1  We thank the reviewers for their comments and time.

2  **Reviewer #1.**  *(Relationship to prior work, particularly (Lyu and Li, 2019).)* We will expand our Related Work section
3  to explicitly separate our technical contribution from prior work, expanding comments found throughout this response.
4  As the reviewer mentions, prior work often *assumes* directional convergence and alignment, but neither indicates a
5  possible proof, nor even provides conclusive evidence. Regarding (Lyu and Li, 2019), they sidestepped directional
6  convergence by using subsequences (cf. their Theorem 4.4), and in fact they provided a pathological example where
7  directional convergence *fails* (cf. their Theorem J.1). Our key technical tool for directional convergence, the notion of
8  o-minimal definability, was only used in (Lyu and Li, 2019) to ensure a nonsmooth chain rule. Alignment, in our exact
9  form, does not appear in prior work, which excites us as it implies both existing and new margin maximization results.

10  *(Overly technical presentation.)* We agree, with hindsight, that our presentation was encumbered with too much focus
11  on technical details, such as how to handle nonsmooth models as mentioned by the reviewer. We will expand intuition
12  and machine learning connections, and move material to the appendices, as the reviewer suggests.

13  *(Minor comments.)* Thank you! We will address these comments in our revisions.

14  **Reviewer #2.**  *(Discrete-time analysis.)* We agree that a discrete-time analysis is essential. We touched upon this
15  in our "Concluding Remarks", but will expand the material; for instance, one can easily adapt our analysis to handle
16  extremely small step sizes, but handling a practical choice is much more challenging. Getting convergence rates is
17  another important open problem, and might be hard even for the gradient flow, despite our present work.

18  *(Non-homogeneous networks.)* We agree this is important, and provide illustrative support, in the form of Figure 2b on
19  page 3 (DenseNet) and Figure 3 in the appendix (ResNet), which does seem to suggest directional convergence holds.

20  *("One might argue that the directional convergence part is relatively easier".)* We will highlight in our revised Related
21  Work section that this question is tricky, and has stymied many mathematicians. As discussed at the end of Section 1.1
22  in our submission, to prove the related gradient conjecture of René Thom, mathematicians had to develop the whole
23  area of o-minimal structures. Even with this powerful technique, as far as we know the existing results on directional
24  convergence all roughly assume piecewise polynomials or real-analytic functions, and therefore cannot analyze a
25  nonsmooth function composed with the exp/logistic loss, which is more relevant to the deep learning community.

26  *("Sub-differential and continuous dynamics are a big field".)* We surveyed and cited recent work, e.g., (Davis et al.,
27  2020); we had to prove many new results (e.g., the unbounded nonsmooth Kurdyka-Łojasiewicz inequalities) to show
28  directional convergence and alignment.

29  **Reviewer #3.**  *(Generalization.)* We agree that generalization is essential, and will include expanded discussion in our
30  revisions. Briefly, on the empirical side, we point the reviewer to the large-scale experiment we cited, by Shallue et
31  al. (2018), which finds that even uncommonly large amounts of training do not seem to hurt generalization. On the
32  theoretical side, we will mention various margin-based generalization bounds, and tie them to our alignment property.

33  *(Relationship to (Lyu and Li, 2019).)* Our analysis of directional convergence is not relevant to this prior work: they did
34  not prove directional convergence but instead must use subsequences. Please refer to lines 2-9 above for more details.

35  *(Non-homogeneity and discrete time.)* We agree these are important; please refer to lines 14-19 above for more details.

36  **Reviewer #5.**  *(Regarding "(1)".)* Our analysis can be extended to many other decreasing losses with an exp tail, but
37  we chose to focus on the exp/logistic loss to highlight the key ideas, and moreover because the logistic loss is one of the
38  most widely-used losses in machine learning. Since our core technical work is on solutions at infinity, losses like the
39  squared loss, which imply a (finite) minimizer, require a different analysis, though some of our lemmas will still apply.

40  *(Regarding "(2)".)* This initialization assumption was first introduced in prior work (Lyu and Li, 2019), and allows us
41  to focus on the late-training regime; we can handle random initialization by first invoking a standard overparameterized
42  analysis as a lemma (these results only hold near random initialization). Regarding "it seems not hard to show
43  the initialization is close to an optimal classifier", firstly we stress this is an orthogonal concern to our directional
44  convergence result, which ensures gradient flow eventually stabilizes. Secondly, "optimal" in our setting is typically
45  "maximum margin", where it seems the late-training phase is essential, and our results handle a few such cases.

46  *(Regarding "(3)".)* The prior work (Ji and Telgarsky, 2018a) only considered deep linear networks, and the proofs
47  there share almost no techniques with the proofs of our main results, in Theorems 3.1 and 4.1. If this was a typo and
48  the reviewer intended (Ji and Telgarsky, 2019), then our potential function $\mathcal{J}$ in Section 4 is indeed based on their
49  dual potential, but their setting is linear and convex, and our nonlinear nonconvex analysis is significantly different.
50  Regarding the relationship to math literature, such as (Kurdyka 2000a, 2006), we had to overcome many obstacles
51  missing from these prior works, such as how to handle the exp/logistic loss with nonsmooth models.