

TRANSFORMERS LEARN LOW SENSITIVITY FUNCTIONS: INVESTIGATIONS AND IMPLICATIONS

Anonymous authors

Paper under double-blind review

ABSTRACT

Transformers achieve state-of-the-art accuracy and robustness across many tasks, but an understanding of their inductive biases and how those biases differ from other neural network architectures remains elusive. In this work, we identify the sensitivity of the model to token-wise random perturbations in the input as a unified metric which explains the inductive bias of transformers across different data modalities and distinguishes them from other architectures. We show that transformers have lower sensitivity than MLPs, CNNs, ConvMixers and LSTMs, across both vision and language tasks. We also show that this low-sensitivity bias has important implications: i) lower sensitivity correlates with improved robustness; it can also be used as an efficient intervention to further improve the robustness of transformers; ii) it corresponds to flatter minima in the loss landscape; and iii) it can serve as a progress measure for grokking. We support these findings with theoretical results showing (weak) spectral bias of transformers in the NTK regime, and improved robustness due to the lower sensitivity.

1 INTRODUCTION

Transformers, originally introduced for language problems (Vaswani et al., 2017), have become a universal backbone across machine learning — including applications such as vision (Dosovitskiy et al., 2021) and protein structure prediction (Jumper et al., 2021). Several recent works have also found that not only do transformers achieve better accuracy, but they are also more robust to various corruptions and changes in the data distribution (Shao et al., 2021; Mahmood et al., 2021; Bhojanapalli et al., 2021; Paul & Chen, 2022). Despite their practical success, relatively little is understood about what distinguishes transformers from other neural network architectures. *If a transformer and an alternative neural network architecture (such as a CNN or an LSTM) are trained to obtain similar training accuracy on a dataset, then how do the models differ in terms of the functions they learn? Equivalently, what inductive biases do transformers have which distinguish them from other architectures?*

Recently, for the setting of Boolean inputs, Bhattamishra et al. (2023b) and Hahn & Rofin (2024) suggest using the notion of *sensitivity* to distinguish transformers from other candidate architectures. The sensitivity of a function measures how likely the output is to change for random changes to the input. Bhattamishra et al. (2023b) and Hahn & Rofin (2024) show that transformers are biased to learn functions with low sensitivity on Boolean inputs. Sensitivity has several desirable properties as a notion of inductive bias. It is closely related to the Fourier representation of the function and various other notions of Boolean function complexity such as the degree of the function and the size of the smallest decision tree which represents the function (O’Donnell, 2014). Low sensitivity functions correspond to low complexity functions based on all these notions of Boolean function complexity, and hence an inductive bias towards low-sensitivity functions is regarded as an instance of ‘simplicity bias’ of the model (Valle-Perez et al., 2019; Bhattamishra et al., 2023b). Sensitivity also has deep connections to well-studied notions of inductive biases such as spectral bias (Rahaman et al., 2019b), which is a bias towards ‘simple’ functions such as low frequency functions in the Fourier space. Sensitivity has also been found to correlate with better generalization for fully-connected networks (Novak et al., 2018). (See Appendix C for a detailed discussion of the related work.)

Our results. Sensitivity is a promising notion of inductive bias, but has mainly been investigated for Boolean functions so far. Given the numerous modalities of data across which transformers are

054 successful in practice, the goal of our work is to examine if appropriate extensions of the notion of
 055 sensitivity for Boolean functions help understand the inductive bias of transformers across varied
 056 data modalities — and if these notions help explain properties of transformers such as their improved
 057 robustness. We now provide an overview of the main claims and results of the paper. We begin our
 058 investigation with the following question:

059 *What are appropriate notions of sensitivity beyond Boolean data?*
 060

061 To provide a concrete starting point where we can understand the properties of sensitivity with theo-
 062 retical analysis, we first consider the Boolean setup and place the low-sensitivity bias of transformers
 063 on a firmer theoretical foundation in that setting (Section 2). Using prior work on spectral bias in
 064 neural networks (Yang & Salman, 2020) we prove that transformers show a low-sensitivity bias on
 065 Boolean functions, and also prove that low sensitivity leads to better robustness. We then consider the
 066 above question, and propose a suitable notion of sensitivity which takes into account the underlying
 067 metric space (Section 3). To investigate if sensitivity is a suitable notion beyond the Boolean case,
 068 we first consider a synthetic dataset where we can tease apart sensitivity from related notions which
 069 can coincide for Boolean functions — such as a preference towards functions that depend on a sparse
 070 set of tokens. We show that transformers prefer to learn low-sensitivity functions (even if they are not
 071 sparse). Subsequently, we examine if this low-sensitivity bias is widely present across different tasks:

072 *Does low-sensitivity serve as a unified notion of simplicity across vision and language tasks, and*
 073 *does it distinguish between transformers and other architectures?*

074 Here, we first conduct experiments on vision datasets. We empirically compare (Vision-)Transformers
 075 with MLPs, CNNs, and ConvMixers, and observe that transformers have lower sensitivity compared
 076 to other candidate architectures (see Section 4). Similarly, we conduct experiments on language tasks
 077 and observe that transformers learn predictors with lower sensitivity than LSTM models. Furthermore,
 078 transformers tend to have uniform sensitivity to all tokens while LSTMs are more sensitive to more
 079 recent tokens (see Section 5). Given this low-sensitivity bias, we next examine its implications:

080 *What are the implications of a bias towards low-sensitivity functions; is it helpful in certain settings?*
 081

082 We study this in three contexts: robustness, properties of the loss landscape, and training dynamics.

- 083 1. **Lower Sensitivity Correlates with Better Robustness:** We show that transformers have lower
 084 sensitivity and are more robust to corruptions when tested on the CIFAR-10-C dataset, compared
 085 to CNNs (Section 6). We also demonstrate that sensitivity is not only predictive of robustness
 086 but also has prescriptive power: We add a regularization term at training time to encourage the
 087 model to have lower sensitivity. Since sensitivity is efficient to measure empirically, this is easy to
 088 accomplish via data augmentation. We find that models explicitly trained to have lower sensitivity
 089 yield even better robustness on CIFAR-10-C. These results show that *low sensitivity correlates*
 090 *with the improved robustness of transformers.*
- 091 2. **Lower Sensitivity Correlates with Flatter Minima:** We explore the connection between sensi-
 092 tivity and a property of the loss landscape that has been found to correlate to good generalization —
 093 the sharpness of the minima. We compare the sharpness of the minima with and without sensi-
 094 tivity regularization, and our results show that *lower sensitivity correlates with flatter minima.* This
 095 indicates that sensitivity could serve as a unified notion for both robustness and generalization.
- 096 3. **Sensitivity Serves as a Progress Measure for Grokking:** We examine if sensitivity can be
 097 used to understand the training dynamics of transformers, specifically from the perspective of the
 098 phenomenon of *grokking* where test accuracy abruptly improves long after the training loss or
 099 accuracy saturates. We consider modular addition, a task where transformers exhibit grokking.
 100 We show that *sensitivity provides a progress measure that decreases even when the training loss*
 101 *does not reduce and is indicative of stages of grokking.*

102 2 SENSITIVITY AND WEAK SPECTRAL BIAS

103 In this section, we theoretically show that transformers with linear attention exhibit (weak) spectral
 104 bias to learn lower-order Fourier coefficients, which in turn implies a bias to learn low-sensitivity
 105 functions. We start with an overview of Fourier analysis on the Boolean cube and sensitivity.

106 **Fourier analysis on the Boolean cube (O’Donnell, 2014).** The space of real-valued functions on
 107 the Boolean cube $\{0,1\}^d$ forms a 2^d -dimensional space. Any such function can be written as a *unique*

108 *multilinear* polynomial. Specifically, the multilinear monomial functions, $\chi_U(\mathbf{x}) := x^U := \prod_{i \in U} x_i$, for

109 each $U \subseteq [d]$, form a Fourier basis of the function space $\{f : \mathbb{B}^d \rightarrow \mathbb{R}\}$, *i.e.*, their inner products

110 satisfy $\mathbb{E}_{\mathbf{x} \sim \mathbb{B}^d} [\chi_U(\mathbf{x})\chi_V(\mathbf{x})] = \mathbb{1}[U=V]$. Consequently, any function $f : \mathbb{B}^d \rightarrow \mathbb{R}$ can be written

111 as $f(\mathbf{x}) = \sum_{U \subseteq [d]} \hat{f}(U)\chi_U(\mathbf{x})$, for a unique set of coefficients $\hat{f}(U)$, $U \subseteq [d]$, where $[d] = \{1, \dots, d\}$.

112

113

114 **Sensitivity in Boolean function analysis.** Sensitivity is a common complexity measure for

115 Boolean functions. Intuitively, it captures the changes in the output of the function, averaged

116 over the neighbours of a particular input. Formally, let $\mathbb{B}^d := \{\pm 1\}^d$ denote the Boolean cube

117 in dimension d . The sensitivity of a Boolean function $f : \mathbb{B}^d \rightarrow \{\pm 1\}$ at input $\mathbf{x} \in \mathbb{B}^d$ is

118 given by $S(f, \mathbf{x}) = \sum_{i=1}^d \mathbb{1}[f(\mathbf{x}) \neq f(\mathbf{x}^{\oplus i})]$, where $\mathbb{1}[\cdot]$ denotes the indicator function and

119 $\mathbf{x}^{\oplus i} = (x_1, \dots, x_{i-1}, -x_i, x_{i+1}, \dots, x_d)$ denotes the sequence obtained after flipping the i^{th} co-

120 ordinate of \mathbf{x} . Note that in the Boolean case, the neighbor of an input can be obtained by flipping

121 a bit, we will define a more general notion later which holds for more complex data. The average

122 sensitivity is measured by averaging $S(f, \mathbf{x})$ across all inputs,

$$123 S(f) = \mathbb{E}_{\mathbf{x} \sim \mathbb{B}^d} [S(f, \mathbf{x})] = \frac{1}{2^d} \sum_{\mathbf{x} \in \mathbb{B}^d} S(f, \mathbf{x}). \quad (1)$$

124

125 Following [Bhattachishra et al. \(2023b\)](#), when comparing inputs of different lengths, we consider the

126 average sensitivity normalized by the input length, $\bar{S}(f) = \frac{1}{d} S(f)$. The sensitivity of a function

127 f is known to be related to the degree $D(f)$ of the multilinear polynomial which represents f

128 ([Huang, 2019](#); [Hatami et al., 2011](#)), and low-degree functions have lower sensitivity. Specifically, a

129 breakthrough result ([Huang, 2019](#)) showed that $D(f) \leq S_{\max}^2(f)$, where $S_{\max}(f) := \max_{\mathbf{x} \in \mathbb{B}^d} S(f, \mathbf{x})$.

130

131 **Attention Layer.** The output of a single-head self-attention layer, parameterized by key, query, value

132 matrices $\mathbf{W}_Q, \mathbf{W}_K \in \mathbb{R}^{\tilde{d} \times d_h}, \mathbf{W}_V \in \mathbb{R}^{\tilde{d} \times d_v}$ for input $\mathbf{X} \in \mathbb{R}^{T \times \tilde{d}}$ with T tokens of \tilde{d} dimension, is

133 $\text{ATTN}(\mathbf{X}; \mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V) := \varphi(\mathbf{X}\mathbf{W}_Q\mathbf{W}_K^T\mathbf{X}^T)\mathbf{X}\mathbf{W}_V$, where $\varphi(\mathbf{X}\mathbf{W}_Q\mathbf{W}_K^T\mathbf{X}^T) \in \mathbb{R}^{T \times T}$ is

134 the attention map with the softmax map $\varphi(\cdot) : \mathbb{R}^T \rightarrow \mathbb{R}^T$ applied row-wise.

135

136 **Main Results.** Consider any model with at least one self-attention layer, where \mathbf{X} is obtained

137 by reshaping $\mathbf{x} \in \mathbb{B}^d$, $d = T\tilde{d}$. Instead of applying the softmax activation, we consider linear

138 attention and apply an identity activation element-wise with a scaling factor of $d^{-1/2}$. The following

139 result shows that the conjugate kernel (CK) or neural tangent kernel (NTK) (see [Appendix B](#) for an

140 overview) induced by transformers with linear attention exhibit a weak form of spectral bias, where

141 the eigenvalues **do not decrease as** the degree of the multi-linear monomials **increases**, separately for

142 even and odd degrees; see [Appendix B](#) for the proof.

143

144 **Proposition 2.1.** *Let K be the CK or NTK of a transformer with linear attention on a Boolean*

145 *cube \mathbb{B}^d . For any $\mathbf{x}, \mathbf{y} \in \mathbb{B}^d$, we can write $K(\mathbf{x}, \mathbf{y}) = \Psi(\langle \mathbf{x}, \mathbf{y} \rangle)$ for some univariate function*

146 *$\Psi : \mathbb{R} \rightarrow \mathbb{R}$. Further, for every $U \subseteq [d]$, χ_U is an eigenfunction of K with eigenvalue*

$$147 \mu_{|U|} := \mathbb{E}_{\mathbf{x} \sim \mathbb{B}^d} [x^U K(\mathbf{x}, \mathbf{1})] = \mathbb{E}_{\mathbf{x} \sim \mathbb{B}^d} \left[x^U \Psi \left(d^{-1} \sum_i x_i \right) \right],$$

148 where $\mathbf{1} := (1, \dots, 1) \in \mathbb{B}^d$, and the eigenvalues μ_k , $k \in [d]$, satisfy

$$149 \mu_0 \geq \mu_2 \geq \dots \geq \mu_{2k} \geq \dots, \quad \mu_1 \geq \mu_3 \geq \dots \geq \mu_{2k+1} \geq \dots$$

150

151 Note that for a given U , the eigenvalue $\mu_{|U|}$ only depends on x^U and $\sum_i x_i$ by definition, and hence,

152 it is invariant under any permutation of $[d]$. Larger eigenvalues for lower-order monomials indicate

153 that simpler features are learned faster. Since low sensitivity implies learning low-degree polynomials,

154 [Proposition 2.1](#) also implies a weak form of low sensitivity bias.

155

156 We now show a connection between lower sensitivity and better robustness. Given a sample $\mathbf{x} \in \mathbb{B}^d$

157 and some $\rho \in (0, 1)$, consider a noisy sample \mathbf{x}' , where $x'_i = x_i$ with probability ρ and uniformly

158 random, otherwise. It can be verified that the pair $(\mathbf{x}, \mathbf{x}')$ has correlation ρ . The following result

159 shows that if f has a lower sensitivity $S_{\max}(f)$, then there is a lower probability $\Pr[f(\mathbf{x}) \neq f(\mathbf{x}')] of$

160 inconsistent predictions on the pair $(\mathbf{x}, \mathbf{x}')_\rho$; see [Appendix B](#) for the proof.

Proposition 2.2. Given ρ -correlated pair $(\mathbf{x}, \mathbf{x}')$, where $\rho \in (0, 1)$, and function $f: \mathbb{B}^d \rightarrow \{\pm 1\}$ with maximum sensitivity $S_{\max}(f)$, $0 \leq \Pr_{(\mathbf{x}, \mathbf{x}')_\rho} [f(\mathbf{x}) \neq f(\mathbf{x}')] \leq 0.5(1 - \rho^{(S_{\max}(f))^2})$.

Together, Propositions 2.1 and 2.2 imply that transformers have low sensitivity and hence, better robustness. In Section 6, we present experimental evidence showing that the low sensitivity of transformers correlates with their improved robustness.

3 MEASURING SENSITIVITY BEYOND BOOLEAN DATA

While sensitivity appears to be a promising metric to understand the inductive biases of transformers, it is only defined for Boolean data. In order to investigate the inductive biases in real-world image and language tasks, we need an equivalent metric for high-dimensional, real-valued data.

We define the sensitivity metric for high-dimensional data, which is an analog of Eq. (1) as follows.

Definition 3.1. Given a model Φ , dataset \mathcal{D} and distribution \mathcal{P} , sensitivity is computed as:

$$\bar{S}(\Phi) = \frac{1}{T} \mathbb{E}_{\substack{\mathbf{X} \sim \mathcal{D} \\ \mathbf{x} \sim \mathcal{P}}} \left[\sum_{\tau=1}^T \mathbb{1}[\text{SIGN}(\Phi(\theta; \mathbf{X})) \neq \text{SIGN}(\Phi(\theta; \mathbf{X}^{\oplus \tau}))] \right], \quad (2)$$

where $\mathbf{X}^{\oplus \tau}$ is obtained by replacing the τ^{th} token in \mathbf{X} with \mathbf{x} .

