
Separation and Bias of Deep Equilibrium Models on Expressivity and Learning Dynamics

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The deep equilibrium model (DEQ) generalizes the conventional feedforward
2 neural network by fixing the same weights for each layer block and extending
3 the number of layers to infinity. This novel model directly finds the fixed points
4 of such a forward process as features for prediction. Despite empirical evidence
5 showcasing its efficacy compared to feedforward neural networks, a theoretical
6 understanding for its separation and bias is still limited. In this paper, we take a
7 step by proposing some separations and studying the bias of DEQ in its expressive
8 power and learning dynamics. The results include: (1) A general separation is
9 proposed, showing the existence of a width- m DEQ that any fully connected neural
10 networks (FNNs) with depth $O(m^\alpha)$ for $\alpha \in (0, 1)$ cannot approximate unless
11 its width is sub-exponential in m ; (2) DEQ with polynomially bounded size and
12 magnitude can efficiently approximate certain steep functions (which has very large
13 derivatives) in L^∞ norm, whereas FNN with bounded depth and exponentially
14 bounded width cannot unless its weights magnitudes are exponentially large; (3)
15 The implicit regularization caused by gradient flow from a diagonal linear DEQ
16 is characterized, with specific examples showing the benefits brought by such
17 regularization. From the overall study, a high-level conjecture from our analysis
18 and empirical validations is that DEQ has potential advantages in learning certain
19 high-frequency components.

20 1 Introduction

21 Implicit deep learning [1], a paradigm that generalizes the recursive principles of traditional explicit
22 models, has gained renewed interest with the advent of novel neural network architectures. Among
23 these, deep equilibrium model (DEQ) [2] stands out as a commonly utilized model. In contrast to
24 explicit neural network that derives features through forward propagation, DEQ computes features
25 directly by solving an equilibrium equation induced by the implicit layer. Since the equilibrium state
26 is also the limit point of the infinitely recursive iterations of the implicit layer, DEQ can be regarded
27 as a new neural network that models the limit of a multi-layer weight-tied neural network with the
28 depth goes to infinity.

29 Nowadays, DEQ has become a popular and widely studied model in the field of machine learning.
30 On the empirical side, competitive performances against explicit feedforward neural networks have
31 been achieved in various real applications such as natural language processing [2], computer vision
32 [3], image generation [4], and solving inverse problems [5]. On the theoretic side, a main research
33 line is to study the well-posedness of DEQ. This line aims to analyze when unique equilibrium can
34 be guaranteed by DEQ and some weight parameterization and initialization techniques have been
35 proposed to ensure the well-posedness [6, 7, 8].

36 However, despite wide studies on DEQ, an understanding of the basic learning theory for its separa-
 37 tion and bias against explicit feedforward neural networks is still limited. For the expressivity, a
 38 preliminary study about the connections between DEQ and fully-connected network (FNN) is pro-
 39 vided in the seminar work [2], where it is shown that every FNN can be reformulated as a large DEQ
 40 under a specific weight re-parameterization, whereas, a deeper study on the provable and quantitative
 41 advantage of DEQ in its expression power is still lacking. Besides, there is another research line that
 42 studies the learning properties of DEQ using the so-called neural tangent kernel (NTK) view [9],
 43 originating from analyzing FNNs [10, 11]. It is shown [12, 13] that under suitable initialization, the
 44 dynamic of over-parameterized DEQ can be approximated by a linear kernel model, therefore global
 45 convergence of Gradient Descent algorithm and possible generalization can be achieved under some
 46 regimes. However, it is still not known whether DEQ has potential advantages over FNNs, even in
 47 such simplified settings. A study on the separation and bias of DEQ over FNN can provide us with
 48 clear and intuitive suggestions about when DEQ is preferred in practice, thus it is strongly desired. In
 49 this paper, we initialize the study by analyzing its expressive power and learning dynamics. The main
 50 results are sketched as follows.

- 51 1. We first propose a general separation showing that there exists a width- m DEQ which cannot
 52 be approximated to a constant accuracy by an FNN with depth $O(m^\alpha)$ for $\alpha \in (0, 1)$ unless
 53 its width is $\exp(\Omega(m^{1-\alpha}))$. This is achieved by comparing the the number of linear regions
 54 that the two networks can generate. Based on the result, we further prove that a width- m
 55 DEQ can generate at most 2^m linear regions, which has provable advantages than FNNs.
- 56 2. We then propose another separation, where a steep function in $[0, 1]^d$ being the solution
 57 to fixed point equation is considered as the target function. We show that a DEQ with
 58 size and magnitude bounded by $O(\varepsilon^{-1})$ can approximate this function to $O(\varepsilon)$ -accuracy in
 59 L^∞ norm, whereas an FNN with bounded depth and exponentially bounded width cannot
 60 unless its weights is $\exp(\Omega(d))$. For the technical contribution, we manage to show that
 61 an approximation of the fixed point mapping by the implicit layer can also guarantee the
 62 approximation the solution defined by the fixed point equation even if the Lipschitz constant
 63 of the fixed point mapping is very close to 1 by a new observation as shown in Lemma 3.
- 64 3. Finally, we study the bias of DEQ from the perspective of learning dynamics. We propose a
 65 general characterization of regularization for gradient flow in an overparameterized setting.
 66 We further analyze the dynamics of both gradient flow and gradient descent, showing that
 67 under mild conditions, convergence is guaranteed, and the model tends to produce ‘dense’
 68 features. Then we offer a concrete example on a specific Out-of-Distribution (OOD) task,
 69 demonstrating that this bias can help reduce the OOD error.

70 Finally, we conduct experiments to validate our theoretical results. From the overall study, a high-level
 71 conjecture is that DEQ has potential advantages in learning certain high-frequency components.

72 **Notations.** We use standard notation $O(\cdot)$ and $\Omega(\cdot)$ to hide constants. We use σ to denote the ReLU
 73 function, i.e., $\sigma(x) = \max(0, x)$, and we use $\text{sgn}(\cdot)$ to denote the sign function. We use $\text{diag}(\cdot)$
 74 to transform a vector into a diagonal matrix with the vector’s elements on the diagonal. We denote by
 75 $\|\cdot\|_p$ the ℓ^p vector norm or the subordinate matrix norm, and by $\|h\|_{L^p(K)}$ the L^p -norm of a function
 76 h on a compact set K . For a vector or vector-valued function \mathbf{v} , we denote v_i the i -th entry of the
 77 vector or the function. For a function $u : \mathbb{R} \rightarrow \mathbb{R}$, we denote $u^{\circ n}$ the n -fold composition of u .

78 2 Related Works

79 In this section we briefly review the literature that are most related to us.

80 **Theoretical Studies on DEQs.** Theoretical research on DEQs has primarily focused on ensuring
 81 their well-posedness [6, 7]. To guarantee well-posedness, different strategies are proposed, including
 82 new parameterizations of DEQ [6, 7], regularization [14], special initialization [8]. Another research
 83 line delve into the learning properties of DEQ. The expressivity of DEQ is preliminarily studied in
 84 [2]. Additionally, some recent works [15, 16, 13] couple the dynamics of over-parameterized DEQs
 85 with a linear kernel using the NTK method. They manage to prove the global convergence and study
 86 the generalization [16]. Nevertheless, an in-depth study on the potential or quantifiable advantage of
 87 DEQ over FNN is still lacking.

88 **Separations on Expressivity of Neural Networks.** The separation on expressivity of neural
 89 networks is a fundamental study characterizing functions that can be approximated efficiently by one
 90 type of neural architecture but not by another. These architectures include FNNs [17, 18], CNNs [19],
 91 RNNs [20], etc. Since DEQ can be viewed as an infinitely deep weight-tied neural networks, depth
 92 separation [21] is most relevant to our study. A key study by Telgarsky [22] constructs a saw-tooth
 93 function that have many oscillations to give a separation, which further inspires a series of separations
 94 [23, 24, 25]. In addition to depth, some recent works study the separation regarding the overall
 95 number of neurons in networks [26] or the magnitude of parameters [27] of FNNs. In this paper, the
 96 first separation result is also inspired by Telgarsky’s construction, whereas we focus on the separation
 97 between DEQ and FNN and provide a more refined analysis regarding the networks’ depth. The
 98 second separation is new.

99 **Implicit Bias of Learning Dynamics on Neural Networks.** The implicit bias of learning dynamics
 100 plays a key role in determining what particular optima can be found by the algorithms when there
 101 are multiple optima. A series of papers study the implicit regularization of gradient-based methods,
 102 showing that under varying settings, these algorithms bias towards solutions with specific properties
 103 [28, 29], such as norm minimization [30], sparsity [31] and low complexity [32, 33, 34]. Due to
 104 the theoretical barrier in analyzing nonlinear neural networks [35], most existing works focus on
 105 simplified models such as random feature models [36, 30], networks with quadratic activations [37]
 106 and diagonal linear networks [31]. This paper follows similar strategies and analyzes the implicit
 107 bias of a simplified diagonal linear DEQ from learning dynamics.

108 3 Preliminaries of DEQ

109 The DEQ is an implicit-depth model [2] that employs the same weights in each layer block of a
 110 feedforward neural network and extends the number of layer to infinity. The layer blocks used in
 111 DEQ can be fully connected, convolutional, or Transformer blocks, resulting in different variants of
 112 deep equilibrium networks. In this paper, we consider a vanilla DEQ with ReLU activation as the
 113 generalization of an FNN. Specifically, an L -layer FNN from \mathbb{R}^d to \mathbb{R}^s can be expressed as

$$114 \mathbf{z}^1 = \mathbf{x}; \quad \mathbf{z}^{i+1} = \sigma(\mathbf{W}_i \mathbf{z}^i + \mathbf{b}_i), \quad 1 \leq i \leq L-2; \quad \mathbf{y} = \mathbf{W}_L \mathbf{z}^{L-1}, \quad (1)$$

114 where $\mathbf{x} \in \mathbb{R}^d$ and $\mathbf{y} \in \mathbb{R}^s$. In DEQ, each \mathbf{W}_i and \mathbf{b}_i in Eq. (1) is replaced by the same weight \mathbf{W} and
 115 bias \mathbf{b} , and a linear transform of the input $\mathbf{U} \mathbf{x}$ is added to each layer, i.e., $\mathbf{z}^l = \sigma(\mathbf{W} \mathbf{z}^{l-1} + \mathbf{U} \mathbf{x} + \mathbf{b})$
 116 for all l . By extending the layer l to infinity, the feature and the prediction of this DEQ can be
 117 expressed as

$$118 \begin{aligned} \mathbf{z} &= \sigma(\mathbf{W} \mathbf{z} + \mathbf{U} \mathbf{x} + \mathbf{b}), \\ \mathbf{y} &= \mathbf{A} \mathbf{z}, \end{aligned} \quad (2)$$

118 where $\mathbf{W} \in \mathbb{R}^{m \times m}$, $\mathbf{U} \in \mathbb{R}^{m \times d}$, $\mathbf{b} \in \mathbb{R}^m$, and $\mathbf{A} \in \mathbb{R}^{s \times m}$. We call $\sigma(\mathbf{W} \mathbf{z} + \mathbf{U} \mathbf{x} + \mathbf{b})$ the implicit
 119 layer and m the width of DEQ. In this paper, we mainly consider $s = 1$, i.e., DEQ as a scalar function
 120 on \mathbb{R}^d .

121 In [2], the authors show that every FNN can be reformulated as a large DEQ with specific weight
 122 reparameterization. Specifically, the depth- L FNN described in Eq. (1) is equivalent to a DEQ in the
 123 form of Eq.(2) with

$$124 \mathbf{A} = (\mathbf{0}, \quad \dots, \quad \mathbf{I}), \quad \mathbf{W} = \begin{pmatrix} \mathbf{0} & & & & \\ \mathbf{W}_2 & \mathbf{0} & & & \\ & \mathbf{W}_3 & \mathbf{0} & & \\ & & \ddots & \ddots & \\ & & & \mathbf{W}_{L-1} & \mathbf{0} \end{pmatrix}, \quad \mathbf{U} = \begin{pmatrix} \mathbf{W}_1 \\ \mathbf{0} \\ \vdots \\ \mathbf{0} \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \\ \vdots \\ \mathbf{b}_{L-1} \end{pmatrix}. \quad (3)$$

124 4 Separation on the Expressivity of DEQ

125 In this section, we focus on the separations on the expressivity of DEQ.

126 **4.1 General Separation over FNNs**

127 The following theorem states a general separation between DEQ and FNN from the size of networks.
 128 The motivation behind the theorem is a common observation that functions with many linear pieces
 129 are typically hard to be approximated by functions having fewer linear pieces.

Theorem 1. *Let $m \in \mathbb{N}^+$. Assume that $L \leq m^\alpha$ for some $0 < \alpha < 1$. Then there exists a function $N_d : [0, 1]^d \rightarrow \mathbb{R}$ computed by a width- m ReLU-DEQ, such that for any function N_f computed by a depth- L ReLU-FNN with width at most $2^{m^{1-\alpha}-2}$, it holds that*

$$\int_{[0,1]^d} |N_d(\mathbf{x}) - N_f(\mathbf{x})| d\mathbf{x} \geq \frac{1}{16}.$$

130 The proof involves quantifying the number of linear regions¹ generated by a DEQ compared to an
 131 FNN. Specifically, we show in the proof that there exists a DEQ producing 2^m linear pieces whereas
 132 no-so-deep FNNs, i.e., FNNs with depth $O(m^\alpha)$ cannot generate such a large number of linear
 133 regions unless the width is sub-exponentially large.

134 Moreover, the example of the hard-to-approximate DEQ enables us to derive an exact bound on the
 135 number of linear regions that a DEQ can generate. This result is of independent interest and is stated
 136 in the proposition below.

137 **Proposition 1.** *Let $m > 0$. A width- m DEQ has at most 2^m linear regions in the input space.
 138 Moreover, this upper bound is attainable, i.e., there exists a width- m DEQ that computes a function
 139 with 2^m linear regions on \mathbb{R}^d .*

140 **Remark 1.** *As a comparison, the work of [38] analyzes ReLU-FNNs. It shows that for a ReLU-FNN
 141 with a total of \tilde{N} neurons of arbitrary depth, the maximal number of linear regions is bounded above
 142 by $2^{\tilde{N}}$. To the best of our knowledge, it is yet to be determined whether this bound is achievable.
 143 Consequently, width- m DEQs can potentially generate a larger number of linear regions compared
 144 to FNNs with m neurons, as DEQs have been shown to achieve their upper bound.*

145 Theorem 1 shows that there exists a width- m DEQ that is hard to be approximated by FNN with
 146 depth $O(m^\alpha)$. This theorem along with Proposition 1 reveals that, although DEQ computes features
 147 by solving an equilibrium function induced by a shallow implicit layer, its complexity in terms of
 148 expressing linear regions of DEQ can be larger than that of not-so-deep FNN.

149 **4.2 Separation on Certain Steep Functions**

150 In this section, we present another separation concerning both the size and parameter magnitude
 151 of neural networks, which more explicitly reveals the bias and potential advantages of DEQ on
 152 expressivity. The separation is based on the observation that the fixed point of a DEQ can be rewritten
 153 as the solution to an optimization problem under certain conditions.

154 To be specific, consider a simple quadratic optimization problem with the optimization variable
 155 $\mathbf{z} \in \mathbb{R}^m$ and a parameter $\mathbf{x} \in \mathbb{R}^d$:

$$\min_{\mathbf{z}} \frac{1}{2} \mathbf{z}^T \mathbf{A}(\mathbf{x}) \mathbf{z} + \mathbf{b}^T(\mathbf{x}) \mathbf{z} + \mathbf{c}, \quad (4)$$

where $\mathbf{A}(\mathbf{x})$ is a positive definite matrix parameterized by \mathbf{x} and $\eta \mathbf{I} \succ \mathbf{A}(\mathbf{x}) \succ \mathbf{0}$ for some $\eta > 0$. Approximating $\mathbf{z} = \mathbf{z}(\mathbf{x})$, i.e., the optimum as a function of the parameter \mathbf{x} , serves useful primitives in various applications. Directly approximating $\mathbf{z}(\mathbf{x})$ by FNN requires the approximation of $\mathbf{z}(\mathbf{x}) = -\mathbf{A}(\mathbf{x})^{-1} \mathbf{b}(\mathbf{x})$. On the other hand, from the optimality condition, $\mathbf{z}(\mathbf{x})$ is implicitly defined through fixed point equation

$$\mathbf{z} = \mathbf{z} - \frac{1}{\eta} (\mathbf{A}(\mathbf{x}) \mathbf{z} + \mathbf{b}(\mathbf{x})).$$

¹We follow the definition of linear regions in [38]: For any piecewise linear function $F : \mathbb{R}^{n_0} \rightarrow \mathbb{R}$, a linear region of the function is a subset $D \subset \mathbb{R}^{n_0}$ satisfying 1) F is linear on D ; 2) If F is linear on some set $\tilde{D} \supset D$, then $\tilde{D} = D$.

156 Hence, approximating $\mathbf{z}(\mathbf{x})$ by DEQ may only require the approximation of the fixed point mapping
 157 $\mathbf{z} - \frac{1}{\eta} (\mathbf{A}(\mathbf{x}) \mathbf{z} + \mathbf{b}(\mathbf{x}))$ by the implicit layer. To some extent, the approximation problem is ‘altered’
 158 due to the model difference, which possibly leads to distinctive division in approximation.

