

---

# On the Failure of a Universal Linear Representation Hypothesis in Deep Neural Networks

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 The Linear Representation Hypothesis (LRH) posits that semantic concepts in deep  
2 neural networks are encoded as linear directions in activation space, recoverable via  
3 linear probes. We provide a comprehensive theoretical analysis demonstrating that  
4 LRH cannot be universally true. Our contributions include: (1) sharp combinatorial  
5 bounds using VC-dimension theory showing exponential gaps between concept  
6 complexity and linear decodability, (2) circuit complexity arguments proving  
7 that depth-separated functions cannot have linear intermediate representations  
8 without exponential dimension, (3) explicit algebraic constructions of non-linear  
9 concept families with detailed complexity analysis, (4) measure-theoretic results  
10 on the generic failure of linear representations, and (5) information-theoretic lower  
11 bounds on representation dimension. We complement theory with reproducible  
12 experiments and discuss precise conditions under which LRH holds approximately.

## 1 Introduction and Motivation

13 The interpretability of deep neural networks remains a central challenge in machine learning [14, 7, 9].  
14 A fundamental assumption underlying many interpretability methods is the **Linear Representation**  
15 **Hypothesis (LRH)**, defined as follows. For semantic concepts learned by deep networks, internal  
16 activations encode these concepts as linear directions: there exists a vector  $\mathbf{v}$  such that the concept  
17 value is approximately  $\text{sign}(\mathbf{v}^\top \phi(\mathbf{x}))$  where  $\phi(\mathbf{x})$  represents intermediate activations.  
18 This hypothesis motivates numerous interpretability techniques including linear probing [1, 3],  
19 activation patching [22], and concept bottleneck models [13]. However, the theoretical foundations  
20 of LRH remain poorly understood.

22 **Our Contributions** We provide the first comprehensive theoretical analysis of the limitations of  
23 the linear representation hypothesis through multiple complementary lenses:

- 24 1. *Combinatorial Analysis*: Sharp VC-dimension bounds showing exponential separation  
25 between concept complexity and linear representability
- 26 2. *Circuit Complexity*: Rigorous depth-separation arguments proving incompatibility with  
27 known lower bounds
- 28 3. *Algebraic Constructions*: Explicit families of functions demonstrating non-linear interme-  
29 diate encodings
- 30 4. *Information Theory*: Lower bounds on representation dimension for linear concept recovery
- 31 5. *Measure Theory*: Generic failure results for random concept-representation pairs
- 32 6. *Empirical Validation*: Reproducible experiments confirming theoretical predictions

33 Our analysis reveals fundamental limitations while identifying precise conditions enabling approxi-  
34 mate linear decodability.

35 **2 Mathematical Framework and Definitions**

36 **2.1 Basic Setup**

37 Let  $\mathcal{X}$  denote an input space (typically  $\mathbb{R}^m$  or  $\{0, 1\}^m$ ) equipped with a probability measure  $\mu$ .  
 38 We consider concepts as measurable functions  $c : \mathcal{X} \rightarrow \{0, 1\}$  and neural networks with the  
 39 decomposition:  $f(\mathbf{x}) = g(\phi(\mathbf{x}))$  where  $\phi : \mathcal{X} \rightarrow \mathbb{R}^d$  represents intermediate activations and  
 40  $g : \mathbb{R}^d \rightarrow \mathbb{R}$  captures subsequent computation.

41 **Definition 1** (Linear Threshold Functions). *The class of linear threshold functions on  $\mathbb{R}^d$  is:  $\mathcal{L}_d =$   
 42  $\{h(\mathbf{z}) = \text{sign}(\mathbf{w}^\top \mathbf{z} - \theta) : \mathbf{w} \in \mathbb{R}^d, \theta \in \mathbb{R}\}$*

43 **Definition 2** (Exact Linear Representability). *A concept  $c$  is exactly linearly representable in  $\phi$  if:  
 44  $\exists h \in \mathcal{L}_d$  such that  $c(\mathbf{x}) = h(\phi(\mathbf{x}))$  for all  $\mathbf{x} \in \mathcal{X}$*

45 **Definition 3** (Statistical Linear Representability). *For distribution  $\mathcal{D}$  and  $\varepsilon \geq 0$ , concept  $c$  is  
 46  $(\varepsilon, \mathcal{D})$ -linearly representable if:  $\inf_{h \in \mathcal{L}_d} \mathbb{P}_{\mathbf{x} \sim \mathcal{D}}[h(\phi(\mathbf{x})) \neq c(\mathbf{x})] \leq \varepsilon$*

47 **Definition 4** (Universal Linear Representation Hypothesis). *The strong form of LRH states: For any  
 48 semantic concept  $c$  arising in natural tasks, there exists a layer  $\ell$  in typical trained networks such  
 49 that  $c$  is  $(0.1, \mathcal{D})$ -linearly representable in the  $\ell$ -th layer representation  $\phi_\ell$ .*

50 **2.2 Complexity-Theoretic Framework**

51 We analyze concept complexity through multiple measures:

52 **Definition 5** (VC-Dimension of Concept Class). *For concept class  $\mathcal{C}$ , the VC-dimension  $\text{VC}(\mathcal{C})$  is  
 53 the largest  $m$  such that some set of  $m$  points can be shattered by  $\mathcal{C}$ .*

54 **Definition 6** (Linear Separation Complexity). *For finite set  $S = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subset \mathcal{X}$  with concept  $c$ ,  
 55 the linear separation complexity is:  $\text{LSC}_\phi(S, c) = \min\{d : \exists h \in \mathcal{L}_d \text{ s.t. } h(\phi(\mathbf{x}_i)) = c(\mathbf{x}_i) \forall i\}$*

56 **3 Combinatorial Impossibility: Sharp VC-Dimension Analysis**

57 Our first main result establishes fundamental combinatorial limitations using VC-dimension theory.

58 **3.1 Classical Foundation**

59 **Lemma 1** (Cover's Dichotomy Theorem [6]). *For  $n$  points in general position in  $\mathbb{R}^d$ , the number of  
 60 distinct dichotomies realizable by hyperplanes is:  $C(n, d) = 2 \sum_{i=0}^d \binom{n-1}{i}$  when  $n > d + 1$ , and  
 61  $C(n, d) = 2^n$  otherwise.*

62 **Theorem 1** (Sharp Combinatorial Impossibility). *Let  $S = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subset \mathcal{X}$  with representations  
 63  $\phi(\mathbf{x}_i) \in \mathbb{R}^d$  in general position. Then: (i) **Exponential Gap**: For  $n > d + 1$ , the fraction of  
 64 linearly separable concepts is:  $\rho(n, d) = \frac{C(n, d)}{2^n} \leq \frac{2^{d+1} n^d}{d! \cdot 2^n} = O\left(\frac{n^d}{2^n}\right)$ . (ii) **Sharp Threshold**: When  
 65  $d = o(n/\log n)$ , we have  $\rho(n, d) = o(1)$ . (iii) **Lower Bound**: For any  $\varepsilon > 0$ , if  $d < (1 - \varepsilon) \log_2 n$ ,  
 66 then  $\rho(n, d) < 2^{-\varepsilon n}$ .*

