
Hidden Progress in Deep Learning: SGD Learns Parities Near the Computational Limit

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 There is mounting empirical evidence of *emergent phenomena* in the capabilities
2 of deep learning methods as we scale up datasets, model sizes, and training times.
3 While there are some accounts of how these resources modulate statistical capacity,
4 far less is known about their effect on the *computational* problem of model training.
5 This work conducts such an exploration through the lens of learning k -sparse
6 parities of n bits, a canonical family of problems which pose theoretical compu-
7 tational barriers. In this setting, we find that neural networks exhibit surprising
8 phase transitions when scaling up dataset size and running time. In particular, we
9 demonstrate empirically that with standard training, a variety of architectures learn
10 sparse parities with $n^{O(k)}$ examples, with loss (and error) curves abruptly dropping
11 after $n^{O(k)}$ iterations. These positive results nearly match known SQ lower bounds,
12 even without an explicit sparsity-promoting prior. We elucidate the mechanisms of
13 these phenomena with a theoretical analysis: we find that the phase transition in
14 performance is *not* due to SGD “stumbling in the dark” until it finds the hidden set
15 of features (a natural algorithm which also runs in $n^{O(k)}$ time); instead, we show
16 that SGD gradually amplifies a *Fourier gap* in the population gradient.

17 1 Introduction

18 Neural networks perform better with more resources (data, model size, training time), but different
19 tasks exhibit qualitatively different dependencies of performance on resources. In particular, while
20 many learning tasks exhibit continuous improvement in performance with increasing resources, other
21 cases show *discontinuous* improvement, where a capability emerges at a certain threshold. Through a
22 statistical lens, it is well-understood that larger models, trained with more data, can fit more complex
23 and expressive functions. However, far less is known about the analogous *computational* question:
24 how does the scaling of these resources influence the success of gradient-based optimization for
25 training these models?

26 These emergent *phase transitions* cannot be explained via statistical capacity alone. In many cases
27 we see a phase transition even when the amount of data remains fixed, with only model size or
28 training time increasing. A timely example is the emergence of reasoning and few-shot learning
29 capabilities when scaling up language models (56, 16, 18, 36). Power et al. (55) give examples
30 exhibiting discontinuous improvements in population accuracy (“grokking”) when running time
31 increases, while data and model size remain fixed.

32 In this work, we analyze the computational aspect of scaling in deep learning in a simple synthetic
33 setting which already exhibits discontinuous improvements. Specifically, we consider the *sparse*
34 *parity problem*— where the label is the parity (XOR) of $k \ll n$ bits in a random length- n binary string.
35 This is a canonical problem which is computationally difficult for a range of algorithms, including
36 gradient-based (41) and streaming (44) algorithms. We focus on analyzing the resource measure of

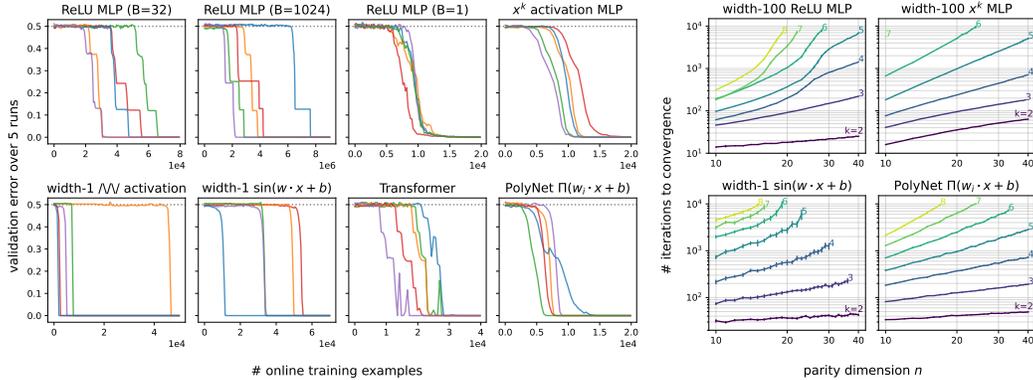


Figure 1: Main empirical findings at a glance. A variety of neural networks, with standard training and initialization, can solve the (n, k) -parity learning problem, with a number of iterations scaling as $n^{O(k)}$. *Left*: Training curves under various algorithmic choices (architecture, batch size, learning rate) on the $(n = 50, k = 3)$ -parity problem. *Right*: Median convergence times for small (n, k) .

37 *training time*, and demonstrate that the loss curves for sparse parities display a phase transition across
 38 a variety of architectures and hyperparameters (see Figure 1, left). A striking observation is that SGD
 39 robustly finds the sparse subset (and hence, effectively, reaches 0 error) with a variety of activation
 40 functions and initialization schemes, even with *no* over-parameterization. This robust convergence is
 41 the starting point for our investigation.

42 Perhaps the most natural hypothesis to explain SGD’s success in learning parities would be that it
 43 simply “stumbles in the dark”, essentially performing random search for the unknown target (e.g.
 44 via stochastic gradient Langevin dynamics). If that were the case, we might expect to observe a
 45 convergence time of $2^{\Omega(n)}$, like a naive search over parameters or subsets of indices. However,
 46 Figure 1 (*right*), already provides some evidence against this “random search” hypothesis: the
 47 convergence time is closer to an $n^{O(k)}$ scaling. Notably, such a convergence rate implies that SGD is
 48 closer to achieving the *optimal computation time* among a natural class of algorithms, namely SQ
 49 (statistical query) algorithms.

50 Through an extensive empirical analysis of the scaling behavior of a variety of models, as well as
 51 theoretical analysis, we give strong evidence against the “stumbling in the dark” viewpoint. Instead,
 52 there is a *hidden progress measure* under which SGD is steadily improving. Furthermore, and
 53 perhaps surprisingly, we show that SGD achieves a computational runtime much closer to the optimal
 54 SQ lower bound than simply doing (non-sparse) parameter search. More generally, our theoretical
 55 and empirical investigations reveal a number of notable phenomena regarding the dependence of
 56 SGD’s performance on resources, and we provide further theoretical and empirical evidence of phase
 57 transitions with data, model size, and training time.

58 1.1 Our contributions

59 **SGD efficiently learns sparse parities, in theory and practice.** It is known from SQ lower bounds
 60 that (noisy) gradient descent on *any* architecture requires at least $n^{\Omega(k)}$ computational steps to
 61 learn k -sparse n -dimensional parities (for background, see Appendix A). However, with standard
 62 architectures and initialization, which do not explicitly encode a sparsity prior, one may expect SGD’s
 63 performance to be much worse, on the order of $2^{\Omega(n)}$. We give extensive empirical evidence that this
 64 is not the case.

65 **Empirical Finding 1.** *For all small instances $(n \leq 30, k \leq 4)$ of the sparse parity problem,*
 66 *architectures $\mathcal{A} \in \{2\text{-layer MLPs, Transformers}^1, \text{sinusoidal (and other non-standard) neurons,}$*
 67 *PolyNets²*, *initializations in $\{\text{uniform, Gaussian, Bernoulli}\}$, and batch sizes $1 \leq B \leq 1024$, SGD*
 68 *on \mathcal{A} solves the (n, k) -sparse parity problem (w.p. ≥ 0.2) within at most $c \cdot n^{\alpha k}$ steps, for small*
 69 *architecture-dependent constants c, α .*

¹With a smaller range of hyperparameters.