An important consideration here is to define \mathcal{P} . While one can replace a token with a randomly selected token to measure sensitivity, this may not ensure that the new token lies in a neighbourhood of the original token. Capturing how the output changes with local perturbations according to the metric of the underlying space is an important aspect of the sensitivity definition for Boolean functions (as discussed in Gopalan et al. (2016), for e.g.), and appropriate extensions of sensitivity beyond Boolean functions should capture this property. For input spaces such as natural images or text embeddings, there is more structure in the tokens, and a randomly selected token can lie far from the original token’s neighborhood. Therefore, instead of replacing a token with a random token, we inject small perturbations into the token to evaluate sensitivity. This allows us to control the size of the neighbourhood by selecting the strength of the noise perturbation.

Formally, for each token e_τ of an input \mathbf{X} , let $\mathbf{x} := e_\tau + \xi$ be a perturbed token, where $\xi \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ is an isotropic Gaussian with variance σ^2 . We measure sensitivity by replacing e_τ with \mathbf{x} as per Definition 3.1, with \mathcal{P} as $\mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$. For image data, each token e_τ corresponds to different patches (see Section 4 for further details), while for language data, the tokens correspond to embeddings of sub-words (see Section 5 for more details). Figure 1 illustrates the measurement for images.

The important characteristics of this metric are that it is a unified notion of complexity across vision and language tasks, and as we will see later, it distinguishes transformers from various other architectures. For instance, we compare the sensitivity measured with token-wise perturbations as mentioned above, with random perturbations across the input in Appendix A.3 and find that the gap in the latter metric is not as large as the proposed metric.

We also note that while for Boolean data, sensitivity aligns with other notions of complexity, such as sparsity, this may or may not be the case in settings with high-dimensional or real-valued data. In the following section, we present experimental results for a self-attention model, in a synthetic data setting to demonstrate this. Specifically, we show that in the synthetic dataset, the related notion of using a sparse set of input tokens may or may not align with low-sensitivity, but the model learns the low-sensitivity function in both cases.

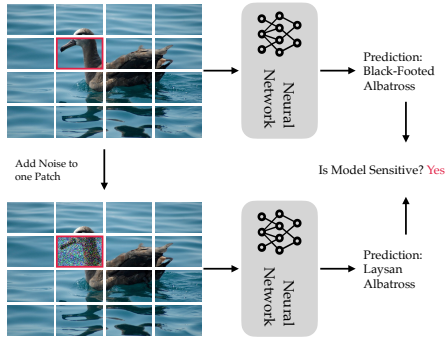


Figure 1: **Measuring Sensitivity in Vision Tasks.** A patch is first selected to add Gaussian noise corruptions. Then the original image and the corrupted image are fed into the same neural network to make predictions. If the predictions are inconsistent, then the neural network is sensitive to this patch. The process is repeated for every patch to measure the overall sensitivity.

3.1 EXPERIMENTS ON SYNTHETIC DATA

We construct a synthetic dataset to examine the inductive bias of a single-layer self-attention model. We show that in the presence of two solutions with the same predictive power but different sensitivity values, this model learns the low-sensitivity function. We begin by describing the experimental setup and then discuss our results.

Setup. We compose a single-head self-attention layer with a linear head $U \in \mathbb{R}^{\tilde{d} \times \tilde{d}}$ to obtain the final prediction, and write the full model as

$$\Phi(\theta; \mathbf{X}) := \langle U, \varphi(\mathbf{X} \mathbf{W}_Q \mathbf{W}_K^\top \mathbf{X}^\top) \mathbf{X} \mathbf{W}_V \rangle, \quad (3)$$

where $\theta := \text{concat}(\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V, U)$. We consider this model for the experiments in this section, with all the parameters initialized randomly at a small scale. Next, we describe the process to generate the dataset. We first define the vocabulary as follows:

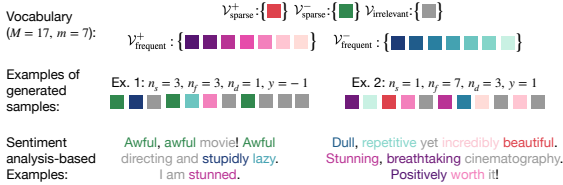


Figure 2: Visualization of the synthetic data generation process (see Section 3.1 for details). For simplicity, we represent each d -dimensional token with a square. Middle row: In each case, given a label y , we randomly sample $T = 11$ tokens, with n_s tokens from $\mathcal{V}_{\text{sparse}}^y$, $\lfloor (n_f + n_d)/2 \rfloor$ tokens from $\mathcal{V}_{\text{frequent}}^y$, $n_f - \lfloor (n_f + n_d)/2 \rfloor$ tokens from $\mathcal{V}_{\text{frequent}}^{-y}$ and the remaining tokens from $\mathcal{V}_{\text{irrelevant}}$. Note that in the first example, since $n_s = 3$ and $n_d = 1$, a predictor that relies (only) on the sparse tokens is less sensitive compared to the one that relies on the frequent tokens. On the other hand, in the second example, since $n_s = 1$ and $n_d = 3$, the predictor that relies on the frequent tokens is less sensitive. Bottom row: We include two sentiment analysis-based examples to illustrate the synthetic data samples in the second row, using the same colors as the first two rows.

Table 1: Comparison of sensitivity values for models that use only sparse or frequent tokens for the settings considered in Figure 3.

(n_s, n_f, n_d, m)	Using sparse tokens	Using frequent tokens
$(3, 5, 1, 16)$; Fig. 3 left col.	0	0.2878
$(1, 17, 7, 20)$; Fig. 3 right col.	0.0339	0

Definition 3.2 (Synthetic Vocabulary). Consider a vocabulary of M distinct tokens $\mathcal{V} := \{e_1, \dots, e_M\}$, where $e_i \in \{0, 1\}^d$ denotes the i^{th} basis vector. We define smaller subsets of *sparse* tokens and larger subsets of *frequent* tokens for each label $y = \pm 1$, as well as a subset of *irrelevant* tokens:

$$\mathcal{V}_{\text{sparse}}^+ := \{e_1\}, \mathcal{V}_{\text{sparse}}^- := \{e_2\}, \mathcal{V}_{\text{irrelevant}} := \{e_{2m+3}, \dots, e_M\}$$

$$\mathcal{V}_{\text{frequent}}^+ := \{e_3, e_5, \dots, e_{2m+1}\}, \mathcal{V}_{\text{frequent}}^- := \{e_4, e_6, \dots, e_{2m+2}\}.$$

Let T denote the sequence length of each data point, n_f and n_s denote the number of frequent and sparse tokens, respectively, such that $n_s < n_f < \min(m, T - n_s)$, and n_d be a parameter satisfying $n_d \leq n_f$. Next, we describe the process of generating the dataset \mathcal{D} ; see Fig. 2 for an example.

Definition 3.3 (Dataset Generation). Consider the vocabulary in Definition 3.2. To generate a data point (\mathbf{X}, y) , we first sample the label $y \in \{\pm 1\}$ uniformly at random. We divide the indices $[T]$ into three sets $\mathcal{I}_{\text{frequent}}, \mathcal{I}_{\text{sparse}}$ and $\mathcal{I}_{\text{irrelevant}}$, and sample each set as follows:

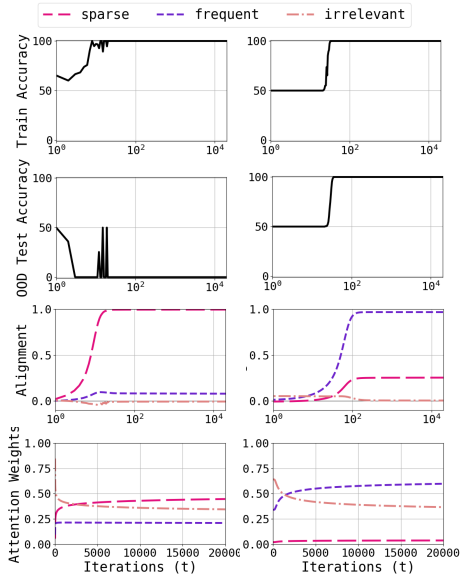


Figure 3: Train and test dynamics for a single-layer self-attention model (Eq. (3)) using the synthetic data visualized in Fig. 2; see Section 3.1 for details. **Left column:** the predictor that uses *sparse* tokens has lower sensitivity (Ex. 1 in Figure 2), **Right column:** the predictor that uses *frequent* tokens has lower sensitivity (Ex. 2 in Figure 2); see Appendix A.1 for more examples.

- $\mathcal{I}_{\text{frequent}}$ is composed of $\lfloor (n_f + n_d)/2 \rfloor$ tokens uniformly sampled from $\mathcal{V}_{\text{frequent}}^y$ and $n_f - \lfloor (n_f + n_d)/2 \rfloor$ tokens uniformly sampled from $\mathcal{V}_{\text{frequent}}^{-y}$.
- $\mathcal{I}_{\text{sparse}}$ contains n_s tokens uniformly sampled from $\mathcal{V}_{\text{sparse}}^y$.
- The remaining $T - n_f - n_s$ tokens in $\mathcal{I}_{\text{irrelevant}}$ are uniformly sampled from $\mathcal{V}_{\text{irrelevant}}$.

To determine if the tokens in $\mathcal{V}_{\text{sparse}}$ or those in $\mathcal{V}_{\text{frequent}}$ have a more significant impact on the model’s predictions, we adapt the test set generation process by altering the second step in Definition 3.3: we sample the sparse tokens from $\mathcal{V}_{\text{sparse}}^{-y}$ instead of $\mathcal{V}_{\text{sparse}}^y$. If this modification leads to a noticeable drop in the test accuracy, it suggests that the model relies on the sparse feature(s) for its predictions.

We consider two other metrics to examine the role of the attention head and the linear predictor. Define three vectors: $\mathbf{v}_{\text{sp}} := \mathbf{e}_1 - \mathbf{e}_2$, $\mathbf{v}_{\text{freq}} := \sum_{i \in \mathcal{V}_{\text{frequent}}^+} \mathbf{e}_i - \sum_{i \in \mathcal{V}_{\text{frequent}}^-} \mathbf{e}_i$, $\mathbf{v}_{\text{irrel}} := \sum_{i \in \mathcal{V}_{\text{irrelevant}}} \mathbf{e}_i$. We plot the average

alignment (cosine similarity) between the rows of UW_V^T and these vectors to see what tokens the prediction head relies on. Similarly, we plot the sum of the softmax scores for the three types of tokens to see which tokens are selected by the attention mechanism. We set \mathcal{P} in Definition 3.1 as the uniform distribution over \mathcal{V} for computing sensitivity in our experiments.

Results. Figure 3 shows the train and test dynamics of the model in Eq. (3) using synthetic datasets generated by following the process in Definition 3.3 (details in Table 1). We consider two cases: in the first case (left column), using the sparse token leads to a function with lower sensitivity, whereas in the second case (right column), using the frequent tokens leads to lower sensitivity (see Table 1 for a comparison of the sensitivity values). We observe that in the first case, the OOD test accuracy drops to 0, the alignment with \mathbf{v}_{sp} is close to 1 and the attention weights on the sparse tokens are the highest. These results show that the model relies on the sparse token in this case. On the other hand, in the second case, the test accuracy remains high, the alignment with \mathbf{v}_{freq} is close to 1 and the attention weights on the frequent tokens are the highest, which shows that the model relies on the frequent tokens. These results show that the model exhibits a low-sensitivity bias. Note that in both cases, the model can learn a function that relies on a sparse set of inputs (using the sparse tokens), however, it uses these tokens only when doing so leads to lower sensitivity.

4 INVESTIGATIONS ON VISION TASKS

In this section, we test whether our notion of sensitivity captures the inductive bias of transformers on vision tasks. We consider Vision Transformers (ViT, Dosovitskiy et al., 2021) which regard images as a sequence of patches instead of a tensor of pixels.

Definition 4.1 (Tokenization for Vision Transformers). Let $\mathbf{X} \in \mathbb{R}^{n_h \times n_w \times n_c}$ be the image with height n_h , width n_w , and number of channels n_c . A tokenization of \mathbf{X} is a sequence of T image patches $\{\mathbf{e}_1, \dots, \mathbf{e}_T\}$ where each token \mathbf{e}_i represents an image patch of dimension $d = n_w n_h n_c / T$.

Sensitivity is measured on the *training* set because our goal is to understand the simplicity bias of the model at training time, to see if it prefers to learn certain simple classes of functions on the training data. Since different models could have different generalization capabilities, the sensitivity on test data might not reflect the model’s preference for low-sensitivity functions at training time. Further, since the choice of optimization algorithm could in principle introduce its own bias and our goal is to understand the bias of the architecture, we train both the models with the same optimization algorithm, namely SGD; see Fig. 11 in the App. for a comparison with Adam.

We consider three datasets in this section (see Appendix D for details), namely CIFAR-10 (Krizhevsky, 2009), ImageNet-1k (Russakovsky et al., 2015) and Fashion-MNIST (Xiao et al., 2017). We use $\sigma^2 = 1, 15$ and 5, respectively. We use a variant of the ViT architecture for small-scale datasets proposed in Lee et al. (2021), referred to as ViT-small here onwards; see Appendix D for more details and Appendix A.3 for additional results where we show that varying model depths and number of heads does not affect the sensitivity of ViT models. We also compare the sensitivity of the ViT-small model and a ResNet-18 (He et al., 2016) CNN on the SVHN dataset with $\sigma^2 = 1$ in Appendix A.3, which leads to the same conclusion.

Transformers learn lower sensitivity functions than CNNs. Figure 4 shows the train accuracies as well as the sensitivity comparison between two ViTs: the ViT-small model and a ViT-simple

model (Beyer et al., 2022), two CNNs: a ResNet-18 and a DenseNet-121 (Huang et al., 2016), and a ConvMixer model (Trockman & Kolter, 2022a). Note that the train accuracies are comparable for all architectures, which allows for a fair comparison of sensitivity. We observe that the ViTs have significantly lower sensitivity compared to the CNNs and the ConvMixer model. At the end of the training, the sensitivity values are 0.3673 for DenseNet-121, 0.0827 for ResNet-18, 0.0829 for ConvMixer, 0.0050 for ViT-small and 0.0014 for ViT-simple.

Do Transformers have lower sensitivity than CNNs because these models process inputs differently?

ViTs process inputs as a sequence of patches whereas CNNs do not, and hence a natural question to ask is if the difference in sensitivity between the two architectures is due to this difference in processing the inputs as opposed to differences in the architecture. To investigate this, we compare ViTs with ConvMixer (Trockman & Kolter, 2022b). Similar to ViTs, ConvMixer processes the input data in a patch-wise manner, but has two key differences: it does not use the self-attention mechanism, which is the core component of transformers, and it relies on convolutions for the feedforward part as well. The higher sensitivity of the ConvMixer model indicates that the low sensitivity simplicity bias of the transformers is not because they process inputs patch-wise, but rather a result of other components of the architecture.

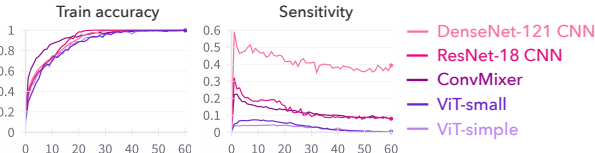


Figure 4: **Sensitivity on CIFAR-10.** Comparison of the sensitivity of two CNNs, two ViTs, and ConvMixer trained on the CIFAR-10 dataset, as a function of training epochs. For a fair comparison, the figure also shows the train accuracies (see App. Fig. 13 for full train dynamics). All models have similar accuracies but the ViTs have significantly lower sensitivity.

Do these observations generalize to pre-trained models? To study this, we consider the ImageNet-1k dataset (Russakovsky et al., 2015). We compare the sensitivity values of pre-trained ConvNext (Liu et al., 2022) and ViT/L-16 (Dosovitskiy et al., 2021) models. For comparable accuracies, ViT/L-16 has a sensitivity of 0.0191, which is lower than that of ConvNext at 0.0342. This shows that the observations on small-scale models studied in this section transfer to large-scale pretrained models.