159 Now, we construct a workable instance. The objective function of our central interest is a special case
 160 of Eq.(4) given by:

$$\min_z (1 + \delta - x_1)z^2 - \delta x_1 z, \quad \mathbf{x} \in [0, 1]^d, \quad (5)$$

161 where $\delta = 2^{-d}$. The solution function is calculated as

$$g(\mathbf{x}) = \frac{\delta x_1}{2(1 + \delta - x_1)}, \quad \mathbf{x} \in [0, 1]^d, \quad (6)$$

162 and it can also be determined by the following fixed point equation

$$z = \tilde{g}(z, \mathbf{x}) := (x_1 - \delta)z + \frac{1}{2}\delta x_1. \quad (7)$$

163 Note that $g(\mathbf{x})$ has very large derivative when x_1 is near 1. It can be regarded as a continuous version
 164 of the common indicator function of the first entry $\frac{1}{2}\mathbf{1}_{x_1=1}(\mathbf{x})$. The separation is presented as follows.

165 **Theorem 2.** *Let $g(\mathbf{x})$ be defined as in Eq.(6) for $\mathbf{x} \in [0, 1]^d$ and $\frac{1}{4} \geq \varepsilon > 0$.*

A. *For any function $N_{fnn}(x)$ implemented by an FNN with depth L and width k where $L \leq C$
 and $k \leq 2^{\frac{d}{2C}}$ for some constant $C = O(1)$. If*

$$\|N_{fnn}(\mathbf{x}) - g(\mathbf{x})\|_{L^\infty([0,1]^d)} \leq \frac{1}{16},$$

*then there exists a weight parameter W_{ij} of the FNN for $1 \leq i \leq L$ and $1 \leq j \leq k$, such
 that*

$$|W_{ij}| \geq 2^{\frac{d}{2C}}.$$

B. *There exists a function N_{deq} implemented by a DEQ with width bounded by $5\varepsilon^{-1}$ and
 weights bounded $2\varepsilon^{-1}$, such that*

$$\|N_{deq}(\mathbf{x}) - g(\mathbf{x})\|_{L^\infty([0,1]^d)} \leq \varepsilon.$$

166 **Remark 2.** *The inapproximability result of FNN in Theorem 2 is stated from the perspective of*
 167 *weight magnitude, which holds practical significance. Exponentially large weight often results in*
 168 *exponential iterations of optimization algorithms in learning with this model, as also noted in [39].*
 169 *Additionally, neural networks in practice typically have small weights due to techniques such as*
 170 *(standard) small initialization, normalization, and gradient clipping.*

171 In Theorem 2, the inapproximability of FNNs is relatively simple: Direct calculation shows that
 172 the derivative of the target function $g(\mathbf{x})$ is exponentially large when $x_1 > 1 - \delta$. To approximate
 173 $g(\mathbf{x})$ in L^∞ norm requires FNNs to have large derivative in certain region, resulting in exponentially
 174 large weight for FNNs with bounded depth. On the other hand, the proof of the approximability of
 175 DEQs is more technical. While \tilde{g} in Eq. (7) seems more benign, it is not clear how to construct the
 176 approximation using the implicit layer in Eq. (2) that resembles an 1-layer FNN with very limited
 177 expressive power. Moreover, even if we manage to approximate \tilde{g} in Eq. (7), it will not necessarily
 178 imply a good approximation between the fixed point of DEQ and the solution of $z = \tilde{g}(z, \mathbf{x})$, i.e., the
 179 target function due to the Lipschitz constant of \tilde{g} with respect to z being very close to 1 when x_1 is
 180 around 1 according to Eq. (7). We provide a proof sketch of this result in Section 4.3.

181 Further insights and implications can be gleaned from Theorem 2. First, it suggests that DEQ may
 182 excel in approximating functions induced by fixed-point iterations. In other words, DEQ may be
 183 better suited for representing algorithms. Second, Theorem 2 implies that functions with large
 184 derivative, or high-frequency components, may be approximated more efficiently by DEQ, as the
 185 function to be approximated by the implicit layer can have much smaller derivative.

186 4.3 Proof Sketch of B. in Theorem 2

187 As discussed in Section 4.2, we want to approximate \tilde{g} using the implicit layer of DEQ. Due to the
 188 limited expressive power of the implicit layer, we propose an equivalent reparameterization of DEQ.

189 **Lemma 1.** Consider a revised DEQ defined as

$$\begin{aligned} \mathbf{z} &= \mathbf{V}\sigma(\mathbf{W}\mathbf{z} + \mathbf{U}\mathbf{x} + \mathbf{b}), \\ \mathbf{y} &= \mathbf{B}\mathbf{z}, \end{aligned} \quad (8)$$

190 where $\mathbf{x} \in \mathbb{R}^d$, $\mathbf{z} \in \mathbb{R}^m$, $\mathbf{W} \in \mathbb{R}^{q \times m}$, $\mathbf{U} \in \mathbb{R}^{q \times d}$, $\mathbf{V} \in \mathbb{R}^{m \times q}$, $\mathbf{B} \in \mathbb{R}^{p \times m}$ and $\|\mathbf{W}\mathbf{V}\|_2 \leq 1$. Then
191 any revised DEQ can be represented by a vanilla DEQ defined as in Eq. (2) with width q .

Lemma 1 enables us to approximate $\tilde{g}(z, \mathbf{x})$ using the revised implicit layer, denoted by $\tilde{h}(z, \mathbf{x})$. Then the crux of the proof centered in bounding the error between the equilibria of two fixed-point equations. To begin, for every \mathbf{x} we denote $\hat{u}(z) = z - \tilde{g}(z, \mathbf{x})$, $\hat{v}(z) = z - \tilde{h}(z, \mathbf{x})$ and consider $|\hat{u}^{\circ 2}(z) - \hat{v}^{\circ 2}(z)|$. Suppose that $\hat{u}(z)$ is $L_{\hat{u}}$ -Lipschitz, then we have

$$|\hat{u}^{\circ 2}(z) - \hat{v}^{\circ 2}(z)| \leq |\hat{u}^{\circ 2}(z) - \hat{u} \circ \hat{v}(z)| + |\hat{u} \circ \hat{v}(z) - \hat{v}^{\circ 2}(z)| \leq (L_{\hat{u}} + 1)|\hat{u}(z) - \hat{v}(z)|.$$

192 Thus if $L_{\hat{u}} < 1$, by recursion, we can bound distance the between the infinitely composition of $\hat{u}(z)$
193 and $\hat{v}(z)$, from which the error of the two fixed points can be bounded.

Lemma 2. Let $\Omega \subset \mathbb{R}$ be a compact set, and $u(z, \mathbf{x}), v(z, \mathbf{x}) : \Omega \times [0, 1]^d \rightarrow \Omega$ be two functions. Assume that for all $\mathbf{x} \in [0, 1]^d$, $u(\cdot, \mathbf{x})$ and $v(\cdot, \mathbf{x})$ are Lipschitz continuous with Lipschitz constant $L_u, L_v < 1$, respectively. Then for any $\mathbf{x} \in [0, 1]^d$, it holds

$$|z_u - z_v| \leq \min\{(1 - L_u)^{-1}, (1 - L_v)^{-1}\} \cdot |u(z, \mathbf{x}) - v(z, \mathbf{x})|$$

194 for all $\forall(z, \mathbf{x}) \in \Omega \times [0, 1]^d$, where z_u and z_v are the fixed point of $z = u(z, \mathbf{x})$ and $z = v(z, \mathbf{x})$,
195 respectively.

196 In our case, $u(z, \mathbf{x})$ and $v(z, \mathbf{x})$ in this Lemma represent $\tilde{g}(z, \mathbf{x})$ and $\tilde{h}(z, \mathbf{x})$, respectively. When
197 $x < 1 - \text{poly}(d)^{-1}$, by calculating $\frac{\partial \tilde{g}(z, \mathbf{x})}{\partial z}$, we have $(1 - L_{\tilde{g}})^{-1} < \text{poly}(d)$. Leveraging this and
198 Lemma 2, we just need $\|\tilde{h} - \tilde{g}\|_{\infty} \leq \text{poly}(d)^{-1}$ to achieve a final accuracy of $O(\varepsilon)$. However,
199 when $x \geq 1 - \delta$, we only have $(1 - L_{\tilde{g}})^{-1} < \exp(\Omega(d))$, which may necessitate an exponential
200 width for the implicit layer to achieve $O(\varepsilon)$ accuracy. In fact, $\tilde{h}(z, \mathbf{x}) = x_i z$ gives an example
201 that even assuming $\|\tilde{h} - \tilde{g}\|_{\infty} \leq \exp(\Omega(d))^{-1}$ is not sufficient to achieve $O(\varepsilon)$ accuracy since
202 $z_{\tilde{h}}(1) - z_{\tilde{g}}(1) = \frac{1}{2}$. So it seems difficult to bound the error without a specific structure of \tilde{h} . To
203 overcome the issue, we observe a *novel* property that enables us to effectively bound the error.

Lemma 3. Let $\xi > 0$. Under the conditions in Lemma 2, if for any interval $T \subset \Omega$ with $\text{diam}(T) > \xi$, $u(z, \mathbf{x}) = v(z, \mathbf{x})$ has a zero in T for all \mathbf{x} , then it holds that

$$|z_u(\mathbf{x}) - z_v(\mathbf{x})| \leq \xi, \quad \forall \mathbf{x} \in [0, 1]^d.$$

204 The intuition behind Lemma 3 is that if for any \mathbf{x} , $z = u(z, \mathbf{x})$ and $z = v(z, \mathbf{x})$ as two monotone
205 univariate functions w.r.t. z can take the same value at frequent intervals, then their zeros will also be
206 close to each other. By using this Lemma, it suffices to construct such $\tilde{h}(z, \mathbf{x})$ that equals $\tilde{g}(z, \mathbf{x})$ at
207 frequent interval of length $O(\varepsilon)$ for every \mathbf{x} .

208 5 The Bias on Learning Dynamics of DEQ

209 In this section, we study the implicit bias of a simplified linear diagonal DEQ and present a concrete
210 example illustrating how such an implicit bias may lead to improve generalization.

211 We begin by considering the model:

$$f(\mathbf{w}, \mathbf{x}) = \sum_{i=1}^d \frac{1}{1 - w_i} x_i := \langle \boldsymbol{\beta}, \mathbf{x} \rangle, \quad \beta_i = \frac{1}{1 - w_i}. \quad (9)$$

212 The model can be regarded as a diagonal linear DEQ in Eq. (2) with $\mathbf{W} = \text{diag}(w_1, w_2, \dots, w_d)$,
213 $\mathbf{U} = \mathbf{I}_d$, $\mathbf{b} = \mathbf{0}$ and $\mathbf{A} = (1, 1, \dots, 1)^T \in \mathbb{R}^d$. Our primary focus lies in minimizing the expected
214 square loss:

$$\min_{\mathbf{w}} L(\mathbf{w}) := \frac{1}{2} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [(y - f(\mathbf{w}, \mathbf{x}))^2]. \quad (10)$$

We are given access to a set of i.i.d. training examples $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$, and we denote the (half) square loss on these examples by $\hat{L}(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (y_i - f(\mathbf{w}, x_i))^2$. Moreover, let

$$\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)^T, \quad \mu_{\min} = \lambda_{\min}(\mathbf{X}\mathbf{X}^T), \quad \mu_{\max} = \lambda_{\max}(\mathbf{X}\mathbf{X}^T).$$

215 We mainly consider the dynamics of gradient flow (GF) and gradient descent (GD) with fixed stepsize
216 η on minimizing $\hat{L}(\mathbf{w})$, expressed as follows

$$(GF) \quad \dot{\mathbf{w}}(t) = -\nabla_{\mathbf{w}} \hat{L}(\mathbf{w}(t)); \quad (GD) \quad \mathbf{w}^{k+1} = \mathbf{w}^k - \eta \nabla_{\mathbf{w}} \hat{L}(\mathbf{w}^k). \quad (11)$$

217 The main theorem below gives a general characterization of the bias of diagonal linear DEQ in the
218 overparameterized regime. The proof is based on the technique proposed in [29].

219 **Theorem 3.** *Let β_i in Eq. (9) be initialed as $\beta_i(0) > 0$ for all i . Suppose that gradient flow for the
220 parameterization problem in Eq. (10) converges to some $\hat{\beta}$ satisfying $\mathbf{X}\hat{\beta} = \mathbf{y}$, then*

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} Q(\beta), \quad \text{s.t. } \mathbf{X}\beta = \mathbf{y}, \quad (12)$$

221 where $Q(\beta) = \sum_{i=1}^d q(\beta_i)$ and $q(x) = \frac{1}{2x^2} + \beta_i(0)^{-3}x$.

222 The theorem implies that the bias of the (simplified) DEQ significantly differs from that of conven-
223 tional linear models and two-layer linear network which tends to give a minimum ℓ_2 -norm interpolator
224 [40]. Specifically, the predictor $\hat{\beta}$ hardly admits parameters of small magnitude due to the penalty
225 term $\frac{1}{2} \sum_{i=1}^d \frac{1}{\beta_i}$. Meanwhile, the predictor can endure parameters of greater magnitude as the penalty
226 $q(x)$ increase almost linearly when x is large.

227 We then study the implicit bias from the learning dynamics of GF and GD. We show that when
228 $\mu_{\min} > 0$, under mild conditions, the convergence of both algorithms is guaranteed. Moreover, in this
229 case, a positive lower bound of the ℓ_{∞} norm of the iterates can be derived, indicating that the model
230 inclines to produce ‘dense’ features in learning process.

Assumption 1. *Denote by β_0 the initialization of β of the model. There exists an optima $\hat{\beta}^*$, i.e.,
 $\mathbf{X}\hat{\beta}^* = \mathbf{y}$ and a constant $c > 0$, such that*

$$\|\hat{\beta}^*\|_{\infty} - \|\hat{\beta}^* - \beta_0\|_2 \geq c > 0.$$

231 **Theorem 4.** *Let $\{\beta(t)\}$ be the process following GF in Eq. (11) and $\{\beta^k\}$ the iterates following GD
232 in Eq. (11). Assume that $\mu_{\min} > 0$ and the initialization $\beta(0)$ and β^0 satisfy Assumption 1 with an
233 optima $\hat{\beta}^*$*

234 A. *$\{\beta(t)\}$ converges to an optima β_f^{∞} with $\|\beta_f^{\infty}\|_{\infty} \geq c$. Moreover, for any $t \geq 0$, we have
235 $c \leq \|\beta(t)\|_{\infty} \leq \|\hat{\beta}^*\|_{\infty} + \|\hat{\beta}^* - \beta_0\|_2$.*

236 B. *If there exists a constant $C > 0$ such that $\|\beta^k\|_{\infty} \leq C$ for all k , then $\{\beta^k\}$ converges to an
237 optima β_d^{∞} with $\|\beta_d^{\infty}\|_{\infty} \geq c$. Moreover, for any $k \geq 0$, we have $c \leq \|\beta^k\|_{\infty} \leq C$.*

238 **Remark 3.** *The assumption in Theorem 4 that $\|\beta^k\|_{\infty}$ is uniformly bounded can be removed if we
239 manually incorporate a constrain on β and optimize the problem using projected gradient descent.
240 In practice, certain reparameterization tricks [6, 7] are proposed to ensure that $\mathbf{I} - \mathbf{W} \succeq m\mathbf{I}$ for
241 some $m > 0$, thus corresponding to the aforementioned assumption.*

242 Based on our results above, we now provide a concrete example to show the advantages brought by
243 the bias of DEQ in out-of-distribution (OOD) tasks. This is motivated by the fact that diversifying
244 spurious features can improve OOD generalization [41]. Specifically, we focus on generalization
245 on the unseen domain (GOTU) setting [34], a rather strong case of OOD generalization where part
246 of the distribution domain is unseen at training but used at testing. As an example, we here utilize
247 the setting in Theorem 3.11 in [34]. Consider the sample space $\mathcal{S} = \{-1, 1\}^d$ and a linear boolean
248 function $f : \mathcal{S} \rightarrow \mathbb{R}$ defined as

$$f(\mathbf{x}) = \hat{f}(\emptyset) + \sum_{i=1}^d \hat{f}(\{i\})x_i, \quad (13)$$

where $\hat{f}(\{i\}) = \mathbb{E}_{\mathbf{x} \sim \mathcal{U}} [x_i f(\mathbf{x})]$ and $\sim \mathcal{U}$ refers to uniform sampling from \mathcal{U} . In training, the k -th component of every accessible sample is fixed as 1, i.e., the unseen domain is $\mathcal{U} = \{\mathbf{x} \in \{\pm 1\}^d : x_k = -1\}$. Denote by $\tilde{f}_{\mathcal{S} \setminus \mathcal{U}}$ the function learned on $\mathcal{S} \setminus \mathcal{U}$. The GOTU error is defined as the generalization completely on the unseen domain, i.e.,

$$GOTU(f, \tilde{f}, \mathcal{U}) = \mathbb{E}_{X \sim \mathcal{U}} [l(\tilde{f}_{\mathcal{S} \setminus \mathcal{U}}(X), f(X))],$$

249 where l is the quadratic loss function. It is shown in [34] that learning this function with diagonal
 250 linear network results in a GOTU error of $4\hat{f}(\{k\})^2 + O(\varepsilon)$ for an infinitesimal ε . On the other
 251 hand, the following proposition shows that under mild conditions, learning such function with
 252 DEQ achieves smaller GOTU error, where we consider DEQ in Eq. (9) with a bias term, i.e.,
 253 $f(\mathbf{w}, \mathbf{x}) = \sum_{i=1}^d \frac{1}{1-w_i} x_i + b$.

