture of Gaussians.

756 In contrast to existing adversarial training methods such as [26] and [51] which only use the embed-
 757 dings, our algorithm gives the practitioner more control over which hidden layer (or distance metric)
 758 to enforce smoothness. In this way, if the practitioner has some knowledge *a priori* of the attacker’s
 759 strategy, they may choose to optimize against the most suitable metric. We present a brief discussion
 760 of the various tradeoffs when choosing α, β in the following section as well as a detailed empirical
 761 analysis in the following subsections. λ is the weight we put on the regularizer in relation to the loss
 762 function \mathcal{L} . We provide a detailed ablation study of the effects of λ in the following subsections.
 763 We surprisingly find that even for $\lambda \ll 1$ we can achieve significant edge in terms of robustness
 764 over existing methods. This is in contrast to virtual adversarial training methods such as [36] which
 765 requires applying a λ -value magnitudes larger. Moreover, for larger λ , we find that the accuracy of
 766 the model is not compromised. This provides some empirical support for Theorem 2.2.

Algorithm 1 Estimating knowledge continuity.

Input: A batch of N data points $\{(x_i, y_i)\}_{i=1}^N$, $M \leq N$, neural network f with n hidden layers, and some $k \in [n]$
Output: An estimation of $\mathbb{E} \sigma_f^k(x, y)$.
Subsample M indices n_1, \dots, n_m uniformly at random from $[N]$ without replacement
 $\sigma_f^k \leftarrow 0$
Losses $\leftarrow \{\mathcal{L}(f(x_{n_i}), y_{n_i})\}_{i=1}^M$
for $(i, j) \in [M] \times [M]$ **do**
 Dist $\leftarrow d_k(f^k(x_{n_i}), f^k(x_{n_j}))$
 $\sigma_f^k \leftarrow \sigma_f^k + |\text{Losses}_i - \text{Losses}_j| / \text{DIST}$
end for
return σ_f^k

Algorithm 2 Regularization of knowledge continuity.

Input: $\alpha, \beta, M, \lambda > 0$.
A neural network f with n hidden layers, loss function \mathcal{L} , and batch $\{(x_i, y_i)\}_{i=1}^N$.
Output: Loss with added knowledge continuity regularization score.
 $X \sim \text{Beta}(\alpha, \beta)$
 $\sigma_f^k \leftarrow (\text{Alg. 1})(f, M, k \leftarrow \max(\lfloor Xn \rfloor, 1))$
return $\frac{1}{N} \sum_{i=1}^N \mathcal{L}(f(x_i), y_i) + \frac{1}{M^2} \lambda \sigma_f^k$

767 I.2 Theoretical Guarantees of Knowledge Continuity Regulation

768 In this subsection, we demonstrate that Alg. 1 is indeed an unbiased estimator for knowledge
769 continuity and also provide some bounds on the rate of convergence of this estimation.

770 **Proposition I.1** (Alg. 1 is an Unbiased Estimator). *Assuming that each data point in the batch,*
771 *$\{(x_i, y_i)\}_{i=1}^N \sim \mathcal{D}_{\mathcal{X}, \mathcal{Y}}$, is sampled i.i.d., then Alg. 1 is an unbiased estimator for $\mathbb{E} \sigma_f^k(x, y)$.*

772 *Proof.* Let $\hat{\theta}$ be the random variable representing the output of Alg. 1. It suffices to show that

$$\mathbb{E}[\hat{\theta}] = \mathbb{E} \sigma_f^k(x, y),$$

773 where the expectation on the left-hand side is taken over the set of all batches. By the definition of
774 Alg. 1,

$$\mathbb{E}[\hat{\theta}] = \mathbb{E} \left(\sum_{i=1}^M \sum_{j=1}^M \frac{1}{M^2} \frac{\Delta \mathcal{L}_f^{(x_{n_j}, y_{n_j})x_{n_i}, y_{n_i}}}{d_k(f^k(x_{n_i}), f^k(x_{n_j}))} \right), \quad (\text{I.1})$$

