

1 We thank the reviewers for the useful feedback. We will add the accidentally missing legend to Fig. 2. (orange line is  
2 the accuracy after removal of memorized examples and green line is accuracy after removal of the same number of  
3 randomly chosen examples). Responses to specific comments are below:

#### 4 R1

5 **more experiments to strengthen the validation of the theory** ▷ We welcome suggestions of other experiments. The  
6 correlation between memorization and influence is addressed by our experiments in Fig 2 since there we measure the  
7 cumulative influence of examples with mem. value above some threshold. In particular, it implies that memorization  
8 and overall influence are positively correlated.

#### 9 R2

10 **the author propose a “closely-related” statistic, called  $\text{infl}_m$ , by keeping  $m$  random subsets instead of removing  
11 one example** ▷ This is not an accurate description of the estimator. We estimate the effect of removal of a single  
12 example from a random subsample of the original dataset.

13 **some explanations missing: relationship between  $\text{infl}$  and  $\text{infl}_m$ , and between  $\text{stddev}$  and time complexity** ▷ The  
14 relationship is explained in Lines 86-89.  $\text{infl}_m$  is not equal to leave-one-out influence but the relationship between  
15 them is that of using a smaller random subset instead of the entire dataset and then taking the expectation. Note that  
16 this is also the classical jackknifing approach in statistics. The relationship between std. dev and time complexity of  
17 estimating the  $\text{infl}_m$  for all training examples at the same time (which is what we need) is stated in Lemma 2.1.

18 **need rudimentary experiment to show effectiveness of the proposed method, compared with naive monte carlo  
19 estimation** ▷ Our estimator is formally equivalent to the naive way to do it. The point of our algorithm is that it  
20 estimates all the values at the same time.

21 **paper reorganization and title change** ▷ Our primary contribution is validating and clarifying an explanation for  
22 a fundamental phenomenon in machine learning. The title and organization are aimed at making this clear. This  
23 suggestion appears to be based on a different view of our contributions with which we respectfully disagree. However  
24 we welcome and will definitely consider concrete suggestions about the title and organization.

#### 25 R3

26 **clarification on the long tail theory. Assume we define a more general measure...** ▷ We do not see how the  
27 proposed definition captures the intuitive notion of memorization since the value is large even if a single example  
28 out of the  $k$  that were removed is not fit by the model. More generally, inference based on sets of examples is what  
29 distinguishes the traditional view of learning from memorization.

30 **Near duplicated examples are dataset artifacts, not demonstration of long tail** ▷ Indeed very high-influence pairs  
31 usually come from artifacts of data collection. However a large fraction of high-influence pairs that have somewhat  
32 lower values (in the 0.15-0.3) range do **not** look like such artifacts. So while memorization may be unnaturally important  
33 for CIFAR and ImageNet due to these artifacts it would still be important without them.

34 **Extend the analysis in pp.7 to last 2, 3, or more layers?** ▷ Thank you for the interesting suggestions. We are definitely  
35 planning additional experiments related to this work (and hope that others will do them too).

36 **why  $m = 0.7n$  not  $0.5n$ ?** ▷ Larger  $m$  makes the value closer to the original leave-one-out estimator and better at  
37 estimating marginal utility since the models become closer to the one computed on the entire dataset.  $m=0.7n$  is both  
38 quite close in accuracy to full-dataset models (unlike  $0.5n$ ) and is sufficiently efficient (efficiency drops linearly as  
39 fraction approaches 1).

#### 40 R4

41 **The only concern I have is that the applications...** ▷ First, by far our main goal is understanding of memorization,  
42 a fundamental question about ML that has been puzzling the research community since the “Understanding Deep  
43 Learning ...” work Zhang et. al. The development of the influence estimator is just a potential bonus and thus we do  
44 not provide a detailed comparison with existing methods. In terms of efficiency note that our method simultaneously  
45 estimates the influence of all training examples on all datapoints. We are not aware of any method that can do that  
46 more efficiently and provide results of comparable quality. That said, we agree that efficiency is a concern for these  
47 applications. We believe that it is possible to develop more efficient estimators of comparable accuracy but leave it for  
48 future work. To stimulate this work we have already made the values of our estimator on CIFAR-100 and ImageNet  
49 publicly available.

50 **Randomness from mini-batch ordering [Toneva et al, ICLR2019]** ▷ The definition of influence/memorization  
51 contains expectation over the randomness of the algorithm. So our estimator measures expected memorization over all  
52 possible choices of minibatches. Also note that despite the use of a related “forgetting” word, the notion is completely  
53 unrelated to memorization that we study. We will clarify that in the related work section.