

1 We thank all the reviewers for the time reading our paper! We will fix all the minor issues, and below we only address
2 the main concerns.

- 3 • **R2:** “The function class it learns seems to be identical to the function classes learned by Allen-Zhu et al. in prior
4 works, which are known to be learnable in polynomial-time via other methods.”

5 We’d like to emphasize that our main focus of the paper is “what function class can *recurrent network (RNN)* learn”.
6 Prior to our work, the only function class known to be learnable are linear functions; but even for this linear class, it
7 could require 2^L samples if known generalization bounds are used, since $\|W\|_2$ is larger than 1 (roughly 2 in our
8 setting). We believe it is (much) more interesting if someone could identify a non-linear class and prove that it is
9 learnable by RNN. This is already a step forward.

- 10 • **R2:** What happens if the inputs are not norm 1? What complexity bounds are obtained in this case? Also, how does
11 the complexity scale with the Lipschitzness/boundedness of the loss function (currently assumed to be 1)?

12 If the inputs are not norm 1 but with norm $< C$, then we can wlog. pad the input (adding $\sqrt{1 - C^2}$ to last
13 coordinate) to make it norm exactly C . Then, the concept class $\phi(\langle w, x \rangle)$ becomes $\phi(C\langle w, x/C \rangle)$ and one can
14 define $\phi'(z) := \phi(Cz)$. Now, the complexity changes from $\mathfrak{C}(\phi, R)$ to $\mathfrak{C}(\phi', R)$: how much it affects the complexity
15 depends on what ϕ is. If ϕ is degree k then this is at most C^k . If $\phi(z) = \sin(z)$ and $\phi'(z) = \sin(Cz)$, this changes
16 the complexity by a factor $(1/\varepsilon)^C$. This is polynomial only if C is a constant, but should be somewhat forgiven:
17 learning $\sin(10\langle w, x \rangle)$ is indeed somewhat “exponentially harder” than learning $\sin(\langle w, x \rangle)$.

18 As for the Lipschitzness/boundedness of the loss function, our final result (time/sample complexities) only scale
19 polynomially with them. We will add a short section explaining this.

- 20 • **R2:** Are there hardness results? is $L^{O(\log 1/\varepsilon)}$ the best possible?

21 There’s no hardness result, but we conjecture $L^{O(\log 1/\varepsilon)}$ is the best possible for vanilla RNN. If more structures
22 such as memory units are added, it may be possible to get $poly(L)$. That’s an interesting future direction.

- 23 • **R2:** My main concern is how much does this paper overlap with several other works analyzing SGD on overparam-
24 eterized networks, most notably Allen-Zhu et al.

25 Allen-Zhu et al. [1] consider two (or three) layer feedforward NN, which uses one (or two) hidden weight matrix to
26 learn one target function. Here in RNN, there’s only **one** weight matrix shared across the entire time horizon to
27 learn L target functions at different input positions.

28 Using “one weight” to learn “one target function” is traditional,
29 but using “one weight” to efficiently learn “ L different target functions” is substantially more difficult,

30 especially when the layer information L is not given as the input to the network. For example, our theorem implies
31 that an RNN can **distinguish** the sequences “AAAB” from “AABA”, since the order of A and B are different.
32 This requires the RNN to keep track, **using one weight matrix**, of the position information of the symbols in the
33 sequence. This is indeed is more difficult.

34 Allen-Zhu et al. [2,3] are about trainability only, and gives no generalization guarantee.

- 35 • **R2:** the main technical difference from previous work seems to be Lemma 5.2b.

36 **No, this is not true.** Due to space, we only present one of our main contributions in 8 pages. As we emphasized
37 on line 269-273 of page 8. Our technical lemmas are in Appendices B/C/D/E/F (regarding RNN at initialization
38 and stability). Although B+F have reused some prior work, Appendices C+D+E are **completely new** what-so-ever.
39 Furthermore, C+D+E give technical lemmas that are (certainly) of independent interests.

- 40 • **R4** also questions our overlap with Allen-Zhu et al [2]. We refer R4 to our answers above.

- 41 • **R4:** Can you give examples about the proposed concept class? Intuitions about \mathfrak{C}_s and \mathfrak{C}_e ?

42 **Counting is an example.** One can define ϕ such that $\phi(a) = 1$ and $\phi(b) = -1$ (this can be achieved by a quadratic
43 function with constant complexity), hence the final prediction (target function) is $\sum_i \phi(x_i)$. If the sequence is
44 $x = a^n b^m$ such that $n = m$, then $\sum_i \phi(x_i) = 0$, otherwise it is non-zero. So this concept class is learnable by our
45 theorem. We will add more examples in the next version. We will also give more intuitions about \mathfrak{C}_s and \mathfrak{C}_e .

- 46 • **R6:** “Robust Large Margin Deep Neural Networks” gives generalization independent of number of neurons

47 First, “Robust...” seems to be about feedforward NN and not RNN. Second, as we have articulated in lines 35-55,
48 there are indeed RNN generalization works [31,9] that do not depend on number of neurons, but they have to depend
49 exponentially on input length.

- 50 • **R6:** what will happen if other parameters are learned? Our same result will hold (almost no proof changes) if
51 A, B, W are all trained together. We will add a paragraph to explain it.