$$= \sum_{i=1}^M \sum_{j=1}^M \frac{1}{M^2} \mathbb{E} \left(\frac{\Delta \mathcal{L}_f^{(x_{n_j}, y_{n_j})x_{n_i}, y_{n_i}}}{d_k(f^k(x_{n_i}), f^k(x_{n_j}))} \right), \quad (\text{I.2})$$

$$= \mathbb{E} \sigma_f^k(x, y). \quad (\text{I.3})$$

775 The second equality follows from the linearity of expectation. ■

776 Next, we state a proposition that gives some theoretical guidance for choosing the hyperparameter M
777 in Alg. 1. In practice, one should choose M to be the batch size as to ensure accurate estimation of
778 the knowledge discontinuity score. We recognize, however, that if $N \gg 1$, choosing $M = N$ may be
779 intractable. We multiplicatively bound the error of the unbiased estimator with respect to M and the
780 overall variance of the δ -knowledge discontinuity.

781 **Definition 8.** *A random variable $\hat{\theta}$ is an (ε, δ) -multiplicative estimator of a random variable θ if*

$$\mathbb{P}[\hat{\theta} \notin (1 \pm \varepsilon)\theta] \leq 1 - \delta.$$

782 The next result is a well-known result from [44] with applications found in [13] and [25].

783 **Theorem I.2** (Median of Means). *Given $\varepsilon, \delta > 0$, and an unbiased estimator $\theta, \hat{\theta}$. We can achieve an*
 784 *(ε, δ) -multiplicative estimator of θ with K independent samples of $\hat{\theta}$ where*

$$K = O\left(\frac{\text{Var}(\hat{\theta})}{(\varepsilon \mathbb{E}\hat{\theta})^2} \ln \frac{1}{\delta}\right),$$

785 where $\text{Var}(\hat{\theta})$ is the variance of the estimator $\hat{\theta}$.

786 **Proposition I.3.** *Suppose $\varepsilon, \delta, \delta^l > 0$, then we can achieve an (ε, δ) -multiplicative estimator of the*
 787 *δ^l -knowledge discontinuity in layer j with $M = \Theta(K)$ using Alg. 1 where*

$$K = O\left(\frac{\delta^l \text{Var}(KD)}{(\varepsilon \mathbb{E}[\Delta \mathcal{L}(f; x, y)])^2} \ln \frac{1}{\delta}\right),$$

788 where $\Delta \mathcal{L}(f; x, y)$ difference in loss of f on any two data points sampled from \mathcal{D} and KD is the
 789 random variable that represents the δ^l -knowledge discontinuities across \mathcal{D} .

790 *Proof.* Consider a variation of the algorithm where we only draw a pair of points. In other words, fix
 791 $M = 2$. Denote the two data points we are considering to be $(x_1, y_1), (x_2, y_2)$. Then, let

$$X := \begin{cases} \frac{|\mathcal{L}(f; x_1, y_1) - \mathcal{L}(f; x_2, y_2)|}{d_j(h_j(x), h_j(x'))}, & \text{if } \|h_j(x_1) - h_j(x_2)\| < \delta^l, \\ 0 & \text{o/w.} \end{cases}$$

792 Since we've already shown that X is an unbiased estimator (see Prop. I.1) of the δ^l -knowledge
 793 discontinuities, it remains to find the variance and squared expectation and apply the Median of
 794 Means theorem (see Theorem I.2). First, we lower bound $(\mathbb{E}X)^2$:

$$\begin{aligned} (\mathbb{E}X)^2 &= \left(\frac{1}{2} \int_{\mathcal{D}|_X} \mathbb{E}_{(x', y') \sim \mathcal{D}|_{V_x}} \left[\frac{|\mathcal{L}(f; x, y) - \mathcal{L}(f; x', y')|}{d_j(h_j(x), h_j(x'))} \right] d\mu_X \right)^2, \\ &\hspace{15em} \text{(from Prop. I.1)} \\ &= \frac{1}{4} \int_{\mathcal{D}|_X \times \mathcal{D}|_X} \mathbb{E}_{(x', y') \sim \mathcal{D}|_{V_{x_1}}} \left[\frac{|\mathcal{L}(f; x_1, y_1) - \mathcal{L}(f; x', y')|}{d_j(h_j(x), h_j(x'))} \right] \\ &\quad \mathbb{E}_{(x', y') \sim \mathcal{D}|_{V_{x_2}}} \left[\frac{|\mathcal{L}(f; x_2, y_2) - \mathcal{L}(f; x', y')|}{d_j(h_j(x), h_j(x'))} \right] d\mu_X(x_1) d\mu_X(x_2), \\ &\geq \frac{1}{4\delta^2} \int_{\mathcal{D}|_X \times \mathcal{D}|_X} \mathbb{E}_{(x', y') \sim \mathcal{D}|_{V_{x_1}}} [|\mathcal{L}(f; x_1, y_1) - \mathcal{L}(f; x', y')|] \cdot \\ &\quad \mathbb{E}_{(x', y') \sim \mathcal{D}|_{V_{x_2}}} [|\mathcal{L}(f; x_2, y_2) - \mathcal{L}(f; x', y')|] d\mu_X(x_1) d\mu_X(x_2), \\ &\hspace{15em} \text{(since } d_j(h_j(x), h_j(x')) \text{)} \\ &\geq \frac{1}{4\delta^2} \int_{\mathcal{D}|_X} \left(\mathbb{E}_{(x', y') \sim \mathcal{D}|_{V_x}} [|\mathcal{L}(f; x, y) - \mathcal{L}(f; x', y')|] \right)^2 d\mu_X(x), \\ &\hspace{15em} \text{(only consider terms where } x_1 = x_2 \text{)} \\ &\geq \frac{1}{4\delta^2} \int_{\mathcal{D}|_X} \left(\int_{\mathcal{D}|_{V_x}} |\mathcal{L}(f; x, y) - \mathcal{L}(f; x', y')|^2 \right) d\mu_X(x), \\ &\hspace{15em} \text{(only consider terms in the product that agree)} \\ &\geq \frac{1}{4\delta^2} \int_{\mathcal{D}|_X \times \mathcal{D}|_X} |\mathcal{L}(f; x, y) - \mathcal{L}(f; x', y')|^2 \chi_\delta(x, x') d\mu(x) d\mu(x') \end{aligned}$$

795 which follows from Tonelli's theorem and $\chi_\delta(x, x') = 1$ if and only if $d_j(h_j(x), h_j(x')) < \delta$ and 0
 796 otherwise. Then, by symmetry, this is equivalent to