Proposition 2. *Let $f(\mathbf{x})$ be defined as in Eq. (13). Assume that*

$$\hat{f}(\{i\}) > 0, \quad \forall 1 \leq i \leq d, \quad \hat{f}(\{k\}) > 1, \quad |\hat{f}(\emptyset)| \leq 2|\hat{f}(\{k\})|.$$

Consider learning f using gradient flow on population loss² on a linear diagonal DEQ with bias initialized by $w_i(0) = b(0) = 0$ for all i with unseen domain $\mathcal{U} = \{\mathbf{x} \in \{\pm 1\}^d : x_k = -1\}$. Then the loss converge to 0, and it holds for the generalization error on the unseen that

$$GOTU \leq 4 \left(\hat{f}(\{k\}) - \left(4 + 3\hat{f}(\{k\})\right)^{-\frac{1}{3}} \right)^2 < 4\hat{f}(\{k\})^2.$$

254 In this setting, the function x_k has a higher frequency component (i.e., degree) compared to the
 255 constant function 1. Consequently, the inductive bias of DEQ enables the model to capture some
 256 information about the high-frequency components. We further conduct experiments to study the
 257 potential advantages of DEQ in learning high-frequency components in Appendix B.2.

258 6 Experiments

259 In this section, we conduct experiments on FNNs and DEQs based on our theoretical results. We first
 260 evaluate the expressivity of both networks on the functions proposed in our two separation results.
 261 Then we experiment on specific OOD tasks. An additional experiment on audio representation is
 262 provided in Appendix B.2.

263 **Piecewise functions.** We first verify the results in Section 4.1. The target function is designed as a
 264 saw-tooth function, as defined in Lemma 4 in Appendix A.1, which can be exactly computed by a
 265 DEQ. We set the number of linear regions of the saw-tooth function to 2^5 and 2^{10} and experiments
 266 on other sawtooth functions can be seen in Appendix B.1. According to Proposition 1, a DEQ with
 267 width 5 and 10 can compute the above functions exactly. Following the standard setting, all models
 268 are trained using ℓ_2 loss with AdamW optimizer [42], with a learning rate of 5e-4, weight decay of
 269 1e-4 and a cosine annealing scheduler for 1000 iterations.

270 Figure 1(a) and Figure 1(d) show that DEQ can achieve nearly zero test loss, demonstrating the
 271 saw-tooth function with 2^m linear regions can be computed by DEQ. On the other hand, a not-so-deep
 272 and not-so-wide FNN fails to achieve test loss as low as DEQ, thus verifying the separation results
 273 between FNN and DEQ.

274 **Solution to quadratic optimization problem.** We then validate the ability of DEQ to approximate
 275 the solution function to certain optimization problems. We empirically demonstrate that DEQ can
 276 approximate such function better than an FNN with a similar number of parameters. We consider
 277 the objective function $g(\mathbf{x})$ defined in Eq. (6), with the input dimension d 10 and 20, and thus δ in
 278 target function being 2^{-10} and 2^{-20} . The input space is $\mathbf{x} \in [0, 1]^d$ with the sampling distribution
 279 $p(\mathbf{x}) = \frac{1}{2(1-\delta)}$ for $0 < x_1 < 1 - \delta$ and $p(\mathbf{x}) = \frac{1}{2\delta}$ for $1 - \delta < x_1 < 1$.

280 In experiment, we train FNN and DEQ models using the ℓ_2 loss. Following the standard setting,
 281 we employ mini-batch SGD optimizer with learning rate of 0.005, weight decay of 1e-4 and cosine
 282 annealing scheduler for all models. To verify results under different network parameters, we adjust

²It is identical to the setting in Theorem 3.11, [34]. Note that optimizing the population loss in generalization cannot reduce the OOD error.

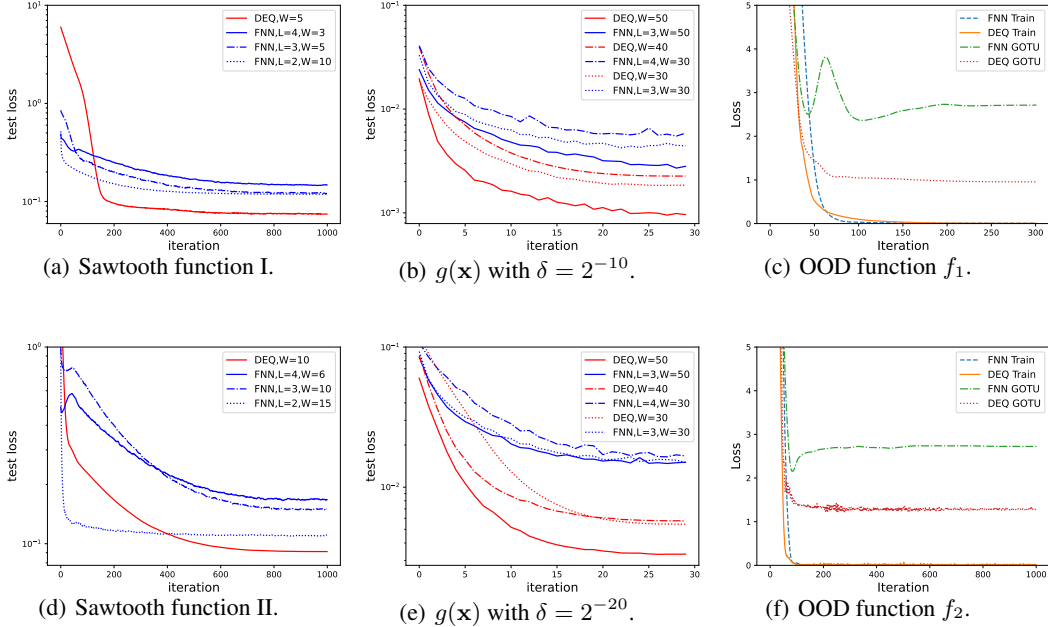


Figure 1: Test losses of FNN and DEQ networks with various width W and depth L . (a) and (b) apply Sawtooth function I and II with 2^5 and 2^{10} linear regions, respectively. (c) and (d) apply function $g(\mathbf{x})$ defined in Eq. (5) with $\delta = 2^{-10}$ and $\delta = 2^{-20}$, respectively. (e) Train loss and the GOTU error of FNN and DEQ on the boolean function f_1 , f_2 and unseen domain given by Eq. (14) and Eq. (15).

283 the layer number and hidden dimension of FNN and the layer width of DEQ while keeping the total
 284 number of parameters of both networks similar.

285 As shown in Figure 1(b) and Figure 1(e), for different network parameters and target functions,
 286 DEQ consistently achieves a lower test loss than FNN, demonstrating the superiority of DEQ to
 287 approximate and represent functions as solutions to certain optimization problems.

288 **Out-of-Distribution tasks.** We further perform experiments on the implicit bias of DEQ to verify
 289 the advantage of DEQ on OOD tasks. We consider 2 linear boolean functions $f : \mathcal{S} \rightarrow \mathbb{R}$ in the form
 290 of Eq. (13) and unseen domains $\mathcal{U} \subset \{\pm 1\}^d$. The first function is an example of the mean function
 291 and the second function is a part of DTFT. Experiments on other OOD functions can be found in
 292 Appendix B.1.

$$f_1(x) = 1.25x_0 + 1.25x_1 + 1.25x_2 + \dots + 1.25x_{10}, \quad \mathcal{U} = \{\mathbf{x} \in \{\pm 1\}^{10} : x_2 = -1\}, \quad (14)$$

$$f_2(x) = \sum_{n=0}^9 \sin\left(\frac{\pi * n}{10}\right)x_n, \quad \mathcal{U} = \{\mathbf{x} \in \{\pm 1\}^{10} : x_1 = -1\} \quad (15)$$

293 For each experiment, we generate all binary sequences in $\{\pm 1\}^d \setminus \mathcal{U}$ for training. We employ AdamW
 294 optimizer with ℓ_2 loss and a cosine annealing scheduler. We can observe in Figure1(c) that the GOTU
 295 error of f_1 is below the threshold of generalization error based on the Proposition 2. As shown in
 296 Figure1(c) and Figure1(f), the training loss converges to 0 and the generalization error on the unseen
 297 domain is bounded, which empirically demonstrates the advantage of DEQ on OOD tasks.
 298

299 7 Conclusions

300 In this paper, we provide two separations of DEQ and FNN and analyze the bias of DEQ through the
 301 lens of learning dynamics. Our theoretical results provably show the advantage of DEQ over FNN in
 302 specific problems and quantify certain learning properties of DEQ. Overall, we conjecture that DEQ
 303 may be advantageous in learning certain high-frequency components.

304 **References**

- 305 [1] L. El Ghaoui, F. Gu, B. Travacca, A. Askari, and A. Tsai, “Implicit deep learning,” *SIAM*
306 *Journal on Mathematics of Data Science*, 2021.
- 307 [2] S. Bai, J. Z. Kolter, and V. Koltun, “Deep equilibrium models,” *Advances in Neural Information*
308 *Processing Systems*, 2019.
- 309 [3] S. Bai, V. Koltun, and J. Z. Kolter, “Multiscale deep equilibrium models,” *Advances in Neural*
310 *Information Processing Systems*, 2020.
- 311 [4] A. Pokle, Z. Geng, and J. Z. Kolter, “Deep equilibrium approaches to diffusion models,”
312 *Advances in Neural Information Processing Systems*, 2022.
- 313 [5] D. Gilton, G. Ongie, and R. Willett, “Deep equilibrium architectures for inverse problems in
314 imaging,” *IEEE Transactions on Computational Imaging*, 2021.
- 315 [6] E. Winston and J. Z. Kolter, “Monotone operator equilibrium networks,” *Advances in Neural*
316 *Information Processing Systems*, 2020.
- 317 [7] M. Revay, R. Wang, and I. R. Manchester, “Lipschitz bounded equilibrium networks,” *arXiv*
318 *preprint arXiv:2010.01732*, 2020.
- 319 [8] A. Agarwala and S. S. Schoenholz, “Deep equilibrium networks are sensitive to initialization
320 statistics,” in *International Conference on Machine Learning*, 2022.
- 321 [9] A. Jacot, F. Gabriel, and C. Hongler, “Neural tangent kernel: Convergence and generalization in
322 neural networks,” *Advances in Neural Information Processing Systems*, 2018.
- 323 [10] S. Du, J. Lee, H. Li, L. Wang, and X. Zhai, “Gradient descent finds global minima of deep
324 neural networks,” in *International Conference on Machine Learning*, 2019.
- 325 [11] Z. Allen-Zhu, Y. Li, and Z. Song, “A convergence theory for deep learning via over-
326 parameterization,” in *International Conference on Machine Learning*, 2019.
- 327 [12] T. Gao, H. Liu, J. Liu, H. Rajan, and H. Gao, “A global convergence theory for deep relu implicit
328 networks via over-parameterization,” in *International Conference on Learning Representations*,
329 2022.
- 330 [13] Z. Ling, X. Xie, Q. Wang, Z. Zhang, and Z. Lin, “Global convergence of over-parameterized
331 deep equilibrium models,” in *International Conference on Artificial Intelligence and Statistics*,
332 2023.
- 333 [14] S. Bai, V. Koltun, and J. Z. Kolter, “Stabilizing equilibrium models by jacobian regularization,”
334 in *International Conference on Machine Learning*, M. Meila and T. Zhang, Eds., 2021.
- 335 [15] K. Kawaguchi, “On the theory of implicit deep learning: Global convergence with implicit
336 layers,” in *International Conference on Learning Representations*, 2021.
- 337 [16] T. Gao and H. Gao, “On the optimization and generalization of overparameterized implicit
338 neural networks,” *arXiv preprint arXiv:2209.15562*, 2022.
- 339 [17] R. Eldan and O. Shamir, “The power of depth for feedforward neural networks,” in *Conference*
340 *on Learning Theory*, 2016.
- 341 [18] M. Telgarsky, “Representation benefits of deep feedforward networks,” *arXiv preprint*
342 *arXiv:1509.08101*, 2015.
- 343 [19] Z. Wang and L. Wu, “Theoretical analysis of the inductive biases in deep convolutional networks,”
344 *Advances in Neural Information Processing Systems*, 2024.
- 345 [20] M. Emami, M. Sahraee-Ardakan, P. Pandit, S. Rangan, and A. K. Fletcher, “Implicit bias of
346 linear rnns,” in *International Conference on Machine Learning*, 2021.
- 347 [21] A. Daniely, “Depth separation for neural networks,” in *Conference on Learning Theory*, 2017.

- 348 [22] M. Telgarsky, “Benefits of depth in neural networks,” in *Conference on Learning Theory*, 2016.
- 349 [23] I. Safran and O. Shamir, “Depth-width tradeoffs in approximating natural functions with neural
350 networks,” in *International Conference on Machine Learning*, 2017.
- 351 [24] V. Chatziafratis, S. G. Nagarajan, and I. Panageas, “Better depth-width trade-offs for neural
352 networks through the lens of dynamical systems,” in *International Conference on Machine
353 Learning*, 2020.
- 354 [25] E. Malach, G. Yehudai, S. Shalev-Schwartz, and O. Shamir, “The connection between approx-
355 imation, depth separation and learnability in neural networks,” in *Conference on Learning
356 Theory*, 2021.
- 357 [26] G. Vardi, D. Reichman, T. Pitassi, and O. Shamir, “Size and depth separation in approximating
358 benign functions with neural networks,” in *Conference on Learning Theory*, 2021.
- 359 [27] G. Vardi and O. Shamir, “Neural networks with small weights and depth-separation barriers,”
360 *Advances in Neural Information Processing Systems*, 2020.
- 361 [28] S. Gunasekar, J. Lee, D. Soudry, and N. Srebro, “Characterizing implicit bias in terms of
362 optimization geometry,” in *International Conference on Machine Learning*, 2018.
- 363 [29] B. Woodworth, S. Gunasekar, J. D. Lee, E. Moroshko, P. Savarese, I. Golan, D. Soudry, and
364 N. Srebro, “Kernel and rich regimes in overparametrized models,” in *Conference on Learning
365 Theory*, 2020.
- 366 [30] P. L. Bartlett, A. Montanari, and A. Rakhlin, “Deep learning: a statistical viewpoint,” *Acta
367 Numerica*, 2021.
- 368 [31] E. Moroshko, B. E. Woodworth, S. Gunasekar, J. D. Lee, N. Srebro, and D. Soudry, “Implicit
369 bias in deep linear classification: Initialization scale vs training accuracy,” *Advances in Neural
370 Information Processing Systems*, 2020.
- 371 [32] Y. Cao, Z. Fang, Y. Wu, D. Zhou, and Q. Gu, “Towards understanding the spectral bias of deep
372 learning,” in *International Joint Conference on Artificial Intelligence*, Z. Zhou, Ed., 2021.
- 373 [33] N. Rahaman, A. Baratin, D. Arpit, F. Draxler, M. Lin, F. Hamprecht, Y. Bengio, and A. Courville,
374 “On the spectral bias of neural networks,” in *International Conference on Machine Learning*.
375 PMLR, 2019.
- 376 [34] E. Abbe, S. Bengio, A. Lotfi, and K. Rizk, “Generalization on the unseen, logic reasoning and
377 degree curriculum,” in *International Conference on Machine Learning*, 2023.
- 378 [35] G. Vardi and O. Shamir, “Implicit regularization in relu networks with the square loss,” in
379 *Conference on Learning Theory*, 2021.
- 380 [36] A. Jacot, B. Simsek, F. Spadaro, C. Hongler, and F. Gabriel, “Implicit regularization of random
381 feature models,” in *International Conference on Machine Learning*, 2020.
- 382 [37] Y. Li, T. Ma, and H. Zhang, “Algorithmic regularization in over-parameterized matrix sensing
383 and neural networks with quadratic activations,” in *Conference On Learning Theory*, 2018.
- 384 [38] G. F. Montufar, R. Pascanu, K. Cho, and Y. Bengio, “On the number of linear regions of deep
385 neural networks,” *Advances in Neural Information Processing Systems*, 2014.
- 386 [39] G. Yehudai and O. Shamir, “On the power and limitations of random features for understanding
387 neural networks,” *Advances in Neural Information Processing Systems*, 2019.
- 388 [40] A. V. Varre, M.-L. Vladarean, L. Pillaud-Vivien, and N. Flammarion, “On the spectral bias of
389 two-layer linear networks,” *Advances in Neural Information Processing Systems*, 2024.
- 390 [41] Y. Lin, L. Tan, Y. Hao, H. Wong, H. Dong, W. Zhang, Y. Yang, and T. Zhang, “Spurious feature
391 diversification improves out-of-distribution generalization,” *arXiv preprint arXiv:2309.17230*,
392 2023.