67 *Proof.* See Appendix A. □

68 **Corollary 1** (Dimension-Dependent Impossibility). *For representation dim  $d$  and dataset size  $n$ : if  
 69  $d = O(\log n)$ , then  $(1 - o(1))$ -fraction of concepts are not linearly separable, if  $d = O(\sqrt{n})$ , then  
 70  $(1 - 2^{-\Omega(\sqrt{n})})$ -fraction are not linearly separable, if  $d = \Theta(n)$ , then  $O(1)$ -fraction may be linearly  
 71 separable.*

72 **3.2 Refined Analysis for Structured Data**

73 Real datasets often have structure beyond general position. We analyze this case:

74 **Definition 7** (Effective Dimension). *For dataset  $S$  with representations  $\{\phi(\mathbf{x}_i)\}$ , the effective  
 75 dimension is:  $d_{\text{eff}}(S) = \dim(\{\phi(\mathbf{x}_i) : i = 1, \dots, n\})$*

76 **Theorem 2** (Structured Data Analysis). *Let  $S$  have effective dimension  $d_{\text{eff}}$  and condition number  $\kappa$ .  
77 Then  $\rho(n, d_{\text{eff}}) \geq \rho_{\text{worst}}(n, d_{\text{eff}}) \cdot \left(1 - O\left(\frac{\log \kappa}{d_{\text{eff}}}\right)\right)$ ,  $\rho_{\text{worst}}$  is the worst-case bound from Theorem 1.*

78 *Proof.* The condition number  $\kappa$  measures how close the data is to being linearly dependent. Using  
79 perturbation theory for linear separability [2], the number of separable dichotomies decreases by at  
80 most  $O(\log \kappa/d_{\text{eff}})$  factor when data becomes ill-conditioned.  $\square$

## 81 4 Circuit Complexity and Depth Separation

82 Our second main contribution leverages circuit complexity theory to prove that universal linear  
83 representability contradicts established depth-separation results.

### 84 4.1 Circuit Complexity Background

85 **Definition 8** (Boolean Circuit Complexity). *For Boolean function  $f : \{0, 1\}^n \rightarrow \{0, 1\}$ :  $C(f) =$   
86  $\min$ imum circuit size computing  $f$ ,  $C_d(f) = \min$ imum size of depth- $d$  circuit computing  $f$ ,  $f$  has  
87 depth separation if  $C_{d-1}(f) = \omega(C_d(f))$ .*

88 **Theorem 3** (Classical Depth Separation [11]). *The parity function  $\text{PARITY}_n(\mathbf{x}) = \bigoplus_{i=1}^n x_i$  satisfies:  
89  $C_{\lceil \log n \rceil}(\text{PARITY}_n) = O(n)$ ,  $C_2(\text{PARITY}_n) = \Omega(2^n/n)$ .*

### 90 4.2 Main Depth-Separation Result

91 **Theorem 4** (Depth-Separation Impossibility). *Let  $\mathcal{F}_n$  be the class of Boolean functions on  $n$  variables  
92 requiring depth  $\Omega(\log n)$  for polynomial-size circuits. If every  $f \in \mathcal{F}_n$  were linearly representable  
93 in some intermediate layer  $\phi : \{0, 1\}^n \rightarrow \mathbb{R}^d$  with  $d = \text{poly}(n)$ , then there exist depth-2 circuits of  
94 polynomial size computing all functions in  $\mathcal{F}_n$ , contradicting known lower bounds.*

95 *Proof.* **Step 1: Structure of Linear Representations.** Suppose  $f \in \mathcal{F}_n$  has linear representation:  
96  $f(\mathbf{x}) = \text{sign}(\mathbf{w}^\top \phi(\mathbf{x}) - \theta)$  where  $\phi$  computes polynomial-size depth- $k$  circuits.

97 **Step 2: Circuit Construction.** We can compute  $f$  via the following depth- $(k+2)$  circuit: **Layers**  
98 **1-k:** Compute  $\phi(\mathbf{x}) = (\phi_1(\mathbf{x}), \dots, \phi_d(\mathbf{x}))$ , **Layer k+1:** Compute  $\sum_{i=1}^d w_i \phi_i(\mathbf{x})$  using addition tree,  
99 **Layer k+2:** Apply threshold to get  $f(\mathbf{x})$ .

100 **Step 3: Depth Reduction.** If  $k = O(1)$  (constant depth intermediate layers), this gives depth- $O(1)$   
101 circuits for all  $f \in \mathcal{F}_n$ . But  $\mathcal{F}_n$  contains functions requiring  $\Omega(\log n)$  depth, yielding contradiction.

102 **Step 4: Polynomial Size Analysis.** Each  $\phi_i$  computed by polynomial-size circuits, and weighted  
103 sum requires  $O(d \log d)$  additional gates. Total size remains polynomial if  $d = \text{poly}(n)$ .

104 Therefore, either: intermediate layers have super-polynomial depth:  $k = \omega(\log n)$ , representation  
105 dimension is exponential:  $d = 2^{\Omega(n)}$ , or functions in  $\mathcal{F}_n$  are not linearly representable. Since practical  
106 networks use  $k = O(\log n)$  depth and  $d = \text{poly}(n)$  dimensions, the third option must hold.  $\square$

107 **Corollary 2** (Specific Function Classes). *The following function classes cannot have polynomial-  
108 dimensional linear representations in shallow intermediate layers: **Parity:**  $\text{PARITY}_n$  requires  
109  $d = \Omega(2^n)$  or depth  $\Omega(\log n)$ , **Majority:** Threshold-of-thresholds functions [16], **Iterated Products:**  
110 Functions defined by  $f(x_1, \dots, x_n) = x_1 x_2 \cdots x_n$  over finite fields, **Recursive Compositions:**  
111 Functions built by composing simple operations  $\Omega(\log n)$  times.*

### 112 4.3 Quantitative Depth-Dimension Trade-offs

113 **Theorem 5** (Depth-Dimension Trade-off). *For function family  $\mathcal{F}_n$  with optimal depth  $D^*$  and  
114 size  $S^*$ , any linear representation in intermediate layer of depth  $k < D^*/2$  requires dimension:  
115  $d \geq \frac{S^*}{2^{D^*-k}} \cdot \Omega\left(\frac{2^n}{n^{O(1)}}\right)$*

116 *Proof.* Using results from [18] on average-case complexity, functions requiring large circuits at  
117 optimal depth must have exponentially larger circuits when depth is significantly reduced. The bound  
118 follows from translating circuit size lower bounds into representation dimension requirements.  $\square$

119 **5 Explicit Algebraic Constructions**

120 We provide analysis of explicit function families showing non-linear intermediate representations.

121 **5.1 Parity Functions: Complete Analysis**

122 **Theorem 6** (Parity Impossibility - Comprehensive). *Let  $\text{PARITY}_k : \{0, 1\}^k \rightarrow \{0, 1\}$  be the  $k$ -bit parity function. For any representation  $\phi : \{0, 1\}^k \rightarrow \mathbb{R}^d$ :*

124 **(i) Exact Case:** If  $\text{PARITY}_k$  is exactly linearly representable in  $\phi$ , then either:  $d \geq 2^{k-1}$ , or  $\phi$  explicitly computes parity information. **(ii) Approximate Case:** For  $\varepsilon < 1/4$ , if  $\text{PARITY}_k$  is  $\varepsilon$ -linearly representable, then  $d \geq \Omega\left(\frac{k}{\log(1/(4\varepsilon))}\right)$  **(iii) Noise Robustness:** Under  $p$ -biased noise with  $p < 1/2 - \gamma$ , linear representability requires:  $d \geq \frac{1}{2\gamma^2} \log\left(\frac{1}{8\varepsilon}\right)$

128 *Proof.* **Part (i) - Exact Case:** Parity partitions  $\{0, 1\}^k$  into sets  $S_0 = \{\mathbf{x} : \text{PARITY}_k(\mathbf{x}) = 0\}$  and  $S_1 = \{\mathbf{x} : \text{PARITY}_k(\mathbf{x}) = 1\}$ , each of size  $2^{k-1}$ . If  $d < 2^{k-1}$  and  $\phi$  is injective on each  $S_i$ , then by pigeonhole principle, some hyperplane separating  $\phi(S_0)$  from  $\phi(S_1)$  exists only if the images are linearly separable. However,  $S_0$  and  $S_1$  have a complex geometric relationship: every element of  $S_0$  differs from some element of  $S_1$  in exactly one bit position. This creates a "checkerboard" pattern that cannot be linearly separated unless  $\phi$  preserves this structure explicitly.