$$= \frac{\mathbb{E}[\Delta \mathcal{L}(f; x, y)]^2}{4\delta^2}.$$

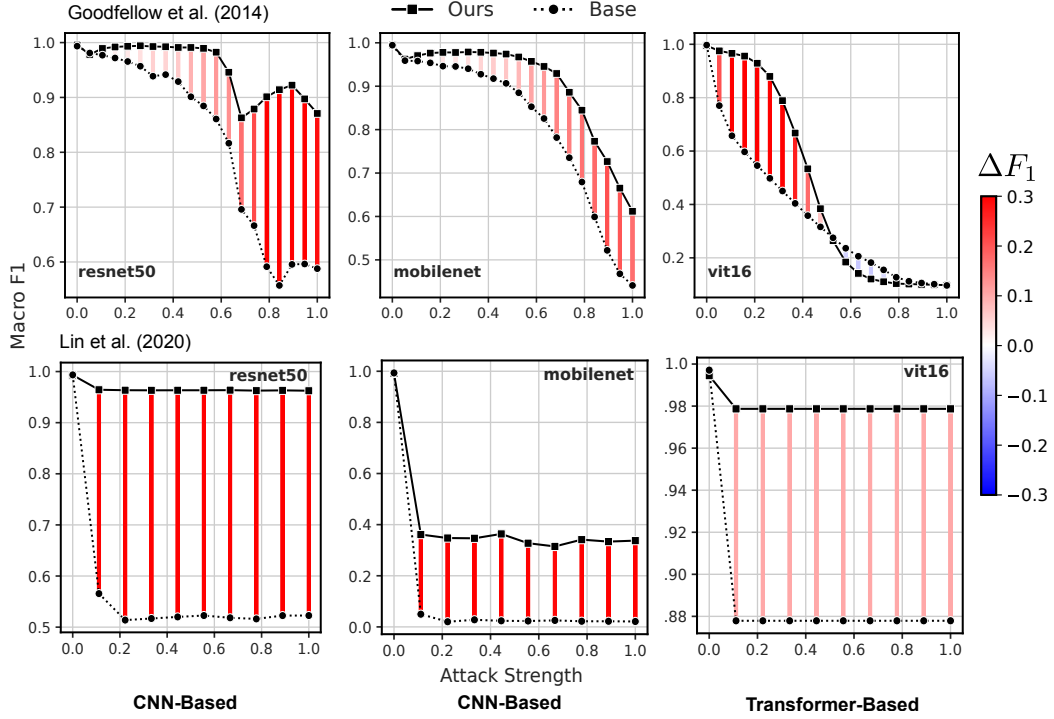


Figure 7: Regulating knowledge continuity on a host of vision models (ResNet50, MobileNetV2, and ViT16). Base models are trained with cross-entropy loss. KCREg (Our) models are finetuned with the additional regularization objective described in Alg. 2. Two adversarial attacks are then performed: the fast-gradient sign method from [19], and an iterative attack SI-NI-FGSM from [35]. We see that regulating knowledge continuity consistently improves/stabilizes robustness. Performance is measured using F1 and the attack strength corresponds to the maximum perturbation magnitude in L2 allowed. Since the pixel values of the images are bounded between $[0, 1]$, we also constrain the attack strength to be between $[0, 1]$.

796 The last equality follows from the fact that Now, we bound the variance of the estimator by above:

$$\begin{aligned}
 \text{Var } X &= \mathbb{E}X^2 - (\mathbb{E}X)^2, \\
 &= \int_{\mathcal{D}|_X \times \mathcal{D}|_X} \frac{|\mathcal{L}(f; x_1, y_1) - \mathcal{L}(f; x_2, y_2)|^2}{d_j(h_j(x), h_j(x'))^2} d\mu(x_1)d\mu(x_2) - (\mathbb{E}KD)^2 \\
 &= \mathbb{E}KD^2 - (\mathbb{E}KD)^2 = \text{Var}(KD).
 \end{aligned}
 \tag{from Prop. I.1}$$

797 Thus, combining both expressions with Theorem. I.2 we yield the desired result. ■

798 I.3 Regulating Knowledge Continuity “In the Wild”

799 We compare our regularization algorithm with several state-of-the-art adversarial and virtual adver-
 800 sarial training algorithms. These results are presented in Table. 1. Additional experiments on MNIST
 801 are also performed. These are presented in Fig. 7.

802 I.4 Ablation Studies

803 Herein, we present ablation studies for the crucial hyperparameters in our regularization algorithm,
 804 Alg. 2: λ which is the weight we assign the knowledge continuity regulation loss and (α, β) which
 805 determines the sampling behavior of the index of the hidden representation space.

Table 1: Comparison of our knowledge continuity algorithm to existing works across various model families and adversarial attack methods. TF, BA, ANLI denote adversarial attacks [29], [34], and [45], respectively. Regulating knowledge continuity to improve robustness is superior across almost all tasks and attacks.

Arch.	Method	IMDB	IMDB _{TF}	IMDB _{BA}	ANLI _{R1}	ANLI _{R2}	ANLI _{R3}
BERT ~110M params	Base	93.6	47.9	45.2	44.5	45.6	33.8
	TF	93.3	69.2	62.5	✗	✗	✗
	ALUM	93.5	56.9	47.8	45.2	46.7	46.3
	KCReg (ours)	94.8	75.1	84.9	45.6	46.9	45.3
GPT2 ~1.5B params	Base	93.6	63.9	54.9	42.7	44.9	43.4
	TF	92.0	64.5	51.3	✗	✗	✗
	ALUM	94.9	49.4	27.5	43.8	45.2	44.6
	KCReg (ours)	94.9	87.8	90.6	47.1	48.1	44.7
T5 ~220M params	Base	93.7	53.9	39.3	46.1	44.7	46.0
	TF	96.8	77.8	60.6	✗	✗	✗
	ALUM	95.1	67.1	51.9	44.5	44.8	44.4
	KCReg (ours)	94.9	89.3	91.3	48.2	45.0	44.3

