

1 We thank the reviewers for their feedback. We would like to note that due to space constraints, we relegated some  
2 theoretical and experimental details to the Appendix. We tried to ensure that the main paper, together with appendix,  
3 is self-contained. In Appendix B, we describe the greedy algorithms of Huang et al and Bengio et al. Appendix D  
4 contains a discussion of the theoretical results. Appendix J contains the experimental setup, hyper-parameter selection.  
5 In the final version of the paper, we will move additional clarifying details from the appendix to the main paper, as  
6 detailed in our responses below. Regarding XGBoost experiments, line 308 has a typo. We actually used decision trees  
7 (not decision stumps) in our experiments. The list of hyper-parameters tuned for XGBoost can be found in Table 7 in  
8 the appendix. We used trees with depth up to 20 in our experiments.

9 **Reviewer 1.** *Algorithms 2, 3.* We believe having both the algorithms can provide a unified view of various boosting  
10 algorithms: many classical boosting algorithms can be derived using Algorithm 3; whereas, greedy algorithms for DNN  
11 training can be derived using Algorithm 2. Due to space constraints, we didn't add experimental results comparing both  
12 the algorithms. In our experiments, we noticed Algorithm 3 has marginally worse performance than Algorithm 2.

13 **Reviewer 2.** *Novelty.* As pointed out in Section 3.1 and Appendix J.1, existing greedy techniques have certain  
14 drawbacks which cause them to perform poorly. One of the contributions of our work is to fix these drawbacks using  
15 better greedy techniques. While the resulting architecture might resemble ResNets, DenseNets, they have certain crucial  
16 differences which cause them to perform better.

17 *Training NNs with boosting has already been proposed.* Proposition 3.1 proves the non-obvious result that all these  
18 techniques are actually equivalent to greedy layer-by-layer technique of Bengio et al and hence face the drawbacks  
19 stated in Section 3.1. Moreover, our proposed techniques strictly improve upon the existing techniques.

20 *Goal of our work.* We would like to emphasize that the reason behind studying greedy algorithms is not to speed-up  
21 optimization of deep networks. Rather, the main reason is to come up with algorithms which have strong theoretical  
22 guarantees and at the same time have good empirical performance. Our generalized framework satisfies both the criteria  
23 (also see our response to Reviewer 3). Moreover, greedy algorithms come with certain additional advantages such as  
24 low memory requirements. For details on computation time, see our response to Reviewer 4.

25 *Empirical Evaluation.* Details about the hyper-parameters such as number of boosting iterations can be found in  
26 Appendix J.2. We used a hold-out set to choose the number of boosting iterations, with an upper limit of 15 iterations.  
27 For all the datasets, the optimal iteration picked using this process is less than 15. Lines 315-317 specify the types of  
28 weak feature transformers used. In Appendix J.1, we provide more experimental results showing why the proposed  
29 techniques work better than StdCompBoost. We believe our results clearly show that the proposed techniques achieve  
30 the above stated goals. On almost all the datasets, our techniques are better than StdCompBoost and other additive  
31 boosting techniques. On some datasets, we even match the performance of end-to-end training. Finally, thanks for  
32 references [1,2]; we will add more experimental results using datasets from these repositories.

33 **Reviewer 3.** *Related Work.* Thanks for pointing out [a,b], which we were unaware of. In [a], the authors propose a slight  
34 variant of AdaBoost in which the initial distribution is allowed to be any probability distribution, and the classification  
35 errors are measured using a general cost function. In [b], the authors consider an additive boosting procedure where  
36 different kinds of weak hypothesis classes are used in each iteration. Our work differs from [a,b] in a number of aspects,  
37 the key aspect being the way we aggregate weak classifiers. [a,b] consider additive combinations of weak classifiers,  
38 whereas we consider more general combinations. Moreover, unlike [a,b], the complexity of weak hypothesis classes in  
39 our algorithms grows with iteration.

40 *Goal of our work.* Our framework is a first step towards bridging the gap between classical methods such as boosting  
41 and modern deep learning methods. Classical boosting comes with strong theoretical guarantees, but doesn't match the  
42 performance of DNNs. Whereas, DNNs have good empirical performance, but don't come up with strong theoretical  
43 guarantees. Our proposed framework allows to derive learning algorithms that come with strong theoretical guarantees  
44 and improve the performance of existing greedy algorithms (see conclusion).

45 *Comparison with end-to-end.* In our experiments we did compare with ResNets: end-to-end training corresponds to  
46 ResNets trained using SGD. We would like to emphasize that our goal is not to beat end-to-end trained deep networks  
47 (see our claims in abstract lines 16-20 and introduction lines 65-70). While matching the performance of end-to-end  
48 training is the ideal goal, in this work we set out with a different goal: we aim to improve the performance of existing  
49 greedy techniques and take a step towards bridging the performance between greedy algorithms and end-to-end training.

50 **Reviewer 4.** *Theoretical guarantees.* For the theoretical guarantees to hold, we only need the feature transformers that  
51 are fit at each iteration to satisfy certain weak learning condition (Definition 4.1). This condition doesn't require the  
52 shallow networks to be trained optimally. This is similar to the weak learning condition in traditional boosting, which  
53 requires the weak learners added at each iteration to have better than random performance.

54 *Computation.* In terms of training time, we didn't notice any significant advantage for greedy techniques over end-to-  
55 end training. However, greedy techniques consume very little memory compared to end-to-end training. So greedy  
56 techniques can be very useful for training huge models on GPUs with very little memory. That being said, greedy  
57 techniques has several other advantages: they are easy to optimize and are theoretically much easier to analyze, as one  
58 only needs to understand optimization of shallow networks.