Transformers learn lower sensitivity functions than MLPs. Next, we consider the Fashion-MNIST dataset and compare the sensitivity of ViT-small, a 3-hidden-layer CNN, an MLP with LeakyReLU activation and an MLP with sigmoid activation (see Fig. 15 in the Appendix for the training curves). At the end of training, the sensitivity values are 0.0559 for the MLP with LeakyReLU, 0.0505 for the MLP with sigmoid, 0.0453 for the CNN and 0.0098 for the ViT.

Thus, transformers learn lower sensitivity functions compared to MLPs, ConvMixers, and CNNs.

5 INVESTIGATIONS ON LANGUAGE TASKS

In this section, we investigate the sensitivity of transformers on natural language tasks, where each datapoint is a sequence of tokens. Similar to the comparison of ViTs with MLPs and CNNs in Section 4, we compare a RoBERTa (Liu et al., 2019) transformer model with LSTMs (Hochreiter & Schmidhuber, 1997), an alternative auto-regressive model, in this section. Recall that we consider a transformer with linear attention for the results in Section 2. Aligning with this setup, we also consider a RoBERTa model with ReLU activation in the attention layer (*i.e.*, replacing $\varphi(\cdot)$ in Eq. (3) with $\text{ReLU}(\cdot)$) for our experiments.

We use the usual RoBERTa-like tokenization procedure to process inputs for all the models so that they are represented as $\langle s \rangle e_1, \dots, e_T \langle /s \rangle$ where each e_j represents tokens that are usually subwords and $\langle s \rangle$ represents the classification (CLS) token, T the sequence length, and $\langle /s \rangle$ the separator token. We denote $e_0 = \langle s \rangle$ and $e_{T+1} = \langle /s \rangle$. For each token e_j , a token embedding $h_E(\cdot): [M] \rightarrow \mathbb{R}^d$ is trained during the process, where M denotes the vocabulary size. For transformers, we also train a separate positional encoder $h_P(\cdot): [N] \rightarrow \mathbb{R}^d$, where N denotes the maximum sequence length. We denote $e_j^{\text{LSTM}} = h_E^{\text{LSTM}}(e_j)$ and $e_j^{\text{RoBERTa}} = h_E^{\text{RoBERTa}}(e_j) + h_P^{\text{RoBERTa}}(j)$ as the embedding tokens of LSTM and RoBERTa, respectively. We omit the superscript for convenience.

To control the relative magnitude of noise, the embeddings $e_\tau \leftarrow \text{LayerNorm}(e_\tau)$ are first layer-normalized (Ba et al., 2016) before the additive Gaussian corruption. To better control possible

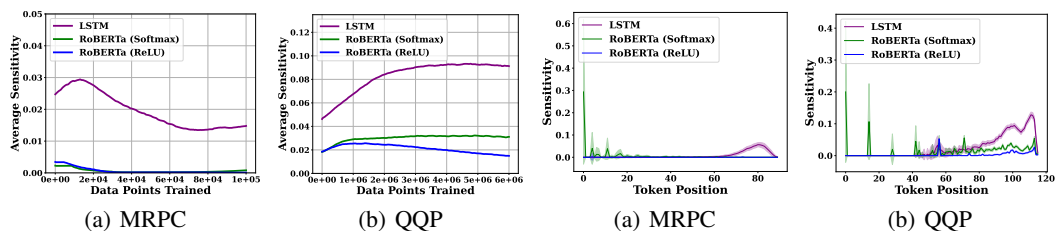


Figure 5: **Sensitivity over Datapoints Trained.** Figure 6: **Sensitivity over Token Position.** On both datasets, the Transformer-based model both datasets, LSTMs are more sensitive to later tokens than early ones, while RoBERTa displays much lower sensitivity compared to LSTMs during the entire training process. RoBERTa with ReLU activation has lower sensitivity compared to its Softmax counterpart at later stages of training. RoBERTa with ReLU activation has lower sensitivity compared to its Softmax counterpart at later early bumps in early tokens which come from the CLS token $\langle s \rangle$.

confounders, we limit both LSTM and RoBERTa to having the same number of layers. Both models are trained from scratch, *without* any pretraining on larger corpora, to ensure fair comparisons.

We consider two binary classification datasets, MRPC (Dolan & Brockett, 2005) and QQP (Iyer et al., 2017) (see Appendix D for details), which are relatively easy to learn without pretraining (Kovaleva et al., 2019). Empirically, we set $\sigma^2 = 15$ (results with $\sigma^2 = 4$ in App. A.4 yield similar observations as the results in this section). Similar to Section 4, we measure sensitivity on the train set (results on the validation set in App. A.4 yield similar observations). We include results with different depth values for RoBERTa as well as using GPT-2 in App. A.4 and they lead to similar conclusions.

Transformers learn lower sensitivity functions than LSTMs. As shown in Figure 5, both RoBERTa models have lower sensitivity than LSTMs on both datasets, regardless of the number of datapoints trained. Even at initialization with random weights, LSTMs are more sensitive. At the end of training, the sensitivity values on the MRPC dataset are 0.15, 0.002 and 0.001 for the LSTM, the RoBERTa model with softmax activation and the RoBERTa with ReLU activation, respectively. On the QQP dataset, LSTM, RoBERTa-softmax and RoBERTa-ReLU have sensitivity values of 0.09, 0.03 and 0.02, respectively. Interestingly, RoBERTa with ReLU activation also has lower sensitivity than its softmax counterpart. This may be because softmax attention encourages sparsity because of which the model can be more sensitive to a particular token; see Ex. 2 in Fig. 2 and bottom row of Fig. 3 for an example where sparsity can lead to higher sensitivity.

LSTMs are more sensitive to later tokens. In Fig. 6, we plot sensitivity over the token positions. We observe that LSTMs exhibit larger sensitivity towards the end of the sequence, i.e. at later token positions. In contrast, transformers are relatively uniform. Similar observations were made by (Fu et al., 2023) for a linear regression setting: LSTMs do more local updates and only remember the most recent observations, whereas transformers preserve global information and have longer memory.

Transformers are sensitive to the CLS token. In Fig. 6, we also observe that the RoBERTa model with softmax activation has frequent bumps in the sensitivity values at early token positions. This is because different sequences have different lengths and while computing sensitivity versus token positions, we align all the sequences to the right. These bumps at early token positions indeed correspond to the starting token after the tokenization procedure, the CLS token $\langle s \rangle$. This aligns with the observation of Jawahar et al. (2019) that the CLS token gathers all global information. Perturbing the CLS token corrupts the aggregation and results in high sensitivity. We also observe that RoBERTa with ReLU activation seems less sensitive to the CLS token compared to its softmax counterpart.

6 IMPLICATIONS OF LOW SENSITIVITY BIAS

We saw in Section 4 that transformers learn lower sensitivity functions than CNNs. In this section, we first compare the test performance of these models on the CIFAR-10-C dataset and show that transformers are more robust than CNNs. Next, we add a regularization term while training the transformer, to encourage lower sensitivity. The results demonstrate that lower sensitivity leads to improved robustness. We then explore the connection between sensitivity and the flatness of the minima. Our results show that lower sensitivity leads to flatter minima. Finally, we examine if

sensitivity can be used to understand the training dynamics of transformers, where we find sensitivity to be a suitable progress measure for certain grokking instances.

6.1 LOWER SENSITIVITY LEADS TO IMPROVED ROBUSTNESS

The CIFAR-10-C dataset (Hendrycks & Dietterich, 2019) was developed to benchmark the performance of various NNs on object recognition tasks under common corruptions that are not confusing to humans. Images from the test set of CIFAR-10 are corrupted with 14 types of algorithmically generated corruptions from blur, noise, weather, and digital categories (see Fig. 1 in Hendrycks & Dietterich (2019) for examples).

Fig. 7 compares the performance of two CNNs: ResNet-18 and DenseNet-121 with two ViTs: ViT-small and ViT-simple on various corruptions from the CIFAR-10-C dataset, at the end of training. We observe that the ViTs have lower sensitivity and better test performance on almost all corruptions compared to the CNNs, which have a higher sensitivity. Since the definition of sensitivity involves the addition of noise and ViTs have lower sensitivity, one can expect to be robust to various noise corruptions. However, the ViTs also have better test performance on several corruptions from weather and digital categories, which are significantly different from noise corruptions. This is consistent with the observations in Mahmood et al. (2021); Bhojanapalli et al. (2021).

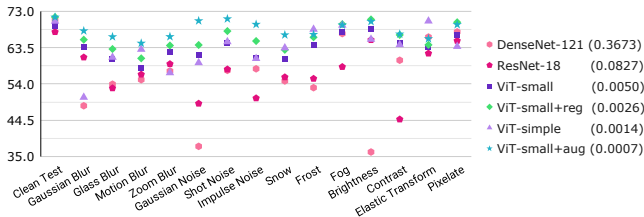


Figure 7: Comparison of the test accuracies on CIFAR-10 and on various corruptions from the CIFAR-10-C dataset (see Section 6 for details) of various models trained on the CIFAR-10 dataset, at the last training epoch (see App. Fig. 17 for a comparison of the accuracies as a function of training epochs.). We observe that the Vision Transformer models ViT-small and ViT-simple exhibit lower sensitivity and higher robustness to corruptions compared to the CNN models DenseNet-121 and ResNet-18. Additionally, encouraging lower sensitivity while training through regularization (ViT-small-reg) and data augmentation (ViT-small-aug) leads to improved robustness (see Section 6 for details).

Next, we conduct an experiment to investigate the role of low sensitivity in the robustness of transformers. We add a regularization term while training the model to explicitly encourage it to have lower sensitivity. If this model is more robust, then we can disentangle the role of low sensitivity from the role of the architecture and establish a concrete connection between lower sensitivity and improved robustness. To add the regularization, we use the fact that sensitivity can be estimated efficiently via sampling and consider two methods. In the first method (augmentation), we augment the training set by injecting the images with Gaussian noise (mean 0, variance 0.1) while preserving the label, and train the ViT on the augmented training set. In the second method (regularization), we add a mean squared error term using the model outputs for the original image and the image with Gaussian noise (mean 0, variance 1) injected into a randomly selected patch.

Fig. 7 also shows the test performance of ViT-small trained with augmentation and regularization methods on various corruptions from CIFAR-10-C. We observe that ViTs trained with these methods exhibit lower sensitivity compared to vanilla training. This is accompanied by an improved test performance on various corruptions, particularly on the noise and blur categories. **As encouraging lower sensitivity improves robustness**, the inductive bias of transformers to learn functions of lower sensitivity **could explain their better** robustness (to common corruptions) compared to CNNs.

6.2 LOWER SENSITIVITY LEADS TO FLATTER MINIMA

In this section, we investigate the connection between low sensitivity and flat minima. Consider a linear model $\Phi(\theta; x) = \theta^\top x$. Measuring sensitivity involves perturbing the input by some Δx . Prediction on the perturbed input is equivalent to perturbing the weight vector with $\Delta \theta = \frac{\theta^\top \Delta x}{\|x\|_2} x$, as

$$\Phi(\theta; x + \Delta x) = \theta^\top (x + \Delta x) = \Phi(\theta; x) + \theta^\top \Delta x = \Phi(\theta; x) + \Delta \theta^\top x = \Phi(\theta + \Delta \theta; x).$$

This draws a natural connection between sensitivity, which is measured with perturbation in the input space, and flatness of minima, which is measured with perturbation in the weight space (Keskar et al.,

Below, we investigate whether such a connection extends to more complex architectures such as transformers. Given model Φ and train set \mathcal{D} , we consider two metrics to measure the flatness of the minimum, based on the model outputs and model predictions, respectively,

$$\text{ShOp} := \mathbb{E}_{\mathbf{x} \sim \mathcal{D}, \boldsymbol{\xi} \sim \mathcal{N}(0, \sigma^2 \mathbf{I})} |\Phi(\boldsymbol{\theta}; \mathbf{x}) - \Phi(\boldsymbol{\theta} + \boldsymbol{\xi}; \mathbf{x})|, \quad \text{ShPred} := \mathbb{E}_{\mathbf{x} \sim \mathcal{D}, \boldsymbol{\xi} \sim \mathcal{N}(0, \sigma^2 \mathbf{I})} \mathbb{1}[f(\boldsymbol{\theta}; \mathbf{x}) \neq f(\boldsymbol{\theta} + \boldsymbol{\xi}; \mathbf{x})],$$

where $f(\boldsymbol{\theta}; \mathbf{x}) = \mathbb{1}[\Phi(\boldsymbol{\theta}; \mathbf{x}) \geq 0]$. Intuitively, for flatter minima, the model output and hence its prediction would remain relatively invariant to small perturbations in the model parameters.

Table 2 shows a comparison of these metrics for the ViT-small model trained with and without the sensitivity regularization at the end of training. Both metrics indicate that lower sensitivity corresponds to a flatter minimum. It is widely believed that flatter minima correlate with better generalization (Jiang* et al., 2020; Keskar et al., 2017; Neyshabur et al., 2017), though they may not always be correlated (Andriushchenko et al., 2023). Our results indicate that low-sensitivity correlates with improved generalization and investigating this connection for other settings can be an interesting direction for future work.

Table 2: Comparison of two sharpness metrics at the end of training the ViT-small model on the CIFAR-10 dataset with and without the sensitivity regularization. Lower values correspond to flatter minima; see text for discussion.

Setting	ShOp	ShPred
ViT-small + vanilla training	39.166	0.5346
ViT-small + sensitivity regularization	9.025	0.3982

6.3 SENSITIVITY AS A PROGRESS MEASURE FOR GROKING

In this section, we investigate if the sensitivity notion could serve as a progress measure for grokking (Nanda et al., 2023; Chen et al., 2024). We train an one-layer Transformer model on the modular addition task $a + b \bmod 113$. When evaluating sensitivity, we

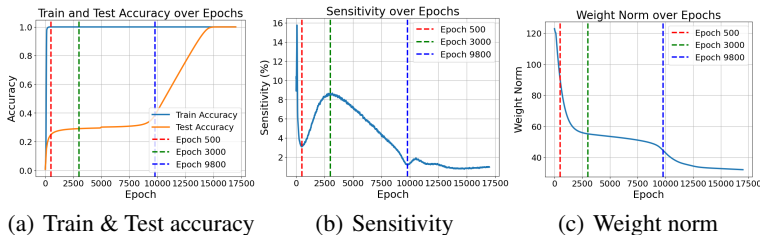


Figure 8: Sensitivity measures progress on modular addition task $a + b \bmod 113$ and indicates different stages of grokking.

add a random Gaussian noise with $\sigma = 0.1$ to the number embeddings. As shown in Figure 8, the test accuracy stays low from epoch 500 to 9,800 while the training accuracy saturates. However, sensitivity values continue to decrease smoothly starting from epoch 3000, and hence it provides a measure of the *hidden progress* (Barak et al., 2022) which the model makes even though the loss does not change, and indicates stages of grokking. In contrast, the weight norm is not a progress measure since it has the same flat curve as test accuracy during epoch 3,000 to 9,800.

As discussed by Nanda et al. (2023), grokking occurs when the model learns to use Fourier features to solve the task. A further bump in sensitivity after the grokked phase at epoch 9800 suggests that the model initially learns less robust Fourier features. At this stage, a small random noise could slightly disrupt the model’s performance. Over time, the Fourier basis becomes more robust. See Appendix A.5 for further discussion and results for more settings.

7 CONCLUSION

In this work, we investigate how the notion of sensitivity, which has shown promise in understanding inductive biases of Transformers on Boolean functions in prior work, can be extended to more realistic settings involving real-valued data. Our results show that transformers learn functions that have low sensitivity to small token-wise input perturbations, compared to other architectures, across vision and language tasks. We corroborate these observations with theoretical results, showing that transformers exhibit spectral bias and lower sensitivity corresponds to better robustness. We also demonstrate three important implications of this low-sensitivity bias: it correlates with improved robustness, flatter minima in the loss landscape, and serves as a progress measure that offers insights about the training dynamics. Investigating sensitivity as a progress measure in more settings can be an interesting direction for future work.