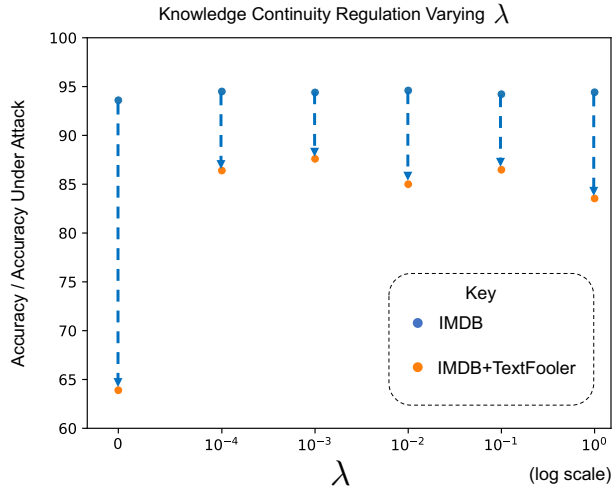


Figure 8: The accuracy of the model (both not under/under adversarial attack) on the IMDB dataset versus varying the weight given to the knowledge continuity regularization term (λ).

806 **Ablation Study of λ .** The weight given to the regularizer (λ) is ablated over, with the results shown
807 in Fig. 8. For any positive λ , there is an immediate large improvement in adversarial robustness. Next,
808 as λ is systematically increased by factors of 10, we do not see a significant change in the accuracy
809 (not under attack). This corroborates Theorem 2.2, as it demonstrates that regulating knowledge
810 discontinuities (no matter how strongly) is not at odds with minimizing the empirical risk of our
811 model. On the other hand, we also do not see a significant increase in adversarial robustness as
812 λ increases. This may imply that we have reached the threshold of adversarial robustness under
813 TextFooler [29]. Specifically, the adversarial attacks generated by TextFooler may not be valid in
814 that they have flipped the ground-truth label. Therefore, we believe that a good λ for this particular
815 application should lie somewhere between 0 and 1×10^{-4} .

816 **Ablation Study of (α, β)** In this subsection, we briefly discuss how the α, β hyperparameters which
817 determine the shape of the Beta distribution in Alg. 2 affect the final performance and robustness of
818 our model on the IMDB dataset. Recall that the shape of the Beta distribution determines the index
819 of the hidden layers we are using to compute the knowledge continuity. Thus, they are crucial in
820 determining the behavior of our regularizer.

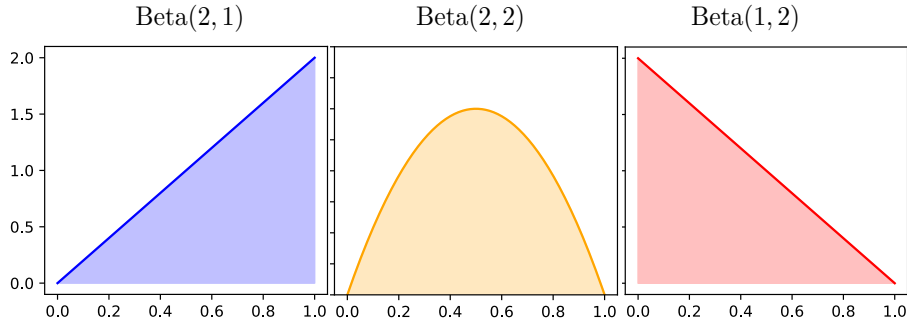


Figure 9: The Beta distributions that we ablated over with the probability density function of their parameterizations shown.