²A non-standard architecture introduced in this work; see Section 3 for the definition.

70 For a subset of these architectures, we performed more extensive training, to show scaling behaviors
71 of the computation time;³ see Figure 1 (*right*) and Figure 10 in the appendix.

72 The above results indicate that SGD succeeds at solving parities much faster than the $2^{\Omega(n)}$ steps that
73 would be required by random search without a sparse prior. This suggests that despite flat loss and
74 accuracy curves before the phase transition, SGD makes progress “under the hood”. We show that
75 this is indeed the case by coming up with a *progress measure* that continuously improves throughout
76 training. We can explain the loss and accuracy curves using this progress measure, by showing that
77 they typically rise sharply once the measure passes a certain threshold.

78 **Theoretical Analysis.** Our empirical results suggest that, in a number of computational steps
79 matching the SQ limit, SGD is able to solve the parity problem and identify the influential coordinates,
80 without an explicit sparse prior. We give a theoretical analysis which validates this claim.

81 **Informal Theorem 2.** *On 2-layer MLPs of width $2^{\Theta(k)}$, and with batch size $n^{O(k)}$, SGD converges*
82 *with high probability to a solution with at most ϵ error on the (n, k) -parity problem in at most*
83 *$2^{O(k)} \cdot \text{poly}(1/\epsilon)$ iterations.*

84 We also present a stronger analysis for an idealized architecture (which we call the *disjoint-PolyNet*),
85 which allows for any batch size, and captures the phase transitions observed in the error curves.

86 **Informal Theorem 3.** *On disjoint-PolyNets, SGD (with any batch size $B \geq 1$) converges with high*
87 *probability to a solution with at most ϵ error on the (n, k) -parity problem in at most $n^{O(k)} \cdot \log(1/\epsilon)$*
88 *iterations. Continuous-time gradient flow exhibits a phase transition: it spends a $1 - o(1)$ fraction of*
89 *its time before convergence with error $\geq 49\%$.*

90 In addition to refuting the alternative “random search” hypothesis, our work also poses a counterex-
91 ample to other models for the computational mechanisms of deep learning; for instance, it provides a
92 setting where deep neural nets successfully learn a concept to 100% accuracy while the corresponding
93 Neural Tangent Kernel (NTK) only achieves trivial performance, hence showing the importance of
94 *feature learning*. We also construct a counterexample to the “*deep only works if shallow is good*”
95 principle of (48), demonstrating a case where a deep network can get near-perfect accuracy even
96 when greedy layerwise training (e.g. Belilovsky et al. (13)) cannot beat trivial performance. By
97 providing positive theory and empirics which elude these simplified explanations of SGD, we hope to
98 point the way to a more complete understanding of learning dynamics in the challenging cases where
99 no apparent progress is made for extended periods of time.

100 1.2 Related work

101 We present the most directly related work on feature learning, and learning parities with neural nets.
102 A broader discussion can be found in Appendix A.3.

103 **SGD and feature learning.** Theoretical analysis of gradient descent over neural networks is
104 notoriously hard, due to the non-convex nature of the optimization problem. That said, it has been
105 established that in some settings, the dynamics of GD keep the weights close to their initialization,
106 thus behaving like convex optimization over the Neural Tangent Kernel (see, for example, (38, 7, 22)).
107 In contrast, it has been shown that in various tasks, moving away from the fixed features of the NTK
108 is essential for the success of neural networks trained with GD (for example (71, 6, 69) and the review
109 in (50)). These results demonstrate that feature learning is an important part of the GD optimization
110 process. Our work also focuses on a setting where feature learning is essential for the target task. In
111 our theoretical analysis, we show that the initial population gradient encodes the relevant features for
112 the problem. The importance of the first gradient step for feature learning has been recently studied
113 in (12).

114 **Learning parities with neural networks.** The problem of learning parities using neural networks
115 has been investigated in prior works from various perspectives. It has been demonstrated that parities
116 are hard for gradient-based algorithms, using similar arguments as in the SQ analysis (63, 1). One
117 possible approach for overcoming the computational hardness is to make favorable assumptions
118 on the input distribution. Indeed, recent works show that under various assumptions on the input
119 distribution, neural networks can be efficiently trained to learn parities (XORs) (20, 64, 27, 50). In

³While our focus is on the performance as a function of *training time*, we also performed some experiments
on performance as a function of *model size*, see Section 5.

120 contrast to these results, this work takes the approach of intentionally focusing on a hard benchmark
 121 task, without assuming that the distribution has some favorable (namely, non-uniform) structure.
 122 This setting allows us to probe the performance of deep learning at a known computational limit.
 123 Notably, the work of (8) provides analysis for learning polynomials (and in particular, parities) under
 124 the uniform distribution. However, their main results require a network of size $n^{O(k)}$ (i.e., extremely
 125 overparameterized network), and provides only partial theoretical and empirical evidence for the
 126 success of smaller networks. Studying a related subject, some works have shown that neural networks
 127 display a spectral bias, learning to fit low-frequency coefficients before high-frequency ones (57, 17).

128 2 Preliminaries

129 We define necessary notation here; see Appendix A for more background and technical ingredients.

130 **Sparse parities.** For integer $n \geq 1$ and non-empty $S \subseteq [n]$, the (n, S) -parity function $\chi_S :$
 131 $\{\pm 1\}^n \rightarrow \{\pm 1\}$ is defined as $\chi_S(x) = \prod_{i \in S} x_i$. We define the (n, S) -parity distribution \mathcal{D}_S as the
 132 distribution over $(x, y = \chi_S(x))$ ⁴ where x is drawn from the uniform distribution over $\{\pm 1\}^n$, which
 133 we denote by $\text{Unif}(\{\pm 1\}^n)$. We define the (n, k) -parity learning problem as the task of recovering
 134 S using samples from \mathcal{D}_S , where S is chosen at random in $\binom{[n]}{k}$. Statistically, it is possible to do
 135 so using $\Theta(\log \binom{n}{k}) \approx k \log n$ samples. However, in the *statistical query* (SQ) model (41), (which
 136 has been shown to encapsulate gradient-based methods such as GD or SGD (2)), this task requires
 137 $\Omega(n^k)$ queries (assuming constant level of noise). While learning noiseless parities can be solved by
 138 Gaussian elimination using $O(n)$ samples, learning sparse *noisy* parities, even at a very small noise
 139 level (i.e., $o(1)$ or $n^{-\delta}$) is believed to inherently require $n^{\Omega(k)}$ computational steps and samples, or
 140 exponential computation with $n^{o(k)}$ samples. This was first explicitly conjectured by Alekhnovich
 141 (5), and has been the basis for several cryptographic schemes (e.g., (37, 9, 10, 15)).