134 **Part (ii) - Approximate Case:** Using VC-dimension bounds for linear threshold functions, any  $\varepsilon$ -approximation requires the representation to capture at least  $(1 - 2\varepsilon)$  fraction of the  $2^k$  dichotomy correctly. From fat-shattering dimension analysis [2], this requires:  $d \geq \Omega\left(\frac{\varepsilon^{-2} \log(2^k/\varepsilon)}{k}\right) = \Omega\left(\frac{k}{\log(1/(4\varepsilon))}\right)$ . **Part (iii) - Noise Analysis:** Under  $p$ -biased noise, each bit flips with probability  $p$ . The effective signal-to-noise ratio becomes  $(1 - 2p)^k \geq (2\gamma)^k$ . Using concentration inequalities, maintaining  $\varepsilon$ -accuracy requires:  $d \geq \frac{1}{2\gamma^2} \log\left(\frac{1}{8\varepsilon}\right)$   $\square$

140 **5.2 XOR and Generalized Parity**

141 **Example 1** (XOR Geometric Analysis). *For XOR on  $\{(0, 0), (0, 1), (1, 0), (1, 1)\}$  with labels  $(0, 1, 1, 0)$ , the convex hulls are:  $\text{conv}(\{(0, 0), (1, 1)\}) = \{t(1, 1) : t \in [0, 1]\}$  and  $\text{conv}(\{(0, 1), (1, 0)\}) = \{(s, 1-s) : s \in [0, 1]\}$ . These intersect at  $(1/2, 1/2)$ , proving no linear separator exists. Any representation  $\phi$  making XOR linearly separable must map these four points to  $\mathbb{R}^d$  such that positive and negative examples become linearly separable.*

146 **Minimal Dimension:** The minimum  $d$  for which some  $\phi : \{0, 1\}^2 \rightarrow \mathbb{R}^d$  makes XOR linearly separable is  $d = 2$ . For example:  $\phi(x_1, x_2) = (x_1 + x_2, x_1 - x_2)$  achieves linear separation via  $w_1 z_1 + w_2 z_2 \geq \theta$  with appropriate weights.

149 **Theorem 7** (Generalized Parity Functions). Define the  $k$ -subset parity family:  $\text{PARITY}_{S,k}(\mathbf{x}) = \bigoplus_{i \in S} x_i$  where  $S \subset [k]$ . **(i) Family Complexity:** The VC-dimension of this family is  $\text{VC} = k$ . **(ii) Linear Representation:** For uniform distribution over  $\{0, 1\}^k$ , the expected dimension required for  $\varepsilon$ -linear representation of random  $\text{PARITY}_{S,k}$  is:  $\mathbb{E}[d] = \Theta\left(\frac{k}{\varepsilon^2 \log(1/\varepsilon)}\right)$ . **(iii) Worst-Case:** There exists choice of  $S$  requiring:  $d = \Omega\left(\frac{2^k}{\sqrt{k}}\right)$

154 **5.3 Polynomial and Rational Functions**

155 **Theorem 8** (Polynomial Function Analysis). Consider polynomial concepts  $p : \mathbb{R}^n \rightarrow \{0, 1\}$  defined by  $p(\mathbf{x}) = \mathbf{1}\{Q(\mathbf{x}) \geq 0\}$  where  $Q$  is degree- $d$  polynomial.

157 **(i) Linear Representability:** If  $Q$  has  $m$  monomials and  $p$  is linearly representable in  $\phi : \mathbb{R}^n \rightarrow \mathbb{R}^D$ , then either  $D \geq m - 1$ , or  $\phi$  computes polynomial features explicitly.

159 **(ii) Specific Cases:** **Quadratic:**  $x_1^2 + x_2^2 \geq 1$  requires  $D \geq 2$ , **Product:**  $x_1 x_2 \geq 0$  requires  $D \geq 1$  but representation must compute products, **High-degree:** Degree- $d$  polynomials require  $D = O(n^d)$  in worst case.

## 162 6 Information-Theoretic Analysis

163 We now provide information-theoretic lower bounds on representation dimension.

### 164 6.1 Mutual Information Bounds

165 **Definition 9** (Concept Information Content). *For concept  $c$  and representation  $\phi$ , define:  $I(c; \phi) = \sup_{\text{linear } h} I(c; h(\phi))$  where  $I(\cdot; \cdot)$  denotes mutual information.*

167 **Theorem 9** (Information-Theoretic Lower Bound). *For concept  $c$  with Shannon entropy  $H(c) = h$  bits, any representation  $\phi : \mathcal{X} \rightarrow \mathbb{R}^d$  achieving  $\varepsilon$ -linear representability must satisfy:  $d \geq \frac{h - H(\varepsilon)}{2 \log(1 + \text{SNR})}$  where SNR is the effective signal-to-noise ratio of  $\phi$ .*

170 *Proof.* Using the data processing inequality, the mutual information between  $c$  and any linear function  
171 of  $\phi$  is bounded by the channel capacity of the representation.

172 For  $d$ -dimensional representation with bounded coordinates, the capacity is approximately  $d \log(1 +$   
173 SNR). The  $\varepsilon$ -error requires preserving at least  $h - H(\varepsilon)$  bits of information about the concept.  
174 Combining these:  $d \log(1 + \text{SNR}) \geq h - H(\varepsilon)$ .  $\square$

175 **Corollary 3** (High-Entropy Concepts). *For concepts with near-maximal entropy ( $H(c) \approx \log |\mathcal{X}|$ ),*

176 *linear representation requires:  $d = \Omega\left(\frac{\log |\mathcal{X}|}{\log \text{SNR}}\right)$*

### 177 6.2 Rate-Distortion Analysis

178 **Theorem 10** (Rate-Distortion for Linear Concepts). *Define the linear rate-distortion function:  
179  $R_{\text{linear}}(D) = \inf_{p(\hat{\mathbf{z}}|\mathbf{z}): \mathbb{E}[\|\mathbf{z} - \hat{\mathbf{z}}\|^2] \leq D} I(\mathbf{z}; \hat{\mathbf{z}})$  where the infimum is over linear encoders of  $d$ -  
180 dimensional representations. For Gaussian representations,  $R_{\text{linear}}(D) = \frac{d}{2} \log\left(\frac{\sigma^2}{D}\right)$  where  $\sigma^2$   
181 is the variance. Concepts requiring high-rate encoding cannot be linearly recovered from low-  
182 dimensional representations.*

## 183 7 Measure-Theoretic Results

184 We analyze the generic behavior of linear representability using measure theory.

### 185 7.1 Random Representation Analysis

186 **Theorem 11** (Generic Failure of Linear Representability). *Let  $\phi$  be a random Gaussian representa-  
187 tion:  $\phi(\mathbf{x}) \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_d)$  independently for each  $\mathbf{x}$ . For concept  $c$  and dataset  $S = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ :*

188 **(i) Probability of Linear Separability:**  $\mathbb{P}[c \text{ linearly separable in } \phi] =$   
189  $\frac{2^{d+1}}{\pi^{d/2} \Gamma(d/2+1)} \int_0^\infty r^d e^{-r^2/2\sigma^2} \Phi(r)^{n-d-1} dr$  where  $\Phi$  is the standard Gaussian CDF.