- 393 [42] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *arXiv preprint*
394 *arXiv:1412.6980*, 2014.
- 395 [43] Z. Huang, S. Bai, and J. Z. Kolter, “(Implicit)²: Implicit layers for implicit representations,”
396 *Advances in Neural Information Processing Systems*, 2021.

397 **A Proofs**

398 **A.1 Proofs in Subsection 4.1**

399 In following technical lemma, we show that there exists a width- m ReLU-DEQ computing a function
400 with 2^m linear regions.

Lemma 4. *Let $m \in \mathbb{N}^+$. For all $m \geq 1$, consider the following function on $[0, 1]^d$:*

$$\phi^{(m)}(\mathbf{x}) = \begin{cases} 2^m x_1 - 2i, & x_1 \in \left[\frac{2i}{2^m}, \frac{2i+1}{2^m}\right], \quad 0 \leq i \leq 2^{m-1} - 1, \\ -2^m x_1 + 2i + 2, & x_1 \in \left[\frac{2i+1}{2^m}, \frac{2i+2}{2^m}\right], \quad 0 \leq i \leq 2^{m-1} - 1. \end{cases}$$

401 *Then there exists a DEQ with width m that exactly computes $-\phi^{(m)}(\mathbf{x}) + 2^m x_1$ on $[0, 1]^d$. Moreover,*
402 *the DEQ has 2^m linear regions on $[0, 1]^d$.*

403 *Proof.* Since $2^m x_1$ is a linear function with respect to z_1 , by definition, $-\phi^{(m)}(\mathbf{x}) + 2^m x_1$ has 2
404 linear regions on $\left[\frac{2i}{2^m}, \frac{2i+2}{2^m}\right] \times [0, 1]^{d-1}$ for all $0 \leq i \leq 2^{m-1} - 1$. Thus it has 2^m linear regions on
405 $[0, 1]^d$. It suffices to show that existence of a DEQ computing $-\phi^{(m)}(\mathbf{x}) + 2^m x_1$.

Consider a width- m DEQ with weight matrices as follows:

$$\mathbf{A}^T = \begin{pmatrix} -2^m \\ -2^{m-1} \\ \vdots \\ -2 \end{pmatrix}, \mathbf{W} = \begin{pmatrix} 0 & & & & \\ -4 & 0 & & & \\ -8 & -4 & 0 & & \\ \vdots & \vdots & \vdots & \ddots & \\ -2^m & -2^{m-1} & -2^{m-2} & \dots & 0 \end{pmatrix}, U_1 = \begin{pmatrix} 2 \\ 4 \\ \vdots \\ 2^m \end{pmatrix}, \mathbf{b} = \begin{pmatrix} -1 \\ -1 \\ \vdots \\ -1 \end{pmatrix},$$

406 where U_1 denotes the first column of \mathbf{U} and $\mathbf{U} = (U_1 \ \mathbf{0})$. When $m = 1$, $\mathbf{W} = \mathbf{0}$ and other matrices
407 follow the above expressions. Direct calculations show that the fixed point \mathbf{z} satisfy

$$z_1(\mathbf{x}) = \sigma(2x_1 - 1), \quad z_t(\mathbf{x}) = \sigma\left(-\sum_{i=1}^{t-1} 2^{t-i+1} z_i(\mathbf{x}) + 2^t x_1 - 1\right), \quad \forall 2 \leq t \leq m. \quad (16)$$

408 Note that $\{\phi^{(m)}(\mathbf{x})\}$ admits a recursive expression:

$$\phi^{(m+1)}(\mathbf{x}) = 2\phi^{(m)}(\mathbf{x}) - 2\sigma(2\phi^{(m)}(\mathbf{x}) - 1), \quad \forall m \geq 0, \quad (17)$$

409 for $\phi^{(0)}(\mathbf{x}) := x_1$. We now show by induction that $z_t(\mathbf{x}) = \sigma(2\phi^{(t-1)}(x) - 1)$ for all $1 \leq t \leq m$.
410 When $t = 1$, it is true immediately from Eq. (16) and 17. Assume it is true for some $t < m$, then by
411 Eq. (16) we have

$$\begin{aligned} z_{t+1}(\mathbf{x}) &= \sigma\left(\sum_{i=1}^t -2^{t-i+2} z_i(x) + 2^{t+1} x_1 - 1\right) \\ &= \sigma\left(\sum_{i=1}^t -2^{t-i+2} \sigma(2\phi^{(i-1)}(x) - 1) + 2^{t+1} x_1 - 1\right) \\ &= \sigma\left(\sum_{i=1}^t -2^{t-i+2} \left(\phi^{(i-1)}(x) - \frac{\phi^{(i)}(x)}{2}\right) + 2^{t+1} x_1 - 1\right) \\ &= \sigma(-2^{t+1} \phi^{(0)}(x) + 2\phi^{(t)}(x) + 2^{t+1} x_1 - 1) = \sigma(2\phi^{(t)}(x) - 1), \end{aligned}$$

412 where we use the induction in the second line, Eq. (17) in the third line, and $\phi^{(0)}(\mathbf{x}) = x_1$ in the last
413 line. Thus the induction holds.

414 Using the induction and Eq. (17) for the DEQ, we have

$$\begin{aligned}
\mathbf{A} \mathbf{z}(\mathbf{x}) &= \sum_{i=1}^m -2^{m+1-i} z_i(\mathbf{x}) \\
&= \sum_{i=1}^m -2^{m+1-i} \sigma(2\phi^{(i-1)}(\mathbf{x}) - 1) \\
&= \sum_{i=1}^m -2^{m+1-i} \left(\phi^{(i-1)}(\mathbf{x}) - \frac{\phi^{(i)}(\mathbf{x})}{2} \right) \\
&= -2^m \phi^{(0)}(\mathbf{x}) + \phi^{(m)}(\mathbf{x}) = \phi^{(m)}(\mathbf{x}) - 2^m x_1,
\end{aligned}$$

415 and the lemma follows. \square

416 In prove the theorem, we also need the following lemma which is proved in [18].

417 **Lemma 5** (Lemma 2.1 in [18]). *Let $k \in \mathbb{N}^+$, $L \geq 2$ and $\rho(x) : \mathbb{R} \rightarrow \mathbb{R}$ be a piecewise affine linear*
418 *function with p pieces. Then every $f : \mathbb{R} \rightarrow \mathbb{R}$ implemented by an FNN with depth L , width k and*
419 *activation function ρ has at most $(pk)^{L-1}$ linear regions.*

420 Note that in Lemma 4, the function computed by DEQ is a variant of the saw-tooth function that
421 has many linear regions. On the other hand, Lemma 5 provides an upper bound on the number of
422 linear regions generated by FNN. Combining these two lemmas and using a technique similar to that
423 in Theorem 1.1, [22], we are able to prove Theorem 1.

Proof of Theorem 1. Let $N_d(\mathbf{x})$ be the DEQ in Lemma 4 that computes $2^m x_1 - \phi^{(m)}(\mathbf{x})$ and denote
the width of the FNN that computes $N_f(x)$ by k . For any $\mathbf{y} \in [0, 1]^{d-1}$, define $p_{\mathbf{y}}(x) : [0, 1] \rightarrow$
 $[0, 1]^d$ as $p_{\mathbf{y}} = (x_1, \mathbf{y})$. Then for $N_f \circ p_{\mathbf{y}}(x)$, by Lemma 5, the number of linear regions is upper
bounded by

$$(pk)^{L-1} \leq 2^{(m^{1-\alpha}-1)(L-1)} \leq 2^{m-2},$$

424 where $p = 2$ denotes the number of linear pieces of ReLU activation function. Therefore, $N_f \circ$
425 $p_{\mathbf{y}}(x) - 2^m x$ has at most 2^{m-2} linear regions on $[0, 1]$.

Note that $\phi^{(m)}(\mathbf{x})$ only depends on the first entry of \mathbf{x} , for simplicity, we define $\varphi^{(m)}(x) : \mathbb{R} \rightarrow \mathbb{R}$
as $\varphi^{(m)}(x) = \phi^{(m)} \circ p_{\mathbf{y}}(\mathbf{x})$. Now we claim that there exists at least $3 \cdot 2^{m-3} - 2$ small intervals
 $\{T_l\}_{l=1}^{2^{m-1}}$ with $\text{diam}(T_l) = 2^{-m}$, such that for any \mathbf{y} , it holds

$$\text{sgn} \left(\varphi^{(m)}(x) - \frac{1}{2} \right) \neq \text{sgn} \left(N_f \circ p_{\mathbf{y}}(x) - 2^m x - \frac{1}{2} \right), \quad \forall x \in T_l, \quad \forall l.$$

For simplicity, denote $\tilde{\varphi}(x) = \varphi^{(m)}(x) - \frac{1}{2}$ and $\tilde{N}_f(x) = N_f \circ p_{\mathbf{y}}(x) - 2^m x - \frac{1}{2}$. Denote \mathcal{P}_ϕ and \mathcal{P}_N
the partitions of $[0, 1]$ into intervals so that $\text{sgn}(\tilde{\varphi}(x) - \frac{1}{2})$ and $\text{sgn}(\tilde{N}_f(x) - 2^m x)$ remains
constant within each interval, respectively. Let \mathcal{I}_ϕ be the set of all intervals partitioned by \mathcal{P}_ϕ and
 \mathcal{I}_N be the set of all intervals partitioned by \mathcal{P}_N . By definition, $|\mathcal{I}_\phi| = 2^m + 1$. Since $\tilde{N}_f(x)$ has at
most 2^{m-2} linear regions, the number of the boundary points of the intervals in \mathcal{I}_N is upper bounded
 $2^{m-2} + 1$. So there are at least $3 \cdot 2^{m-2}$ intervals in \mathcal{I}_ϕ that do not intersect with any boundary points
of intervals, i.e., lie completely in an interval in \mathcal{I}_N . Denote this set of intervals by \mathcal{I}'_ϕ . On the other
hand, for every $J \in \mathcal{I}_N$ that contains i_J intervals in \mathcal{I}'_ϕ , there will be $\frac{i_J+1}{2}$ intervals when i_J is odd
and $\frac{i_J}{2}$ intervals when i_J is even, on which $\text{sgn}(\tilde{\varphi}(x)) = \text{sgn}(\tilde{N}_f(x))$. Note that $\sum_{J \in \mathcal{I}_N} i_J = 2^m + 1$.
Therefore, among the sets in \mathcal{I}'_ϕ , the number of sets on which $\text{sgn}(\tilde{\varphi}(x)) \neq \text{sgn}(\tilde{N}_f(x))$ is at least

$$3 \cdot 2^{m-2} - \sum_{J \in \mathcal{I}_N} \frac{i_J + 1}{2} \geq 2^{m-3}.$$

426 Note that except for two intervals, every $T \in \mathcal{I}'_\phi$ can be represented as $[\frac{4i+1}{2^{m+1}}, \frac{4i+3}{2^{m+1}}]$ or $[\frac{4i-1}{2^{m+1}}, \frac{4i+1}{2^{m+1}}]$
427 for some i , thus $\text{diam}(T_l) = 2^{-m}$, which proves the claim. Moreover, on each T_l , direct calculations
428 show $\int_{T_l} |\phi^{(m)}(x) - \frac{1}{2}| dx \geq 2^{-m-2}$.

429 Therefore, by using the claim, we have

$$\begin{aligned}
& \int_{[0,1]^d} |N_d(\mathbf{x}) - N_f(\mathbf{x})| d\mathbf{x} \\
&= \int_{[0,1]^{d-1}} \int_{[0,1]} \left| 2^m x_1 - \phi^{(m)}(x_1) - N_f \circ p_{\mathbf{y}}(x_1) \right| dx_1 d\mathbf{y} \\
&\geq \int_{[0,1]^{d-1}} \int_{\cup_i T_i} \left| 2^m x_1 - \phi^{(m)}(x_1) - N_f \circ p_{\mathbf{y}}(x_1) \right| dx_1 d\mathbf{y} \\
&\geq \int_{[0,1]^{d-1}} \int_{\cup_i T_i} \left| \frac{1}{2} - \phi^{(m)}(x_1) \right| dx_1 d\mathbf{y} \\
&\geq \int_{[0,1]^{d-1}} |T_i| \cdot 2^{-m-2} d\mathbf{y} \\
&\geq (3 \cdot 2^{m-3} - 2) \cdot 2^{-m-2} \geq \frac{1}{16}.
\end{aligned}$$

430

□

431 Next we turn to proof Proposition 1. We use $\text{Diag}(\cdot)$ to extract the diagonal elements of a matrix into
432 a vector. The proof of Proposition 1 relies on following explicit expression of ReLU DEQ.

433 **Lemma 6.** *Let $\mathbf{W}, \mathbf{U}, \mathbf{b}$ be the weights of a DEQ with $\|\mathbf{W}\|_2 < 1$. Then for any $\mathbf{x} \in \mathbb{R}^d$, there
434 exists a diagonal matrix $\mathbf{D} \in \mathbb{R}^{d \times d}$ whose diagonal entries are either 1 or 0, such that*

$$\text{sgn}(\text{diag}((\mathbf{I} - \mathbf{W}\mathbf{D})^{-1})(\mathbf{U}\mathbf{x} + \mathbf{b})) = \text{Diag}(\mathbf{D}). \quad (18)$$

435 Moreover, fix \mathbf{D} , for all \mathbf{x} that Eq. (18) holds, we have

$$\mathbf{z}(\mathbf{x}) = (\mathbf{I} - \mathbf{D}\mathbf{W})^{-1} \mathbf{D}(\mathbf{U}\mathbf{x} + \mathbf{b}). \quad (19)$$

436 *Proof.* Recall that the fixed point $\mathbf{z}(\mathbf{x})$ satisfies

$$\mathbf{z} = \sigma(\mathbf{W}\mathbf{z} + \mathbf{U}\mathbf{x} + \mathbf{b}). \quad (20)$$

437 For each z_i , if the i -th entry of $(\mathbf{W}\mathbf{z} + \mathbf{U}\mathbf{x} + \mathbf{b})$ is smaller than 0, then $z_i = 0$. Without loss of
438 generality, we assume that the first t ($t \leq m$) entries of $(\mathbf{W}\mathbf{z} + \mathbf{U}\mathbf{x} + \mathbf{b})$ are greater than 0, and the
439 rest $m - t$ entries are smaller than 0. Denote by $\mathbf{v} = \mathbf{U}\mathbf{x} + \mathbf{b}$ and the corresponding block matrices
440 $\mathbf{z}, \mathbf{W}, \mathbf{v}$ by

$$\mathbf{z} = \begin{pmatrix} \tilde{\mathbf{z}} \\ \mathbf{0} \end{pmatrix}, \mathbf{W} = \begin{pmatrix} \mathbf{W}_{11} & \mathbf{W}_{12} \\ \mathbf{W}_{21} & \mathbf{W}_{22} \end{pmatrix}, \mathbf{v} = \begin{pmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \end{pmatrix}, \quad (21)$$

441 where $\tilde{\mathbf{z}} \in \mathbb{R}^t$, $\mathbf{W}_{11} \in \mathbb{R}^{t \times t}$, and $\mathbf{v}_1 \in \mathbb{R}^t$. Then, Eq.(20) is equivalent to

$$\tilde{\mathbf{z}} = \mathbf{W}_{11}\tilde{\mathbf{z}} + \mathbf{v}_1, \quad \mathbf{W}_{21}\tilde{\mathbf{z}} + \mathbf{v}_2 \leq 0, \quad \tilde{\mathbf{z}} > 0. \quad (22)$$

Now we define $\mathbf{D} = \begin{pmatrix} \mathbf{I}_t & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix}$ and show that it is the desired matrix. Note that $\|\mathbf{W}\|_2 < 1$ and
 $\|\mathbf{D}\|_2 = 1$, we have

$$\|\mathbf{W}_{11}\|_2 = \|\mathbf{W}\mathbf{D}\|_2 \leq \|\mathbf{W}\|_2 \|\mathbf{D}\|_2 < 1,$$

442 showing that $\mathbf{I}_t - \mathbf{W}_{11}$ is invertible. Thus Eq.(22) gives

$$\tilde{\mathbf{z}} = (\mathbf{I}_t - \mathbf{W}_{11})^{-1} \mathbf{v}_1 > 0, \quad \mathbf{W}_{21}(\mathbf{I}_t - \mathbf{W}_{11})^{-1} \mathbf{v}_1 + \mathbf{v}_2 \leq 0. \quad (23)$$

443 Additionally, by simple calculation, we have

$$\begin{aligned}
(\mathbf{I} - \mathbf{D}\mathbf{W})^{-1} &= \begin{pmatrix} (\mathbf{I}_t - \mathbf{W}_{11})^{-1} & (\mathbf{I}_t - \mathbf{W}_{11})^{-1} \mathbf{W}_{12} \\ \mathbf{0} & \mathbf{I} \end{pmatrix}, \\
(\mathbf{I} - \mathbf{W}\mathbf{D})^{-1} &= \begin{pmatrix} (\mathbf{I}_t - \mathbf{W}_{11})^{-1} & \mathbf{0} \\ \mathbf{W}_{21}(\mathbf{I}_t - \mathbf{W}_{11})^{-1} & \mathbf{I} \end{pmatrix}.
\end{aligned} \quad (24)$$

444 Combining Eq. (21), (23) and (24), we have

$$\begin{aligned} (\mathbf{I} - \mathbf{W}\mathbf{D})^{-1}(\mathbf{U}\mathbf{x} + \mathbf{b}) &= \begin{pmatrix} \mathbf{W}_{11}\tilde{\mathbf{z}} + \mathbf{v}_1 \\ \mathbf{W}_{21}(\mathbf{I}_t - \mathbf{W}_{11})^{-1}\mathbf{v}_1 + \mathbf{v}_2 \end{pmatrix}, \\ \mathbf{z} &= \begin{pmatrix} (\mathbf{I}_t - \mathbf{W}_{11})^{-1}\mathbf{v}_1 \\ \mathbf{0} \end{pmatrix} = (\mathbf{I} - \mathbf{D}\mathbf{W})^{-1}\mathbf{D}(\mathbf{U}\mathbf{x} + \mathbf{b}). \end{aligned}$$

445 Finally, since the output of the implicit layer is unique, in the sense of permuting the entries of \mathbf{D} ,
446 there always exists a matrix \mathbf{D} such that the Lemma follows. \square

447 Note that there are at most 2^m diagonal matrix whose diagonal entries are either 1 or 0, the upper
448 bound of the number of linear regions is 2^m . Thus Proposition 1 follows straightforwardly from 6
449 and 4.