142 **Notation for neural networks and training.** Our main results are presented in the online learning
 143 setting, with a stream of i.i.d. batches of examples. At each iteration $t = 1, \dots, T$, a learning
 144 algorithm receives a batch of B examples $\{(x_{t,i}, y_{t,i})\}_{i=1}^B$ drawn i.i.d. from \mathcal{D}_S , then outputs a
 145 classifier $\hat{y}_t : \{\pm 1\}^n \rightarrow \{\pm 1\}$. We say that the algorithm solves the parity task in t steps with ϵ error,
 146 if with probability at least $1 - \epsilon$ over both training and internal randomness, as well as $(x, y) \sim \mathcal{D}_S$,
 147 $\hat{y}_t(x) = y$. We will focus on the case that $\hat{y}_t = \text{sign}(f(x; \theta_t))$ for some parameters θ_t in a continuous
 148 domain Θ and for a continuous function $f : \{\pm 1\}^n \times \Theta \rightarrow \mathbb{R}$.⁵ A ubiquitous online learning
 149 algorithm is gradient descent (GD). For a choice of loss function $\ell : \{\pm 1\} \times \mathbb{R} \rightarrow \mathbb{R}$, initialization
 150 θ_0 (that are chosen from some distribution), learning rate schedule $\{\eta_t\}_{t=1}^T \subseteq \mathbb{R}$ and weight-decay
 151 schedule $\{\lambda_t\}_{t=1}^T \subseteq \mathbb{R}$, GD refers the standard iterative update rule using the regularized, empirical
 152 loss function, which is a function of *architecture* f . The learning rate η_t can also be a vector (e.g.,
 153 allowing different rate schedules for different layers).

154 3 Empirical findings

155 3.1 SGD on neural networks learns sparse parities

156 The central phenomenon of study in this work is the empirical observation that neural networks, with
 157 standard initialization and training, can solve the (n, k) -parity problem in a number of iterations
 158 scaling as $n^{O(k)}$ on small instances. We observed robust positive results for randomly-initialized
 159 SGD on the following architectures, indexed by Roman numerals:

- 160 • *2-layer MLPs:* ReLU ($\sigma(z) = (z)_+$) or polynomial ($\sigma(z) = z^k$) activation, in a wide variety of
 161 width regimes $r \geq k$. Settings (i), (ii), (iii) (resp. (iv), (v), (vi)) use $r = \{10, 100, 1000\}$ ReLU
 162 (resp. polynomial) activations. We also consider $r = k$ (exceptional settings (*i), (*ii)), the
 163 minimum width for representing a k -wise parity for both activations.
- 164 • *Single neurons:* Next, we consider non-standard activation functions σ which allow a one-neuron
 165 architecture $f(x; w) = \sigma(w^\top x)$ to realize k -wise parities. The constructions stem from letting

⁴Our theoretical analyses and experiments can tolerate noisy parities, that is, random flipping of the label. For ease of presentation, we state the non-noisy setting.

⁵When $f(x; \theta) = 0$ in practice (e.g. with sign initialization), we break the tie arbitrarily. We ensure in the theoretical analysis that this does not happen.

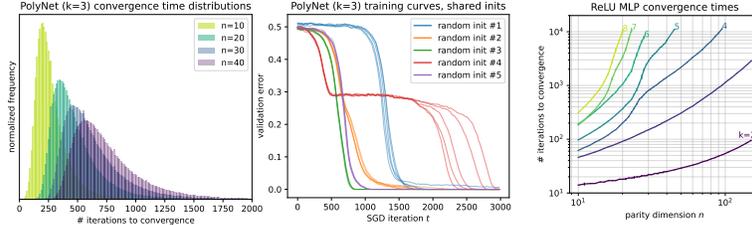


Figure 2: Black-box observations about the training dynamics. *Left*: Histograms of convergence times over 10^6 random trials, with heavy upper tails but no observed successes near $t = 0$ (unlike random search). *Center*: Optimization path (and convergence time) depend heavily on initialization, not the randomness of SGD; $B = 128, \eta = 0.01$ shown here. *Right*: The power-law exponent (α such that $t_c \propto n^\alpha$) eventually degrades on larger problem instances.

166 $w^* = \sum_{i \in S} e_i$, and constructing $\sigma(\cdot)$ to interpolate (the appropriate scaling of) $\frac{k-w^* \top x}{2} \bmod 2$
 167 with a piecewise linear k -zigzag activation (vii), or a degree- k polynomial (viii). Going a step
 168 further, a single ∞ -zigzag (ix) or sinusoidal (x) neuron can represent *all* k -wise parities. We also
 169 removed the second trainable layer (setting $u = 1$), obtaining settings (xi), (xii), (xiii), (xiv). We
 170 found that wider architectures with these activations also trained successfully.
 171 • *Transformers*: Motivated by recent theoretical and empirical work on the ability of self-attention
 172 to learn sparse functions and parities (47, 23, 32), we consider a simplified specialization of the
 173 Transformer architecture to this sequence classification problem. This is the less-robust setting
 174 (*iii); the architecture and optimizer are described in Appendix D.1.3.
 175 • *PolyNets*: Our final setting (xv) is the PolyNet, a slightly modified version of the parity machine
 176 architecture. Parity machines have been studied extensively in the statistical mechanics of ML
 177 literature (see the related work section) as well as in a line of work on ‘neural cryptography’ (59). A
 178 parity machine outputs the sign of the product of k linear functions of the input. A PolyNet simply
 179 outputs the product itself. Both architectures can clearly realize k -sparse parities. The PolyNet
 180 architecture was originally motivated by the search for an idealized setting where an end-to-end
 181 optimization trajectory analysis is tractable (see Section 4.1); we found in these experiments that
 182 this architecture trains very stably and sample-efficiently.

183 **Robust space of positive results.** All of the networks listed above were observed to successfully learn
 184 sparse parities in a variety of settings. We summarize our findings as follows: for all combinations
 185 of $n \in \{10, 20, 30\}$, $k \in \{2, 3, 4\}$, batch sizes $B \in \{1, 2, 4, \dots, 1024\}$, initializations {uniform,
 186 Gaussian, Bernoulli}, loss functions {hinge, square, cross entropy}, and architecture configurations
 187 {(i), (ii), ..., (xv)}, SGD solved the parity problem (with 100% accuracy, validated on a batch
 188 of 2^{13} samples) in at least 20% of 25 random trials, for at least one choice of learning rate $\eta \in$
 189 $\{0.001, 0.01, 0.1, 1\}$. The models converged in $t_c \leq c \cdot n^{\alpha k} \leq 10^5$ steps, for small architecture-
 190 dependent constants c, α (see Appendix C). Figure 1 (*left*) shows some representative training curves.

191 **Less robust configurations.** Settings (*i) and (*ii), where the MLP just barely represents a k -sparse
 192 parity, and the Transformer setting (*iii), are less robust to small batch sizes. In these settings, the
 193 same positive results as above only held for sufficiently large batch sizes: $B \geq 16$. Also, setting (*iii)
 194 used the Adam optimizer; see Appendix D.1.3 for details.

195 **Phase transitions in training curves.** For almost all of the architectures, we find that that the training
 196 curves exhibit phase transitions in terms of running time (and thus, in the online learning setting,
 197 dataset size as well): long durations of seemingly no progress, followed by periods of rapid decrease
 198 in the validation error. Strikingly, for architectures (v) and (vi), this plateau is absent: the error in the
 199 initial phase appears to decrease with a linear slope. See Appendix C.8 for more plots.

200 3.2 Random search or hidden progress?