REFERENCES

- 540
541
542 Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What learning
543 algorithm is in-context learning? investigations with linear models. In *The Eleventh International
544 Conference on Learning Representations*, 2023. 27
- 545 Maksym Andriushchenko, Francesco Croce, Maximilian Mueller, Matthias Hein, and Nicolas Flam-
546 marion. A modern look at the relationship between sharpness and generalization. In *International
547 Conference on Machine Learning*, 2023. URL [https://api.semanticscholar.org/
548 CorpusID:256846369](https://api.semanticscholar.org/CorpusID:256846369). 10
- 549 Sanjeev Arora, Nadav Cohen, Wei Hu, and Yuping Luo. Implicit regularization in deep matrix
550 factorization. In *Advances in Neural Information Processing Systems*, volume 32, 2019. 28
551
- 552 Devansh Arpit, Stanisław Jastrzębski, Nicolas Ballas, David Krueger, Emmanuel Bengio, Maxinder S
553 Kanwal, Tegan Maharaj, Asja Fischer, Aaron Courville, Yoshua Bengio, et al. A closer look at
554 memorization in deep networks. In *International conference on machine learning*, pp. 233–242.
555 PMLR, 2017. 27
- 556 Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer normalization, 2016. 7
557
- 558 Yu Bai, Fan Chen, Huan Wang, Caiming Xiong, and Song Mei. Transformers as statisticians:
559 Provable in-context learning with in-context algorithm selection. *arXiv preprint arXiv:2306.04637*,
560 2023. 27
- 561 Boaz Barak, Benjamin Edelman, Surbhi Goel, Sham Kakade, Eran Malach, and Cyril Zhang. Hidden
562 progress in deep learning: Sgd learns parities near the computational limit. *Advances in Neural
563 Information Processing Systems*, 35:21750–21764, 2022. 10
564
- 565 Ronen Basri, David Jacobs, Yoni Kasten, and Shira Kritchman. The convergence rate of neural
566 networks for learned functions of different frequencies, 2019. 28
- 567 Lucas Beyer, Xiaohua Zhai, and Alexander Kolesnikov. Better plain vit baselines for imagenet-1k,
568 2022. 7
569
- 570 Satwik Bhattamishra, Arkil Patel, Phil Blunsom, and Varun Kanade. Understanding in-context
571 learning in transformers and llms by learning to learn discrete functions, 2023a. 27
- 572 Satwik Bhattamishra, Arkil Patel, Varun Kanade, and Phil Blunsom. Simplicity bias in transformers
573 and their ability to learn sparse Boolean functions. In *Proceedings of the 61st Annual Meeting of
574 the Association for Computational Linguistics*, 2023b. 1, 3, 28
575
- 576 Srinadh Bhojanapalli, Ayan Chakrabarti, Daniel Glasner, Daliang Li, Thomas Unterthiner, and
577 Andreas Veit. Understanding robustness of transformers for image classification. In *Proceedings of
578 the IEEE/CVF international conference on computer vision*, pp. 10231–10241, 2021. 1, 9, 27, 28
- 579 Alberto Bietti and Julien Mairal. *On the inductive bias of neural tangent kernels*. Curran Associates
580 Inc., Red Hook, NY, USA, 2019. 28
581
- 582 Christopher M. Bishop. Training with noise is equivalent to tikhonov regularization. *Neural
583 Computation*, 7:108–116, 1995. 28
- 584 Guy Blanc, Neha Gupta, Gregory Valiant, and Paul Valiant. Implicit regularization for deep neural
585 networks driven by an ornstein-uhlenbeck like process. In *Conference on learning theory*, pp.
586 483–513. PMLR, 2020. 28
587
- 588 Simone Bombari and Marco Mondelli. Towards understanding the word sensitivity of attention
589 layers: A study via random features, 2024. 28
- 590 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
591 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel
592 Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler,
593 Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray,
Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever,

- 594 and Dario Amodei. Language models are few-shot learners. In H. Larochelle, M. Ranzato,
595 R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neural Information Processing Systems*,
596 volume 33, pp. 1877–1901. Curran Associates, Inc., 2020. 27
- 597 Yuan Cao, Zhiying Fang, Yue Wu, Ding-Xuan Zhou, and Quanquan Gu. Towards understanding the
598 spectral bias of deep learning. In *IJCAI*, 2021. 28
- 600 Angelica Chen, Ravid Shwartz-Ziv, Kyunghyun Cho, Matthew L Leavitt, and Naomi Saphra. Sudden
601 drops in the loss: Syntax acquisition, phase transitions, and simplicity bias in MLMs. In *The Twelfth*
602 *International Conference on Learning Representations*, 2024. URL [https://openreview.net/](https://openreview.net/forum?id=M05PiKHELW)
603 [forum?id=M05PiKHELW](https://openreview.net/forum?id=M05PiKHELW). 10
- 604 Ting-Rui Chiang. On a benefit of mask language modeling: Robustness to simplicity bias, 2021. 28
- 605 Lenaic Chizat and Francis Bach. Implicit bias of gradient descent for wide two-layer neural networks
606 trained with the logistic loss. In *Conference on Learning Theory*, pp. 1305–1338. PMLR, 2020. 28
- 609 Arthur Conmy, Augustine N. Mavor-Parker, Aengus Lynch, Stefan Heimersheim, and Adrià Garriga-
610 Alonso. Towards automated circuit discovery for mechanistic interpretability. In *Neural Information*
611 *Processing Systems*, 2023. 27
- 612 Ekin Dogus Cubuk, Barret Zoph, Dandelion Mané, Vijay Vasudevan, and Quoc V. Le. Autoaugment:
613 Learning augmentation strategies from data. *2019 IEEE/CVF Conference on Computer Vision and*
614 *Pattern Recognition (CVPR)*, pp. 113–123, 2019. 28
- 616 Alexander G. de G. Matthews, Mark Rowland, Jiri Hron, Richard E. Turner, and Zoubin Ghahramani.
617 Gaussian process behaviour in wide deep neural networks, 2018. 26
- 618 William B. Dolan and Chris Brockett. Automatically constructing a corpus of sentential paraphrases.
619 In *Proceedings of the Third International Workshop on Paraphrasing (IWP2005)*, 2005. URL
620 <https://aclanthology.org/I05-5002>. 8, 30
- 622 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
623 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit,
624 and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale.
625 In *International Conference on Learning Representations*, 2021. URL [https://openreview.](https://openreview.net/forum?id=YicbFdNTTy)
626 [net/forum?id=YicbFdNTTy](https://openreview.net/forum?id=YicbFdNTTy). 1, 6, 7
- 627 Spencer Frei, Gal Vardi, Peter L Bartlett, Nathan Srebro, and Wei Hu. Implicit bias in leaky relu
628 networks trained on high-dimensional data. *arXiv preprint arXiv:2210.07082*, 2022. 28
- 629 Deqing Fu, Tian-Qi Chen, Robin Jia, and Vatsal Sharan. Transformers learn higher-order optimization
630 methods for in-context learning: A study with linear models. *arXiv*, abs/2310.17086, 2023. 8, 27
- 632 Shivam Garg, Dimitris Tsipras, Percy Liang, and Gregory Valiant. What can transformers learn
633 in-context? a case study of simple function classes. *ArXiv*, abs/2208.01066, 2022. 27
- 634 Adrià Garriga-Alonso, Carl Edward Rasmussen, and Laurence Aitchison. Deep convolutional
635 networks as shallow gaussian processes. In *International Conference on Learning Representations*,
636 2019. URL <https://openreview.net/forum?id=Bklfsi0cKm>. 26
- 638 Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A. Wichmann, and
639 Wieland Brendel. Imagenet-trained CNNs are biased towards texture; increasing shape bias
640 improves accuracy and robustness. In *International Conference on Learning Representations*,
641 2019. URL <https://openreview.net/forum?id=Bygh9j09KX>. 27, 28
- 642 Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel, Matthias
643 Bethge, and Felix A Wichmann. Shortcut learning in deep neural networks. *Nature Machine*
644 *Intelligence*, 2(11):665–673, 2020. 27
- 645 Soumya Suvra Ghosal, Yifei Ming, and Yixuan Li. Are vision transformers robust to spurious
646 correlations?, 2022. 27

- 648 Parikshit Gopalan, Noam Nisan, Rocco A Servedio, Kunal Talwar, and Avi Wigderson. Smooth
649 boolean functions are easy: Efficient algorithms for low-sensitivity functions. In *Proceedings of*
650 *the 2016 ACM Conference on Innovations in Theoretical Computer Science*, pp. 59–70, 2016. 4
- 651
- 652 Tianyu Guo, Wei Hu, Song Mei, Huan Wang, Caiming Xiong, Silvio Savarese, and Yu Bai. How
653 do transformers learn in-context beyond simple functions? a case study on learning with repre-
654 sentations. In *The Twelfth International Conference on Learning Representations*, 2024. URL
655 <https://openreview.net/forum?id=ikwEDva1JZ>. 27
- 656 Suchin Gururangan, Swabha Swayamdipta, Omer Levy, Roy Schwartz, Samuel Bowman, and
657 Noah A. Smith. Annotation artifacts in natural language inference data. In Marilyn Walker,
658 Heng Ji, and Amanda Stent (eds.), *Proceedings of the 2018 Conference of the North American*
659 *Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume*
660 *2 (Short Papers)*, pp. 107–112, New Orleans, Louisiana, June 2018. Association for Computational
661 Linguistics. doi: 10.18653/v1/N18-2017. URL <https://aclanthology.org/N18-2017>. 28
- 662 Michael Hahn and Mark Rofin. Why are sensitive functions hard for transformers? In Lun-Wei
663 Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the*
664 *Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 14973–15008, Bangkok,
665 Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.
666 acl-long.800. URL <https://aclanthology.org/2024.acl-long.800>. 1
- 667 Michael Hanna, Ollie Liu, and Alexandre Variengien. How does GPT-2 compute greater-than?:
668 Interpreting mathematical abilities in a pre-trained language model. In *Thirty-seventh Conference*
669 *on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=p4PckNQR8k>. 27
- 670
- 671 Jeff Z HaoChen, Colin Wei, Jason Lee, and Tengyu Ma. Shape matters: Understanding the implicit
672 bias of the noise covariance. In *Conference on Learning Theory*, pp. 2315–2357. PMLR, 2021. 28
- 673
- 674 Michael Hassid, Hao Peng, Daniel Rotem, Jungo Kasai, Ivan Montero, Noah A. Smith, and Roy
675 Schwartz. How much does attention actually attend? questioning the importance of attention in
676 pretrained transformers, 2022. 27
- 677
- 678 Pooya Hatami, Raghav Kulkarni, and Denis Pankratov. Variations on the sensitivity conjecture.
679 *Theory of Computing*, pp. 1–27, 2011. 3
- 680 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep Residual Learning for Im-
681 age Recognition. In *Proceedings of 2016 IEEE Conference on Computer Vision and Pattern*
682 *Recognition, CVPR ’16*, pp. 770–778. IEEE, June 2016. doi: 10.1109/CVPR.2016.90. URL
683 <http://ieeexplore.ieee.org/document/7780459>. 6
- 684
- 685 Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common
686 corruptions and perturbations. In *International Conference on Learning Representations*, 2019.
687 URL <https://openreview.net/forum?id=HJz6tiCqYm>. 9
- 688 Sepp Hochreiter and Jürgen Schmidhuber. Long Short-Term Memory. *Neural Computation*, 9
689 (8):1735–1780, 11 1997. ISSN 0899-7667. doi: 10.1162/neco.1997.9.8.1735. URL <https://doi.org/10.1162/neco.1997.9.8.1735>. 7
- 690
- 691 Jiri Hron, Yasaman Bahri, Jascha Sohl-Dickstein, and Roman Novak. Infinite attention: Nngp and
692 ntk for deep attention networks. In *Proceedings of the 37th International Conference on Machine*
693 *Learning, ICML’20*. JMLR.org, 2020. 26
- 694
- 695 Gao Huang, Zhuang Liu, and Kilian Q. Weinberger. Densely connected convolutional networks.
696 *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2261–2269,
697 2016. URL <https://api.semanticscholar.org/CorpusID:9433631>. 7
- 698
- 699 Hao Huang. Induced subgraphs of hypercubes and a proof of the sensitivity conjecture. *Annals of*
700 *Mathematics*, 190(3):949–955, 2019. 3, 27
- 701
- 701 Minyoung Huh, Hossein Mobahi, Richard Zhang, Brian Cheung, Pulkit Agrawal, and Phillip Isola.
The low-rank simplicity bias in deep networks. *Trans. Mach. Learn. Res.*, 2023, 2021. 27, 28

- 702 Shankar Iyer, Nikhil Dandekar, and Kornel Csernai. First Quora Dataset
703 Release: Question Pairs, 2017. URL [https://quoradata.quora.com/
704 First-Quora-Dataset-Release-Question-Pairs](https://quoradata.quora.com/First-Quora-Dataset-Release-Question-Pairs). Online. 8, 30
705
- 706 Ganesh Jawahar, Benoît Sagot, and Djamé Seddah. What does BERT learn about the structure
707 of language? In Anna Korhonen, David Traum, and Lluís Màrquez (eds.), *Proceedings of the
708 57th Annual Meeting of the Association for Computational Linguistics*, pp. 3651–3657, Florence,
709 Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1356. URL
710 <https://aclanthology.org/P19-1356>. 8
- 711 Ziwei Ji and Matus Telgarsky. Risk and parameter convergence of logistic regression. *arXiv preprint
712 arXiv:1803.07300*, 2018. 28
713
- 714 Ziwei Ji and Matus Telgarsky. Directional convergence and alignment in deep learning. *Advances in
715 Neural Information Processing Systems*, 33:17176–17186, 2020. 28
716
- 717 Ziwei Ji and Matus Telgarsky. Characterizing the implicit bias via a primal-dual analysis. In
718 *Algorithmic Learning Theory*, pp. 772–804. PMLR, 2021. 28
- 719 Ziwei Ji, Nathan Srebro, and Matus Telgarsky. Fast margin maximization via dual acceleration. In
720 *International Conference on Machine Learning*, pp. 4860–4869. PMLR, 2021. 28
721
- 722 Yiding Jiang*, Behnam Neyshabur*, Hossein Mobahi, Dilip Krishnan, and Samy Bengio. Fantasic
723 generalization measures and where to find them. In *International Conference on Learning
724 Representations*, 2020. URL <https://openreview.net/forum?id=SJgIPJBFvH>. 10
725
- 726 John M. Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ron-
727 neberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Zidek, Anna Potapenko, Alex Bridg-
728 land, Clemens Meyer, Simon A A Kohl, Andy Ballard, Andrew Cowie, Bernardino Romera-
729 Paredes, Stanislav Nikolov, Rishub Jain, Jonas Adler, Trevor Back, Stig Petersen, David A. Reiman,
730 Ellen Clancy, Michal Zielinski, Martin Steinegger, Michalina Pacholska, Tamas Berghammer, Se-
731 bastian Bodenstern, David Silver, Oriol Vinyals, Andrew W. Senior, Koray Kavukcuoglu, Pushmeet
732 Kohli, and Demis Hassabis. Highly accurate protein structure prediction with alphafold. *Nature*,
733 596:583 – 589, 2021. URL <https://api.semanticscholar.org/CorpusID:235959867>. 1
734
- 735 Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping Tak Peter
736 Tang. On large-batch training for deep learning: Generalization gap and sharp minima. In
737 *International Conference on Learning Representations*, 2017. URL [https://openreview.net/
738 forum?id=H1oyRlYgg](https://openreview.net/forum?id=H1oyRlYgg). 9, 10
- 739 P. Kirichenko, Pavel Izmailov, and Andrew Gordon Wilson. Last layer re-training is sufficient for
740 robustness to spurious correlations. *ArXiv*, abs/2204.02937, 2022. 28
741
- 742 Goro Kobayashi, Tatsuki Kuribayashi, Sho Yokoi, and Kentaro Inui. Analyzing feed-forward blocks
743 in transformers through the lens of attention maps. In *The Twelfth International Conference on
744 Learning Representations*, 2024. URL <https://openreview.net/forum?id=mYWsyTuiRp>. 27
745
- 746 Yiwen Kou, Zixiang Chen, and Quanquan Gu. Implicit bias of gradient descent for two-layer relu
747 and leaky relu networks on nearly-orthogonal data, 2023. 28
748
- 749 Olga Kovaleva, Alexey Romanov, Anna Rogers, and Anna Rumshisky. Revealing the dark secrets of
750 BERT. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan (eds.), *Proceedings of the 2019
751 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint
752 Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 4365–4374, Hong Kong,
753 China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1445.
754 URL <https://aclanthology.org/D19-1445>. 8, 30
755
- 756 Alex Krizhevsky. Learning multiple layers of features from tiny images. pp. 32–33, 2009. URL
757 <https://www.cs.toronto.edu/~kriz/learning-features-2009-TR.pdf>. 6, 29
- 758 Yann LeCun and Corinna Cortes. The mnist database of handwritten digits. 2005. 19