450 A.2 Proofs in Subsection 4.2

451 A.2.1 Inapproximability of FNNs

452 The goal of this section is to prove the following proposition, which is an extended version of the
453 inapproximability result in Theorem 2.

454 **Assumption 2.** *The activation function $\tilde{\sigma}$ is of $C^0(\mathbb{R})$ and continuous differentiable except for at
455 most finitely many points. And there exists an absolute constant $C_{\tilde{\sigma}} > 0$, such that $|\tilde{\sigma}'(x)| \leq C_{\tilde{\sigma}}$ for
456 all x on which $\tilde{\sigma}$ is differentiable.*

Proposition 3 (Inapproximability of FNN). *Let $N_{\text{fnn}}(x)$ be computed by an FNN with depth L , width
 k , and an activation function $\tilde{\sigma}$ satisfying Assumption 2 on $\mathbf{x} \in [0, 1]^d$. Let $g(\mathbf{x})$ be defined as in
Eq.(6), and $\frac{1}{4} \geq \varepsilon > 0$. If $\|N_{\text{fnn}}(\mathbf{x}) - g(\mathbf{x})\|_{L^\infty([0,1]^d)} \leq \varepsilon$, then there exists a weight parameter W_{ij}
of the FNN for $1 \leq i \leq L$ and $1 \leq j \leq k$, such that*

$$|W_{ij}| \geq \frac{1}{C_{\tilde{\sigma}}k} \cdot 2^{\frac{d-4}{L}}.$$

457 *Proof.* By assumption, $N_{\text{fnn}}(\mathbf{x})$ is of $C(\mathbb{R}^d)$ and continuous differentiable except for at most finitely
458 many points, then by the intermediate value theorem, we have

$$\max_{x \in [0,1]^d} \left| \frac{\partial N_{\text{fnn}}(\mathbf{x})}{\partial x_1} \right| \geq \left| \frac{g_1(1) - g_1(1 - \delta)}{\delta} \right| \geq \frac{\frac{1}{2} - \frac{\delta(1-\delta)}{4\delta} - 2 \cdot \frac{1}{16}}{\delta} \geq \frac{1}{8\delta} - 1 \geq 2^{d-4}, \quad (25)$$

where $\frac{\partial N_{\text{fnn}}(\mathbf{x})}{\partial x_1}$ refers to the subgradient on the non-differentiable points. Additionally, by definition,

$$N_{\text{fnn}}(\mathbf{x}) = \mathbf{W}_L \tilde{\sigma}(\mathbf{W}_{L-1} \tilde{\sigma}(\cdots \tilde{\sigma}(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) \cdots)) + \mathbf{b}_{L-1}$$

459 Then direct calculation gives

$$\nabla N_{\text{fnn}}(\mathbf{x}) = \mathbf{W}_1^T \mathbf{D}_1 \cdots \mathbf{D}_{L-1} \mathbf{W}_L^T, \quad (26)$$

where $\mathbf{D}_l = \text{diag}(\tilde{\sigma}'(W_l \tilde{\sigma}(\cdots \tilde{\sigma}(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) \cdots)) + b_l)$ for $1 \leq l \leq L - 1$. By Assumption 2, it
holds that

$$\|\mathbf{D}_l\|_\infty \leq C_{\tilde{\sigma}}, \quad \forall 1 \leq l \leq L - 1.$$

Then combining Eq. (25) and (26), we have

$$2^{d-4} \leq \|\nabla N_{\text{fnn}}(\mathbf{x})\|_\infty \leq \prod_{i=1}^L \|\mathbf{D}_i \mathbf{W}_i\|_\infty \leq C_{\tilde{\sigma}}^L \cdot \prod_{i=1}^L \|\mathbf{W}_i\|_\infty.$$

Therefore, there exists at least one \mathbf{W}_i for $1 \leq i \leq L$, such that

$$\|\mathbf{W}_i\|_\infty \geq C_{\tilde{\sigma}}^{-1} 2^{\frac{d-4}{L}}.$$

Finally, by the definition of $\|\cdot\|_\infty$, there exists an entry W_{ij} with $1 \leq j \leq k$, such that

$$|W_{ij}| \geq \frac{1}{C_{\tilde{\sigma}}k} \cdot 2^{\frac{d-4}{L}}.$$

460 \square

461 **Remark 4.** Assumption 2 is mild and one can verify that most commonly used activation functions
 462 such as ReLU, GeLU, sigmoid and tanh satisfy the assumption.

To prove the inapproximability of FNNs in Theorem 2, we take $C_{\bar{\sigma}} = 1$ in Proposition 3 as $|\sigma'(x)| \leq 1$ and derive

$$|W_{ij}| \geq k^{-1} \cdot 2^{\frac{d-4}{L}} \leq 2^{-\frac{d}{2\bar{\sigma}} + \frac{d-4}{\bar{\sigma}}} \geq 2^{-\frac{d}{2\bar{\sigma}}},$$

463 which finishes the proof.

464 A.2.2 Approximability of DEQs

465 This section centers around the approximability result of DEQs. We restate the approximability result
 466 of Theorem 2 as the following proposition.

Proposition 4 (Approximability of DEQ). *Let $g(\mathbf{x})$ be defined as in Eq.(6) on $[0, 1]^d$. $\forall \frac{1}{4} \geq \varepsilon > 0$, there exists a DEQ N_{deq} with width bounded by $5\varepsilon^{-1}$ and weights bounded by $2\varepsilon^{-1}$, such that*

$$\|N_{deq}(\mathbf{x}) - g(\mathbf{x})\|_{L^\infty([0,1]^d)} \leq \varepsilon.$$

467 The proof of the proposition requires some intermediate steps regarding the constructing approxima-
 468 tion by DEQ and bounding the fixed-points' error. For simplicity, in the rest of the section, for any
 469 function f which is continuous differentiable except for at most finitely many points, we denote f'
 470 the derivative of f on the differentiable points, and the subgradient of f on the non-differentiable
 471 points.

472 The next lemma considers approximating the square function using a 2-layer FNN.

Lemma 7. *For any $N \in \mathbb{N}^+$, there exists a function ϕ implemented by a 2-layer ReLU FNN with width $2N$ such that*

$$|\phi(x) - x^2| \leq 4N^{-2}, \quad |\phi'(x)| \leq 2 - \frac{1}{N}, \quad \forall x \in [-1, 1].$$

Proof. Denote $\frac{1}{N}$ by t for simplicity. Let $\{x_i\}_{i=1}^{2N+1}$ be $2N + 1$ points on \mathbb{R} defined as follows:

$$x_1 = -1, x_2 = -1 + t, \dots, x_{2N} = 1 - t, x_{2N+1} = 1.$$

473 We consider the following function $\phi(x)$ that interpolates x^2 on all $\{x_i\}_{i=1}^{2N+1}$:

$$\phi(x) = \sigma(tx) + \sigma(-tx) + \sum_{i=1}^{N-1} \sigma(2tx - 2it^2) + \sum_{i=1}^{N-1} \sigma(-2tx + 2it^2). \quad (27)$$

It can be seen that $\phi(x)$ can be implemented by a 2-layer ReLU FNN with width $2N$ and weight bounded by $2t$. By the interpolation property of $\phi(x)$, on every $[x_j, x_{j+1}]$, it holds

$$\max_{x \in [x_j, x_{j+1}]} |\phi(x) - x^2| = \phi\left(\frac{x_j + x_{j+1}}{2}\right) - \left(\frac{x_j + x_{j+1}}{2}\right)^2 = \frac{t^2}{4}.$$

Thus we have $|\phi(x) - x^2| \leq 4N^{-2}$ for all $x \in [-1, 1]$. Moreover, since $\phi(x)$ is convex, we have

$$|\phi'(x)| \leq \frac{1 - (1-t)^2}{t} = 2 - t.$$

474

□

475 We now move to prove the equivalence between the revised DEQ and vanilla DEQ.

Proof of Lemma 1. For any $\hat{\mathbf{z}}^0 \in \mathbb{R}^m$, we define a sequence $\{\hat{\mathbf{z}}^k\}$ as

$$\hat{\mathbf{z}}^{k+1} = \sigma(\mathbf{W} \mathbf{V} \hat{\mathbf{z}}^k + \mathbf{U} \mathbf{x} + \mathbf{b}).$$

Since $\|\mathbf{W} \mathbf{V}\|_2 \leq 1$, $\{\hat{\mathbf{z}}^k\}$ converges and the limit $\hat{\mathbf{z}}^*$ is the fixed point of $\hat{\mathbf{z}} = \sigma(\mathbf{W} \mathbf{V} \hat{\mathbf{z}} + \mathbf{U} \mathbf{x} + \mathbf{b})$. Now we set $\mathbf{z}^0 = \mathbf{V} \mathbf{y}^0$ and define another sequence $\{\mathbf{z}^k\}$ as

$$\mathbf{z}^{k+1} = \mathbf{V} \sigma(\mathbf{W} \mathbf{z}^k + \mathbf{U} \mathbf{x} + \mathbf{b}), \quad \forall k \geq 0.$$

It follows immediately by induction that $\mathbf{z}^k = \mathbf{V}\hat{\mathbf{z}}^k$ for all $k \geq 0$. Note that $\{\mathbf{z}^k\}$ converges and the limit \mathbf{z}^* is exactly the fixed point of the revised DEQ in Eq. (8). Therefore, it holds that

$$\mathbf{z}^* = \lim_{k \rightarrow \infty} \mathbf{z}^k = \lim_{k \rightarrow \infty} \mathbf{V}\hat{\mathbf{z}}^k = \mathbf{V}\hat{\mathbf{z}}^*.$$

476 The desired DEQ is constructed as

$$\begin{aligned}\hat{\mathbf{z}} &= \sigma(\mathbf{W}\mathbf{V}\hat{\mathbf{z}} + \mathbf{U}\mathbf{x} + \mathbf{b}), \\ \hat{\mathbf{y}} &= (\mathbf{B}\mathbf{V})\hat{\mathbf{z}}\end{aligned}$$

477

□

478 In the following we turn to bound the error between the equilibria of two fixed-point equations. We
479 start with the proof of Lemma 2.

480 *Proof of Lemma 2.* The existence of z_u and z_v follows from the fixed point theorem since $u(\cdot, \mathbf{x})$ and
481 $v(\cdot, \mathbf{x})$ are contraction mappings. For simplicity, we denote $u_{\mathbf{x}}(z) = u(z, \mathbf{x})$ and $v_{\mathbf{x}}(z) = v(z, \mathbf{x})$.
482 Note that the range of $u_{\mathbf{x}}$ and $v_{\mathbf{x}}$ are in Ω . Then $\forall n \in \mathbb{N}^+$, we have

$$\begin{aligned}\|u_{\mathbf{x}}^{\circ n} - v_{\mathbf{x}}^{\circ n}\| &\leq \left\| u_{\mathbf{x}}^{\circ n} - u_{\mathbf{x}}\left(v_{\mathbf{x}}^{\circ(n-1)}\right) \right\| + \left\| u_{\mathbf{x}}\left(v_{\mathbf{x}}^{\circ(n-1)}\right) - v_{\mathbf{x}}^{\circ n} \right\| \\ &\leq L_u \left\| u_{\mathbf{x}}^{\circ(n-1)} - v_{\mathbf{x}}^{\circ(n-1)} \right\| + \|u_{\mathbf{x}} - v_{\mathbf{x}}\| \\ &\leq L_u \left(\left\| u_{\mathbf{x}}^{\circ(n-2)} - v_{\mathbf{x}}^{\circ(n-2)} \right\| + \|u_{\mathbf{x}} - v_{\mathbf{x}}\| \right) + \|u_{\mathbf{x}} - v_{\mathbf{x}}\| \\ &\leq \dots \\ &\leq (1 + L_u + \dots + L_u^{n-1}) \|u_{\mathbf{x}} - v_{\mathbf{x}}\| \\ &= \frac{1 - L_u^n}{1 - L_u} \|u_{\mathbf{x}} - v_{\mathbf{x}}\|.\end{aligned}$$

483 By definition, $\forall (z, \mathbf{x}) \in \Omega \times [0, 1]^d$, $z_u(\mathbf{x}) = \lim_{n \rightarrow \infty} u_{\mathbf{x}}^{\circ n}(z)$, and $z_v(\mathbf{x}) = \lim_{n \rightarrow \infty} v_{\mathbf{x}}^{\circ n}(z)$. Hence,
484 we have

$$\begin{aligned}|z_u(\mathbf{x}) - z_v(\mathbf{x})| &\leq \lim_{n \rightarrow \infty} |u_{\mathbf{x}}^{\circ n}(z) - v_{\mathbf{x}}^{\circ n}(z)| \\ &\leq \lim_{n \rightarrow \infty} \frac{1 - L_u^n}{1 - L_u} |u(z, \mathbf{x}) - v(z, \mathbf{x})| \\ &\leq \frac{1}{1 - L_u} |u(z, \mathbf{x}) - v(z, \mathbf{x})|.\end{aligned}$$

485 Finally, by the symmetry of u and v , we also have $|z_u(\mathbf{x}) - z_v(\mathbf{x})| \leq \frac{1}{1 - L_v} |u(z, \mathbf{x}) - v(z, \mathbf{x})|$. The
486 proof is finished. □

487 We also need Lemma 3 to bound the error.

488 *Proof of Lemma 3.* We use the intermediate value theorem to proof the lemma. Define $q(z, \mathbf{x}) =$
489 $z - v(z, \mathbf{x})$. The fixed point $z_v(\mathbf{x})$ is unique zero of $q(z, \mathbf{x}) = 0$. Since $v(z, \mathbf{x})$ is L_v Lipschitz with
490 respect to z and $L_v < 1$, $q(z, \mathbf{x})$ is monotonically increasing with respect to z for all \mathbf{x} .

491 Fix z_u , the proof proceeds by discussing the following 2 cases:

- 492 • If $q(z_u, \mathbf{x}) \leq 0$, i.e., $u(z_u, \mathbf{x}) = z_u \leq v(z_u, \mathbf{x})$, we consider $T = [z_u, z_u + \xi] \subset \Omega$.
493 By assumption, there exists $z^* \in T$, such that $u(z^*, \mathbf{x}) = v(z^*, \mathbf{x})$. Note that $q(\cdot, \mathbf{x})$ is
494 monotonically increasing, thus we have $q(z^*, \mathbf{x}) \geq 0$. By the continuity of $q(z, \mathbf{x})$ w.r.t.
495 z and the intermediate value theorem, $q(z, \mathbf{x})$ must have a zero in $[z_u, z_0] \subset T$, which is
496 $z_v(\mathbf{x})$ by definition. Hence, it holds that $|z_u - z_v| \leq \xi$.
- 497 • If $q(z_u, \mathbf{x}) \geq 0$, i.e., $u(z_u, \mathbf{x}) = z_u \geq v(z_u, \mathbf{x})$, we consider $T = [z_u - \xi, z_u] \subset \Omega$. It
498 follows from similar deductions that $|z_u - z_v| \leq \xi$ in this case.

499 We finish the proof. \square

500 With the results above, we begin our formal proof of Proposition 4. The proof is sketched as follows:
 501 First, we consider a fixed point equation $z = \tilde{g}(z, \mathbf{x})$ that induce the target function $g(\mathbf{x})$. We show
 502 that there exists a function $\tilde{h}(z, \mathbf{x}) : \mathbb{R}^{d+1} \rightarrow \mathbb{R}$ computed by a 2-layer FNN with width $O(\varepsilon^{-1})$ that
 503 can approximate $\tilde{g}(z, \mathbf{x})$ in sup-norm to an accuracy of $O(\varepsilon^2)$. Moreover, $z = \tilde{h}(z, \mathbf{x})$ is a well-posed
 504 fixed point equation and induces a revised DEQ. Second, we bound the error between $g(\mathbf{x})$ and $h(\mathbf{x})$,
 505 where $h(\mathbf{x})$ is the fixed point of $z = \tilde{h}(z, \mathbf{x})$. The proof is further divided into two parts: When
 506 $1 - x_1 > \frac{\varepsilon}{2}$, by using Lemma 2, we can bound the error $\|h - g\|$ by $\varepsilon \cdot \|\tilde{h} - \tilde{g}\|$. When $1 - x_1 < \frac{\varepsilon}{2}$,
 507 we show that the conditions of Proposition 4 holds for $\xi = \varepsilon$, thus $\|h - g\|$ is upper bounded by ε .