201 The remainder of this paper seeks to answer the question: “By what mechanism does deep learning
 202 solve these emblematic computationally-hard optimization problems?”

203 A natural hypothesis would be that SGD somehow implicitly performs Monte Carlo random search,
 204 “bouncing around” the loss landscape in the absence of a useful gradient signal. Upon closer inspection,
 205 several empirical observations clash with this hypothesis:

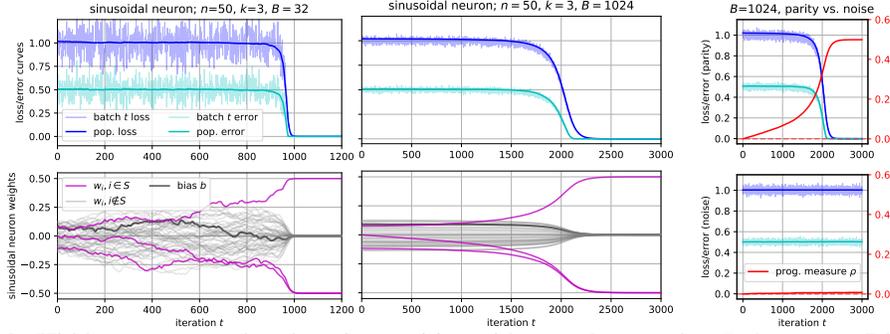


Figure 3: Hidden progress when learning parities with neural networks. *Left, center*: Black-box losses and accuracies exhibit a long plateau and sharp phase transition (top), hiding gradual progress in the SGD iterates (bottom). *Right*: A hidden progress measure which distinguishes gradual feature amplification (top) from training on noise (bottom).

- 206 • *Scaling of convergence times*: Without an explicit sparsity prior in the architecture or initialization,
 207 it is unclear how to account for the $n^{\Theta(k)}$ runtimes observed in experiments. The initializations,
 208 which certainly do not prefer sparse functions⁶, are close to the correct solutions with probability
 209 $2^{-\Omega(n)} \ll n^{-k}$. On larger instances (n, k) , the power-law exponents worsen; see Figure 2 (*right*),
 210 and the discussion in Appendix C.2.
- 211 • *No early convergence*: Over a large number of random trials, no copies of this randomized
 212 algorithm get “lucky” (i.e. solve the problem in significantly fewer than the median number of
 213 iterations); see Figure 2 (*left*). The success times of random exhaustive search would be distributed
 214 as $\text{Geom}(1/\binom{n}{k})$, whose probability mass is highest at $t = 0$ and decreases monotonically with t .
- 215 • *Sensitivity to initialization, not stochastic batches*: Running these training setups over multiple
 216 stochastic batches from a common initialization, we find that loss curves and convergence times
 217 are highly correlated with the architecture’s random initialization; see Figure 2 (*center*).

218 Even these observations, which do not probe the internal state of the algorithm, suggest that exhaustive
 219 search is an insufficient picture of the training dynamics, and a different mechanism is at play.

220 4 Theoretical analyses

221 4.1 Provable emergence of the parity indices in high-precision gradients

222 We now provide a theoretical account for the success of SGD in solving the (n, k) -parity problem.
 223 Our main theoretical observation is that, in many cases, the *population* gradient of the weights at
 224 initialization contains enough “information” for solving the parity problem. That is, given an accurate
 225 enough estimate of the initial gradient (by e.g. computing the gradient over a large enough batch
 226 size), the relevant subset S can be found.

227 As a warm-up example, consider training a single ReLU neuron $f(x; w) = \sigma(w^\top x)$ w.r.t. the
 228 correlation loss $\ell(y, \hat{y}) = -y\hat{y}$, from the initialization $w = [1, \dots, 1]$. While a single neuron cannot
 229 express the parity, we observe that the population gradient can indicate what the correct subset
 230 is: $\mathbb{E}_{(x,y) \sim \mathcal{D}_S} [\nabla_{w_i} \ell(y, f(x; w))] = \mathbb{E}_{(x,y) \sim \mathcal{D}_S} [-yx_i \sigma'(w^\top x)]$ which corresponds to either the
 231 order- $(k-1)$ Fourier coefficient $S \setminus \{i\}$ of the function $x \mapsto \sigma'(w^\top x)$ (if $i \in S$ is a relevant
 232 coordinate) or the order- $(k+1)$ coefficient $S \cup \{i\}$ (if $i \notin S$ is irrelevant). When σ is the ReLU
 233 function and $w = [1, \dots, 1]$, $\sigma'(w^\top x) = \frac{\text{sign}(\sum_i x_i) + 1}{2}$ is just a shifted majority function of x . The
 234 Fourier spectrum of majority is well-understood: for even k , there is a gap between these Fourier
 235 coefficients that is detectable using $n^{O(k)}$ samples.

236 This analysis can be further extended to a ReLU neuron initialized with weights $w \in \{\pm 1\}^n$ and a
 237 small bias. In fact, we can show that taking a single gradient step with large enough batch size on a

⁶Indeed, under all standard architectures and initialization, the probability that a random network is $\Omega(1)$ -
 correlated with a sparse parity would be $2^{-\Omega(n)}$, since with that probability $1 - o(1)$ of the total influence would
 be accounted by the $n - k$ irrelevant features.

238 ReLU network trained with the hinge-loss $\ell(y, \hat{y}) = \max\{1 - y\hat{y}, 0\}$, already finds features that can
 239 solve the parity problem. Using this, we show that training a ReLU MLP with SGD solves the parity
 240 problem:

241 **Theorem 4.** Fix some $\epsilon \in (0, 1)$, let k be some even number, and assume that n is some odd number
 242 satisfying $n \geq \Omega(k^4 \log(nk/\epsilon))$. There exist a symmetric random initialization scheme and learning
 243 rate and weight decay schedules s.t. for every $S \subseteq [n]$ of size k , training a ReLU MLP of size
 244 $r = O(2^k k \log(k/\epsilon))$ with batch size $B = n^{O(k)} \log^2(r/\epsilon)$ for $T = \text{poly}(k, r, 1/\epsilon)$ iterations on
 245 \mathcal{D}_S w.r.t. the hinge loss, finds a network $f(x; \theta_t)$ with expected loss: $\mathbb{E}[\ell(f(x; \theta_t), y)] \leq \epsilon$, where
 246 the expectation is over the randomness of initialization, training and sampling $(x, y) \sim \mathcal{D}_S$.

247 The analysis presented above shows that when the batch size scales with $n^{O(k)}$, SGD over MLP solves
 248 the parity problem. However, in our experimental setting, the number of steps, and not the batch size,
 249 scales with $n^{O(k)}$. While we believe that running SGD with small batch size and small learning rate
 250 essentially amplifies the signal in the population gradient, behaving similarly to performing a large
 251 step over a large batch size, we do not have a complete analysis for training MLPs with SGD in the
 252 small batch size regime. To complement the above result, in Section 4.2 we analyze the trajectory for
 253 a variant of the PolyNet architecture trained with gradient flow.