- 756 Jaehoon Lee, Jascha Sohl-dickstein, Jeffrey Pennington, Roman Novak, Sam Schoenholz, and
757 Yasaman Bahri. Deep neural networks as gaussian processes. In *International Conference on*
758 *Learning Representations*, 2018. URL <https://openreview.net/forum?id=B1EA-M-0Z>. 25,
759 26
- 760 Seung Hoon Lee, Seunghyun Lee, and Byung Cheol Song. Vision transformer for small-size datasets,
761 2021. 6
- 762 Zhiyuan Li, Yuping Luo, and Kaifeng Lyu. Towards resolving the implicit bias of gradient descent
763 for matrix factorization: Greedy low-rank learning, 2021. 28
- 764 Zhiyuan Li, Tianhao Wang, Jason D Lee, and Sanjeev Arora. Implicit bias of gradient descent
765 on reparametrized models: On equivalence to mirror descent. *Advances in Neural Information*
766 *Processing Systems*, 35:34626–34640, 2022. 28
- 767 Sungbin Lim, Ildoo Kim, Taesup Kim, Chiheon Kim, and Sungwoong Kim. Fast autoaugment. In
768 *Neural Information Processing Systems*, 2019. 28
- 769 Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike
770 Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining
771 approach, 2019. 7
- 772 Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining
773 Xie. A convnet for the 2020s. In *2022 IEEE/CVF Conference on Computer Vision and Pattern*
774 *Recognition (CVPR)*, pp. 11966–11976, 2022. doi: 10.1109/CVPR52688.2022.01167. 7
- 775 Kaifeng Lyu and Jian Li. Gradient descent maximizes the margin of homogeneous neural networks.
776 In *International Conference on Learning Representations*, 2020. 28
- 777 Kaifeng Lyu, Zhiyuan Li, Runzhe Wang, and Sanjeev Arora. Gradient descent on two-layer nets:
778 Margin maximization and simplicity bias. In *Neural Information Processing Systems*, 2021. 28
- 779 Kaleel Mahmood, Rigel Mahmood, and Marten Van Dijk. On the robustness of vision transformers
780 to adversarial examples. In *Proceedings of the IEEE/CVF International Conference on Computer*
781 *Vision*, pp. 7838–7847, 2021. 1, 9, 27, 28
- 782 Luke Melas-Kyriazi. Do you even need attention? a stack of feed-forward layers does surprisingly
783 well on imagenet. *ArXiv*, abs/2105.02723, 2021. URL [https://api.semanticscholar.org/](https://api.semanticscholar.org/CorpusID:233864618)
784 [CorpusID:233864618](https://api.semanticscholar.org/CorpusID:233864618). 27
- 785 Kevin Meng, David Bau, Alex J Andonian, and Yonatan Belinkov. Locating and editing factual
786 associations in GPT. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho
787 (eds.), *Advances in Neural Information Processing Systems*, 2022. URL [https://openreview.](https://openreview.net/forum?id=-h6WAS6eE4)
788 [net/forum?id=-h6WAS6eE4](https://openreview.net/forum?id=-h6WAS6eE4). 27
- 789 Depen Morwani, Jatin Batra, Prateek Jain, and Praneeth Netrapalli. Simplicity bias in 1-hidden layer
790 neural networks, 2023. 27, 28
- 791 Mor Shpigel Nacson, Jason Lee, Suriya Gunasekar, Pedro Henrique Pamplona Savarese, Nathan
792 Srebro, and Daniel Soudry. Convergence of gradient descent on separable data. In *The 22nd*
793 *International Conference on Artificial Intelligence and Statistics*, pp. 3420–3428. PMLR, 2019. 28
- 794 Vaishnavh Nagarajan, Anders Andreassen, and Behnam Neyshabur. Understanding the failure modes
795 of out-of-distribution generalization. In *International Conference on Learning Representations*,
796 2021. URL https://openreview.net/forum?id=fSTD6NFIW_b. 28
- 797 Preetum Nakkiran, Gal Kaplun, Dimitris Kalimeris, Tristan Yang, Benjamin L. Edelman, Fred Zhang,
798 and Boaz Barak. *SGD on Neural Networks Learns Functions of Increasing Complexity*. Curran
799 Associates Inc., Red Hook, NY, USA, 2019. 27
- 800 Neel Nanda, Lawrence Chan, Tom Lieberum, Jess Smith, and Jacob Steinhardt. Progress measures
801 for grokking via mechanistic interpretability, 2023. 10, 24, 27

- 810 Muzammal Naseer, Kanchana Ranasinghe, Salman Hameed Khan, Munawar Hayat, Fahad Shahbaz
811 Khan, and Ming-Hsuan Yang. Intriguing properties of vision transformers. In *Neural Informa-*
812 *tion Processing Systems*, 2021. URL [https://api.semanticscholar.org/CorpusID:](https://api.semanticscholar.org/CorpusID:235125781)
813 [235125781](https://api.semanticscholar.org/CorpusID:235125781). 27
- 814 Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y. Ng. Reading
815 digits in natural images with unsupervised feature learning. In *NIPS Workshop on Deep Learn-*
816 *ing and Unsupervised Feature Learning 2011*, 2011. URL [http://ufldl.stanford.edu/](http://ufldl.stanford.edu/housenumbers/nips2011_housenumbers.pdf)
817 [housenumbers/nips2011_housenumbers.pdf](http://ufldl.stanford.edu/housenumbers/nips2011_housenumbers.pdf). 21, 29
- 818 Behnam Neyshabur, Ryota Tomioka, and Nathan Srebro. In search of the real inductive bias:
819 On the role of implicit regularization in deep learning. *CoRR*, abs/1412.6614, 2014. URL
820 <https://api.semanticscholar.org/CorpusID:6021932>. 27
- 821 Behnam Neyshabur, Srinadh Bhojanapalli, David McAllester, and Nathan Srebro. Exploring
822 generalization in deep learning. In *Neural Information Processing Systems*, 2017. URL
823 <https://api.semanticscholar.org/CorpusID:9597660>. 10
- 824 Roman Novak, Yasaman Bahri, Daniel A Abolafia, Jeffrey Pennington, and Jascha Sohl-
825 Dickstein. Sensitivity and generalization in neural networks: an empirical study. *arXiv preprint*
826 *arXiv:1802.08760*, 2018. 1
- 827 Roman Novak, Lechao Xiao, Yasaman Bahri, Jaehoon Lee, Greg Yang, Daniel A. Abolafia, Jeffrey
828 Pennington, and Jascha Sohl-dickstein. Bayesian deep convolutional networks with many channels
829 are gaussian processes. In *International Conference on Learning Representations*, 2019. URL
830 <https://openreview.net/forum?id=B1g30j0qF7>. 26
- 831 Ryan O’Donnell. *Analysis of Boolean Functions*. Cambridge University Press, 2014. doi: 10.1017/
832 CBO9781139814782. 1, 2, 26, 27
- 833 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor
834 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward
835 Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang,
836 Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning
837 library, 2019. 29
- 838 Sayak Paul and Pin-Yu Chen. Vision transformers are robust learners. In *Proceedings of the AAAI*
839 *conference on Artificial Intelligence*, volume 36, pp. 2071–2081, 2022. 1, 27, 28
- 840 Mohammad Pezeshki, Sekouba Kaba, Yoshua Bengio, Aaron C. Courville, Doina Precup, and
841 Guillaume Lajoie. Gradient starvation: A learning proclivity in neural networks. In *Neural*
842 *Information Processing Systems*, 2020. 28
- 843 Mary Phuong and Christoph H Lampert. The inductive bias of re{lu} networks on orthogonally
844 separable data. In *International Conference on Learning Representations*, 2021. URL [https:](https://openreview.net/forum?id=krz7T0xU9Z_)
845 [//openreview.net/forum?id=krz7T0xU9Z_](https://openreview.net/forum?id=krz7T0xU9Z_). 28
- 846 Alethea Power, Yuri Burda, Harri Edwards, Igor Babuschkin, and Vedant Misra. Grokking: General-
847 ization beyond overfitting on small algorithmic datasets, 2022. 27
- 848 Philip Quirke and Fazl Barez. Understanding addition in transformers. In *International Conference*
849 *on Learning Representations (ICLR)*, Vienna, Austria, 2024. 27
- 850 Maithra Raghu, Thomas Unterthiner, Simon Kornblith, Chiyuan Zhang, and Alexey Dosovitskiy.
851 Do vision transformers see like convolutional neural networks? In A. Beygelzimer, Y. Dauphin,
852 P. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems*,
853 2021. URL <https://openreview.net/forum?id=G18FHfMVTzu>. 27
- 854 Nasim Rahaman, Aristide Baratin, Devansh Arpit, Felix Draxler, Min Lin, Fred Hamprecht, Yoshua
855 Bengio, and Aaron Courville. On the spectral bias of neural networks. In Kamalika Chaudhuri
856 and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine*
857 *Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 5301–5310. PMLR,
858 09–15 Jun 2019a. URL <https://proceedings.mlr.press/v97/rahaman19a.html>. 27
- 860

- 864 Nasim Rahaman, Aristide Baratin, Devansh Arpit, Felix Draxler, Min Lin, Fred Hamprecht, Yoshua
865 Bengio, and Aaron Courville. On the spectral bias of neural networks. In *International conference*
866 *on machine learning*, pp. 5301–5310. PMLR, 2019b. 1
- 867
868 Sylvestre-Alvise Rebuffi, Sven Gowal, Dan Andrei Calian, Florian Stimberg, Olivia Wiles, and
869 Timothy Mann. Data augmentation can improve robustness. In *Neural Information Processing*
870 *Systems*, 2021. 28
- 871
872 Herbert E. Robbins. A stochastic approximation method. *Annals of Mathematical Statistics*, 22:
873 400–407, 1951. 28
- 874
875 Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang,
876 Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet
877 Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115
(3):211–252, 2015. doi: 10.1007/s11263-015-0816-y. 6, 7, 29
- 878
879 Shiori Sagawa, Pang Wei Koh, Tatsunori B. Hashimoto, and Percy Liang. Distributionally robust
880 neural networks. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=ryxGuJrFvS>. 28
- 881
882 Harshay Shah, Kaustav Tamuly, Aditi Raghunathan, Prateek Jain, and Praneeth Netrapalli. The
883 pitfalls of simplicity bias in neural networks. *Advances in Neural Information Processing Systems*,
884 33, 2020. 27, 28
- 885
886 Rulin Shao, Zhouxing Shi, Jinfeng Yi, Pin-Yu Chen, and Cho-Jui Hsieh. On the adversarial robustness
887 of visual transformers. *arXiv preprint arXiv:2103.15670*, 1(2), 2021. 1, 27, 28
- 888
889 Lujia Shen, Yuwen Pu, Shouling Ji, Changjiang Li, Xuhong Zhang, Chunpeng Ge, and Ting Wang.
890 Improving the robustness of transformer-based large language models with dynamic attention.
arXiv preprint arXiv:2311.17400, 2023. 28
- 891
892 Zhouxing Shi, Huan Zhang, Kai-Wei Chang, Minlie Huang, and Cho-Jui Hsieh. Robustness verifica-
893 tion for transformers. *arXiv preprint arXiv:2002.06622*, 2020. 28
- 894
895 Zhouxing Shi, Yihan Wang, Huan Zhang, Jinfeng Yi, and Cho-Jui Hsieh. Fast certified robust training
896 with short warmup. *Advances in Neural Information Processing Systems*, 34:18335–18349, 2021.
28
- 897
898 Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Suriya Gunasekar, and Nathan Srebro. The implicit
899 bias of gradient descent on separable data. *The Journal of Machine Learning Research*, 19(1):
900 2822–2878, 2018. 28
- 901
902 Davoud Ataee Tarzanagh, Yingcong Li, Christos Thrampoulidis, and Samet Oymak. Transformers as
support vector machines. *ArXiv*, abs/2308.16898, 2023a. 27
- 903
904 Davoud Ataee Tarzanagh, Yingcong Li, Xuechen Zhang, and Samet Oymak. Max-margin token
905 selection in attention mechanism, 2023b. 27
- 906
907 Rishabh Tiwari and Pradeep Shenoy. Overcoming simplicity bias in deep networks using a feature
sieve, 2023. 28
- 908
909 Asher Trockman and J. Zico Kolter. Patches are all you need? *Transactions on Machine Learning*
910 *Research*, 2023, 2022a. URL <https://api.semanticscholar.org/CorpusID:245633423>.
7, 27
- 911
912 Asher Trockman and J Zico Kolter. Patches are all you need?, 2022b. URL <https://openreview.net/forum?id=TVHS5Y4dNm>. 7
- 913
914 Guillermo Valle-Perez, Chico Q Camargo, and Ard A Louis. Deep learning generalizes because the
915 parameter-function map is biased towards simple functions. *International Conference on Learning*
916 *Representations*, 2019. 1, 27
- 917
Gal Vardi. On the implicit bias in deep-learning algorithms, 2022. 28

- 918 Bhavya Vasudeva, Kameron Shahabi, and Vatsal Sharan. Mitigating simplicity bias in deep learning
919 for improved ood generalization and robustness, 2023. 28
920
- 921 Bhavya Vasudeva, Puneesh Deora, and Christos Thrampoulidis. Implicit bias and fast convergence
922 rates for self-attention, 2024. 27
- 923 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N
924 Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon,
925 U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett
926 (eds.), *Advances in Neural Information Processing Systems*, volume 30. Curran Associates,
927 Inc., 2017. URL [https://proceedings.neurips.cc/paper_files/paper/2017/file/
928 3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf). 1
- 929 Johannes von Oswald, Eyvind Niklasson, E. Randazzo, João Sacramento, Alexander Mordvintsev,
930 Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient descent. In
931 *International Conference on Machine Learning*, 2022. 27
932
- 933 Fali Wang, Zheng Lin, Zhengxiao Liu, Mingyu Zheng, Lei Wang, and Daren Zha. Macrobert:
934 Maximizing certified region of bert to adversarial word substitutions. In *Database Systems for
935 Advanced Applications: 26th International Conference, DASFAA 2021, Taipei, Taiwan, April
936 11–14, 2021, Proceedings, Part II* 26, pp. 253–261. Springer, 2021. 28
- 937 Kevin Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt. Inter-
938 pretability in the wild: a circuit for indirect object identification in gpt-2 small, 2022. 27
939
- 940 Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking
941 machine learning algorithms, 2017. 6, 29
- 942 Zhi-Qin John Xu, Yaoyu Zhang, Tao Luo, Yanyang Xiao, and Zheng Ma. Frequency principle:
943 Fourier analysis sheds light on deep neural networks. *arXiv preprint arXiv:1901.06523*, 2019. 27,
944 28
- 945 Greg Yang. Tensor programs i: Wide feedforward or recurrent neural networks of any architecture
946 are gaussian processes, 2021. 26
947
- 948 Greg Yang and Hadi Salman. A fine-grained spectral perspective on neural networks, 2020. 2, 25, 26,
949 27
- 950 Hongyi Zhang, Moustapha Cisse, Yann N. Dauphin, and David Lopez-Paz. mixup: Beyond empirical
951 risk minimization. In *International Conference on Learning Representations*, 2018. URL [https:
952 //openreview.net/forum?id=r1Ddp1-Rb](https://openreview.net/forum?id=r1Ddp1-Rb). 28
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

972	APPENDIX	
973		
974		
975	A Additional Experiments	19
976	A.1 Synthetic Data and the MNIST Dataset	19
977	A.2 Sensitivity with Random Noise instead of Token-wise Noise	19
978	A.3 Vision Tasks	20
979	A.4 Language Tasks	21
980	A.5 Sensitivity as a Progress Measure for Grokking	24
981	A.6 Sensitivity as a Progress Measure for Learning Sparse Parities	24
982		
983	B Proofs for Section 2	25
984	B.1 Proof of Proposition 2.1	26
985	B.2 Proof of Proposition 2.2	26
986		
987	C Related Work	27
988		
989	D Details of Experimental Settings	29
990		
991	E Limitations	30
992		
993		
994		
995		
996		
997		
998		
999	A ADDITIONAL EXPERIMENTS	

1000 A ADDITIONAL EXPERIMENTS

1001 In this section, we include some additional results to supplement the main experimental results for
 1002 synthetic data as well as the vision and language tasks.