Proof of Proposition 4. Let $g(\mathbf{x})$ be defined as in Eq.(6). Recall that $g(\mathbf{x})$ is the fixed point of the fixed point equation

$$z = \tilde{g}(z, \mathbf{x}) := zx_1 + \delta \left(\frac{x_1}{2} - z \right).$$

508 **Approximate \tilde{g} using 2-layer FNN.** By Lemma 7, $\forall N \in \mathbb{N}^+$, there exist $\mathbf{a} \in \mathbb{R}^{2N}$, $\tilde{\mathbf{b}} \in \mathbb{R}^{2N}$, $\tilde{\mathbf{W}} \in$
 509 \mathbb{R}^{2N} and a function $\phi(x) = \mathbf{a}^T \sigma(\tilde{\mathbf{W}}x + \tilde{\mathbf{b}})$, such that for all $x \in [-1, 1]$, it holds

$$|\phi_N(x) - x^2| \leq 4N^{-2}, \quad |\phi'_N(x)| \leq 2 - \frac{1}{N}. \quad (28)$$

Now, we define

$$\tilde{h}(z, \mathbf{x}) = \frac{1}{2} \left[\phi_N \left(z + \frac{x_1}{2} \right) - \phi_N \left(z - \frac{x_1}{2} \right) \right] + \delta \left(\frac{x_1}{2} - z \right).$$

510 1. $\tilde{h}(z, \mathbf{x})$ can be implemented by a 2-layer ReLU FNN with width $4N + 2$ for $(z, \mathbf{x}) \in$
 511 $[-\delta, \frac{1}{2}] \times [0, 1]^d$. To see this, when $(z, \mathbf{x}) \in [-\delta, \frac{1}{2}] \times [0, 1]^d$, it holds $z + \frac{x_1}{2} \in [0, 1]$,
 512 $z - \frac{x_1}{2} \in [-1, 0]$. Then

$$\begin{aligned} & \begin{pmatrix} \frac{\mathbf{a}^T}{2} & \frac{\mathbf{a}^T}{2} & -\delta & \delta \end{pmatrix} \sigma \left(\begin{pmatrix} \tilde{\mathbf{W}} & \mathbf{0} \\ \mathbf{0} & \tilde{\mathbf{W}} \\ 0 & 1 \\ 0 & -1 \end{pmatrix} \begin{pmatrix} 1 & \frac{1}{2} & \mathbf{0} \\ 1 & -\frac{1}{2} & \mathbf{0} \end{pmatrix} \begin{pmatrix} z \\ x_1 \\ \mathbf{x}_{-1} \end{pmatrix} + \begin{pmatrix} \tilde{\mathbf{b}} \\ \tilde{\mathbf{b}} \\ 0 \\ 0 \end{pmatrix} \right) \\ &= \begin{pmatrix} \frac{\mathbf{a}^T}{2} & \frac{\mathbf{a}^T}{2} & -\delta & \delta \end{pmatrix} \sigma \left(\begin{pmatrix} \tilde{\mathbf{W}} \left(z + \frac{x_1}{2} \right) + \tilde{\mathbf{b}} \\ \tilde{\mathbf{W}} \left(z - \frac{x_1}{2} \right) + \tilde{\mathbf{b}} \\ \left(z - \frac{x_1}{2} \right) \\ - \left(z - \frac{x_1}{2} \right) \end{pmatrix} \right) \\ &= \frac{1}{2} \mathbf{a}^T \sigma \left(\tilde{\mathbf{W}} \left(z + \frac{x_1}{2} \right) + \tilde{\mathbf{b}} \right) + \frac{1}{2} \mathbf{a}^T \sigma \left(\tilde{\mathbf{W}} \left(z - \frac{x_1}{2} \right) + \tilde{\mathbf{b}} \right) \\ &\quad + \delta \left(\sigma \left(-z + \frac{x_1}{2} - \sigma \left(z - \frac{x_1}{2} \right) \right) \right) \\ &= \frac{1}{2} \left[\phi_N \left(z + \frac{x_1}{2} \right) - \phi_N \left(z - \frac{x_1}{2} \right) \right] + \delta \left(\frac{x_1}{2} - z \right) = \tilde{h}(z, \mathbf{x}), \end{aligned} \quad (29)$$

513 where the first line resembles a function implemented by an FNN with width $4N + 2$.

514 2. $\tilde{h}(z, \mathbf{x})$ approximate $\tilde{g}(z, \mathbf{x})$ well on $(z, \mathbf{x}) \in [-\delta, \frac{1}{2}] \times [0, 1]^d$. Since $z + \frac{x_1}{2} \in [0, 1]$ and
 515 $z - \frac{x_1}{2} \in [-1, 0]$, from Eq. (28), we have

$$\begin{aligned} |\tilde{h}(z, \mathbf{x}) - \tilde{g}(z, \mathbf{x})| &= \frac{1}{2} \left[\phi_N \left(z + \frac{x_1}{2} \right) - \left(z + \frac{x_1}{2} \right)^2 \right] - \frac{1}{2} \left[\phi_N \left(z - \frac{x_1}{2} \right) - \left(z - \frac{x_1}{2} \right)^2 \right] \\ &\leq \frac{1}{2} \left| \phi_N \left(z + \frac{x_1}{2} \right) - \left(z + \frac{x_1}{2} \right)^2 \right| + \frac{1}{2} \left| \phi_N \left(z - \frac{x_1}{2} \right) - \left(z - \frac{x_1}{2} \right)^2 \right| \\ &\leq \frac{1}{2} \left(\frac{t^2}{4} + \frac{t^2}{4} \right) = \frac{t^2}{4}. \end{aligned} \quad (30)$$

516 3. The fixed point equation $z = \tilde{h}(z, \mathbf{x})$ is well-posed on $[-\delta, \frac{1}{2}] \times [0, 1]^d$. for the partial
 517 derivative $\frac{\partial \tilde{h}(z, \mathbf{x})}{\partial z}$, we have

$$\begin{aligned} \left| \frac{\partial \tilde{h}(z, \mathbf{x})}{\partial z} \right| &= \frac{1}{2} \left(\phi'_N \left(z + \frac{x_1}{2} \right) - \phi'_N \left(z - \frac{x_1}{2} \right) \right) - \delta \\ &\leq \frac{1}{2} \left(\phi'_N \left(z + \frac{x_1}{2} \right) - \phi'_N \left(z + \frac{x_1}{2} - 1 \right) \right) - \delta \\ &\leq 1 - \delta < 1, \end{aligned}$$

518 where the second line holds because $\phi'(x)$ is monotonically increasing and $x_1 < 1$. There-
 519 fore, the fixed point equation $z = \tilde{h}(z, \mathbf{x})$ has a unique solution for all \mathbf{x} .

4. Note that $\tilde{h}(z, \mathbf{x})$ can be computed by a revised DEQ defined in Eq. (8) with

$$\mathbf{V} = \begin{pmatrix} \frac{\mathbf{a}^T}{2} & \frac{\mathbf{a}^T}{2} & -\delta & \delta \end{pmatrix}, \quad \mathbf{W} = \begin{pmatrix} \tilde{\mathbf{W}} \\ \tilde{\mathbf{W}} \\ 1 \\ -1 \end{pmatrix}, \quad \mathbf{B} = \mathbf{1}.$$

520 And it can be verified that $\|\mathbf{W}\mathbf{V}\|_2 = 1 - t - 2\delta \leq 1$. By Lemma 1, the fixed point of
 521 $z = \tilde{h}(z, \mathbf{x})$ can be computed by a DEQ with width $4N + 2$, which we denote by $N_{\text{deq}}(\mathbf{x})$.
 522 Further calculations shows that the weight of the DEQ is also bounded by $2t$.

523 **Approximate g using the induced DEQ.** We will bound $\|N_{\text{deq}}(\mathbf{x}) - g(\mathbf{x})\|_{L^\infty([0,1]^d)}$ using Lemma
 524 2 and Lemma 3. Let $\Omega = [-\delta, \frac{1}{2}]$ and assume that $t > 10\delta$. It can be easily verified that both the
 525 range of $\tilde{g}(z, \mathbf{x})$ and $\tilde{h}(z, \mathbf{x})$ are in Ω when $(z, \mathbf{x}) \in \Omega \times [0, 1]^d$.

1. When $x_1 \leq 1 - \frac{t}{2}$, by definition, the Lipschitz constant of $\tilde{g}(\cdot, \mathbf{x})$ is upper bounded by
 $\max \left| \frac{\partial \tilde{g}(z, \mathbf{x})}{\partial z} \right|$. Leveraging Lemma 2 and Eq.(30), we have

$$|N_{\text{deq}}(\mathbf{x}) - g(\mathbf{x})| \leq \left| 1 - \frac{\partial \tilde{g}(z, \mathbf{x})}{\partial z} \right|^{-1} |\tilde{h}(z, \mathbf{x}) - \tilde{g}(z, \mathbf{x})| \leq \frac{2}{t} \cdot \frac{t^2}{4} = \frac{t}{2}.$$

2. When $1 > x_1 > 1 - \frac{t}{2}$, if $z + \frac{x_1}{2} = nt$ for some $\frac{N}{2} - 1 \leq n \leq N$, we have

$$z - \frac{x_1}{2} = nt - \frac{x_1}{2} \in \left(\left(n - \frac{N}{2} \right) t, \left(n - \frac{N}{2} - 1 \right) t \right).$$

Note that $\phi_N(x) > x^2$ for all $x \in [0, 1] \setminus t\mathbb{N}$ and $\phi_N(x) = x^2$ for all $x \in [0, 1] \cap t\mathbb{N}$. Thus,
 when $z = nt - \frac{x_1}{2}$, we have

$$\tilde{h}(z, \mathbf{x}) < \frac{1}{2} \left(\left(z + \frac{x_1}{2} \right)^2 - \left(z - \frac{x_1}{2} \right)^2 \right) + \delta \left(\frac{x_1}{2} - z \right) = \tilde{g}(z, \mathbf{x}).$$

526 Note that for every $T \subset \Omega$ with $|T| \leq t$, there exists $z_g \in T$, such that $z_g = nt - \frac{x_1}{2}$ and
 527 thus $\tilde{h}(z_g, \mathbf{x}) < \tilde{g}(z_g, \mathbf{x})$.

On the other hand, if $z = \frac{x_1}{2} - kt$ for some $0 \leq k \leq \frac{N}{2} - 1$, we have

$$z + \frac{x_1}{2} = kt + \frac{x_1}{2} \in \left(\left(-k + \frac{N}{2} - 1 \right) t, \left(-k + \frac{N}{2} \right) t \right).$$

Similarly, we have

$$\tilde{h}(z, \mathbf{x}) > \frac{1}{2} \left(\left(z + \frac{x_1}{2} \right)^2 - \left(z - \frac{x_1}{2} \right)^2 \right) + \delta \left(\frac{x_1}{2} - z \right) = \tilde{g}(z, \mathbf{x}).$$

Note that for every $T \subset \Omega$ with $|T| \leq t$, there exists $z_l \in T$, such that $z_l = \frac{x_1}{2} - kt$ and
 thus $\tilde{h}(z_l, \mathbf{x}) > \tilde{g}(z_l, \mathbf{x})$. From the intermediate value theorem, there exists $z^* \in T$, such
 that $\tilde{h}(z^*, \mathbf{x}) = \tilde{g}(z^*, \mathbf{x})$. Thus it follows from Lemma 3 immediately that

$$|N_{\text{deq}}(\mathbf{x}) - g(\mathbf{x})| \leq t.$$

Additionally, when $x_1 = 1$, by simple calculations, we have $\tilde{h}(\frac{1}{2}, \mathbf{x}) = \tilde{g}(\frac{1}{2}, \mathbf{x}) = \frac{1}{2}$, indicating that $N_{\text{deq}}(\mathbf{x}) = g(\mathbf{x}) = \frac{1}{2}$. Combining all the results above, we have

$$|N_{\text{deq}}(\mathbf{x}) - g(\mathbf{x})| \leq |g(\mathbf{x}) - \bar{z}'| \leq t, \quad \mathbf{x} \in [0, 1]^d.$$

528 By choosing $t = \varepsilon$, we finish the proof.

529

□

530 A.3 Proofs in Section 5

531 We start with the proof of Theorem 3.

532 *Proof of Theorem 3.* Denote $\{\beta(t)\}$ the process that follows the gradient flow dynamics $\dot{\mathbf{w}}(t) =$
 533 $-\nabla_{\mathbf{w}} \hat{L}(\mathbf{w}(t))$ initialized by $\beta(0) > 0$. Recall that the empirical loss is $\frac{1}{2} \|\mathbf{X}\beta - \mathbf{y}\|_2^2$, then the
 534 dynamics of $\{\beta(t)\}$ can be computed as follows:

$$\begin{aligned} \frac{d\beta(t)}{dt} &= \nabla_{\mathbf{w}} \beta(t) \cdot \frac{dw(t)}{dt} \\ &= \nabla_{\mathbf{w}} \beta(t) \cdot \left(-\nabla_{\mathbf{w}} \left(\frac{1}{2} \|\mathbf{X}\beta(t) - \mathbf{y}\|_2^2 \right) \right) \\ &= \nabla_{\mathbf{w}} \beta(t)^2 \cdot \left(-\nabla_{\beta} \left(\frac{1}{2} \|\mathbf{X}\tilde{\beta}(t) - \mathbf{y}\|_2^2 \right) \right) \\ &= -(\mathbf{X}^T \mathbf{r}(t)) \odot \tilde{\beta}(t)^{\odot 4}, \end{aligned} \quad (31)$$

535 where $\mathbf{r}(t) = \mathbf{X}\beta(t) - \mathbf{y}$ denotes the residual. For any $t > 0$, it can be verified easily from Eq. (31)
 536 that

$$-\frac{1}{3}\beta(t)^{\odot -3} + \frac{1}{3}\beta(0)^{\odot -3} = -\mathbf{X} \int_0^t \mathbf{r}(s) ds. \quad (32)$$

537 For simplicity, we denote $\mathbf{v}(t) = \int_0^t \mathbf{r}(s) ds$. Then from Eq. (32), we have

$$\beta(t) = (3\mathbf{X}^T \mathbf{v}(t) + \beta(0)^{\odot -3})^{\odot -\frac{1}{3}} \quad (33)$$

538 By assumption, $\beta(t)$ converges to some $\beta^\infty \in \mathbb{R}^d$ when $t \rightarrow \infty$, thus $\mathbf{v}(t)$ converges to some
 539 $\mathbf{v}^\infty := \int_0^\infty \mathbf{r}(s) ds$. By letting $t \rightarrow \infty$ in Eq. (33), we have

$$\beta^\infty = (3\mathbf{X}^T \mathbf{v}^\infty + \beta(0)^{\odot -3})^{\odot -\frac{1}{3}}. \quad (34)$$

Next we want to show that β^∞ satisfies the KKT condition of the optimization problem in Eq. (12). Given access to the expression of $Q(\beta)$, the KKT optimality conditions can be expressed as

$$\mathbf{X}\beta^* = \mathbf{y}, \quad \nabla Q(\beta^*) = \mathbf{X}\mathbf{v},$$

540 for some $\mathbf{v} \in \mathbb{R}^d$. By the definition of $Q(\beta)$, $\nabla Q(\beta^*) = \mathbf{X}^T \mathbf{v}$ is equivalent to

$$(\mathbf{X}^T \mathbf{v})_i = (\nabla Q(\beta^*))_i = q'(\beta_i^*) = -(\beta_i^*)^{-3} + \beta_i(0)^{-3}, \quad \forall i.$$

541 On the other hand, from Eq. (34), it can be verified that

$$-(\beta_i^\infty)^{-3} + \beta_i(0)^{-3} = -3(\mathbf{X}^T \mathbf{v}^\infty)_i - \beta_i(0)^{-3} + \beta_i(0)^{-3} = -3(\mathbf{X}^T \mathbf{v}^\infty)_i, \quad \forall i.$$

542 Thus it holds that $\nabla Q(\beta^\infty) = -\frac{1}{3} \mathbf{X}(\mathbf{v}^\infty)$. Combining this with the assumption that $\mathbf{X}\beta^\infty = \mathbf{y}$, we
 543 derive that β^∞ satisfies the KKT condition. Moreover, by simple calculation, $Q(\beta)$ is convex, which
 544 make β^∞ an optimum of the problem.