254 We note that, while the above analysis applies for ReLU MLPs with a specific initialization scheme, a
 255 similar feature emergence phenomenon can be observed in a broader set of architectures and setting.
 256 Indeed, for feature emergence to occur, we only require that there is a “gap” between the relevant and
 257 irrelevant Fourier coefficients. Formally, denote $\hat{f}(S) := \mathbb{E}[f(x)\chi_S(x)]$ the Fourier coefficient of f
 258 at S , and observe the following definition:

259 **Definition 1** (Fourier gap). For some function $f : \{\pm 1\}^n \rightarrow \mathbb{R}$ and some subset S of size k , we
 260 say that f has a γ -Fourier gap at S if (1) for every $(k-1)$ -element subset $S' \subseteq S$, it holds that
 261 $|\hat{f}(S')| \geq \gamma$, and (2) for every subset $S' \subseteq [n]$ of size $k+1$ it holds that $|\hat{f}(S')| \leq \gamma/2$.

262 Now, given a network architecture where some neuron has a γ -Fourier gap with respect to the target
 263 subset S , we can generalize the result of the ReLU neuron. That is, we show that the subset S can be
 264 determined by observing an estimate of the population gradient at initialization:

265 **Proposition 5.** Let σ be some activation function and let ℓ be some loss function. Denote $f(x; w) =$
 266 $\sigma(w^\top x)$. Fix some subset $S \subseteq [n]$. Let $g \in \mathbb{R}^n$ be an estimate of the population gradient such
 267 that $\|g - \mathbb{E}_{(x,y) \sim \mathcal{D}_S} [\nabla_w \ell(y, f(x; w))]\|_\infty \leq \gamma/4$. Then, for every w s.t. $\sigma'(w^\top x)$ has a γ -Fourier
 268 gap w.r.t. to S and $\ell'(f(x; w), y) = -y$ for all x , the target subset S is detected by g , namely
 269 $S = \{i \in [n] : |g_i| \geq 3\gamma/4\}$.

270 **Comparison with NTK analysis.** In recent years, many theoretical works have studied the behavior
 271 of SGD on neural networks through the lens of the neural tangent kernel (NTK) (38). It is therefore
 272 important to highlight the fact that the NTK (or, in fact, any kernel) cannot solve the sparse parity
 273 problem. The following result (see (40, 49)) shows that no kernel can achieve small error on the
 274 sparse parity problem, unless the size of the kernel is $n^{\Omega(k)}$:

275 **Theorem 6.** Let $\Psi : \{\pm 1\}^n \rightarrow \mathbb{R}^D$ be some D -dimensional embedding with $\sup_x \|\Psi(x)\|_2 \leq 1$,
 276 and let $R > 0$ be some number. If $DR^2 < \epsilon^2 \cdot \binom{n}{k}$, then there exists some (n, k) -parity distribution
 277 \mathcal{D}_S s.t. $\inf_{\|w\| \leq R} \mathbb{E}_{(x,y) \sim \mathcal{D}_S} [\ell(\Psi(x)^\top w, y)] > 1 - \epsilon$.

278 4.2 Disjoint-PolyNet: an idealized architecture for trajectory analyses

279 In this section, we present an idealized architecture (a version of PolyNets (xv)) that exhibits similar
 280 behavior to MLPs (experimentally) and is technically easier to analyze. More specifically, we consider
 281 the *disjoint-PolyNet* which takes a product over k linear functions where the linear functions are
 282 computed on k disjoint partitions of the input P_1, \dots, P_k with each $P_i = \{(i-1)n' + 1, \dots, in'\}$
 283 with $n' = n/k^7$, that is, $f(x; w_{1:k}) := \prod_{i=1}^k w_i^\top x_{P_i}$. As noted in the related work section, this is
 284 equivalent to a tree parity machine but with real-valued rather than ± 1 output.

285 In order for the disjoint-PolyNet to be able to express the class of parity problems, we assume that
 286 the set S of size k in the (n, k) -parity problem is selected such that exactly one index belongs to each

⁷We assume for simplicity that n is divisible by k .

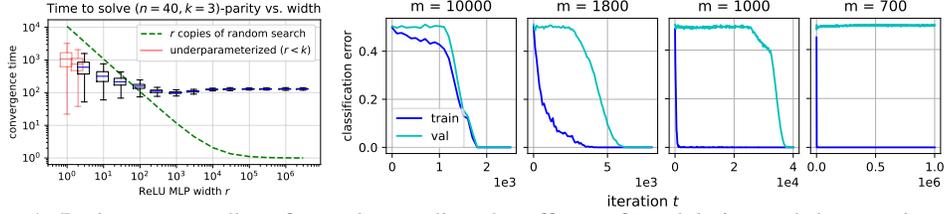


Figure 4: Parity as a sandbox for understanding the effects of model size and dataset size. *Left*: Success times vs. network width r on a fixed $(40, 3)$ -parity task: in accordance with the theory, parallelization experiences diminishing returns (unlike expected success times for random search, shown in green). Underparameterized models ($r = 1, 2$) were considered successful upon reaching 55% accuracy. *Right*: Training curves for an identical setup ($(50, 3)$ -parity task, architecture, and training algorithm), varying only the sample size m . The two center panels display “grokking”: a large gap between the time to zero train error vs. zero test error.

287 disjoint partition, that is, for all $i \in [k]$, $S \cap P_i = 1$. We refer to this problem as the (n, k) -disjoint
 288 parity problem. Note that there are still $(n')^k = (n/k)^k$ different possibilities for set S under this
 289 restriction. For fixed k , these represent a constant portion of the $\binom{n}{k} \approx (ne/k)^k$ (by Stirling’s
 290 approximation) possibilities for S in the general non-disjoint case.

291 Consider training a disjoint-PolyNet w.r.t. the correlation loss. WLOG, let $S = \{1, n' + 1, \dots, (k -$
 292 $1)n' + 1\}$ and $e_1 = (1, 0, \dots, 0)$. The population gradient is non-zero at i iff $i \in S$: $g_i(w_{1:k}) =$
 293 $\mathbb{E}[\nabla_{w_i} \ell(f(x; w_{1:k}), y)] = -\mathbb{E}\left[y \left(\prod_{j \neq i} w_j^\top x_{P_j}\right) x_{P_i}\right] = -\left(\prod_{j \neq i} w_{j,1}\right) e_1$.

294 Now we consider the gradient flow dynamics of disjoint-PolyNet, which provide a mathematically
 295 tractable case study for the trajectory of sparse parity learning. For each $i \in [k]$, in this section
 296 we treat w_i as a function from $\mathbb{R}_{\geq 0} \rightarrow \mathbb{R}^{n'}$ which satisfies the following differential equation:
 297 $\dot{w}_i = -g_i(w_{1:k}(t))$. For clarity of exposition, assume all-ones initialization.⁸ Then all of the relevant
 298 weights $\{w_{i,1} : i \in [k]\}$ follow the same trajectory, which we denote by $v : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}$. By analyzing
 299 the resulting differential equations, we can formally exhibit “phase transition”-like behavior in the
 300 fully deterministic gradient flow setting.

301 **Theorem 7** (Gradient flow on disjoint-PolyNets). *Suppose $k \geq 3$. Let $T(.49)$ be the time it takes*
 302 *for error to fall below .49, and let $T(0)$ be the time it takes to reach zero error. Then $\frac{T(.49)}{T(0)} =$*
 303 $1 - O\left((n')^{1-k/2}\right)$.