821 We finetune {BERT, T5, GPT2} models on the IMDB dataset with the hyperparameters described in
 822 the next subsection. The results are displayed in Table 2. Across all models we observe a decrease in
 823 robustness for $\alpha = 1, \beta = 2$. These values correspond to a right-skewed distribution which places
 824 high sampling probability on the earlier (closer to the input) hidden layers. Intuitively, perturbations
 825 in the early layers should correspond to proportional textual perturbations in the input text. Pure
 826 textual perturbations with respect to some metric like the Levenshtein distance should be only
 827 loosely if not completely (un)correlated with the actual labels of these inputs. Therefore, enforcing
 828 knowledge continuity with respect to this metric should not see increase robustness. Moreover, we
 829 also observe a larger decrease in accuracy (not under attack) with the same parameters. This suggests
 830 that maintaining this sort of knowledge continuity in the earlier layers is harder to converge on and
 831 there may be a “push-and-pull” behavior between optimizing knowledge continuity and accuracy
 832 (not under attack). Surprisingly, we observe no significant difference between the other α, β values
 shown in the table.

Table 2: We train finetune {BERT, T5, GPT2} using knowledge continuity regularization, as described in Alg. 2. We varied the α, β hyperparameters for the Beta distribution as to determine the effect of these parameters on model performance and robustness. The rows of the table are labeled with the format: Model+Reg $_{(\alpha, \beta)}$. The bolded entries of the table correspond to the best performing metrics out of the knowledge continuity regulated models.

Model	IMDB	IMDB _{TF}
BERT _{BASE}	93.6	47.9
BERT _{BASE} +Reg _(2,1)	94.8	75.1
BERT _{BASE} +Reg _(2,2)	89.2	74.1
BERT _{BASE} +Reg _(1,2)	87.0	68.2
GPT2	93.6	63.9
GPT2+Reg _(2,1)	94.6	85.0
GPT2+Reg _(2,2)	94.9	87.8
GPT2+Reg _(1,2)	93.1	84.9
T5 _{BASE}	93.7	53.9
T5 _{BASE} +Reg _(2,1)	95.0	88.9
T5 _{BASE} +Reg _(2,2)	94.9	89.3
T5 _{BASE} +Reg _(1,2)	94.6	88.1

834 We did not formally benchmark other configurations of α, β such as increasing their magnitude to
 835 impose a sharper distribution. During training, we noticed that using these sharper distributions
 836 both significantly slowed the model’s convergence and decreased the model’s accuracy (not under
 837 attack). It could be that though knowledge continuity itself is a *local* property the enforcement of
 838 this *local* property requires change on a *global* scale. In other words, one cannot simply reduce the
 839 knowledge discontinuities or uniformly converge with respect to one layer without participation from
 840 other layers. The extent to which other layers are involved in the regularization of a specific one is an
 841 interesting question that we leave for future research.

842 I.5 Training Details

843 In this section, we describe in detail the training objectives, procedures, algorithms, and hyperparme-
 844 ters that we used in the main text and further experiments done in the appendix.

845 **Brute-Force Adversarial Training.** For all models undergoing adversarial training, we first finetune
 846 the model against the training set. Then, attack it using the TextFooler [29] algorithm with examples
 847 from the training set. After the attacks are concluded, we then incorporate the text of successful
 848 adversarial attacks back into the training set and proceed to finetune again. This procedure iteratively
 849 continues. For the sake of computational efficiency, for all models we applied this procedure once.
 850 The parameters we are using during the adversarial attack is the same hyperparameters we actually
 851 use at test-time. Specifically, we impose a query budget of 300 queries.

852 **Plain Finetuning on IMDB.** The IMDB dataset consist of 50,000 examples with 25,000 for training
 853 and 25,000 for testing. We split the test set 40%-60% to create a validation and test set of 10,000
 854 and 15,000 examples, respectively. Examples were sampled uniformly at random during the splitting
 855 process. Since adversarial attacks were costly, we uniformly subsampled 5,000 examples from this
 856 15,000 to benchmark robustness in the experiments related to the regularizer. However, for the
 857 experiments estimating the knowledge vulnerability score, we performed adversarial attacks on all
 858 15,000 datapoints in the test set. We found no significant difference between robustness estimation
 859 on this 5,000 subsample versus and the entire 15,000 dataset.