545

□

546 *Proof of Theorem 4. Gradient Flow.* We first show that the distance between $\beta(t)$ and $\hat{\beta}^*$ is bounded.
 547 From the dynamic of $\beta(t)$ shown in Eq. (31), we can derive the gradient flow of $\|\beta(t) - \hat{\beta}^*\|_2^2$ as
 548 below:

$$\frac{d}{dt} \|\beta(t) - \hat{\beta}^*\|_2^2 = \left(\frac{d\tilde{\beta}(t)}{dt} \right)^T (\beta(t) - \hat{\beta}^*) = -\left\| \mathbf{X}(\beta(t) - \hat{\beta}^*) \odot \beta(t)^{\odot 2} \right\|_2^2 \leq 0. \quad (35)$$

549 Therefore, $\|\beta(t) - \hat{\beta}^*\|_2^2$ is monotonically non-increasing and upper bounded by $\|\beta(0) - \hat{\beta}^*\|_2^2$ for
 550 all t . By Assumption 1, we then have $\|\beta(t)\|_\infty \geq c > 0$ for all t . To prove the convergence, we
 551 denote $\mathbf{r}(t) = \mathbf{X}\beta(t) - \mathbf{y}$. The gradient flow of $\|\mathbf{r}(t)\|_2^2$ is

$$\frac{d}{dt}\|\mathbf{r}(t)\|_2^2 = \left(\frac{d\tilde{\beta}(t)}{dt}\right)^T \tilde{\mathbf{X}}^T \mathbf{r}(t) = -(\mathbf{r}(t)^T \tilde{\mathbf{X}} \tilde{\mathbf{X}}^T \mathbf{r}(t)) \odot \tilde{\beta}(t)^{\odot 4}. \quad (36)$$

Combining this with the fact that $\mu_{\min} > 0$ and the lower boundedness of $\|\beta(t)\|_\infty$, we then have

$$\frac{d}{dt}\|\mathbf{r}(t)\|_2^2 \leq c^4 \mu_{\min} \|\mathbf{r}(t)\|_2^2,$$

552 which proves the convergence of gradient flow.

Gradient Descent. The proof of gradient descent follows from a similar strategy. We first give an explicit expression of the update on β^k . In the following we denote $\mathbf{r}^k = \mathbf{X}\beta^k - \mathbf{y}$. Recall that the gradient descent iterate on w_i is

$$w_i^{k+1} = w_i^k - \eta \frac{1}{(1-w_i)^2} \tilde{\mathbf{x}}_i \mathbf{r}^k,$$

553 where $\tilde{\mathbf{x}}_i$ denotes the i -th column of $\tilde{\mathbf{X}}$. Then by the definition of β , we have

$$\begin{aligned} \beta_i^{k+1} &= \frac{1}{1-w_i^{k+1}} = \left(\frac{1}{1-w_i^k} - \frac{1}{1-w_i^k + \eta\beta_i^k \tilde{\mathbf{x}}_i^T \mathbf{r}^k} \right) \frac{1-w_i^k}{\eta\beta_i^k \tilde{\mathbf{x}}_i^T \mathbf{r}^k} \\ &= \frac{\beta_i^k - \beta_i^{k+1}}{\eta(\beta_i^k)^3 \tilde{\mathbf{x}}_i^T \mathbf{r}^k}, \end{aligned}$$

554 Equivalently, the update of β can be expressed as

$$\beta^{k+1} = \beta^k - \eta \mathbf{X}^T \mathbf{r}^k \odot \mathbf{u}^k, \quad \mathbf{u}^k := (\beta^k)^{\odot 3} \odot \beta^{k+1}. \quad (37)$$

555 We now show that with an appropriate choice of η , the distance between β^k and $\hat{\beta}^*$ is bounded. By
 556 Eq. (37), we have

$$\begin{aligned} \|\beta^{k+1} - \hat{\beta}^*\|_2^2 - \|\beta^k - \hat{\beta}^*\|_2^2 &= \|\beta^{k+1} - \beta^k\|_2^2 - 2(\beta^k - \hat{\beta}^*)^T (\beta^{k+1} - \beta^k) \\ &= \eta^2 \|\mathbf{X}^T \mathbf{r}^k \odot \mathbf{u}^k\|_2^2 - 2\eta \|\mathbf{r}^k \odot (\mathbf{u}^k)^{\odot \frac{1}{2}}\|_2^2 \\ &\leq \mu_{\max} \eta^2 \sum_{i=1}^n (r_i^k u_i^k)^2 - 2\eta \sum_{i=1}^n (r_i^k)^2 u_i^k. \end{aligned} \quad (38)$$

Assume $\beta^k > 0$ for all k so that $u_i^k > 0$ for all i . Now we set $\eta < \frac{1}{C\mu_{\max}}$. With these conditions, it holds for each i that

$$\mu_{\max} \eta^2 (u_i^k)^2 - 2\eta u_i^k \leq 0.$$

557 Combining this with Eq. (38), we have $\|\beta^{k+1} - \hat{\beta}^*\|_2^2 \leq \|\beta^k - \hat{\beta}^*\|_2^2$. By Assumption 1, it can be
 558 shown that $\|\beta^k\|_\infty \geq c > 0$ for all k . Similar to the proof for gradient flow, we turn to the update of
 559 $\|\mathbf{r}^k\|_2$. Note the the loss function is μ_{\max} -smooth w.r.t. β , thus we have

$$\|\mathbf{r}^{k+1}\|_2^2 \leq \|\mathbf{r}^k\|_2^2 + 2(\mathbf{r}^k)^T \mathbf{X}(\beta^{k+1} - \beta^k) + \mu_{\max} \|\beta^{k+1} - \beta^k\|_2^2.$$

560 Substituting the update of β^k in Eq. (37) into the above equation, we have

$$\begin{aligned} \|\mathbf{r}^{k+1}\|_2^2 &\leq \|\mathbf{r}^k\|_2^2 - 2\eta (\mathbf{r}^k)^T \mathbf{X} (\mathbf{X}^T \mathbf{r}^k \odot (\beta^{k+1})^3 \odot \beta^k) + \eta^2 \mu_{\max} \|\mathbf{X}^T \mathbf{r}^k \odot (\beta^{k+1})^3 \odot \beta^k\|_2^2 \\ &= \|\mathbf{r}^k\|_2^2 - 2\eta \sum_{i=1}^n (l_i^k)^2 u_i^k + \eta^2 \mu_{\max} \sum_{i=1}^n (l_i^k u_i^k)^2, \end{aligned} \quad (39)$$

561 where we denote $\mathbf{X}^T \mathbf{r}^k = \mathbf{l}^k$ for simplicity. For every fixed l_i^k , the quadratic function $f(u_i^k) =$
 562 $-2\eta (l_i^k)^2 u_i^k + \eta^2 \mu_{\max} (l_i^k u_i^k)^2$ attains its minima at $u_i^k = \frac{1}{\eta \mu_{\max}} > C$, from which we know that $f(u_i^k)$
 563 is monotonically decreasing for $u_i^k < C$. Hence, by the fact that $u_i^k > c$, it holds that

$$-2\eta (l_i^k)^2 u_i^k + \eta^2 \mu_{\max} (l_i^k u_i^k)^2 \leq (-2\eta c + \eta^2 \mu_{\max} c^2) (l_i^k)^2 \leq 0, \quad \forall 1 \leq i \leq n, \quad (40)$$

Note that $\sum_{i=1}^n (l_i^k)^2 = \|\mathbf{X}^T \mathbf{r}^k\|_2^2 \leq \mu_{\max} \|\mathbf{r}^k\|_2^2$. Leveraging this and Eq. (39) and Eq. (40), we have

$$\|\mathbf{r}^{k+1}\|_2^2 \leq (1 - (-2\eta c + \eta^2 \mu_{\max} c^2)) \|\mathbf{r}^k\|_2^2.$$

Moreover, to ensure $\beta_i^1 = \frac{\beta_i^0}{1 + \eta(\beta_i^0)^3 \mathbf{x}_i^T \mathbf{r}^0} \geq 0$ for all i , we choose

$$\eta \leq \frac{1}{\|\beta^0\|_\infty^3 \mathbf{x}_i^T \mathbf{r}^0} \leq \frac{1}{C^3 \|\mathbf{X}\|_2 \|\mathbf{r}^0\|_2} \leq \frac{1}{C^4 \mu_{\max} \|\mathbf{r}^0\|_2}.$$

With this choice of η , we can prove by induction that the assumption $\beta^k > 0$ holds for all k . Indeed, $k = 0$ follows immediately from assumption. If $\beta^t > 0$ holds, then from update of β^{t+1} , we have

$$\beta_i^{t+1} = \frac{\beta_i^t}{1 + \eta(\beta_i^t)^3 \mathbf{x}_i^T \mathbf{r}^t} \geq \frac{\beta_i^t}{1 + \eta(\beta_i^t)^3 \|\mathbf{x}_i^T\|_2 \|\mathbf{r}^t\|_2} \geq \frac{\beta_i^t}{1 + \eta C^3 \|\mathbf{X}\|_2 \|\mathbf{r}^0\|_2} \geq 0.$$

564 Thus the induction holds. Finally, we set $\eta = \min \left\{ \frac{2}{C^4 \mu_{\max}}, \frac{1}{C^4 \mu_{\max} \|\mathbf{r}^0\|_2} \right\}$, the theorem follows. \square

565 Next, we move to prove Proposition 2. For completeness, we formally introduce the definition of
566 GOTU in [34] as below.

Definition 1 (Generalization on the Unseen, [34]). *Let $\ell : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ be a loss function and \mathcal{S} be a given sample space. During training, part of \mathcal{S} is not sampled, which we call the unseen domain \mathcal{U} , while in testing, we sample from the full set \mathcal{S} . Let f be the target function and $\tilde{f}_{\mathcal{S} \setminus \mathcal{U}}$ the function learned by a learning algorithm on $\mathcal{S} \setminus \mathcal{U}$. The generalization on the unseen for an algorithm \tilde{f} and target function f is defined as*

$$GOTU(f, \tilde{f}, \mathcal{U}) = \mathbb{E}_{X \sim_{\mathcal{U}} \mathcal{U}} [\ell(\tilde{f}_{\mathcal{S} \setminus \mathcal{U}}(X), f(X))],$$

567 where $\sim_{\mathcal{U}} \mathcal{U}$ refers to uniform sampling from \mathcal{U} .

Proof of Proposition 2. We first give an explicit expression of the expected loss and gradient flow dynamics. Denote

$$\tilde{f}_\beta(\mathbf{x}) = \sum_{i=1}^d \beta_i x_i + b = f(\mathbf{w}, \mathbf{x}) = \sum_{i=1}^d \frac{1}{1 - w_i} x_i + b.$$

By definition, the half ℓ_2 loss on any sample \mathbf{x} is

$$\ell(\tilde{f}_\beta(x), f(\mathbf{x})) = \frac{1}{2} (\tilde{f}_\beta(x) - f(\mathbf{x}))^2 = \frac{1}{2} \left(b - \hat{f}(\emptyset) + \sum_{i=1}^d (\beta_i - \hat{f}(\{i\})) x_i \right)^2$$

568 Denote the distribution on the training set by U_{-k}^{d-1} . Note that $\{1, x_1, \dots, x_d\}$ are orthogonal in the
569 Hilbert space $\mathcal{S} = \{\pm 1\}^d$ equipped with the inner product $\langle g, h \rangle = \mathbb{E}_{\mathbf{x} \sim_{\mathcal{U}} \{\pm 1\}^d} [g(\mathbf{x})h(\mathbf{x})]$. Denote
570 the distribution on the training samples by U_{-k}^{d-1} . By using Parseval's Theorem, the expected loss on
571 the training set can be expressed as:

$$\begin{aligned} \mathbb{E}_{U_{-k}^{d-1}} [\ell(\tilde{f}_\beta(x), f(\mathbf{x}))] &= \frac{1}{2} \mathbb{E}_{U_{-k}^{d-1}} \left[\left(b - \hat{f}(\emptyset) + \sum_{i=1}^d (\beta_i - \hat{f}(\{i\})) x_i \right)^2 \right] \\ &= \frac{1}{2} \left(b - \hat{f}(\emptyset) + \beta_k - \hat{f}(\{k\}) \right)^2 + \frac{1}{2} \sum_{i \neq k}^d (\beta_i - \hat{f}(\{i\}))^2. \end{aligned}$$

572 Then we can derive the gradient flow for β_i and b as below

$$\begin{aligned} \frac{db(t)}{dt} &= -(b(t) - \hat{f}(\emptyset) + \beta_k(t) - \hat{f}(\{k\})), \\ \frac{d\beta_k(t)}{dt} &= -(b(t) - \hat{f}(\emptyset) + \beta_k(t) - \hat{f}(\{k\}))\beta_k(t)^4, \\ \frac{d\beta_i(t)}{dt} &= -(\beta_i(t) - \hat{f}(\{i\}))\beta_i(t)^4, \quad \forall i \neq k. \end{aligned}$$

573 For simplicity, denote $B = \hat{f}(\emptyset) + \hat{f}(\{k\})$. Using the above, we have

$$\begin{aligned} \frac{d}{dt}(b(t) + \beta_k(t) - B)^2 &= -2(b(t) + \beta_k(t) - B)^2(1 + \beta_k(t)^4), \\ \frac{d}{dt}(\beta_i(t) - \hat{f}(\{i\}))^2 &= -2(\beta_i(t) - \hat{f}(\{i\}))^2\beta_i(t)^4, \end{aligned} \quad (41)$$

574 which shows that $|b(t) + \beta_k(t) - B|^2$ and $|\beta_i(t) - \hat{f}(\{i\})|^2$ is monotonically nonincreasing. Since
575 $\beta_i(0)$ and $\hat{f}(\{i\})$ are greater than 0, from the monotonicity we know that $\beta_i(t) > 0$ for all t .
576 Therefore, the convergence of gradient flow follows from Eq. (41) that both $|b(t) + \beta_k(t) - B|^2$ and
577 $|\beta_i(t) - \hat{f}(\{i\})|^2$ decrease linearly.

578 Denote the limit of $\beta_i(t)$ and $b(t)$ by β_i^∞ and b^∞ , respectively. We now turn to estimate the GOTU
579 error.

580 1. When $B > 1$, it holds that $b(0) + \beta_k(0) - B < 0$, thus $b(t)$ and $\beta_k(t)$ is monotonically
581 increasing. Using the fact that $\beta_k(t) > \beta_k(0) = 1$, we know that

$$\frac{d}{dt}(\beta_k(t) - b(t)) = -2(b(t) + \beta_k(t) - B)^2(\beta_k(t)^4 - 1) < 0.$$

582 Combing this with $\beta_k^\infty + b^\infty = B$, it can be verified that $\beta_k^\infty \geq \frac{B+1}{2}$. Then by definition
583 and Parseval's theorem, the GOTU loss is

$$\begin{aligned} GOTU(f, \tilde{f}_\beta, \{x_k = -1\}) &= \left(b^\infty - \hat{f}(\emptyset) - \beta_k^\infty + \hat{f}(\{k\})\right)^2 + \sum_{i \neq k}^d (\beta_i^\infty - \hat{f}(\{i\}))^2 \\ &= 4(\hat{f}(\{k\}) - \beta_k^\infty)^2, \end{aligned}$$

584 where we use the convergence of the flow in the second line. Leveraging the bound of β_k^∞ ,
585 we derive that

$$4(\hat{f}(\{k\}) - \beta_k^\infty)^2 \leq 4\left(\hat{f}(\{k\}) - \frac{B+1}{2}\right)^2. \quad (42)$$

586 By the assumption that $\hat{f}(\emptyset) < 2\hat{f}(\{k\})$, we have $\frac{B+1}{2} < \frac{3\hat{f}(\{k\})+1}{2} < 2\hat{f}(\{k\})$. Leverag-
587 ing this in Eq. (42), we know that

$$GOTU(f, \tilde{f}_\beta, \{x_k = -1\}) \leq (\hat{f}(\{k\}) + 1)^2. \quad (43)$$

2. When $B < 1$, similar to the proof of Theorem 3, we have from the dynamic of $\beta_k(t)$ that

$$\beta_k(t)^{-3} - 1 = 3 \int_0^t (b(s) + \beta_k(s) - B) ds \leq 3(1 - B),$$

588 where we use the monotonicity of $b(s) + \beta_k(s) - B$ and the convergence of the flow.
589 Therefore, it holds that $\beta_k^\infty \geq (3(1 - B) + 1)^{-\frac{1}{3}}$. We can bound the GOTU error as

$$4(\hat{f}(\{k\}) - \beta_k^\infty)^2 \leq 4(\hat{f}(\{k\}) - (3(1 - B) + 1)^{-\frac{1}{3}})^2. \quad (44)$$

590 By using the assumption that $\hat{f}(\emptyset) > -2\hat{f}(\{k\})$, Eq. (44) gives

$$GOTU(f, \tilde{f}_\beta, \{x_k = -1\}) \leq 4\left(\hat{f}(\{k\}) - \left(4 + 3\hat{f}(\{k\})\right)^{-\frac{1}{3}}\right)^2. \quad (45)$$

591 Then the proposition follows from Eq. (43) and (45).

592 □

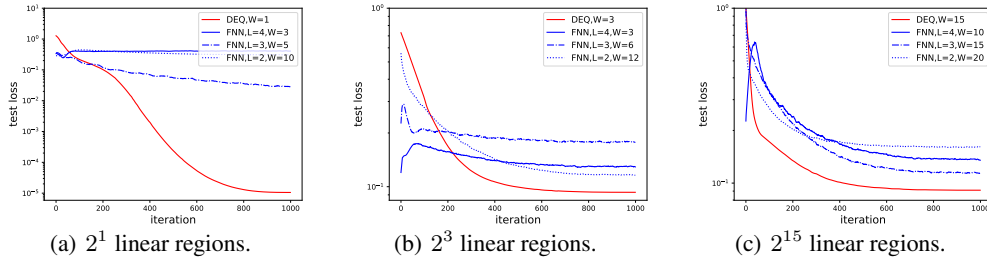


Figure 2: Test loss of FNN and DEQ trained on sawtooth functions with 2^1 , 2^3 , 2^{15} linear regions.