304 In other words, the network takes much longer to reach slightly better than trivial accuracy than it
 305 takes to go from slightly better than trivial to perfect accuracy.

306 We can also analyze the trajectory of disjoint-PolyNets trained with online SGD, confirming that a
 307 neural network trained with batch size 1 SGD can learn k -sparse parities within $n^{O(k)}$ iterations.

308 **Theorem 8** (SGD on disjoint-PolyNets). *Suppose we train a disjoint-PolyNet, initialized as above,*
 309 *with online SGD. Then there exists an adaptive learning rate schedule such that for any $\epsilon > 0$, with*
 310 *probability .99, the error falls below ϵ within $\tilde{O}\left((n')^{(2k-1)} \log(1/\epsilon)\right)$ steps.*

311 Extended versions of these theorems, along with proofs, can be found in Appendix B.3.

312 5 Hidden progress: discussion and additional experiments

313 In this section, we advocate for sparse parities as an idealized testbed for understanding algorithms
 314 and phenomena in modern deep learning. These are accompanied by experimental vignettes which
 315 are auxiliary to the core results from Section 3, with more systematic studies deferred to future work.
 316 Details are given in Appendix D.

317 **A progress measure for parity.** To begin, the black-box experiments in Section 3 suggest that
 318 random search is the incorrect model of SGD’s behavior in this setting. Using the theoretical insight

⁸Results for ± 1 initialization and Gaussian initialization are qualitatively similar and can be found in the Appendix.

319 that amplifying a precise initial population gradient suffices for learning parities, we construct one
320 possible progress measure, which is a function of the sequence of weights $w_0, \dots, w_t \in \mathbb{R}^n$ so far:
321 $\rho(w_{0:t}) := \|w_t - w_0\|_\infty$, using the fact that $w_t - w_0$ is an estimate for the initial population gradient
322 in the linearized setting. Figure 3 shows how gradual weight movement (and thus, progress) can be
323 hidden behind plateauing losses.

324 **Roles of overparameterization vs. oversampling.** An interesting consequence of our analysis is that
325 it illuminates scaling behaviors with respect to a third fundamental resource parameter: *model size*,
326 which we study in terms of network width r . If SGD operated by a “random search” mechanism, one
327 would expect width to provide a parallel speedup. Instead, we SGD sequentially amplifies progress.
328 The sharp lower tails in Figure 2 (*left*) imply that running r identical copies of SGD does *not* give
329 $(1/r) \times$ speedups; more directly, in Figure 4 (*left*), convergence times for sparse parities empirically
330 plateau at large model sizes.

331 It is a significant challenge to understand the interactions between network *depth* and computation,
332 and largely outside the scope of this work. However, in Appendix C.7, we provide a brief note on
333 using parities and polynomial-activation MLPs to construct a simple counterexample to the “*deep*
334 *only works if shallow is good*” principle of (48), demonstrating a case where a deep network can get
335 near-perfect accuracy even when greedy layerwise training (e.g. (13)) cannot beat trivial performance.

336 **Learning and grokking in the finite-sample multi-pass setting.** The main theoretical and empirical
337 results in this work consider online learning algorithms which couple the resources of training time
338 and independent samples. However, due to the computational-statistical gap in parity learning, these
339 positive results are suboptimal in terms of sample efficiency. We find that minibatch SGD (with
340 weight decay) can empirically solve sparse parities, even from a sample of size $m \ll n^k$. For small
341 values of m , we reliably observe the *grokking* phenomenon (55): overfitting for a long time, then
342 a *delayed* phase transition for the generalization error; see the two center panels of Figure 4 (*right*).

343 6 Conclusion

344 This work puts forward parity learning as a stylized test case to explore some of the puzzling features
345 of the role of computational (as opposed to statistical) resources in deep learning. These include
346 discontinuous improvements (a.k.a. *emergent capabilities*), feature learning, and universality of
347 architectures. In particular, we show that deep learning on parities exhibits a phase transition behavior,
348 that it is successfully learned by a variety of deep-net architectures, and that this success cannot be
349 explained as a “random exhaustive search”, nor through frameworks such as the neural tangent kernel
350 or layer-by-layer learning.

351 However, there are more experimental and theoretical questions, even for this simplified case of parity
352 learning. Our focus in this work was on the online learning case, where training time and samples
353 arise in tandem. However, we believe it would be instructive to investigate parity learning when three
354 resources of samples, time, and model size are scaled separately. Some very preliminary findings
355 along these lines are presented in Section 3.

356 Extending our theoretical results to the small-batch setting, as well as to more architectures, is an
357 open problem. Resolving it would require a better understanding of anti-concentration (lower bound
358 on deviation from mean) of Fourier coefficients, a phenomenon that is much less studied than the
359 concentration of these coefficients. We would also want to extend the analysis beyond parities to
360 tasks that are not aligned with the elementary basis such as low-rank tensor recovery.

361 Another important follow-up direction is understanding the extent that these insights extend from
362 parity learning to real-world problems, as well as the extent into which non-synthetic tasks (in,
363 e.g., natural language processing and program synthesis) embed within them parity-like subtasks of
364 exhaustive combinatorial search.

365 **Broader impact.** This work seeks to contribute to the foundational understanding of computational
366 scaling behaviors in deep learning. Our theoretical and empirical analyses are in a heavily-idealized
367 synthetic problem setting. Hence, we see no direct societal impacts of the results in this study.