We train all models using the following hyperparameter and optimizer configurations:

Table 3: Training hyperparameters and optimizer configurations for finetuning models {BERT, GPT2, T5} on IMDB without any form of regularization or adversarial training.

HYPERPARAMETER	VALUE
OPTIMIZER	ADAM
ADAM β_1	0.9
ADAM β_2	0.999
ADAM ϵ	1×10^{-8}
MAX GRADIENT NORM	1.0
LEARNING RATE SCHEDULER	LINEAR
EPOCHS	20
BATCH SIZE	32
LEARNING RATE	5×10^{-5}
WEIGHT DECAY	1×10^{-9}

860

861 **Knowledge Discontinuity Regulation on IMDB.** For enforcing the knowledge discontinuity on
 862 IMDB, we use a constant $\lambda = 1 \times 10^{-2}$ for all models. As shown in Table 2, we varied $\alpha, \beta \in$
 863 $\{1, 2\} \times \{1, 2\}$ and displayed the best models in terms of robustness in Table. 1 in the main text.
 864 We train all models for 50 epochs. Other than that all the other hyperparameters and optimizer
 865 configurations are the same as regular finetuning (see Table 3).

866 **Knowledge Discontinuity Regulation on ANLI.** Optimizing over the ANLI dataset was significantly
 867 harder than on IMDB. As a result, for each model class {BERT, GPT2, T5} we performed a quick
 868 hyperparameter search over λ (1×10^{-4}), the learning rate (5×10^{-5}), and weight decay (1×10^{-9})
 869 fixing the parameterization of the Beta distribution to be the best values on the IMDB dataset. That is,
 870 for T5: $\alpha = 2, \beta = 1$; BERT-Base-Uncased: $\alpha = 2, \beta = 1$; GPT2: $\alpha = 2, \beta = 2$.

871 **ALUM on IMDB and ANLI.** We train all ALUM models for 50 epochs (the same as knowledge
872 discontinuity regularized models). For hyperparameters specific to the ALUM algorithm we choose
873 all of the same ones as its authors, [36], with the exception of α (analogous to the λ in our algorithm,
874 essentially the weight put on the virtual adversarial training loss term). The authors of the original
875 paper choose $\alpha = 10$. We, however, found that this applied to finetuning does not converge at all.
876 Thus, with a rough binary search in the parameter space we found $\alpha = 1 \times 10^{-3}$ to be the best with
877 respect to both performance and robustness.

878 We keep the same hyperparameters on ANLI, however, we impose early stopping during the training
879 process. That is, we choose the best model with respect to its performance on the **dev** set.

880 **J Limitations**

881 The certification guarantees of our definition knowledge continuity is a probabilistic one. Specifically,
882 this randomness is over the data distribution. However, this does not protect against out-of-distribution
883 attacks that plague large language models such as [59, 76]. More work is needed to yield deterministic
884 results that do not become vacuous in discrete settings. As mentioned in Section 2.4, our expres-
885 siveness bounds only apply under little restrictions to the metric decompositions of the estimator f .
886 Though we see some empirical verification for this in Appendix I, it remains unclear whether or not
887 we can tighten these bounds.

888 **K Broader Impacts**

889 This contribution is concerned with robust deep learning models. As deep learning becomes ubiqui-
890 tous as the mode for artificial intelligence, their applications in increasingly critical areas to the lay
891 and corporations alike demand not only both high inferential accuracy and confidence. Robustness ad-
892 dresses this latter point, by making deep learning models more robust, we improve the trustworthiness
893 of their decision-making and protect them against adversaries. More specifically, our contribution
894 unifies separate robustness efforts from continuous and discrete domains.

895 **L Reproducibility**

896 All of our experiments were conducted on four NVIDIA RTX A6000 GPUs as well as four NVIDIA
897 Quadro RTX 6000 GPUs. The rest of our code base including implementations of the algorithms and
898 figures described in the manuscript are attached as supplementary materials.