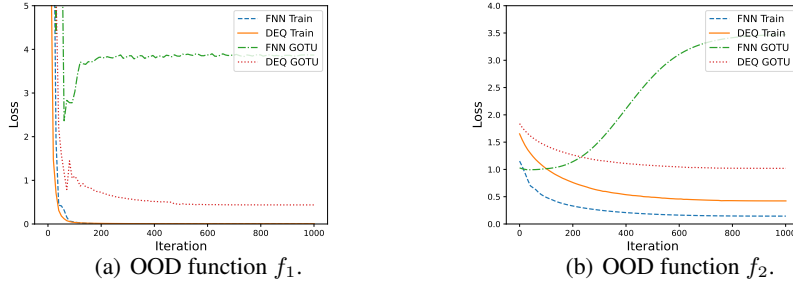


Figure 3: Train and test loss of DEQ and FNN trained on OOD tasks f_1 and f_2 .

593 B Experiment Details

594 B.1 Supplementary Experiments in Section 6

595 In this subsection, We first show how DEQ and FNN perform on various linear regions of sawtooth
 596 function. We report results of other sawtooth functions with less or more linear regions. Figure 2
 597 present the test loss for sawtooth functions with 2^1 , 2^3 and 2^{15} linear regions. For all experiments,
 598 we execute our program on Nvidia GTX 1660 and all the program occupies less than 10M memory
 599 and runs for less than 2 minutes. In consistency with our results in Section 6, we can see that DEQ
 600 outperforms FNN with similar size of network on every sawtooth function in our experiment and the
 601 test loss of DEQ converges closer to zero loss.

602 We next conduct OOD experiments on the following 2 functions and unseen domains. The first
 603 function is a higher-dimensional form of Eq. (14) which is a form of mean function. The second
 604 function is the majority function on 3 bits with the maximum degree 3. The expressions of these
 605 functions are presented below.

$$f_1(x) = 1.25 * x_0 + 1.25 * x_1 + \dots + 1.25 * x_{20}, \quad \mathcal{U} = \{\mathbf{x} \in \{\pm 1\}^{10} : x_1 = -1\},$$

$$f_2(x) = \frac{1}{2}(x_0 + x_1 + x_2 - x_1x_2x_3), \quad \mathcal{U} = \{\mathbf{x} \in \{\pm 1\}^{10} : x_0x_1 = -1\}.$$

606 For all experiments, we generate all binary sequences in $\mathcal{U}^c = \{\pm 1\}^d \setminus \mathcal{U}$ for training. Figure 3(a)
 607 shows that the GOTU error does not increase significantly compared to Figure 1 where the ambient
 608 dimension is 10. In consistency with our results in Section 6, we can learn from Figure 3(a) that when
 609 learning a linear boolean function on population loss on DEQ, the training loss converges to zero
 610 and the generalization error on the unseen is bounded. As is shown in Figure 3(b), when learning
 611 the unlinear boolean function, DEQ can also achieve nearly zero train loss with smaller GOTU error
 612 compared with FNN.

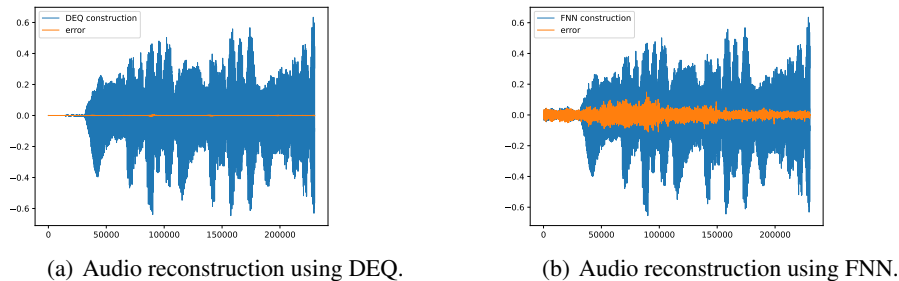


Figure 4: The reconstruction results with DEQ and FNN and the error computed by subtracting the original signal.

613 B.2 Experiment on Audio Representation

614 Inspired by the overall studies, we conduct experiments on a real tasks of audio representation
 615 to verify the potential advantage of DEQ in learning functions with high-frequency component.
 616 We utilize the setting of experiments in [43], where the very-high-frequency audio signals were
 617 represented using a conventional explicit network and an(implicit)² network, which is variant of DEQ
 618 employing a neural block with three layers and specific activation functions such as $\sin(x)$ Although
 619 Huang et al. [43] shows that (implicit)² network outperforms conventional explicit networks in audio
 620 representation [43], revealing the advantage of DEQ to an extend, it is unclear whether the superiority
 621 of the (implicit)² network is attributed solely to the carefully-designed block. In contrast, we apply
 622 DEQ and FNN in their basic forms to represent the audio signal in our experiment to further explore
 623 the potential advantages of DEQ in real scenarios.

624 Following the setting in [43], we train the models to fit a 7-second music piece. We set the width of
 625 DEQ to 20, the layer of FNN to 3 and the hidden dimension of FNN to 20. This setting enables the
 626 model to exactly fit the audio signal based on our experiments.

627 In Figure 4, we show the reconstruction results with DEQ and FNN and the error computed by
 628 subtracting the original signal. We observe that DEQ outperforms FNN with a noticeable error,
 629 verifying the advantages of DEQ in representing high-frequency components.

630 **NeurIPS Paper Checklist**

631 **1. Claims**

632 Question: Do the main claims made in the abstract and introduction accurately reflect the
633 paper's contributions and scope?

634 Answer: [Yes]

635 Justification: See Abstract and Introduction, where enumerates the contributions.

636 Guidelines:

- 637 • The answer NA means that the abstract and introduction do not include the claims
638 made in the paper.
- 639 • The abstract and/or introduction should clearly state the claims made, including the
640 contributions made in the paper and important assumptions and limitations. A No or
641 NA answer to this question will not be perceived well by the reviewers.
- 642 • The claims made should match theoretical and experimental results, and reflect how
643 much the results can be expected to generalize to other settings.
- 644 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
645 are not attained by the paper.

646 **2. Limitations**

647 Question: Does the paper discuss the limitations of the work performed by the authors?

648 Answer: [Yes]

649 Justification: See Related Works and Section 5. We mention that we study a simplified DEQ
650 due to technical issues.

651 Guidelines:

- 652 • The answer NA means that the paper has no limitation while the answer No means that
653 the paper has limitations, but those are not discussed in the paper.
- 654 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 655 • The paper should point out any strong assumptions and how robust the results are to
656 violations of these assumptions (e.g., independence assumptions, noiseless settings,
657 model well-specification, asymptotic approximations only holding locally). The authors
658 should reflect on how these assumptions might be violated in practice and what the
659 implications would be.
- 660 • The authors should reflect on the scope of the claims made, e.g., if the approach was
661 only tested on a few datasets or with a few runs. In general, empirical results often
662 depend on implicit assumptions, which should be articulated.
- 663 • The authors should reflect on the factors that influence the performance of the approach.
664 For example, a facial recognition algorithm may perform poorly when image resolution
665 is low or images are taken in low lighting. Or a speech-to-text system might not be
666 used reliably to provide closed captions for online lectures because it fails to handle
667 technical jargon.
- 668 • The authors should discuss the computational efficiency of the proposed algorithms
669 and how they scale with dataset size.
- 670 • If applicable, the authors should discuss possible limitations of their approach to
671 address problems of privacy and fairness.
- 672 • While the authors might fear that complete honesty about limitations might be used by
673 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
674 limitations that aren't acknowledged in the paper. The authors should use their best
675 judgment and recognize that individual actions in favor of transparency play an impor-
676 tant role in developing norms that preserve the integrity of the community. Reviewers
677 will be specifically instructed to not penalize honesty concerning limitations.

678 **3. Theory Assumptions and Proofs**

679 Question: For each theoretical result, does the paper provide the full set of assumptions and
680 a complete (and correct) proof?

681 Answer: [Yes]

682 Justification: For assumptions, see Section 5. For proofs, see Appendix A.

683 Guidelines:

- 684 • The answer NA means that the paper does not include theoretical results.
- 685 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
- 686 referenced.
- 687 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 688 • The proofs can either appear in the main paper or the supplemental material, but if
- 689 they appear in the supplemental material, the authors are encouraged to provide a short
- 690 proof sketch to provide intuition.
- 691 • Inversely, any informal proof provided in the core of the paper should be complemented
- 692 by formal proofs provided in appendix or supplemental material.
- 693 • Theorems and Lemmas that the proof relies upon should be properly referenced.

694 4. Experimental Result Reproducibility

695 Question: Does the paper fully disclose all the information needed to reproduce the main ex-

696 perimental results of the paper to the extent that it affects the main claims and/or conclusions

697 of the paper (regardless of whether the code and data are provided or not)?

698 Answer: [Yes]

699 Justification: See Section 6. We provide details of our experiment.

700 Guidelines:

- 701 • The answer NA means that the paper does not include experiments.
- 702 • If the paper includes experiments, a No answer to this question will not be perceived
- 703 well by the reviewers: Making the paper reproducible is important, regardless of
- 704 whether the code and data are provided or not.
- 705 • If the contribution is a dataset and/or model, the authors should describe the steps taken
- 706 to make their results reproducible or verifiable.
- 707 • Depending on the contribution, reproducibility can be accomplished in various ways.
- 708 For example, if the contribution is a novel architecture, describing the architecture fully
- 709 might suffice, or if the contribution is a specific model and empirical evaluation, it may
- 710 be necessary to either make it possible for others to replicate the model with the same
- 711 dataset, or provide access to the model. In general, releasing code and data is often
- 712 one good way to accomplish this, but reproducibility can also be provided via detailed
- 713 instructions for how to replicate the results, access to a hosted model (e.g., in the case
- 714 of a large language model), releasing of a model checkpoint, or other means that are
- 715 appropriate to the research performed.
- 716 • While NeurIPS does not require releasing code, the conference does require all submis-
- 717 sions to provide some reasonable avenue for reproducibility, which may depend on the
- 718 nature of the contribution. For example
 - 719 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
 - 720 to reproduce that algorithm.
 - 721 (b) If the contribution is primarily a new model architecture, the paper should describe
 - 722 the architecture clearly and fully.
 - 723 (c) If the contribution is a new model (e.g., a large language model), then there should
 - 724 either be a way to access this model for reproducing the results or a way to reproduce
 - 725 the model (e.g., with an open-source dataset or instructions for how to construct
 - 726 the dataset).
 - 727 (d) We recognize that reproducibility may be tricky in some cases, in which case
 - 728 authors are welcome to describe the particular way they provide for reproducibility.
 - 729 In the case of closed-source models, it may be that access to the model is limited in
 - 730 some way (e.g., to registered users), but it should be possible for other researchers
 - 731 to have some path to reproducing or verifying the results.

732 5. Open access to data and code

733 Question: Does the paper provide open access to the data and code, with sufficient instruc-

734 tions to faithfully reproduce the main experimental results, as described in supplemental

735 material?

736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787

Answer: [Yes]

Justification: We will provide the open access to the code after the paper publication.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: See Section 6 and Appendix B.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [No]

Justification: This paper does not include experiments on large real datasets.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.

- 788
- 789
- 790
- 791
- 792
- 793
- 794
- 795
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
 - For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
 - If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments Compute Resources

797 Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

800 Answer: [Yes]

801 Justification: See Appendix B.1.

802 Guidelines:

- 803
- 804
- 805
- 806
- 807
- 808
- 809
- 810
- The answer NA means that the paper does not include experiments.
 - The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
 - The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
 - The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code Of Ethics

812 Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics [https://neurips.cc/public/EthicsGuidelines?](https://neurips.cc/public/EthicsGuidelines)

814 Answer: [Yes]

815 Justification:

816 Guidelines:

- 817
- 818
- 819
- 820
- 821
- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
 - If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
 - The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader Impacts

823 Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

825 Answer: [NA]

826 Justification: This paper mainly focus on the theory of neural networks.

827 Guidelines:

- 828
- 829
- 830
- 831
- 832
- 833
- 834
- 835
- 836
- 837
- 838
- The answer NA means that there is no societal impact of the work performed.
 - If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
 - Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
 - The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to

839 generate deepfakes for disinformation. On the other hand, it is not needed to point out
840 that a generic algorithm for optimizing neural networks could enable people to train
841 models that generate Deepfakes faster.

- 842 • The authors should consider possible harms that could arise when the technology is
843 being used as intended and functioning correctly, harms that could arise when the
844 technology is being used as intended but gives incorrect results, and harms following
845 from (intentional or unintentional) misuse of the technology.
- 846 • If there are negative societal impacts, the authors could also discuss possible mitigation
847 strategies (e.g., gated release of models, providing defenses in addition to attacks,
848 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
849 feedback over time, improving the efficiency and accessibility of ML).

850 11. Safeguards

851 Question: Does the paper describe safeguards that have been put in place for responsible
852 release of data or models that have a high risk for misuse (e.g., pretrained language models,
853 image generators, or scraped datasets)?

854 Answer: [NA]

855 Justification: This paper focuses on theory.

856 Guidelines:

- 857 • The answer NA means that the paper poses no such risks.
- 858 • Released models that have a high risk for misuse or dual-use should be released with
859 necessary safeguards to allow for controlled use of the model, for example by requiring
860 that users adhere to usage guidelines or restrictions to access the model or implementing
861 safety filters.
- 862 • Datasets that have been scraped from the Internet could pose safety risks. The authors
863 should describe how they avoided releasing unsafe images.
- 864 • We recognize that providing effective safeguards is challenging, and many papers do
865 not require this, but we encourage authors to take this into account and make a best
866 faith effort.

867 12. Licenses for existing assets

868 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
869 the paper, properly credited and are the license and terms of use explicitly mentioned and
870 properly respected?

871 Answer: [NA]

872 Justification: See Section 6 and Appendix B. Although experiments in Appendix B.2 utilize
873 another experiment, the code, data and models are created by ourselves.

874 Guidelines:

- 875 • The answer NA means that the paper does not use existing assets.
- 876 • The authors should cite the original paper that produced the code package or dataset.
- 877 • The authors should state which version of the asset is used and, if possible, include a
878 URL.
- 879 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 880 • For scraped data from a particular source (e.g., website), the copyright and terms of
881 service of that source should be provided.
- 882 • If assets are released, the license, copyright information, and terms of use in the
883 package should be provided. For popular datasets, paperswithcode.com/datasets
884 has curated licenses for some datasets. Their licensing guide can help determine the
885 license of a dataset.
- 886 • For existing datasets that are re-packaged, both the original license and the license of
887 the derived asset (if it has changed) should be provided.
- 888 • If this information is not available online, the authors are encouraged to reach out to
889 the asset's creators.

890 13. New Assets

891 Question: Are new assets introduced in the paper well documented and is the documentation
892 provided alongside the assets?

893 Answer: [NA]

894 Justification:

895 Guidelines:

- 896 • The answer NA means that the paper does not release new assets.
- 897 • Researchers should communicate the details of the dataset/code/model as part of their
898 submissions via structured templates. This includes details about training, license,
899 limitations, etc.
- 900 • The paper should discuss whether and how consent was obtained from people whose
901 asset is used.
- 902 • At submission time, remember to anonymize your assets (if applicable). You can either
903 create an anonymized URL or include an anonymized zip file.

904 14. **Crowdsourcing and Research with Human Subjects**

905 Question: For crowdsourcing experiments and research with human subjects, does the paper
906 include the full text of instructions given to participants and screenshots, if applicable, as
907 well as details about compensation (if any)?

908 Answer: [NA]

909 Justification:

910 Guidelines:

- 911 • The answer NA means that the paper does not involve crowdsourcing nor research with
912 human subjects.
- 913 • Including this information in the supplemental material is fine, but if the main contribu-
914 tion of the paper involves human subjects, then as much detail as possible should be
915 included in the main paper.
- 916 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
917 or other labor should be paid at least the minimum wage in the country of the data
918 collector.

919 15. **Institutional Review Board (IRB) Approvals or Equivalent for Research with Human 920 Subjects**

921 Question: Does the paper describe potential risks incurred by study participants, whether
922 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
923 approvals (or an equivalent approval/review based on the requirements of your country or
924 institution) were obtained?

925 Answer: [NA]

926 Justification:

927 Guidelines:

- 928 • The answer NA means that the paper does not involve crowdsourcing nor research with
929 human subjects.
- 930 • Depending on the country in which research is conducted, IRB approval (or equivalent)
931 may be required for any human subjects research. If you obtained IRB approval, you
932 should clearly state this in the paper.
- 933 • We recognize that the procedures for this may vary significantly between institutions
934 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
935 guidelines for their institution.
- 936 • For initial submissions, do not include any information that would break anonymity (if
937 applicable), such as the institution conducting the review.