1 We thank all reviewers for their helpful reviews. Please see our response below.

2 **Correctness of Theorem 1 proof** Thank you R3 for pointing out the mistake in the current proof. The mistake is
3 in the very last step of the proof, where we tried to show $W_1 - W_2^\top = 0$ (lines 593-594). Fortunately, we have the
4 following fix that asserts the correctness of Theorem 1. We hope R3 will update their score based on this revised proof.

5 From Lemma 2 (line 583), $C = \frac{1}{n}(I - W_2 W_1)XX^\top \in \mathbb{R}^{m \times m}$ is positive semi-definite, and $\Lambda^2 = \text{diag}(\lambda_1, \ldots, \lambda_k)$
6 is positive definite. Define $A = W_1 - W_2^\top \in \mathbb{R}^{k \times m}$. We prove below that $A = 0$ follows from line 592 which states

$$\forall v \in \mathbb{R}^k, \quad 0 = v^\top ACA^\top v + v^\top \Lambda^2 AA^\top v \tag{1}$$

7 *Proof.* Since $ACA^\top \succeq 0$, we have $\forall v, v^\top ACA^\top v \geq 0$. Hence, from (1), $\forall v, v^\top \Lambda^2 AA^\top v \leq 0$. Consider setting
8 $v = e_i$, the $i^{th}$ coordinate vector in $\mathbb{R}^k$ ($i^{th}$ entry is 1, and all others are 0). We must have $e_i^\top \Lambda^2 AA^\top e_i = \lambda_i^2 \|A_i\|_2^2 \leq 0$
9 ($A_i$ denotes the $i^{th}$ row of $A$). Since $\lambda_i > 0$, we have $A_i = 0$. Since this holds for all $i = 1, \ldots, k$, we have $A = 0$. $\square$

10 **How to choose the non-uniform regularization parameters (R4)** This is a great
11 question. It's indeed difficult to choose an "optimal" set of $\lambda$ values for the non-
12 uniform $\ell_2$ regularization (see below for justifications). However, note that our goal
13 is the analysis, and that the difficulty in choosing the $\lambda$ values a priori contributes
14 to the argument of weak symmetry breaking by the non-uniform $\ell_2$ regularization.



Figure 1: Optimal $\lambda$ values.

15 As for why it is difficult to choose an "optimal" set of $\lambda$ values, note that optimal
16 values at global minima are not optimal in general. At global minima, the optimal $\lambda$
17 values can be obtained by solving the $\min \max$ optimization in line 509 (an example
18 of such optimized $\lambda$ values for MNIST, $k = 20$ is shown in Figure 1). However, this
19 set of $\lambda$ values concentrate on the larger side, and significantly slow down learning
20 (suboptimal *away from* the global minima). In our experiments, we first make a (rough) estimate of the $k^{th}$ largest data
21 singular value, and linearly space the $\lambda$ values from 0 to the estimated $\sigma_k$ value (see Appendix G for details).
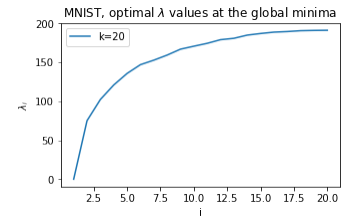
22 **Practical utility of the analyzed algorithms (R1, R3)** A big part of our mo-
23 tivation was to understand the slowness of learning neural net representations,
24 as opposed to reducing loss — a topic of immense practical importance, as
25 evidenced by the recent flurry of interest on early training effects. Linear au-
26 toencoders are one of the few examples where we can determine exactly what
27 representation *ought* to be learned, which makes them a particularly useful
28 model system for understanding convergence of representations.



Figure 2: Mini-batch experiment.

29 **Probabilistic interpretations of the non-uniform $\ell_2$ regularization and the
30 nested dropout?** As pointed out by R2, the probabilistic interpretation of non-
31 uniform $\ell_2$ regularization is a straightforward generalization of that studied
32 in Kunin et al. In particular, we can assign a diagonal Gaussian prior to the
33 weights whose precision parameters equal the $\lambda_i$ value for the corresponding
34 latent dimension. There is a well-known Bayesian interpretation of dropout
35 [Gal and Ghahramani, 2016], and extending this to nested dropout is an interesting direction for future work.
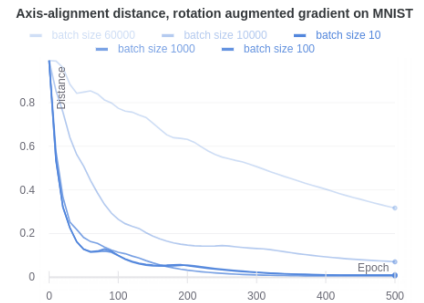
36 **Relation to additional prior work** Compared to Wager et al. [2013] (R1), which discusses the connection of
37 dropout and adaptive $\ell_2$ regularization in generalized linear models, we work with a different class of models (linear
38 autoencoders) and study a different type of dropout (nested dropout in the latent units). Nevertheless, it is an interesting
39 future direction to investigate the connection between the deterministic nested dropout and common types of regularizers.

40 Concurrent work [Oftadeh et al., 2020] (pointed out by R3) addresses the identifiability issue in linear autoencoders by
41 proposing a new loss function. The loss function proposed in Oftadeh et al. is a special case of deterministic nested
42 dropout (section 5 of our paper) where the prior $p_B(\cdot)$ is a uniform distribution. Oftadeh et al. show that the local
43 minima correspond to ordered, axis-aligned representations, but do not analyze the speed of convergence under the new
44 loss. Hence, our analysis provides additional insight into their method. We will include the above discussions, as well
45 as the additional citations regarding the connections between linear VAEs and pPCA (R3) in the revised paper.

46 **Mini-batch training for the rotation augmented gradient? (R1)** We did additional experiments on MNIST using
47 the rotation augmented gradient, with various batch sizes and $k = 20$ (Figure 2). The results show that the rotation
48 augmented gradient works well with mini-batches. Larger batches improve per-epoch convergence up to a point of
49 diminishing returns, similarly to standard models and algorithms (e.g. Shallue et al., JMLR 2019).

50 **Writing & clarity** We thank R2 and R4 for the writing & clarity suggestions, and will address them in the revised
51 paper. Note that the preconditioning of the Adam optimizer is not compatible with the rotation augmented gradient, so
52 only the SGD is relevant for this algorithm (see Appendix G for more experimental details).