

1 Thanks for the helpful feedback. We got a clear sense of where more clarification would be helpful.

2 **General Discussion** Our work is part of the following larger and important discussion within the NeurIPS community:  
3 To what solution do neural nets (trained w. GD) converge to in the overparameterized setting? This discussion was  
4 reignited by two recent NeurIPS papers (Gunasekar et al. 2017) and (Vaskevicius et al. 2019). They showed that in a  
5 certain context, continuous GD updates converge to sparse solutions. In a recent COLT paper (Amid, Warmuth 2020)  
6 this was related to a particular MD algorithm: if the edges of a linear neuron are doubled (as in Fig. 2), then continuous  
7 GD on this network simulates the unnormalized exponentiated gradient algorithm (EGU). EGU is the classical MD  
8 algorithm (with the log link) and there is a long history of results showing that it does well when the solution is sparse.  
9 Previously it was thought that GD cannot take advantage of the sparsity of the solution.

10 In this paper we show that the reach of GD is even much farther: Any CMD update can be simulated by another CMD  
11 provided a certain reparameterization function exists. We give many reparameterization examples in the paper and in  
12 particular, we develop a family of MD updates (using the  $\log_\tau$  link) that lie between GD and EGU using vanilla GD  
13 (by means of a reparameterization).

14 What is the surprising insight? The structure of the network is key! Much research has been done in exploring multi  
15 layer nets with fully connected (FC) layers. In particular, wide FC layers have been investigated. For example the  
16 NeurIPS paper, Arora et al. “On Exact Computation with an Infinitely Wide Neural Net.” NeurIPS 19, shows that an  
17 infinitely wide neural net behaves like a kernel method.

18 Our research leads us in a different direction in that a linear neuron in which the edges are doubled (Fig. 2) can exploit  
19 sparsity. In particular, combined with some well known bounds for any kernel method, this shows that GD on the thin  
20 network of Fig. 2. can learn certain sparse problems exponentially faster than any kernel method. We believe that this  
21 discussion is central to neural net research and important to the NeurIPS community.

22 We will add more context along the above to the final version. We will also smoothen the transitions and give more  
23 extensive definitions as suggested by the reviewers. The discussion of duality is important for the optimization context.  
24 We will relegate much of that to the appendix in exchange for giving more motivation for the main results of the paper.

25 **Reviewer #2** • “given that this work is being submitted to NeurIPS, I don’t see enough of a relationship to artificial  
26 neural networks. . . Describing the specifics of using this work to help the neural processing world feels as though it  
27 would be more appropriate.”

28 Please note that many similar theoretical work such as (Gunasekar et al. 2017) and (Vaskevicius et al. 2019)  
29 have appeared before in NeurIPS. That is, NeurIPS has always been one of the top venues for publishing results on  
30 learning theory and optimization. The current submission is a theoretical work, but also has practical implications for  
31 understanding the training dynamics of deep neural networks. Therefore, we would like to encourage Reviewer #2 to  
32 cast their judgment solely based on the quality of the paper and significance of the theoretical results and defer the  
33 concerns about the relevance of the work to the NeurIPS community to the ACs.

34 **Reviewer #4** • *Lack of definitions:* We will add the missing definitions and improve the flow.

35 • *Meaning of ‘independent variables:* We refer to independence by means of elementary calculus. For independent  
36 variables, a variable in an equation may have its value freely chosen without considering values of the other variables.

37 • *Motivation of the dual form:* The dual form of MD is extremely important for efficient implementation of updates  
38 and analyzing the worst case regret bounds. We will make the motivation clear. We will also improve the write up to  
39 ease the transition between the sections.

40 • “Why should we care about reparameterized CMD as GD?”

41 For a long time, many MD updates such as EGU had been considered fundamentally different than GD (multiplicative  
42 vs. additive). Therefore, the EGU update was considered to be unrealizable by additive updates such as backprop. We  
43 show that EGU (and many other CMD) updates are realizable “exactly” via backprop on a reparameterized network.  
44 Following (Amid, Warmuth 2020), we believe that the discretized updates will closely mimic the original updates.  
45 This has important implications for implementing these updates via GD in the existing platforms (such as TF).

46 **Reviewer #5** • “the authors show that CMD can be presented as minimizing a trade-off between the loss and a  
47 Bregman momentum. . . I am unsure about the significance of this contribution.”

48 Prior to this work, CMD was motivated as the limit point of MD when the step-size goes to zero. Alternatively, we  
49 motivate CMD as the minimization of a functional objective. This novel view of CMD is useful for deriving worst case  
50 regret bounds in the continuous-time, which we defer to future work. This also allows a novel approach for deriving  
51 MD updates by directly discretizing the functional objective (please see Appendix C for more details).