1003 A.1 SYNTHETIC DATA AND THE MNIST DATASET

1004 In this section, we present some additional results for the low-sensitivity bias of a single-layer
 1005 self-attention model (Eq. (3)) on the synthetic dataset generated based on Definition 3.3, visualized in
 1006 Fig. 2. Similar to the results in Section 3.1, we consider three data settings where using the sparse
 1007 token leads to a function with lower sensitivity (Fig. 9, top row) and three settings where using
 1008 the frequent token leads to lower sensitivity (Fig. 9, bottom row). The exact data settings and a
 1009 comparison of the sensitivity values for each setting are shown in Table 3. These results yield similar
 1010 conclusions as in Section 3.1: in both cases, the model uses tokens which leads to a lower sensitivity
 1011 function.

1012 Continuing from the synthetic data, we now consider a slightly more complicated dataset, namely
 1013 MNIST (LeCun & Cortes, 2005). The MNIST dataset consists of $70k$ black-and-white images of
 1014 handwritten digits of resolution 28×28 . There are $60k$ images in the training set and $10k$ images in
 1015 the test set. It is released under the CC BY-SA 3.0 license. We compare the sensitivity of a ViT-small
 1016 model with an MLP on a binary digit classification task (< 5 or ≥ 5). In our experiments, each image
 1017 is divided into $T = 16$ patches of size 7×7 for the ViT-small model. For the MLP, the inputs are
 1018 vectorized as usual. With this setting, we measure the sensitivity of the two models using patch token
 1019 replacement as per Definition 3.1. As shown in Figure 10, when achieving the same training accuracy,
 1020 the ViT shows lower sensitivity compared to the MLP.

1021 A.2 SENSITIVITY WITH RANDOM NOISE INSTEAD OF TOKEN-WISE NOISE

1022 In this section, we consider changing the way we compute sensitivity, to see if the resulting metric
 1023 also distinguishes transformers from other architectures. Instead of token-wise perturbations, we add
 1024
 1025

Table 3: Comparison of sensitivity values for models that use only sparse or frequent tokens for the settings considered in Fig. 9.

Data Setting (n_f, m)	Top row in Fig. 3 ($n_s = 3, n_d = 1$)			Bottom row in Fig. 3 ($n_s = 1, n_d = 7$)		
	(3, 6)	(5, 16)	(7, 28)	(7, 10)	(17, 20)	(32, 36)
Using sparse tokens	0	0	0	0.0339	0.0339	0.0339
Using frequent tokens	0.1315	0.2878	0.4502	0	0	0

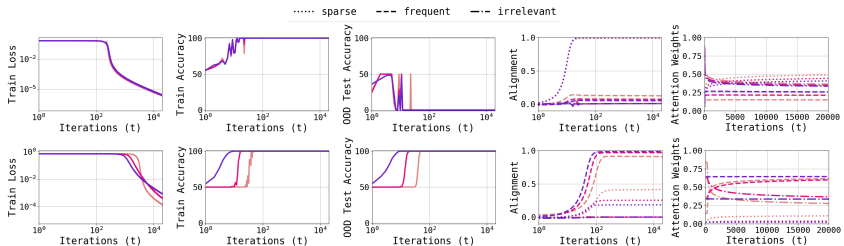


Figure 9: Train and test dynamics for a single-layer self-attention model (Eq. (3)) using the synthetic data visualized in Fig. 2; see Section 3.1 for details. The top row corresponds to the cases where the predictor that uses sparse tokens has lower sensitivity, while the bottom row corresponds to the cases where using the frequent tokens leads to lower sensitivity. The precise data settings for this figure, as well as a comparison of sensitivity values, are shown in Table 3.

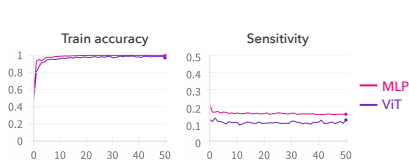


Figure 10: **Sensitivity on MNIST.** ViT and MLP get similar accuracy, but the ViT has lower sensitivity.

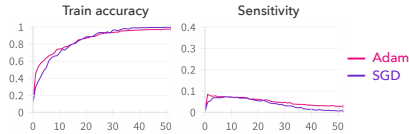


Figure 11: **Sensitivity using SGD and Adam.** Comparison of train accuracies and sensitivity values of the ViT-small model trained on the CIFAR-10 dataset using SGD and Adam optimizers.

Gaussian noise across the entire input with a smaller variance so that the transformer’s sensitivity in this case is similar to the sensitivity with token-wise noise.

Table 4: Comparison of sensitivity values measured with random and patch-wise noise for various model-dataset settings. Token-wise perturbations lead to a larger gap between the sensitivity of transformer-based models compared to other architectures.

Model and dataset	Random noise	Token-wise noise
ResNet-18 on CIFAR-10	0.0172	0.0827
ViT-small on CIFAR-10	0.0082	0.0050
LSTM on QQP	0.11	0.09
RoBERTa on QQP	0.05	0.03

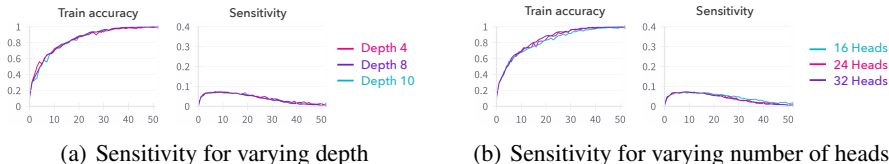
In Table 4, we compare the sensitivity values at the end of training for ResNet-18 and ViT-small on the CIFAR-10 dataset (variance 0.025) and LSTM and RoBERTa on the QQP dataset (variance 0.5). We find that for random perturbations, the difference between sensitivity values is much smaller for CIFAR-10 and similar for the QQP dataset, compared to patch-wise perturbations. These results suggest that measuring sensitivity with patch-wise noise is indeed the metric that we should consider since it distinguishes transformers from other architectures with a larger gap.

A.3 VISION TASKS

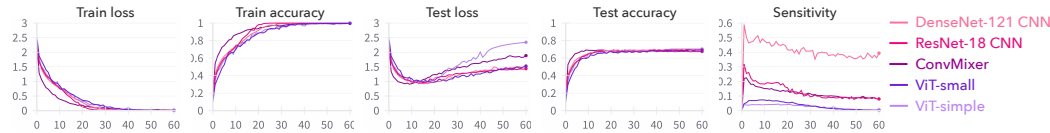
Effect of Depth, Number of Heads and the Optimization Algorithm. In Fig. 11, we compare the sensitivity values of a ViT-small model trained on CIFAR-10 dataset with SGD and Adam

1080 optimization algorithms. Although the model trained with Adam has a slightly higher sensitivity, the
 1081 sensitivity values for both the models are quite similar. This indicates that the low-sensitivity bias is
 1082 quite robust to the choice of the optimization algorithm.

1083 In Fig. 12, we compare the sensitivity values of a ViT-small model with different depth and number
 1084 of attention heads, when trained on the CIFAR-10 dataset. Note that for our main results, we use
 1085 a model with depth 8 and 32 heads. We observe that the train accuracies and the sensitivity values
 1086 remain the same across the different model settings. This indicates that the low-sensitivity bias is
 1087 quite robust to the model setting.

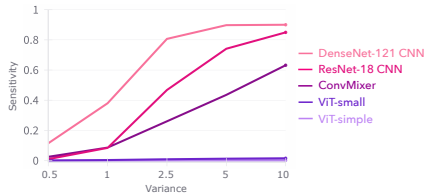


1094 (a) Sensitivity for varying depth (b) Sensitivity for varying number of heads
 1095
 1096 **Figure 12: Sensitivity for Various Model Settings.** Comparison of train accuracies and sensitivity
 1097 values on the CIFAR-10 dataset when varying the depth and number of heads of the ViT-small model.
 1098 We observe that for the same train accuracy, the sensitivity values remain very similar for different
 1099 model settings.



1102
 1103
 1104
 1105
 1106 **Figure 13: Comparison of the sensitivity of two CNNs, two ViTs, and ConvMixer trained on the**
 1107 **CIFAR-10 dataset, as a function of training epochs. For a fair comparison, the figure also shows the**
 1108 **train and test accuracies and loss values (cross-entropy loss). All models have similar accuracies but**
 1109 **the ViTs have significantly lower sensitivity than the other models.**

1111 **Effect of Variance.** In Fig. 14, we compare
 1112 the effect of the variance σ^2 used while evalu-
 1113 ating sensitivity for different models trained
 1114 on the CIFAR-10 dataset. We observe that the
 1115 ViTs have significantly lower sensitivity than the
 1116 other models and the difference becomes starker
 1117 as the variance level increases.



1118 **Results on SVHN Dataset.** Fig. 16 shows the
 1119 training accuracy and sensitivity of a ResNet-18
 1120 and a ViT-small model trained on SVHN dataset
 1121 (Netzer et al., 2011). At the end of training,
 1122 the sensitivity values are: 0.0516 for ResNet-18
 1123 and 0.0147 for ViT-small. Similar to the
 1124 observations for CIFAR-10, we see that the ViT
 1125 has a significantly lower sensitivity.

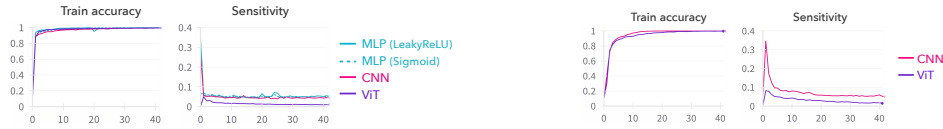
1126 **Figure 14: Sensitivity for Different Variances.**
 1127 Comparison of the sensitivity of two CNNs, two
 1128 ViTs, and ConvMixer trained on the CIFAR-10
 1129 dataset, as a function of different variance levels,
 1130 at the end of training. The ViTs have significantly
 1131 lower sensitivity at any variance and the difference
 1132 grows as variance increases.

1126 **Additional Results on CIFAR-10-C.** Fig. 17 and Fig. 18 show the test performance on various
 1127 corruptions from the CIFAR-10-C dataset with severity level 2 and 1, respectively. We observe that
 1128 CNNs have lower test accuracies on corrupted images compared to ViTs. Further, encouraging lower
 1129 sensitivity in the ViT leads to better robustness.

1130 A.4 LANGUAGE TASKS

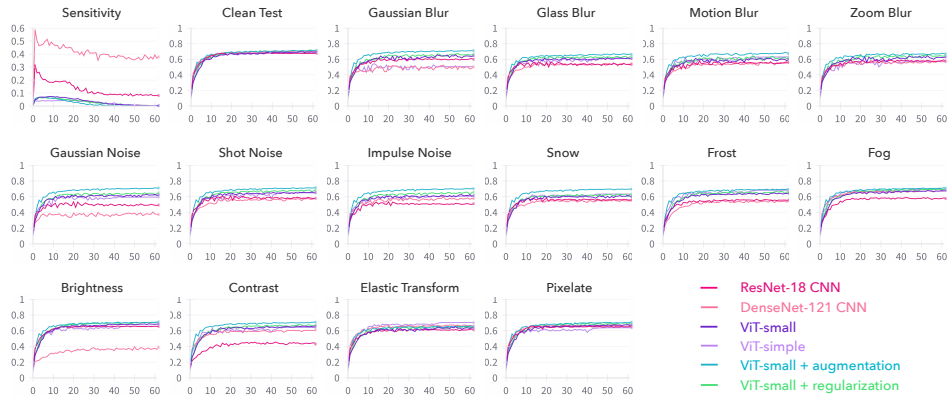
1131 **Sensitivity Measured with Variance $\sigma^2 = 4$.** Alternative to the main experiments with $\sigma^2 = 15$, we
 1132 also evaluate sensitivity with a different corruption strength $\sigma^2 = 4$ on the QQP dataset, as shown in

1134
1135
1136
1137
1138



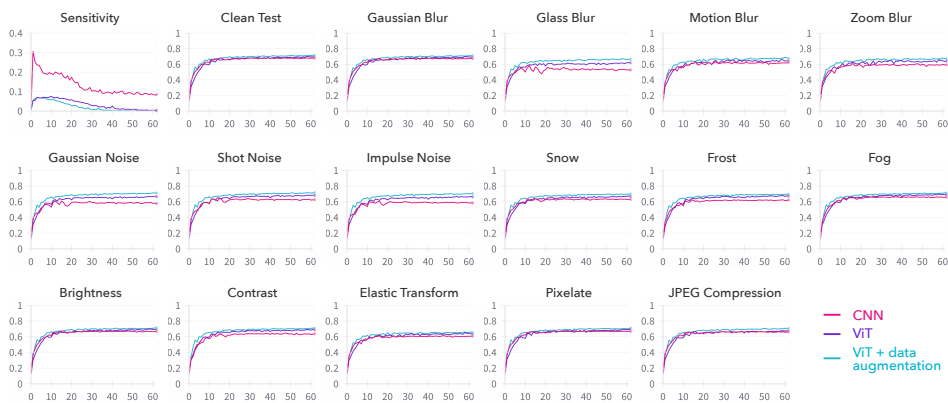
1139 **Figure 15: Sensitivity on Fashion-MNIST.** Comparison of sensitivity of a ViT with a CNN, an MLP with
1140 LeakyReLU activation, and an MLP with sigmoid activation, as a function of training epochs. All the models
1141 have similar accuracies but the ViT has significantly lower sensitivity.
1142 **Figure 16: Sensitivity on SVHN.** Comparison of sensitivity of a ResNet-18 CNN and
1143 a ViT-small trained on SVHN dataset, as a function of training epochs. Both the models
1144 have similar accuracies but the ViT has significantly lower sensitivity.

1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158



1159 **Figure 17: Comparison of the sensitivity, test accuracy on CIFAR-10 and test accuracies on various**
1160 **corruptions from the CIFAR-10-C dataset (see Section 6 for details) of two CNNs and two ViTs**
1161 **trained on the CIFAR-10 dataset, as a function of the training epochs. We also compare with ViT-**
1162 **small trained with data augmentation/regularization, which encourage low sensitivity (see Section 6**
1163 **for discussion).**

1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177



1178 **Figure 18: Comparison of the sensitivity, test accuracy on CIFAR-10, and test accuracies on various**
1179 **corruptions from the CIFAR-10-C dataset (see Section 6 for details) of a ResNet-18 CNN and a**
1180 **ViT-small model trained on the CIFAR-10 dataset, as a function of the training epochs. We also**
1181 **compare with ViT-small trained with data augmentation, which acts as a regularizer to encourage**
1182 **low sensitivity (see Section 6 for discussion). Here, we use severity level 1, while in Figure 17, we**
1183 **considered severity level 2.**

1184
1185
1186
1187

Fig. 19. We observe the same trends as in Figures 5 and 6: RoBERTa models have lower sensitivity
than the LSTM and the LSTM is more sensitive to more recent tokens. These results indicate that the
low-sensitivity bias is robust to the choice of corruption strength.

1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241

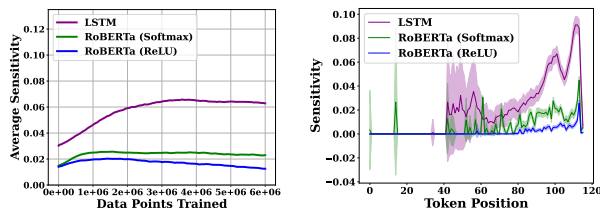


Figure 19: **Sensitivity on the QQP Dataset with Variance $\sigma^2 = 4$.** Results with alternative variance yield observations that are consistent with the setup in the main text. (Left) LSTM has higher sensitivity than the RoBERTa models. (Right) Softmax activation for RoBERTa induces higher sensitivity towards the CLS token.

Sensitivity Measured on the Validation Set. We also test sensitivity on the validation set for the two language datasets. As seen in Figure 20, the results on the validation set are consistent with those on the train set in Fig. 5.

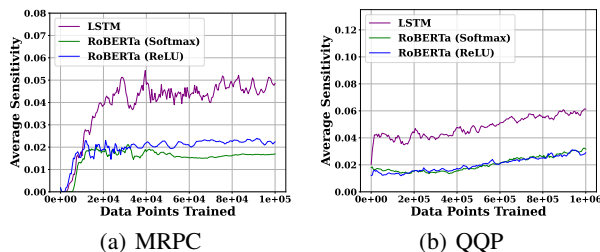


Figure 20: **Sensitivity on the Validation Sets.** Similar to the observation in Figure 5, the RoBERTa models have lower sensitivity than the LSTM for both the datasets. However, the difference between RoBERTa-ReLU and RoBERTa-softmax is less marginal on the validation set compared to the training set.