368 **References**

- 369 [1] Emmanuel Abbe and Colin Sandon. Poly-time universality and limitations of deep learning.
370 *arXiv preprint arXiv:2001.02992*, 2020.
- 371 [2] Emmanuel Abbe, Prithish Kamath, Eran Malach, Colin Sandon, and Nathan Srebro. On the
372 power of differentiable learning versus pac and sq learning. *Advances in Neural Information*
373 *Processing Systems*, 34, 2021.
- 374 [3] Emmanuel Abbe, Elisabetta Cornacchia, Jan Hązła, and Christopher Marquis. An initial
375 alignment between neural network and target is needed for gradient descent to learn. *arXiv*
376 *preprint arXiv:2202.12846*, 2022.
- 377 [4] Naman Agarwal, Rohan Anil, Elad Hazan, Tomer Koren, and Cyril Zhang. Disentangling
378 adaptive gradient methods from learning rates. *arXiv preprint arXiv:2002.11803*, 2020.
- 379 [5] Michael Alekhnovich. More on average case vs approximation complexity. In *44th Annual*
380 *IEEE Symposium on Foundations of Computer Science, 2003. Proceedings.*, pages 298–307.
381 IEEE, 2003.
- 382 [6] Zeyuan Allen-Zhu and Yuanzhi Li. What can resnet learn efficiently, going beyond kernels?
383 *Advances in Neural Information Processing Systems*, 32, 2019.
- 384 [7] Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via
385 over-parameterization. In *International Conference on Machine Learning*, pages 242–252.
386 PMLR, 2019.
- 387 [8] Alexandr Andoni, Rina Panigrahy, Gregory Valiant, and Li Zhang. Learning polynomials with
388 neural networks. In *International conference on machine learning*, pages 1908–1916. PMLR,
389 2014.
- 390 [9] Benny Applebaum, David Cash, Chris Peikert, and Amit Sahai. Fast cryptographic primitives
391 and circular-secure encryption based on hard learning problems. In *Annual International*
392 *Cryptology Conference*, pages 595–618. Springer, 2009.
- 393 [10] Benny Applebaum, Boaz Barak, and Avi Wigderson. Public-key cryptography from different
394 assumptions. In *Proceedings of the forty-second ACM symposium on Theory of computing*,
395 pages 171–180, 2010.
- 396 [11] Gerard Ben Arous, Reza Gheissari, and Aukosh Jagannath. Online stochastic gradient descent
397 on non-convex losses from high-dimensional inference. *J. Mach. Learn. Res.*, 22:106–1, 2021.
- 398 [12] Jimmy Ba, Murat A Erdogdu, Taiji Suzuki, Zhichao Wang, Denny Wu, and Greg Yang. High-
399 dimensional asymptotics of feature learning: How one gradient step improves the representation.
400 *arXiv preprint arXiv:2205.01445*, 2022.
- 401 [13] Eugene Belilovsky, Michael Eickenberg, and Edouard Oyallon. Greedy layerwise learning can
402 scale to imagenet. In *International conference on machine learning*, pages 583–593. PMLR,
403 2019.
- 404 [14] Andrew C Berry. The accuracy of the gaussian approximation to the sum of independent
405 variates. *Transactions of the american mathematical society*, 49(1):122–136, 1941.
- 406 [15] Andrej Bogdanov, Manuel Sabin, and Prashant Nalini Vasudevan. Xor codes and sparse learning
407 parity with noise. *Proceedings of the Thirtieth Annual ACM-SIAM Symposium on Discrete*
408 *Algorithms*, pages 986–1004, 2019.
- 409 [16] Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal,
410 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
411 few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020.
- 412 [17] Yuan Cao, Zhiying Fang, Yue Wu, Ding-Xuan Zhou, and Quanquan Gu. Towards understanding
413 the spectral bias of deep learning. *arXiv preprint arXiv:1912.01198*, 2019.
- 414 [18] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam
415 Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm:
416 Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022.
- 417 [19] Alexandru Damian, Jason Lee, and Mahdi Soltanolkotabi. Neural networks can learn represen-
418 tations with gradient descent. In *Conference on Learning Theory*, pages 5413–5452. PMLR,
419 2022.

- 420 [20] Amit Daniely and Eran Malach. Learning parities with neural networks. *Advances in Neural*
421 *Information Processing Systems*, 33:20356–20365, 2020.
- 422 [21] Ilias Diakonikolas, Surbhi Goel, Sushrut Karmalkar, Adam R Klivans, and Mahdi Soltanolkotabi.
423 Approximation schemes for relu regression. In *Conference on Learning Theory*, pages 1452–
424 1485. PMLR, 2020.
- 425 [22] Simon S Du, Xiyu Zhai, Barnabas Poczos, and Aarti Singh. Gradient descent provably optimizes
426 over-parameterized neural networks. *arXiv preprint arXiv:1810.02054*, 2018.
- 427 [23] Benjamin L Edelman, Surbhi Goel, Sham Kakade, and Cyril Zhang. Inductive biases and
428 variable creation in self-attention mechanisms. *arXiv preprint arXiv:2110.10090*, 2021.
- 429 [24] Andreas Engel and Christian Van den Broeck. *Statistical mechanics of learning*. Cambridge
430 University Press, 2001.
- 431 [25] Carl-Gustav Esseen. On the liapunov limit error in the theory of probability. *Ark. Mat. Astr.*
432 *Fys.*, 28:1–19, 1942.
- 433 [26] Spencer Frei, Yuan Cao, and Quanquan Gu. Agnostic learning of a single neuron with gradient
434 descent. *Advances in Neural Information Processing Systems*, 33:5417–5428, 2020.
- 435 [27] Spencer Frei, Niladri S Chatterji, and Peter L Bartlett. Random feature amplification: Feature
436 learning and generalization in neural networks. *arXiv preprint arXiv:2202.07626*, 2022.
- 437 [28] Elizabeth Gardner and Bernard Derrida. Three unfinished works on the optimal storage capacity
438 of networks. *Journal of Physics A: Mathematical and General*, 22(12):1983, 1989.
- 439 [29] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedfor-
440 ward neural networks. In *Proceedings of the thirteenth international conference on artificial*
441 *intelligence and statistics*, pages 249–256. JMLR Workshop and Conference Proceedings, 2010.
- 442 [30] Oded Goldreich and Leonid A Levin. A hard-core predicate for all one-way functions. In
443 *Proceedings of the twenty-first annual ACM symposium on Theory of computing*, pages 25–32,
444 1989.
- 445 [31] Sebastian Goldt, Madhu Advani, Andrew M Saxe, Florent Krzakala, and Lenka Zdeborová.
446 Dynamics of stochastic gradient descent for two-layer neural networks in the teacher-student
447 setup. *Advances in neural information processing systems*, 32, 2019.
- 448 [32] Michael Hahn. Theoretical limitations of self-attention in neural sequence models. *Transactions*
449 *of the Association for Computational Linguistics*, 8:156–171, 2020.
- 450 [33] D Hansel, G Mato, and C Meunier. Memorization without generalization in a multilayered
451 neural network. *EPL (Europhysics Letters)*, 20(5):471, 1992.
- 452 [34] Godfrey Harold Hardy, John Edensor Littlewood, George Pólya, György Pólya, et al. *Inequalities*.
453 Cambridge university press, 1952.
- 454 [35] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers:
455 Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE*
456 *international conference on computer vision*, pages 1026–1034, 2015.
- 457 [36] Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza
458 Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al.
459 Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022.
- 460 [37] Yuval Ishai, Eyal Kushilevitz, Rafail Ostrovsky, and Amit Sahai. Cryptography with constant
461 computational overhead. In *Proceedings of the fortieth annual ACM symposium on Theory of*
462 *computing*, pages 433–442, 2008.
- 463 [38] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and
464 generalization in neural networks. *Advances in neural information processing systems*, 31,
465 2018.
- 466 [39] Y Kabashima. Perfect loss of generalization due to noise in $k=2$ parity machines. *Journal of*
467 *Physics A: Mathematical and General*, 27(6):1917, 1994.
- 468 [40] Pritish Kamath, Omar Montasser, and Nathan Srebro. Approximate is good enough: Probabilis-
469 tic variants of dimensional and margin complexity. In *Conference on Learning Theory*, pages
470 2236–2262. PMLR, 2020.