190 **(ii) Phase Transition:** There exists critical ratio  $\alpha_c \approx 0.833$  such that: If  $d/n > \alpha_c$ , then  
191  $\mathbb{P}[\text{separable}] \rightarrow 1$  as  $n \rightarrow \infty$ . If  $d/n < \alpha_c$ , then  $\mathbb{P}[\text{separable}] \rightarrow 0$  as  $n \rightarrow \infty$ .

192 **(iii) Concentration:** The transition is sharp: for any  $\delta > 0$ ,  $\mathbb{P}\left[\left|\frac{d}{n} - \alpha_c\right| > \delta\right] \leq 2 \exp(-cn\delta^2)$  for  
193 some constant  $c > 0$ .

194 *Proof.* **Part (i):** For random Gaussian features, linear separability depends on the geometric arrange-  
195 ment of projected points. Using results from [10] on Gaussian processes, the probability involves  
196 integration over the distribution of inter-point distances in the projected space.

197 **Part (ii):** The phase transition follows from the asymptotic analysis of random matrix theory. When  
198  $d/n > \alpha_c$ , the feature space has sufficient dimensionality to separate most point configurations. The  
199 critical value  $\alpha_c$  comes from the solution to:  $\int_0^{\alpha_c} \sqrt{2\pi t} e^{-1/(2t)} dt = 1$

200 **Part (iii):** Concentration follows from Talagrand's inequality applied to the separability function,  
201 which has bounded differences property.  $\square$

202 **7.2 Robustness to Perturbations**

203 **Theorem 12** (Stability of Non-Linear Concepts). *Let  $c$  be a concept that is not linearly separable  
204 in representation  $\phi$ . For perturbation  $\tilde{\phi}(\mathbf{x}) = \phi(\mathbf{x}) + \epsilon(\mathbf{x})$  where  $\|\epsilon(\mathbf{x})\| \leq \delta$ : (i) **Perturbation  
205 Bound:** If  $c$  requires margin  $\gamma > 0$  for linear separation in  $\phi$ , then it remains non-linearly separable  
206 in  $\tilde{\phi}$  provided  $\delta < \gamma/2$ . (ii) **Generic Robustness:** For random perturbations  $\epsilon(\mathbf{x}) \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ :  
207  $\mathbb{P}[c \text{ becomes linearly separable in } \tilde{\phi}] \leq \exp\left(-\frac{\gamma^2}{8\sigma^2}\right)$*

208 **8 Advanced Constructions and Extensions**

209 **8.1 Hierarchical Concept Families**

210 **Definition 10** (Hierarchical Parity). *Define  $k$ -level hierarchical parity:  $H\text{PARITY}_k(\mathbf{x}) =$   
211  $\text{PARITY}_{k-1}(\text{PARITY}_1(\mathbf{x}^{(1)}), \dots, \text{PARITY}_1(\mathbf{x}^{(2^{k-1}))}, \mathbf{x}$  is partitioned into blocks  $\mathbf{x}^{(i)}$  of size  $2^{k-1}$ .*

212 **Theorem 13** (Hierarchical Impossibility). *For  $k$ -level hierarchical parity on  $n = 2^k$  variables:*

213 *(i) Depth Requirement:* Any circuit computing  $H\text{PARITY}_k$  requires depth  $\geq k$ .

214 *(ii) Linear Representation:* Any intermediate representation at depth  $< k$  making  $H\text{PARITY}_k$  linearly  
215 decodable requires dimension:  $d \geq 2^{2^{k-1}-1}$

216 *(iii) Network Implication:* Deep networks computing hierarchical concepts cannot have linear  
217 intermediate representations of polynomial dimension.

218 *Proof.* The proof proceeds by induction on  $k$ . The base case  $k = 1$  reduces to standard parity. For the  
219 inductive step, suppose level  $(k - 1)$  requires the stated dimension. Then level  $k$  requires computing  
220 parity over  $2^{k-1}$  intermediate results, each requiring exponential representation, leading to doubly  
221 exponential dimension requirement.  $\square$

222 **8.2 Continuous and Mixed-Type Concepts**

223 **Theorem 14** (Continuous Concept Analysis). *For continuous concepts  $c : \mathbb{R}^n \rightarrow [0, 1]$  and represen-  
224 tations  $\phi : \mathbb{R}^n \rightarrow \mathbb{R}^d$ :*

225 *(i) Lipschitz Constraint:* If  $c$  is  $L$ -Lipschitz and linearly representable in  $\phi$  with  $\varepsilon$ -error, then:  
226  $d \geq \frac{L^2 \text{diam}(\mathcal{X})^2}{4\varepsilon^2}$  where  $\text{diam}(\mathcal{X})$  is the diameter of the input space.

227 *(ii) Smooth Concept Classes:* For  $C^r$  concepts with bounded derivatives up to order  $r$ , the dimension  
228 requirement scales as:  $d = \Omega\left(\left(\frac{1}{\varepsilon}\right)^{n/r}\right)$

229 *(iii) Oscillatory Concepts:* Concepts with high-frequency components require exponentially large  
230 representations for linear decodability.

231 **9 Experimental Validation and Reproducible Protocols**

232 We provide comprehensive experimental protocols to validate our theoretical predictions. Please see  
233 Appendix B for synthetic, naturalistic and language model experiments.

234 **10 Conditions for Approximate Linear Representability**

235 Despite our negative results, we identify precise conditions enabling approximate LRH.

236 **10.1 Architectural Inductive Biases**

237 **Theorem 15** (Architectural Bias for Linearity). *Consider networks with architectural constraints  
238 encouraging linear concept development:*

- 239 (i) **Bottleneck Architectures:** Networks with severe dimension reduction layers force concepts into  
 240 linear subspaces. If layer  $\ell$  has dimension  $d \ll$  input complexity, then concepts become approximately  
 241 linear with error:  $\varepsilon \leq O\left(\sqrt{\frac{\text{concept complexity}}{d}}\right)$
- 242 (ii) **Attention Mechanisms:** Self-attention with linear value projections naturally creates linear  
 243 concept combinations:  $\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right)\mathbf{V}$  encourages linear decodability of  
 244 concepts encoded in value vectors.
- 245 (iii) **Skip Connections:** Residual connections  $\mathbf{h}_{l+1} = \mathbf{h}_l + f(\mathbf{h}_l)$  preserve linear information across  
 246 layers, enabling linear concept tracking.