Sensitivity Measured with the GPT-2 Model. Here, we ablate the effect of the transformer architecture on the sensitivity values. We compare a GPT-based model with the two BERT-based models used in our main experiments. The key difference lies in the construction of the attention masks: for GPT models, each token only observes the tokens that appear before it, whereas BERT models are bidirectional, therefore each token observes all the tokens in the sequence. In Fig. 21, we observe that the GPT-2 model has higher sensitivity compared to the RoBERTa models, but the sensitivity is significantly lower than the LSTM. The GPT-2 model is also relatively more sensitive to more recent tokens compared to the RoBERTa models, while also being sensitive to some CLS tokens.

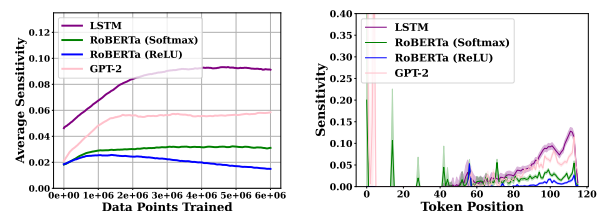


Figure 21: **Sensitivity of GPT-2 on the QQP Dataset.** (Left) We find that the RoBERTa models tend to have lower sensitivity than GPT-2, and all Transformer models have lower sensitivity than LSTM. (Right) The sensitivity per token of GPT-2 is more similar to LSTMs, which is possibly due to their shared auto-regressive design.

1242
 1243
 1244
 1245
 1246
 1247
 1248
 1249
 1250
 1251
 1252
 1253
 1254
 1255
 1256
 1257
 1258
 1259
 1260
 1261
 1262
 1263
 1264
 1265
 1266
 1267
 1268
 1269
 1270
 1271
 1272
 1273
 1274
 1275
 1276
 1277
 1278
 1279
 1280
 1281
 1282
 1283
 1284
 1285
 1286
 1287
 1288
 1289
 1290
 1291
 1292
 1293
 1294
 1295

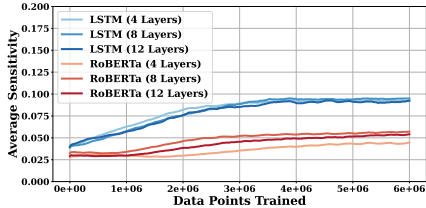


Figure 22: **Sensitivity for Different Model Depths.** We vary the model depths of LSTM and RoBERTa on the QQP datasets and observe that LSTM models tend to have the same sensitivity throughout the entire training. RoBERTa model with 4 layers has slightly lower sensitivity with its 8-layer or 12-layer variants. RoBERTa, regardless of depths, have lower sensitivity than LSTMs.

A.5 SENSITIVITY AS A PROGRESS MEASURE FOR GROKING

Grokking is a phenomenon in which a neural network suddenly and drastically improves its generalization ability after a long period of training, even though it initially overfits to the training data (Nanda et al., 2023). During grokking, the model transitions from memorizing the training data to learning a more general solution, allowing it to perform well on unseen data. This often happens after many training steps, during which time the test accuracy remains low despite perfect training accuracy. The transition is abrupt, making grokking seem like an emergent behavior where the model, after much training, "figures out" the correct approach to the task.

In more technical terms, this shift occurs as the network amplifies structured mechanisms that enable generalization and removes components that only lead to memorization. Grokking has been observed in models trained with regularization techniques like weight decay on algorithmic tasks, where the model learns an underlying structure that generalizes well beyond the training data. In (Nanda et al., 2023), the authors demonstrate that for the modular addition task, the model initially memorizes the training data, leading to low test accuracy. Later, the model discovers how to use trigonometric functions to solve the task. However, the emergence of grokking is difficult to measure in practice, and it is not guaranteed that grokking will occur or that the model will find the correct approach to the task. Nevertheless, as shown in Figure 23, we can clearly observe a significant change in sensitivity between epochs 3000 and 9800, during the saturation of the training loss. This suggests that the model is discovering a more robust solution to the task. We propose that sensitivity can serve as a useful metric for assessing whether the model is learning to grok.

Similar to Nanda et al. (2023)’s classification on phases: memorization, circuit formation, and cleanup, we claim that sensitivity also provides a progress measure, with an extra phase of noise reduction:

Memorization: From epoch 0 to 500, sensitivity drops significantly while training accuracy saturates to 100 and test accuracy increases but cannot saturate.

Circuit Formation: From epoch 500 to 3,000, the sensitivity goes up again but test accuracy remains low and flat. The dramatic fall in the weight norms suggests that circuit formation likely happens due to weight decay.

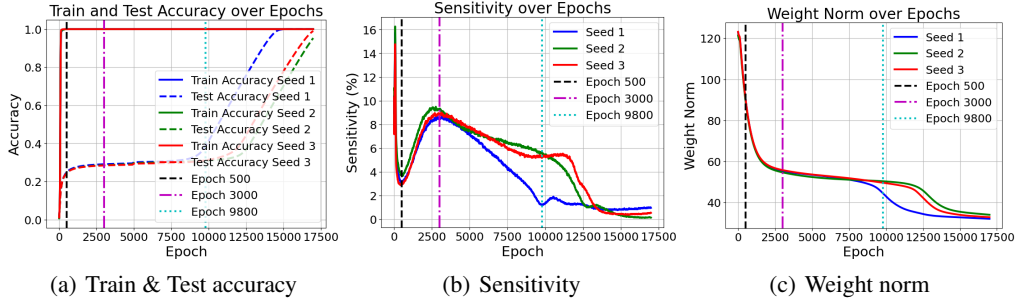
Clean up: From epoch 3,000 to 9,800, sensitivity starts to decrease, and the progress measure indicates that the model starts to learn to use Fourier features.

Noise Reduction: From epoch 9,800 onwards, sensitivity has an upward and then downward trend, and this is not caught by other mechanistic interpretability measures in Nanda et al. (2023). This trend in sensitivity is due to further Fourier noise reduction, where initially, slightly perturbations to the number embeddings could drift the model performance where later on, model learns a more robust Fourier basis.

A.6 SENSITIVITY AS A PROGRESS MEASURE FOR LEARNING SPARSE PARITIES

In this section, we investigate if sensitivity also acts as a progress measure when training transformers on the sparse parity task, with input $\mathbf{x} \in \{\pm 1\}^d$ and label $y = \prod_{i \in S} x_i$, where $S \subset [d]$ with sparsity level $p := |S| < d$.

1296
1297
1298
1299
1300
1301
1302
1303
1304



1305
1306
1307
1308
1309

Figure 23: Sensitivity measures progress on modulo addition task $a + b \pmod{113}$ and indicates different stages of grokking.

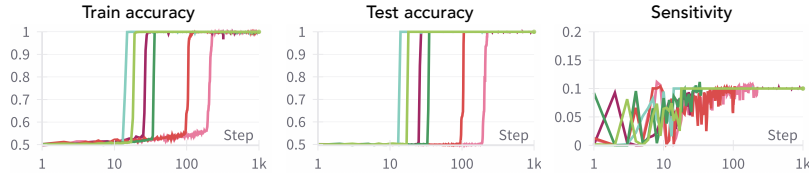
1310
1311
1312
1313
1314
1315
1316
1317

We train a two-layer transformer model on this task with 4 heads using Adam optimizer for 1000 epochs. We evaluate the sensitivity metric based on Eq. (1), using 10^5 samples to estimate the expectation over the Boolean cube. Fig. 24 shows the results on six different settings in terms of (number of train samples, batch size, embedding dimension, learning rate, seed):

- Sea green: (50k, 250, 32, 0.001, 123)
- Parrot green: (50k, 250, 32, 0.001, 42)
- Dark green: (50k, 250, 32, 0.001, 0)
- Maroon: (10k, 100, 64, 0.001, 42)
- Coral: (5k, 25, 32, 0.0001, 42)
- Pink: (5k, 50, 32, 0.0001, 42).

1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328

We observe that in all the settings, the train and test accuracy remain close to 50% for the initial few epochs, while sensitivity increases, which indicates progress in learning the sparse parity task. When the model learns the sparse parity function, the sensitivity converges to $\frac{p}{d} = 0.1$, which coincides with train and test accuracy going to 100%.



1329
1330
1331
1332
1333
1334

Figure 24: Train accuracy, test accuracy, and sensitivity as a function of training epochs, when training a two-layer transformer on the sparse parity task with dimension 40 and sparsity level 4. Sensitivity increases while train and test accuracy are close to 50%, and converges to 0.1 when the model learns the sparse parity function, which coincides with train and test accuracy going to 100%.

B PROOFS FOR SECTION 2

1335
1336
1337

We give a brief overview of the CK and NTK here and refer the reader to (Lee et al., 2018; Yang & Salman, 2020) for more details.

1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349

Consider a model with L layers and widths $\{d_l\}_{l=1}^L$ and an input \mathbf{x} . Let $g^l(\mathbf{x})$ denote the output of the l^{th} layer scaled by $d_l^{-1/2}$. Suppose we randomly initialize weights from the Gaussian distribution $\mathcal{N}(0, 1)$. It can be shown that in the infinite width limit when $\min_{l \in [L]} d^l \rightarrow \infty$, each element of $g^l(\mathbf{x})$ is a Gaussian process (GP) with zero mean and kernel function K^l . The kernel K^L corresponding to the last layer of the model is the CK. In other words, it is the kernel induced by the embedding $\mathbf{x} \mapsto g^{L-1}(\mathbf{x})$ when the model is initialized randomly. On the other hand, NTK corresponds to training the entire model instead of just the last layer. Intuitively, when the model parameters θ stay close to initialization θ_0 , the residual $g^L(\mathbf{x}; \theta) - g^L(\mathbf{x}; \theta_0)$ behaves like a linear model with features given by the gradient at random initialization, $\nabla_{\theta} g^L(\mathbf{x}, \theta_0)$, and the NTK is the kernel of this linear model. The spectra of these kernels provide insights about the implicit prior of a randomly initialized

model as well as the implicit bias of training using gradient descent (Yang & Salman, 2020). The closer these spectra are to the spectrum of the target function, the better we can expect training using gradient descent to generalize.

B.1 PROOF OF PROPOSITION 2.1

(Hron et al., 2020) show that the self-attention layer with linear attention and $d^{-1/2}$ scaling converges in distribution to $\mathcal{GP}(0, K)$ in the infinite width limit, *i.e.* when the number of heads d^H become large. For any layer $l \in [L]$, let \tilde{K}^l denote the kernel induced by the intermediate transformation when applying some nonlinearity ϕ to the output of the previous layer $l - 1$. Let $f_{i,j}^l := \{f_{i,j}^l(\mathbf{x}) : \mathbf{x} \in \mathcal{X}, i \in [T]\}$, where \mathcal{X} denotes the input space of \mathbf{x} . They show the following result for NNs with at least one linear attention layer, in the infinite width limit.

Theorem B.1 (Theorem 3 in (Hron et al., 2020)). *Let $l \in [L]$, and ϕ be such that $|\phi(x)| \leq c + m|x|$ for some $c, m \in \mathbb{R}^+$. Assume g^{l-1} converges in distribution to $g^{l-1} \sim GP(0, K^{l-1})$, such that $g_{i,j}^{l-1}$ and $g_{i,k}^{l-1}$ are independent for any $j \neq k$. Then as $\min\{d^{l,H}, d^l\} \rightarrow \infty$, g^l converges in distribution to $g^l \sim GP(0, K^l)$ with $g_{i,k}^l$ and $g_{i,\ell}^l$ independent for any $k \neq \ell$, and*

$$K^l(\mathbf{x}, \mathbf{x}') = \mathbb{E}[g^l(\mathbf{x})g^l(\mathbf{x}')] = \sum_{i,j=1}^{\bar{d}} (\tilde{K}_{ij}^l(\mathbf{x}, \mathbf{x}'))^2 \tilde{K}_{ab}^l(\mathbf{x}, \mathbf{x}').$$

Similar results are also known for several non-linearities and other layers such as convolutional, dense, average pooling (Lee et al., 2018; de G. Matthews et al., 2018; Garriga-Alonso et al., 2019; Novak et al., 2019; Yang, 2021), as well as residual, positional encoding and layer normalization (Hron et al., 2020; Yang, 2021).

Consequently, any model composed of these layers, such as a transformer with linear attention, also converges to a Gaussian process. This follows using an induction-based argument. It can easily be shown that the induced kernel takes the form

$$K(\mathbf{x}, \mathbf{y}) = \Psi \left(\frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\|\mathbf{x}\| \|\mathbf{y}\|}, \frac{\|\mathbf{x}\|^2}{d}, \frac{\|\mathbf{y}\|^2}{d} \right),$$

for some function $\Psi : \mathbb{R}^3 \rightarrow \mathbb{R}$. In addition, since $\mathbf{x}, \mathbf{y} \in \mathbb{B}^d$, they have the same norm, and Ψ can be treated as a univariate function that only depends on $c = d^{-1} \langle \mathbf{x}, \mathbf{y} \rangle$, *i.e.* $\Psi(c, 1, 1) = \Psi(c)$.

Using this property and the following result, it follows that the kernel induced by a transformer with linear attention is diagonalized by the Fourier basis $\{\chi_U\}_{U \subseteq [d]}$.

Theorem B.2 (Theorem 3.2 in (Yang & Salman, 2020)). *On the d -dimensional boolean cube \mathbb{B}^d , for every $U \subseteq [d]$, χ_U is an eigenfunction of K with eigenvalue*

$$\mu_{|U|} := \mathbb{E}_{\mathbf{x} \sim \mathbb{B}^d} [x^U K(\mathbf{x}, \mathbf{1})] = \mathbb{E}_{\mathbf{x} \sim \mathbb{B}^d} \left[x^U \Psi \left(d^{-1} \sum_i x_i \right) \right],$$

where $\mathbf{1} := (1, \dots, 1) \in \mathbb{B}^d$. This definition of $\mu_{|U|}$ does not depend on the choice S , only on the cardinality of S . These are all of the eigenfunctions of K by dimensionality considerations.

Further, using the following result, it follows that transformers (with linear attention) exhibit weak spectral simplicity bias.

Theorem B.3 (Theorem 4.1 in (Yang & Salman, 2020)). *Let K be the CK or NTK of an MLP on a boolean cube \mathbb{B}^d . Then the eigenvalues $\mu_k, k = 0, \dots, d$, satisfy*

$$\begin{aligned} \mu_0 &\geq \mu_2 \geq \dots \geq \mu_{2k} \geq \dots, \\ \mu_1 &\geq \mu_3 \geq \dots \geq \mu_{2k+1} \geq \dots \end{aligned}$$

B.2 PROOF OF PROPOSITION 2.2

First, we introduce the concept of noise stability $Q_\rho(f)$, which measures the correlation between the outputs of a function f for ρ -correlated pair $(\mathbf{x}, \mathbf{x}')$, as $Q_\rho(f) := \mathbb{E}_{(\mathbf{x}, \mathbf{x}')_\rho} f(\mathbf{x})f(\mathbf{x}')$. Note that $Q_\rho(f)$

is related to $R_\rho(f) := \Pr_{(\mathbf{x}, \mathbf{x}')_\rho} [f(\mathbf{x}) \neq f(\mathbf{x}')] as $Q_\rho(f) = 1 - 2R_\rho(f)$ (O'Donnell, 2014).$

Using the following result, we can relate noise stability to the Fourier weight of the function f at different degrees $i \in [d]$.

Theorem B.4 (Theorem 2.49 in (O’Donnell, 2014)). *For function $f : \mathbb{F}^d \rightarrow \{\pm 1\}$, the noise stability for ρ -correlated pair (x, x') satisfies $Q_\rho(f) = \sum_{U \subseteq [d]} \rho^{|U|} \hat{f}(U)^2 = \sum_{i=0}^d \rho^i W^i[f]$, where*

$$W^i[f] := \sum_{U \subseteq [d], |U|=i} \hat{f}(U)^2.$$

Clearly, $Q_\rho(f) \leq 1$ since the minimum degree of f is 0. Next, we use the following important result, which upper bounds the degree of f in terms of its maximum sensitivity, $S_{\max}(f) := \max_{x \in \mathbb{F}^d} S(f, x)$.