- 471 [41] Michael Kearns. Efficient noise-tolerant learning from statistical queries. *Journal of the ACM*
472 (*JACM*), 45(6):983–1006, 1998.
- 473 [42] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*
474 *arXiv:1412.6980*, 2014.
- 475 [43] Adam R Klivans, Ryan O’Donnell, and Rocco A Servedio. Learning intersections and thresholds
476 of halfspaces. *Journal of Computer and System Sciences*, 68(4):808–840, 2004.
- 477 [44] Gillat Kol, Ran Raz, and Avishay Tal. Time-space hardness of learning sparse parities. In
478 *Proceedings of the 49th Annual ACM SIGACT Symposium on Theory of Computing*, pages
479 1067–1080, 2017.
- 480 [45] Eyal Kushilevitz and Yishay Mansour. Learning decision trees using the fourier spectrum. *SIAM*
481 *Journal on Computing*, 22(6):1331–1348, 1993.
- 482 [46] Beatrice Laurent and Pascal Massart. Adaptive estimation of a quadratic functional by model
483 selection. *Annals of Statistics*, pages 1302–1338, 2000.
- 484 [47] Kevin Lu, Aditya Grover, Pieter Abbeel, and Igor Mordatch. Pretrained transformers as universal
485 computation engines. *arXiv preprint arXiv:2103.05247*, 2021.
- 486 [48] Eran Malach and Shai Shalev-Shwartz. Is deeper better only when shallow is good? *Advances*
487 *in Neural Information Processing Systems*, 32, 2019.
- 488 [49] Eran Malach and Shai Shalev-Shwartz. When hardness of approximation meets hardness of
489 learning. *Journal of Machine Learning Research*, 23(91):1–24, 2022.
- 490 [50] Eran Malach, Pritish Kamath, Emmanuel Abbe, and Nathan Srebro. Quantifying the benefit
491 of using differentiable learning over tangent kernels. In *International Conference on Machine*
492 *Learning*, pages 7379–7389. PMLR, 2021.
- 493 [51] GJ Mitchison and RM Durbin. Bounds on the learning capacity of some multi-layer networks.
494 *Biological Cybernetics*, 60(5):345–365, 1989.
- 495 [52] Ryan O’Donnell. *Analysis of boolean functions*. Cambridge University Press, 2014.
- 496 [53] Manfred Opper. Learning and generalization in a two-layer neural network: The role of the
497 vapnik-chervonvenkis dimension. *Physical review letters*, 72(13):2113, 1994.
- 498 [54] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan,
499 Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas
500 Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy,
501 Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-
502 performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-
503 Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*,
504 pages 8024–8035. Curran Associates, Inc., 2019. URL [http://papers.neurips.cc/paper/](http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf)
505 [9015-pytorch-an-imperative-style-high-performance-deep-learning-library.](http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf)
506 [pdf](http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf).
- 507 [55] Alethea Power, Yuri Burda, Harri Edwards, Igor Babuschkin, and Vedant Misra. Grokking: Gen-
508 eralization beyond overfitting on small algorithmic datasets. *arXiv preprint arXiv:2201.02177*,
509 2022.
- 510 [56] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al.
511 Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- 512 [57] Nasim Rahaman, Aristide Baratin, Devansh Arpit, Felix Draxler, Min Lin, Fred Hamprecht,
513 Yoshua Bengio, and Aaron Courville. On the spectral bias of neural networks. In *International*
514 *Conference on Machine Learning*, pages 5301–5310. PMLR, 2019.
- 515 [58] Maria Refinetti, Sebastian Goldt, Florent Krzakala, and Lenka Zdeborová. Classifying high-
516 dimensional gaussian mixtures: Where kernel methods fail and neural networks succeed. In
517 *International Conference on Machine Learning*, pages 8936–8947. PMLR, 2021.
- 518 [59] Michal Rosen-Zvi, Einat Klein, Ido Kanter, and Wolfgang Kinzel. Mutual learning in a tree
519 parity machine and its application to cryptography. *Physical Review E*, 66(6):066135, 2002.
- 520 [60] David Saad and Sara Solla. Dynamics of on-line gradient descent learning for multilayer neural
521 networks. *Advances in neural information processing systems*, 8, 1995.
- 522 [61] David Saad and Sara A Solla. On-line learning in soft committee machines. *Physical Review E*,
523 52(4):4225, 1995.

- 524 [62] Shai Shalev-Shwartz and Shai Ben-David. *Understanding machine learning: From theory to*
525 *algorithms*. Cambridge university press, 2014.
- 526 [63] Shai Shalev-Shwartz, Ohad Shamir, and Shaked Shammah. Failures of gradient-based deep
527 learning. In *International Conference on Machine Learning*, pages 3067–3075. PMLR, 2017.
- 528 [64] Zhenmei Shi, Junyi Wei, and Yingyu Liang. A theoretical analysis on feature learning in
529 neural networks: Emergence from inputs and advantage over fixed features. In *International*
530 *Conference on Learning Representations*, 2021.
- 531 [65] Roberta Simonetti and Nestor Caticha. On-line learning in parity machines. *Journal of Physics*
532 *A: Mathematical and General*, 29(16):4859, 1996.
- 533 [66] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
534 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information*
535 *processing systems*, 30, 2017.
- 536 [67] Martin J Wainwright. *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48.
537 Cambridge University Press, 2019.
- 538 [68] Timothy LH Watkin, Albrecht Rau, and Michael Biehl. The statistical mechanics of learning a
539 rule. *Reviews of Modern Physics*, 65(2):499, 1993.
- 540 [69] Colin Wei, Jason D Lee, Qiang Liu, and Tengyu Ma. Regularization matters: Generalization and
541 optimization of neural nets vs their induced kernel. *Advances in Neural Information Processing*
542 *Systems*, 32, 2019.
- 543 [70] Gilad Yehudai and Shamir Ohad. Learning a single neuron with gradient methods. In *Conference*
544 *on Learning Theory*, pages 3756–3786. PMLR, 2020.
- 545 [71] Gilad Yehudai and Ohad Shamir. On the power and limitations of random features for under-
546 standing neural networks. *Advances in Neural Information Processing Systems*, 32, 2019.
- 547 [72] Jingzhao Zhang, Sai Praneeth Karimireddy, Andreas Veit, Seungyeon Kim, Sashank Reddi,
548 Sanjiv Kumar, and Suvrit Sra. Why are adaptive methods good for attention models? *Advances*
549 *in Neural Information Processing Systems*, 33:15383–15393, 2020.

550 Checklist

- 551 1. For all authors...
- 552 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
553 contributions and scope? [Yes]
- 554 (b) Did you describe the limitations of your work? [Yes]
- 555 (c) Did you discuss any potential negative societal impacts of your work? [Yes]
- 556 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
557 them? [Yes]
- 558 2. If you are including theoretical results...
- 559 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 560 (b) Did you include complete proofs of all theoretical results? [Yes]
- 561 3. If you ran experiments...
- 562 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
563 mental results (either in the supplemental material or as a URL)? [Yes]
- 564 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
565 were chosen)? [Yes]
- 566 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
567 ments multiple times)? [Yes]
- 568 (d) Did you include the total amount of compute and the type of resources used (e.g., type
569 of GPUs, internal cluster, or cloud provider)? [Yes]
- 570 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 571 (a) If your work uses existing assets, did you cite the creators? [N/A]
- 572 (b) Did you mention the license of the assets? [N/A]

- 573 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
574 (d) Did you discuss whether and how consent was obtained from people whose data you're
575 using/curating? [N/A]
576 (e) Did you discuss whether the data you are using/curating contains personally identifiable
577 information or offensive content? [N/A]
578 5. If you used crowdsourcing or conducted research with human subjects...
579 (a) Did you include the full text of instructions given to participants and screenshots, if
580 applicable? [N/A]
581 (b) Did you describe any potential participant risks, with links to Institutional Review
582 Board (IRB) approvals, if applicable? [N/A]
583 (c) Did you include the estimated hourly wage paid to participants and the total amount
584 spent on participant compensation? [N/A]