## 247 10.2 Training Dynamics and Implicit Regularization

- 248 **Theorem 16** (SGD Bias Toward Linearity). *Under certain conditions, SGD implicitly biases repre-  
 249 sentations toward linear concept encodings:*
- 250 (i) **Overparameterized Regime:** In the NTK limit with infinite width, networks learn approximately  
 251 linear functions of features, promoting linear concept decodability.
- 252 (ii) **Early Stopping:** Stopping training before convergence often preserves linear structure that would  
 253 otherwise be optimized away.
- 254 (iii) **Regularization Effects:**  $L_2$  weight decay and dropout create implicit pressure toward simpler,  
 255 more linear representations.
- 256 **Quantitative Bound:** Under these conditions, concepts are  $\varepsilon$ -linearly decodable with:  $\varepsilon \leq$   
 257  $\frac{C\sqrt{\log(d/\delta)}}{\sqrt{m}}$  where  $m$  is the number of training examples,  $d$  is dimension, and  $C$  depends on the  
 258 regularization strength.

## 259 10.3 Data Structure and Natural Concepts

- 260 **Theorem 17** (Natural Data Promotes Linear Concepts). *Real-world datasets often have structure  
 261 promoting approximate linear concept encodings:*
- 262 (i) **Low Intrinsic Dimension:** If data lies near a  $k$ -dimensional manifold with  $k \ll d$ , then concepts  
 263 aligned with the manifold structure become approximately linear.
- 264 (ii) **Hierarchical Structure:** Datasets with natural hierarchies (e.g., ImageNet taxonomy) allow  
 265 concepts at each level to be linearly separated within the corresponding subspace.
- 266 (iii) **Smoothness Assumptions:** If concepts vary smoothly over the data manifold, local linear  
 267 approximations become globally valid:  $|c(\mathbf{x}) - \mathbf{w}^T \phi(\mathbf{x})| \leq L \cdot \text{curvature}(\mathcal{M}) \cdot \text{diameter}(\mathcal{M})^2$

## 268 11 Implications for Interpretability Methods

- 269 Our theoretical analysis has profound implications for interpretation techniques.

### 270 Linear Probe Limitations and Best Practices

- 271 **Theorem 18** (Interpretation of Probe Results). *For linear probe achieving accuracy  $A$  on concept  $c$ :*
- 272 (i) **Lower Bound on Network Usage:** The network uses concept  $c$  with strength at least:  $\text{Usage}(c) \geq$   
 273  $\max\left(0, \frac{2A-1}{\sqrt{d/n}}\right)$  where  $d$  is representation dimension and  $n$  is dataset size.
- 274 (ii) **Upper Bound Limitations:** High probe accuracy does NOT imply: the concept is the primary  
 275 computational mechanism, the concept is causally relevant for network decisions, the representation  
 276 is interpretable or disentangled.
- 277 (iii) **Failure Interpretation:** Low probe accuracy does NOT imply: the network doesn't use the concept,  
 278 the concept is absent from internal computations, the representation lacks relevant information.

279 **Alternative Interpretability Approaches** Given LRH limitations, we recommend complementary  
280 interpretation methods: **Non-linear Probes**: Use polynomial or neural network probes to capture non-  
281 linear concept encodings, **Geometric Analysis**: Study representational geometry using: Persistent  
282 homology [15], Manifold learning techniques, Clustering and density analysis, **Causal Intervention**:  
283 Use activation patching and causal scrubbing [4] to test concept relevance, **Feature Synthesis**:  
284 Generate synthetic inputs maximizing concept activations, **Circuit Analysis**: Identify computational  
285 subgraphs implementing specific functions [17].

286 **Method-Specific Recommendations** **Concept Bottleneck Models** [13]: Our results suggest these  
287 work best for naturally linear concepts. For non-linear concepts, use non-linear bottlenecks or hierar-  
288 chical concept structures. **Activation Patching**: More robust than linear probes since it tests causal  
289 relevance rather than linear decodability. **Gradient-Based Methods**: Saliency maps and integrated  
290 gradients [19] can capture non-linear concept usage patterns. **Mechanistic Interpretability**: Focus  
291 on identifying computational circuits rather than assuming linear concept encodings.

## 292 **12 Open Questions and Future Directions**

293 **Theoretical Questions** **Tighter Bounds**: Can we achieve matching upper and lower bounds for  
294 specific concept families?, **Average-Case Analysis**: Most results are worst-case. What about typical  
295 concept-representation pairs?, **Algorithmic Aspects**: Given a representation, how efficiently can we  
296 determine if a concept is linearly decodable?, **Multi-Task Settings**: How does linear decodability  
297 change when networks learn multiple related concepts?, **Continual Learning**: How does concept  
298 linearity evolve as networks learn new tasks?

299 **Empirical Investigations** **Large-Scale Studies**: Systematic analysis of linear vs. non-linear  
300 concepts across diverse architectures and datasets, **Intervention Experiments**: Test predictions  
301 by directly manipulating network architectures and training procedures, **Cross-Modal Analysis**:  
302 Compare concept linearity across vision, language, and multimodal models, **Scaling Laws**: How do  
303 our results change with model size, data scale, and compute budget?

304 **Methodological Development** **Non-Linear Probes**: Develop principled methods for non-linear  
305 concept detection, **Representational Metrics**: Create measures of concept complexity and repre-  
306 sentational geometry, **Hybrid Approaches**: Combine linear and non-linear interpretation methods  
307 optimally, **Uncertainty Quantification**: Develop confidence intervals for interpretability claims.

## 308 **13 Conclusion**

309 We have provided a comprehensive theoretical analysis demonstrating fundamental limitations of  
310 the Linear Representation Hypothesis. Through combinatorial arguments, circuit complexity theory,  
311 explicit constructions, information-theoretic analysis, and measure-theoretic results, we have shown  
312 that LRH cannot be universally true.

313 Our key contributions are: **Sharp Impossibility Results**: exponential gaps between concept com-  
314 plexity and linear representability, **Circuit-Theoretic Analysis**: proof that depth-separated functions  
315 cannot have linear intermediate representations, **Constructive Examples**: explicit concept fami-  
316 lies demonstrating non-linear encodings, **Information Bounds**: dimension requirements for linear  
317 concept recovery, **Practical Guidelines**: conditions enabling approximate linear representability.

318 **Broader Impact** Our results have significant implications for **Interpretability Research**: methods  
319 beyond linear probes are needed, **AI Safety**: understand when simple interpretability methods fail,  
320 **Network Design**: architectures that promote concept linearity when desired, **Scientific Under-  
321 standing**: theoretical foundations for representational analysis. While linear probes remain valuable tools,  
322 they must be used with awareness of their fundamental limitations. The future of interpretability lies  
323 in developing robust methods that can handle the full complexity of neural network representations,  
324 combining linear and non-linear approaches as appropriate for each specific context. Our theoretical  
325 framework provides the foundation for this next generation of interpretability methods, ensuring they  
326 are built on solid mathematical principles rather than unverified assumptions.