Theorem B.5 (Theorem 1.4 in (Huang, 2019)). *For function $f : \mathbb{F}^d \rightarrow \{\pm 1\}$, the degree $D(f)$ of the multilinear polynomial which represents f satisfies $D(f) \leq (S_{\max}(f))^2$.*

Using this, and the fact that $Q_\rho(f)$ is minimized when the Fourier weight is concentrated on the highest degree term, we get the lower bound $Q_\rho(f) \geq \rho^{(S_{\max}(f))^2}$, since $\rho < 1$ and $S_{\max}(f) \in [1, d]$.

Using the relation between $Q_\rho(f)$ and $R_\rho(f)$ then finishes the proof.

C RELATED WORK

Understanding Transformers. The emergence of transformers as the go-to architecture for many tasks has inspired extensive work on understanding the internal mechanisms of transformers, including reverse-engineering language models (Wang et al., 2022), the grokking phenomenon (Power et al., 2022; Nanda et al., 2023), manipulating attention maps (Hassid et al., 2022; Kobayashi et al., 2024), automated circuit finding (Conmy et al., 2023), arithmetic computations (Hanna et al., 2023; Quirke & Barez, 2024), optimal token selection (Tarzanagh et al., 2023a;b; Vasudeva et al., 2024), and in-context learning (Brown et al., 2020; Garg et al., 2022; Akyürek et al., 2023; von Oswald et al., 2022; Fu et al., 2023; Bhattamishra et al., 2023a; Guo et al., 2024). Several works investigate why vision transformers (ViTs) outperform CNNs (Trockman & Kolter, 2022a; Raghu et al., 2021; Melas-Kyriazi, 2021), as well as other properties of ViTs, such as robustness to (adversarial) perturbations and distribution shifts (Bai et al., 2023; Shao et al., 2021; Mahmood et al., 2021; Bhojanapalli et al., 2021; Naseer et al., 2021; Paul & Chen, 2022; Ghosal et al., 2022). Further, several works on mechanistic interpretability of transformers share a similar recipe of measuring sensitivity — corruption with Gaussian noise (Meng et al., 2022; Conmy et al., 2023) but on hidden states rather than the input space.

Sensitivity and Spectral Bias. Sensitivity is closely related to spectral bias (Yang & Salman, 2020), which is a bias towards ‘simple’ functions in the Fourier space. Simple functions in the Fourier space generally correspond to low-frequency terms when the input space is continuous, and low-degree polynomials when the input space is discrete. Recent work has shown that deep networks prefer to use low-frequency Fourier functions on images (Xu et al., 2019), and low-degree Fourier terms on Boolean functions (Yang & Salman, 2020). We note that in contrast to some other notions of spectral bias, sensitivity also has the advantage that it can be efficiently estimated on data through sampling — in contrast, estimating all the Fourier coefficients requires time exponential in the dimensionality of the data and hence can be computationally prohibitive (Xu et al., 2019).

Simplicity Bias in DL. Several works (Neyshabur et al., 2014; Valle-Perez et al., 2019; Arpit et al., 2017; Geirhos et al., 2020) show that NNs prefer learning ‘simple’ functions over the data. Nakkiran et al. (2019) show that during the early stages of SGD training, the predictions of NNs can be approximated well by linear models. Morwani et al. (2023) show that 1-hidden-layer NNs exhibit simplicity bias to rely on low-dimensional projections of the data, while (Huh et al., 2021) empirically show that deep NNs find solutions with lower rank embeddings. (Shah et al., 2020) create synthetic datasets where features that can be separated by predictors with fewer piece-wise linear components are considered simpler, and show that in the presence of simple and complex features with equal predictive power, NNs rely heavily on simple features. Geirhos et al. (2019) show that trained CNNs rely more on image textures rather than image shapes to make predictions. Rahaman et al. (2019a)

1458 use Fourier analysis tools and show that deep networks are biased towards learning low-frequency
1459 functions, and (Xu et al., 2019; Cao et al., 2021; Bietti & Mairal, 2019; Basri et al., 2019) provide
1460 further theoretical and empirical evidence for this.

1461 **Implicit Biases of Gradient Methods.** Several works study the implicit bias of gradient-based
1462 methods for linear predictors and MLPs. Pioneering work by Soudry et al. (2018); Ji & Telgarsky
1463 (2018) revealed that linear models trained with gradient descent to minimize an exponentially-tailed
1464 loss on linearly separable data converge (in direction) to the max-margin classifier. Following
1465 this, Nacson et al. (2019); Ji & Telgarsky (2021); Ji et al. (2021) derived fast convergence rates
1466 for gradient-based methods in this setting. Recent works show that MLPs trained with gradient
1467 flow/descent converge to a KKT point of the corresponding max-margin problem in the parameter
1468 space, in both finite (Ji & Telgarsky, 2020; Lyu & Li, 2020) and infinite width (Chizat & Bach, 2020)
1469 regimes. Further, Phuong & Lampert (2021); Frei et al. (2022); Kou et al. (2023) have also studied
1470 ReLU/Leaky-ReLU networks trained with gradient descent on nearly orthogonal data. Li et al. (2022)
1471 show that the training path in over-parameterized models can be interpreted as mirror descent applied
1472 to an alternative objective. In regression problems, when minimizing the mean squared error, the
1473 bias manifests in the form of rank minimization (Arora et al., 2019; Li et al., 2021). Additionally,
1474 the implicit bias of other optimization algorithms, such as stochastic gradient descent and adaptive
1475 methods, has also been explored in various studies (Blanc et al., 2020; HaoChen et al., 2021); see the
1476 recent survey (Vardi, 2022) for a detailed summary.

1477 **Robustness.** Several research efforts have been made to investigate the robustness of Transformers.
1478 Shao et al. (2021) showed that Transformers exhibit greater resistance to adversarial attacks compared
1479 to other models. Additionally, Mahmood et al. (2021) highlighted the notably low transferability of
1480 adversarial examples between CNNs and ViTs. Subsequent research (Shen et al., 2023; Bhojanapalli
1481 et al., 2021; Paul & Chen, 2022) expanded this robustness examination to improve transformer-based
1482 language models. Shi et al. (2020) introduced the concept of robustness verification in Transformers.
1483 Various robust training methods have been suggested to enhance the robustness guarantees of models,
1484 often influenced by or stemming from their respective verification techniques. Shi et al. (2021)
1485 expedited the certified robust training process through the use of interval-bound propagation. Wang
1486 et al. (2021) employed randomized smoothing to train BERT, aiming to maximize its certified robust
1487 space. Recent work of Bombari & Mondelli (2024) shows that randomly-initialized attention layers
1488 tend to have higher word-level sensitivity than fully connected layers. In contrast to our work, they
1489 consider word sensitivity, which has been experimentally shown to be similar for transformers and
LSTMs (Bhattamishra et al., 2023b).

1490 **Spurious Correlations.** A common pitfall to the generalization of neural networks is the presence of
1491 spurious correlations (Sagawa et al., 2020). For example, Geirhos et al. (2019) observed that trained
1492 CNNs are biased towards textures rather than shapes to make predictions for object recognition
1493 tasks. Such biases make NNs vulnerable to adversarial attacks. Gururangan et al. (2018) attribute
1494 the reliance of NNs on spurious features to confounding factors in data collection while Shah et al.
1495 (2020) attribute it to a *simplicity bias*. Several works have studied the underlying causes of simplicity
1496 bias (Chiang, 2021; Nagarajan et al., 2021; Morwani et al., 2023; Huh et al., 2021; Lyu et al., 2021)
1497 and multiple methods have been developed to mitigate this bias and improve generalization (Pezeshki
1498 et al., 2020; Kirichenko et al., 2022; Vasudeva et al., 2023; Tiwari & Shenoy, 2023).

1499 **Data Augmentation.** The essence of data augmentation is to impose some notion of regularization.
1500 The simplest design of data augmentation dates back to Robbins (1951) where image manipulation,
1501 e.g., flip, crop, and rotate, was introduced. Bishop (1995) proved that training with Gaussian noise is
1502 equivalent to Tikhonov regularization. We also note this observation is in parallel to our proposition
1503 in Section 6 that training with Gaussian noise promotes low sensitivity. Recently, mixup-based
1504 augmentation methods have been proposed to improve model robustness by merging two images
1505 as well as their labels (Zhang et al., 2018). Several works also use a combination of existing
1506 augmentation techniques (Cubuk et al., 2019; Lim et al., 2019). A common belief is that data
1507 augmentation can improve model robustness (Rebuffi et al., 2021), and this work bridges the method
1508 (augmentation) and the outcome (robustness) with an explanation — simplicity bias towards low
1509 sensitivity.

1512 D DETAILS OF EXPERIMENTAL SETTINGS

1513 We use PyTorch (Paszke et al., 2019) as our code framework and as our implementation of LSTMs.
 1514 PyTorch is licensed under the Modified BSD license.

1515 **Experimental Settings for Synthetic Data Experiments.** We use standard SGD training with batch
 1516 size 100. We consider $T = 50$ and train with 1000 samples and test on 500 samples generated as per
 1517 Definition 3.3. **Datasets, Model Architectures and Experimental Settings for Vision Tasks.** We
 1518 consider the following datasets:

1519 *Fashion-MNIST.* Fashion-MNIST (Xiao et al., 2017) consists of 28×28 grayscale images of
 1520 Zalando’s articles. This is a 10-class classification task with $60k$ training and $10k$ test images. It is
 1521 released under the MIT license.

1522 *CIFAR-10.* The CIFAR-10 dataset (Krizhevsky, 2009) is a well-known object recognition dataset. It
 1523 consists of 32×32 color images in 10 classes, with $6k$ images per class. There are $50k$ training and
 1524 $10k$ test images. It is released under the MIT license.

1525 *SVHN.* Street View House Numbers (SVHN) (Netzer et al., 2011) is a real-world image dataset used
 1526 as a digit classification benchmark. It contains 32×32 RGB images of printed digits (0 to 9) cropped
 1527 from Google Street View images of house number plates. There are $60k$ images in the train set and
 1528 $10k$ images in the test set. It is released under the CC BY 4.0 license.

1529 *ImageNet-1k.* The ImageNet-1k dataset (Russakovsky et al., 2015), also known as the ILSVRC
 1530 (ImageNet Large Scale Visual Recognition Challenge) dataset, is a widely used benchmark dataset in
 1531 computer vision for tasks such as image classification, object detection, and image segmentation.
 1532 There are 1000 different classes, and approximately 1.28 million training images of size 224×224 .
 1533 It is released under a non-commercial research use license.

1534 For all the datasets, we use the ViT-small architecture implementation available at <https://github.com/lucidrains/vit-pytorch>.
 1535 For the ResNet-18 model used in the experiments on CIFAR-10
 1536 and SVHN datasets, we use the implementation available at <https://github.com/kuangliu/pytorch-cifar>.
 1537 Additionally, for the DenseNet-121 model, ConvMixer model and ViT-simple
 1538 model used in the experiments on CIFAR-10, we use the implementations available at https://github.com/huyvnphan/PyTorch_CIFAR10,
 1539 <https://github.com/locuslab/convmixer>
 1540 and <https://github.com/lucidrains/vit-pytorch>, respectively. All of these models are
 1541 released under the MIT license.

1542 All the models are trained with SGD using batch size 50 for MNIST and 100 for the other datasets.
 1543 We use patch size 7 for MNIST and 4 for the other datasets. We estimate the expectation over \mathcal{P} in
 1544 Definition 3.1 by replacing every patch with a noisy patch 5 times, and sample about 30% of the
 1545 training data to evaluate sensitivity.

1546 For the MNIST experiments, we consider a 1-hidden-layer MLP with 100 hidden units and
 1547 LeakyReLU activation. We set depth as 2, number of heads as 1 and the hidden units in the
 1548 MLP as 128 for the ViT. We train both models with a learning rate of 0.01.

1549 For Fashion-MNIST, we set depth as 2, number of heads as 8 and the hidden units in the MLP as
 1550 256 for the ViT. We consider a 2-hidden layer MLP with 512 and 128 hidden units, respectively. The
 1551 CNN consists of two 2D convolutional layers with 32 output channels and kernel size 3 followed by
 1552 a 2D MaxPool layer with both kernel size and stride as 2 and two fully connected layers with 128
 1553 hidden units. We use LeakyReLU activation for the CNN. We use learning rates of 0.1 for the MLP
 1554 with LeakyReLU, 0.5 for the MLP with sigmoid, 0.005 for the CNN and 0.1 for the ViT.

1555 For CIFAR-10, we set depth as 8, number of heads as 32 and the hidden units in the MLP as 256 for
 1556 the ViT-small model, while these values are set as 6, 16, and 512 for the ViT-simple model. For the
 1557 ConvMixer model, we use depth 6, embedding dimension 128 and kernel size 3. The learning rate is
 1558 set as 0.1 for ViT-small, 0.2 for ViT-simple, 0.06 for ConvMixer, 0.001 for ResNet-18 and 0.005 for
 1559 DenseNet-121.

For SVHN, most of the settings are the same as the CIFAR-10 experiments, except we set the hidden units in the MLP as 512 for the ViT-small model and the learning rate is set as 0.0015 for ResNet-18.

For ImageNet-1k, we sample $20k$ samples from the training set and compare the sensitivity values of pre-trained ConvNext (ConvNextV2-Tiny) and ViT/L-16 with $\sigma^2 = 15$. Both achieved the same training and validation accuracy of 85%, and it ensures our sensitivity comparison is fair.

Datasets, Model Architectures and Experimental Settings for Language Tasks. We consider the following two binary classification datasets, which are relatively easy to learn without pretraining (Kovaleva et al., 2019).

MRPC. Microsoft Research Paraphrase Corpus (MRPC) (Dolan & Brockett, 2005) is a corpus that consists of 5801 sentence pairs. Each pair is labeled if it is a paraphrase or not by human annotators. It has 4076 training examples and 1725 validation examples. It is released under the ODC-By or the Microsoft Research license.

QQP. Quora Question Pairs (QQP) (Iyer et al., 2017) dataset is a corpus that consists of over $400k$ question pairs. Each question pair is annotated with a binary value indicating whether the two questions are paraphrases of each other. It has $364k$ training examples and $40k$ validation examples. It is released under the CC BY-SA 2.5 license.

For both RoBERTa models and LSTM models, we keep the same number of layers: 4 layers. We set number of heads as 8 RoBERTa. We use the AdamW optimizer with a learning rate of 0.0001 and weight decay of 0.0001 for all the tasks. We also use a dropout rate of 0.1. We use a batch size of 32 for all the experiments. The used RoBERTa model is released on Huggingface <https://huggingface.co/FacebookAI/roberta-base> with MIT license.

Experimental Settings for Section 6. We set the learning rate as 0.16 and 0.2 when training the ViT-small with regularization and augmentation, respectively. We use a regularization strength of 0.25. The remaining settings are the same as for the other experiments. For computing the sharpness metrics, we approximate the expectation over the Gaussian noise by averaging over 5 repeats and set σ as 0.005.

Experimental Settings for Appendix A.3. For the experiment with the Adam optimizer, we employ a learning rate scheduler to ensure that the accuracy on the train set is similar to the model trained with SGD. The initial learning rate is 0.002 and after every 8 epochs, it is scaled by a factor of 0.5.

For the remaining experiments in this section, we consider the same settings as for the respective main experiment.

Compute Details. Experiments with synthetic data were run on Google Colab. Experiments on vision and language tasks were run on internal clusters using NVIDIA RTX A6000 GPUs with 48GB of VRAM. For the experiments on vision data, we use two GPUs and the runtime for each setting is about 17 hours. Experiments on language tasks use one GPU and the runtime for each experiment is about 24 hours.

E LIMITATIONS

In our theoretical results, we show that transformers exhibit weak spectral bias, similar to other NN architectures. An important direction for future work is to distinguish transformers from other architectures and show that they exhibit a stronger spectral bias.

Additionally, this work focuses on the inductive bias of the transformer architecture. However, other factors such as the data used while pre-training can also effect the biases these models exhibit on downstream tasks. It would be interesting to explore this effect in the future.

Similarly, the choice of the optimization algorithm used for training can also have an effect. In our experiments on the CIFAR-10 dataset (Fig. 11), we see that SGD and Adam are very similar. However, conducting a more thorough comparison, e.g., by considering second-order optimization methods, can be an important direction for future work.