327 **References**

- 328 [1] Alain, G. and Bengio, Y. (2016). Understanding intermediate layers using linear classifier  
329 probes. *arXiv preprint arXiv:1610.01644*.
- 330 [2] Bartlett, P. L. and Shawe-Taylor, J. (1998). Generalization performance of support vector  
331 machines and other pattern classifiers. *Advances in Kernel Methods*, 43–54.
- 332 [3] Belinkov, Y. and Glass, J. (2017). Analysis methods in neural language processing: A survey.  
333 *Transactions of the Association for Computational Linguistics*, 7, 49–72.
- 334 [4] Chan, L., Garriga-Alonso, A., Goldowsky-Dill, N., Greenblatt, R., Nitishinskaya, J., Radhakr-  
335 ishnan, A., Shlegeris, B., and Thomas, N. (2022). Causal scrubbing: a method for rigorously  
336 testing interpretability hypotheses. *arXiv preprint arXiv:2211.00032*.
- 337 [5] Conneau, A., Kruszewski, G., Lample, G., Barrault, L., and Baroni, M. (2018). What you can  
338 cram into a single vector: Probing sentence embeddings for linguistic properties. *Proceedings*  
339 *of the 56th Annual Meeting of the Association for Computational Linguistics*, 2126–2136.
- 340 [6] Cover, T. M. (1965). Geometrical and statistical properties of systems of linear inequalities with  
341 applications in pattern recognition. *IEEE Transactions on Electronic Computers*, (3), 326–334.
- 342 [7] Doshi-Velez, F. and Kim, B. (2017). Towards a rigorous science of interpretable machine  
343 learning. *arXiv preprint arXiv:1702.08608*.
- 344 [8] Eldan, R. and Shamir, O. (2016). The power of depth for feedforward neural networks. *Confer-  
345 ence on Learning Theory*, 907–940.
- 346 [9] Gilpin, L. H., Bau, D., Yuan, B. Z., Bajwa, A., Specter, M., and Kagal, L. (2018). Explaining  
347 explanations: An overview of interpretability of machine learning. *2018 IEEE 5th International  
348 Conference on Data Science and Advanced Analytics*, 80–89.
- 349 [10] Gordon, Y. (1988). On Milman’s inequality and random subspaces which escape through a  
350 mesh in  $\mathbb{R}^n$ . *Geometric Aspects of Functional Analysis*, 84, 84–106.
- 351 [11] Håstad, J. (1986). Almost optimal lower bounds for small depth circuits. *Proceedings of the  
352 eighteenth annual ACM symposium on Theory of computing*, 6–20.
- 353 [12] Hewitt, J. and Manning, C. D. (2019). A structural probe for finding syntax in word representa-  
354 tions. *Proceedings of the 2019 Conference of the North American Chapter of the Association  
355 for Computational Linguistics*, 4129–4138.
- 356 [13] Koh, P. W., Nguyen, T., Tang, Y. S., Mussmann, S., Pierson, E., Kim, B., and Liang, P. (2020).  
357 Concept bottleneck models. *International Conference on Machine Learning*, 5338–5348.
- 358 [14] Lipton, Z. C. (2018). The mythos of model interpretability. *Communications of the ACM*, 61(10),  
359 36–43.
- 360 [15] Naitzat, G., Zhitnikov, A., and Lim, L. H. (2020). Topology of deep neural networks. *The  
361 Journal of Machine Learning Research*, 21(1), 7503–7542.
- 362 [16] O’Donnell, R. and Servedio, R. A. (2003). Learning monotone decision trees in polynomial  
363 time. *SIAM Journal on Computing*, 37(3), 827–844.
- 364 [17] Olah, C., Cammarata, N., Schubert, L., Goh, G., Petrov, M., and Carter, S. (2020). Zoom in: An  
365 introduction to circuits. *Distill*, 5(3), e00024–001.
- 366 [18] Rossman, B. (2008). On the constant-depth complexity of k-clique. *Proceedings of the fortieth  
367 annual ACM symposium on Theory of computing*, 721–730.
- 368 [19] Sundararajan, M., Taly, A., and Yan, Q. (2017). Axiomatic attribution for deep networks.  
369 *International Conference on Machine Learning*, 3319–3328.
- 370 [20] Telgarsky, M. (2016). Benefits of depth in neural networks. *Conference on Learning Theory*,  
371 1517–1539.

372 [21] Vapnik, V. N. and Chervonenkis, A. Y. (1971). On the uniform convergence of relative frequencies  
 373 of events to their probabilities. *Theory of Probability & Its Applications*, 16(2), 264–280.

374 [22] Vig, J., Gehrman, S., Belinkov, Y., Qian, S., Nevo, D., Singer, Y., and Shieber, S. (2020).  
 375 Causal mediation analysis for interpreting neural NLP: The case of gender bias. *arXiv preprint*  
 376 *arXiv:2004.12265*.

## 377 A Combinatorial Impossibility

378 *Proof.* Theorem 1.

379 **Part (i):** From Lemma 1, when  $n > d + 1$ :  $C(n, d) = 2 \sum_{i=0}^d \binom{n-1}{i} \leq 2(d+1) \binom{n-1}{d}$  Using  
 380 the bound  $\binom{n-1}{d} \leq \left(\frac{e(n-1)}{d}\right)^d$ :  $C(n, d) \leq 2(d+1) \left(\frac{e(n-1)}{d}\right)^d \leq 2^{d+1} \frac{n^d}{d!}$ . Therefore:  $\rho(n, d) \leq$   
 381  $\frac{2^{d+1} n^d}{d! \cdot 2^n}$ .

382 **Part (ii):** When  $d = o(n/\log n)$ , using Stirling's approximation  $d! \geq \sqrt{2\pi d}(d/e)^d$ :  $\rho(n, d) \leq$   
 383  $\frac{2^{d+1} n^d}{\sqrt{2\pi d}(d/e)^d \cdot 2^n} = \frac{2^{d+1} (en/d)^d}{\sqrt{2\pi d} \cdot 2^n}$ . Taking logarithms:  $\log \rho(n, d) \leq (d+1) \log 2 + d \log(en/d) -$   
 384  $\frac{1}{2} \log(2\pi d) - n \log 2$ . When  $d = o(n/\log n)$ , the dominant term is  $-n \log 2$ , so  $\rho(n, d) \rightarrow 0$ .

385 **Part (iii):** For  $d < (1-\varepsilon) \log_2 n$  and for sufficiently large  $n$ , we have  $2^d < n^{1-\varepsilon}$ , thus:

$$386 \rho(n, d) \leq \frac{2^{d+1} n^d}{2^n} \leq \frac{2n^{1-\varepsilon+d/\log_2 n}}{2^n} \leq \frac{2n^{2-\varepsilon}}{2^n} < 2^{-\varepsilon n} \quad \square$$

## 387 B Experimental Validation and Reproducible Protocols

### 388 B.1 Synthetic Experiments

---

#### Algorithm 1 Comprehensive Parity Experiment

---

```

1: Parameters:  $k \in \{5, 10, 15, 20\}$ ,  $N_{\text{train}} = 50000$ ,  $N_{\text{test}} = 10000$ 
2: Architecture: MLP with layers  $[k, 512, 256, 128, 1]$ , ReLU activations
3: Training: Adam optimizer, lr =  $10^{-3}$ , batch size 256, 1000 epochs
4: for each  $k$  do
5:   Generate dataset:  $\mathbf{x}_i \sim \text{Bernoulli}(0.5)^k$ ,  $y_i = \text{PARITY}_k(\mathbf{x}_i)$ 
6:   Train network until convergence (loss <  $10^{-6}$ )
7:   Evaluate full network accuracy  $A_{\text{full}}$ 
8:   for each hidden layer  $\ell \in \{1, 2, 3\}$  do
9:     Extract representations  $\{\phi_\ell(\mathbf{x}_i)\}$ 
10:    Train linear probe:  $\mathbb{R}^{d_\ell} \rightarrow \{0, 1\}$  using logistic regression
11:    Evaluate probe accuracy  $A_{\text{probe}}^{(\ell)}$ 
12:    Compute linear separability score:  $S_\ell = A_{\text{probe}}^{(\ell)} / A_{\text{full}}$ 
13:   end for
14:   Record dimension vs. performance trade-offs
15: end for
16: Expected Results:  $A_{\text{full}} > 0.99$ ,  $A_{\text{probe}}^{(\ell)} \approx 0.5$  for  $\ell < 3$ 
```

---

---

**Algorithm 2** Depth-Separation Validation

---

- 1: **Function Class:** Iterated multiplication  $f(\mathbf{x}) = \prod_{i=1}^k x_i \bmod 2$
- 2: **Networks:** Compare depth-2 vs depth- $\lceil \log k \rceil$  architectures
- 3: **for** each depth  $d \in \{2, 3, \lceil \log k \rceil\}$  **do**
- 4:   Train network with maximum width  $W = 1000$
- 5:   Measure: (1) Training convergence, (2) Test accuracy, (3) Linear probe performance
- 6:   Record computational requirements and representational properties
- 7: **end for**
- 8: **Prediction:** Shallow networks fail; deep networks succeed but with non-linear concepts

---

389 **B.2 Naturalistic Experiments**

---

**Algorithm 3** Vision Task with Compositional Concepts

---

- 1: **Dataset:** CLEVR-style synthetic scenes with compositional attributes
- 2: **Task:** Classify presence of "red cube AND blue sphere" (compositional concept)
- 3: **Network:** ResNet-18 pretrained on ImageNet, fine-tuned on task
- 4: Train to high accuracy on main task
- 5: **for** each layer  $\ell$  in  $\{3, 6, 9, 12, 15, 18\}$  **do**
- 6:   Extract features  $\phi_\ell(\mathbf{x}) \in \mathbb{R}^{d_\ell}$
- 7:   Test linear probes for:
  - Primitive concepts: "red", "cube", "blue", "sphere"
  - Compositional concept: "red cube AND blue sphere"
  - Control concepts: unrelated scene properties
- 8:   Measure probe accuracy vs. layer depth
- 9: **end for**
- 10: **Expected Pattern:** Primitive concepts become linear earlier; compositional concepts require deeper layers

---

390 **B.3 Language Model Experiments**

---

**Algorithm 4** Syntactic vs. Semantic Concept Probing

---

- 1: **Model:** GPT-2 or BERT-base
- 2: **Tasks:** POS tagging (syntactic) vs. sentiment analysis (semantic)
- 3: **Concepts:** Extract intermediate representations for both task types
- 4: **for** each layer  $\ell$  **do**
- 5:   Train linear probes for:
  - Syntactic: POS tags, dependency relations, grammatical number
  - Semantic: Sentiment, topic classification, factual knowledge
  - Complex: Coreference resolution, logical inference
- 6:   Compare probe accuracies across concept types and layers
- 7: **end for**
- 8: **Hypothesis:** Syntactic concepts are more linearly decodable; semantic/logical concepts require non-linear processing

---

391 **Agents4Science AI Involvement Checklist**

392 This checklist is designed to allow you to explain the role of AI in your research. This is important for  
393 understanding broadly how researchers use AI and how this impacts the quality and characteristics  
394 of the research. **Do not remove the checklist! Papers not including the checklist will be desk**  
395 **rejected.** You will give a score for each of the categories that define the role of AI in each part of the  
396 scientific process. The scores are as follows:

- 397 • **[A] Human-generated:** Humans generated 95% or more of the research, with AI being of  
398 minimal involvement.
- 399 • **[B] Mostly human, assisted by AI:** The research was a collaboration between humans and  
400 AI models, but humans produced the majority (>50%) of the research.
- 401 • **[C] Mostly AI, assisted by human:** The research task was a collaboration between humans  
402 and AI models, but AI produced the majority (>50%) of the research.
- 403 • **[D] AI-generated:** AI performed over 95% of the research. This may involve minimal  
404 human involvement, such as prompting or high-level guidance during the research process,  
405 but the majority of the ideas and work came from the AI.

406 These categories leave room for interpretation, so we ask that the authors also include a brief  
407 explanation elaborating on how AI was involved in the tasks for each category. Please keep your  
408 explanation to less than 150 words.

- 409 1. **Hypothesis development:** Hypothesis development includes the process by which you  
410 came to explore this research topic and research question. This can involve the background  
411 research performed by either researchers or by AI. This can also involve whether the idea  
412 was proposed by researchers or by AI.

413 Answer: **[A]**

414 Explanation: The hypothesis that the linear representation hypothesis does not universally  
415 hold in neural networks is entirely human generated.

- 416 2. **Experimental design and implementation:** This category includes design of experiments  
417 that are used to test the hypotheses, coding and implementation of computational methods,  
418 and the execution of these experiments.

419 Answer: **[D]**

420 Explanation: The AI did the theoretical exploration of the research hypotheses from var-  
421 ious angles (combinatorial arguments, circuit complexity theory, explicit constructions,  
422 information-theoretic analysis, and measure-theoretic results), ultimately showing that LRH  
423 cannot be universally true.

- 424 3. **Analysis of data and interpretation of results:** This category encompasses any process to  
425 organize and process data for the experiments in the paper. It also includes interpretations of  
426 the results of the study.

427 Answer: **[C]**

428 Explanation: It was mostly done by AI, with the human verifying and fact-checking the  
429 claims made.

- 430 4. **Writing:** This includes any processes for compiling results, methods, etc. into the final  
431 paper form. This can involve not only writing of the main text but also figure-making,  
432 improving layout of the manuscript, and formulation of narrative.

433 Answer: **[C]**

434 Explanation: It was mostly the AI doing the writing, with human involvement in terms of  
435 prompting or high-level guidance during the research process.

- 436 5. **Observed AI Limitations:** What limitations have you found when using AI as a partner or  
437 lead author?

438 Description: Difficulty in steering these models and generating texts with given constraints.

439 **Agents4Science Paper Checklist**

440 The checklist is designed to encourage best practices for responsible machine learning research,  
441 addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove  
442 the checklist: **Papers not including the checklist will be desk rejected.** The checklist should  
443 follow the references and follow the (optional) supplemental material. The checklist does NOT count  
444 towards the page limit.

445 Please read the checklist guidelines carefully for information on how to answer these questions. For  
446 each question in the checklist:

- 447 • You should answer [Yes] , [No] , or [NA] .
- 448 • [NA] means either that the question is Not Applicable for that particular paper or the  
449 relevant information is Not Available.
- 450 • Please provide a short (1–2 sentence) justification right after your answer (even for NA).

451 **The checklist answers are an integral part of your paper submission.** They are visible to the  
452 reviewers and area chairs. You will be asked to also include it (after eventual revisions) with the final  
453 version of your paper, and its final version will be published with the paper.

454 The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation.  
455 While "[Yes]" is generally preferable to "[No]", it is perfectly acceptable to answer "[No]" provided  
456 a proper justification is given. In general, answering "[No]" or "[NA]" is not grounds for rejection.  
457 While the questions are phrased in a binary way, we acknowledge that the true answer is often more  
458 nuanced, so please just use your best judgment and write a justification to elaborate. All supporting  
459 evidence can appear either in the main paper or the supplemental material, provided in appendix.  
460 If you answer [Yes] to a question, in the justification please point to the section(s) where related  
461 material for the question can be found.

462 **1. Claims**

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

465 Answer: [Yes]

466 Justification: The abstract and introduction argue against the universality of the linear  
467 representation hypothesis, while the claims made by the paper theoretically support these  
468 arguments with theretically proven results.

469 Guidelines:

- 470 • The answer NA means that the abstract and introduction do not include the claims  
471 made in the paper.
- 472 • The abstract and/or introduction should clearly state the claims made, including the  
473 contributions made in the paper and important assumptions and limitations. A No or  
474 NA answer to this question will not be perceived well by the reviewers.
- 475 • The claims made should match theoretical and experimental results, and reflect how  
476 much the results can be expected to generalize to other settings.
- 477 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
478 are not attained by the paper.

479 **2. Limitations**

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

481 Answer: [Yes]

482 Justification: The paper is primarily concerned with limitations of the linear representation  
483 hypotheses and linear probing. We provide concrete recommendations for alternative  
484 interpretability approaches, discuss remaining open questions and future directions.

485 Guidelines:

- 486 • The answer NA means that the paper has no limitation while the answer No means that  
487 the paper has limitations, but those are not discussed in the paper.

- 488           • The authors are encouraged to create a separate "Limitations" section in their paper.  
 489           • The paper should point out any strong assumptions and how robust the results are to  
 490           violations of these assumptions (e.g., independence assumptions, noiseless settings,  
 491           model well-specification, asymptotic approximations only holding locally). The authors  
 492           should reflect on how these assumptions might be violated in practice and what the  
 493           implications would be.  
 494           • The authors should reflect on the scope of the claims made, e.g., if the approach was  
 495           only tested on a few datasets or with a few runs. In general, empirical results often  
 496           depend on implicit assumptions, which should be articulated.  
 497           • The authors should reflect on the factors that influence the performance of the approach.  
 498           For example, a facial recognition algorithm may perform poorly when image resolution  
 499           is low or images are taken in low lighting.  
 500           • The authors should discuss the computational efficiency of the proposed algorithms  
 501           and how they scale with dataset size.  
 502           • If applicable, the authors should discuss possible limitations of their approach to  
 503           address problems of privacy and fairness.  
 504           • While the authors might fear that complete honesty about limitations might be used by  
 505           reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
 506           limitations that aren't acknowledged in the paper. Reviewers will be specifically  
 507           instructed to not penalize honesty concerning limitations.

508           **3. Theory assumptions and proofs**

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

511           Answer: [Yes]

512           Justification: For each claim we provide theoretical proofs, discuss assumptions, limitations  
 513           and best practices.

514           Guidelines:

- 515           • The answer NA means that the paper does not include theoretical results.
- 516           • All the theorems, formulas, and proofs in the paper should be numbered and cross-  
 517           referenced.
- 518           • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 519           • The proofs can either appear in the main paper or the supplemental material, but if  
 520           they appear in the supplemental material, the authors are encouraged to provide a short  
 521           proof sketch to provide intuition.

522           **4. Experimental result reproducibility**

523           Question: Does the paper fully disclose all the information needed to reproduce the main ex-  
 524           perimental results of the paper to the extent that it affects the main claims and/or conclusions  
 525           of the paper (regardless of whether the code and data are provided or not)?

526           Answer: [NA].

527           Justification: The paper is focusing on theory and does not include experiments.

528           Guidelines:

- 529           • The answer NA means that the paper does not include experiments.
- 530           • If the paper includes experiments, a No answer to this question will not be perceived  
 531           well by the reviewers: Making the paper reproducible is important.
- 532           • If the contribution is a dataset and/or model, the authors should describe the steps taken  
 533           to make their results reproducible or verifiable.
- 534           • We recognize that reproducibility may be tricky in some cases, in which case authors  
 535           are welcome to describe the particular way they provide for reproducibility. In the case  
 536           of closed-source models, it may be that access to the model is limited in some way  
 537           (e.g., to registered users), but it should be possible for other researchers to have some  
 538           path to reproducing or verifying the results.

539           **5. Open access to data and code**

540 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
541 tions to faithfully reproduce the main experimental results, as described in supplemental  
542 material?

543 Answer: answerNA

544 Justification: The paper is focusing on theory and does not include experiments requiring  
545 code.

546 Guidelines:

- 547 • The answer NA means that paper does not include experiments requiring code.
- 548 • Please see the Agents4Science code and data submission guidelines on the conference  
549 website for more details.
- 550 • While we encourage the release of code and data, we understand that this might not be  
551 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not  
552 including code, unless this is central to the contribution (e.g., for a new open-source  
553 benchmark).
- 554 • The instructions should contain the exact command and environment needed to run to  
555 reproduce the results.
- 556 • At submission time, to preserve anonymity, the authors should release anonymized  
557 versions (if applicable).

## 558 6. Experimental setting/details

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

562 Answer: [NA]

563 Justification: The paper is focusing on theory and does not include experiments.

564 Guidelines:

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

## 570 7. Experiment statistical significance

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

573 Answer: [NA].

574 Justification: The paper is focusing on theory and does not include experiments.

575 Guidelines:

- 576 • The answer NA means that the paper does not include experiments.
- 577 • The authors should answer "Yes" if the results are accompanied by error bars, confi-  
578 dence intervals, or statistical significance tests, at least for the experiments that support  
579 the main claims of the paper.
- 580 • The factors of variability that the error bars are capturing should be clearly stated  
581 (for example, train/test split, initialization, or overall run with given experimental  
582 conditions).

## 583 8. Experiments compute resources

584 Question: For each experiment, does the paper provide sufficient information on the com-  
585 puter resources (type of compute workers, memory, time of execution) needed to reproduce  
586 the experiments?

587 Answer: [NA]

588 Justification: The paper is focusing on theory and does not include experiments.

589 Guidelines:

- 590           • The answer NA means that the paper does not include experiments.  
591           • The paper should indicate the type of compute workers CPU or GPU, internal cluster,  
592            or cloud provider, including relevant memory and storage.  
593           • The paper should provide the amount of compute required for each of the individual  
594            experimental runs as well as estimate the total compute.

595           **9. Code of ethics**

596           Question: Does the research conducted in the paper conform, in every respect, with the  
597           Agents4Science Code of Ethics (see conference website)?

598           Answer: [Yes]

599           Justification: The paper fully obeys the Agents4Science Code of Ethics.

600           Guidelines:

- 601           • The answer NA means that the authors have not reviewed the Agents4Science Code of  
602            Ethics.  
603           • If the authors answer No, they should explain the special circumstances that require a  
604            deviation from the Code of Ethics.

605           **10. Broader impacts**

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

608           Answer: [Yes]

609           Justification: The paper discusses limitations, alternative approaches and implications for  
610           ensuring interpretability methods are built on solid mathematical principles rather than  
611           unverified assumptions.

612           Guidelines:

- 613           • The answer NA means that there is no societal impact of the work performed.  
614           • If the authors answer NA or No, they should explain why their work has no societal  
615            impact or why the paper does not address societal impact.  
616           • Examples of negative societal impacts include potential malicious or unintended uses  
617           (e.g., disinformation, generating fake profiles, surveillance), fairness considerations,  
618           privacy considerations, and security considerations.  
619           • If there are negative societal impacts, the authors could also discuss possible mitigation  
620           strategies.