

1 We thank the reviewers for their positive comments and detailed feedback. Specific concerns are addressed below.

2 **R4: “Motivation unclear; compare to earlier metric learning work”:** In our related work (lines 97-109), we outline
3 the differences between modern contrastive learning methods (including ours) and the papers listed by R4. In summary:
4 our approach can be seen as a generalization of the prior approaches by using multiple positives and negatives per
5 anchor; the use of a temperature parameter which plays the role of a margin; and the use of multiple views of the data.
6 Each of these components contributes to the performance of the supervised contrastive loss, as shown in our ablation
7 experiments. It is well known that distance metric learning approaches such as triplet loss and the loss proposed in
8 [4] have slow convergence and require proper tuning of hard negatives (e.g. see [1], Section 3). Additionally, some of
9 these losses seem to require very small batch sizes to avoid falling into local minima (e.g., [3] cites a batch size of 16
10 for this exact reason). Our approach overcomes these issues as shown by our application to large scale problems such
11 as ImageNet. In Section 3.2.3, we show analytically the reduction to the triplet loss in Section 3.2.3. We will further
12 clarify these points in the final version.

13 **R4: “Comparison to metric learning methods”:** We tested N-pairs [2] in our framework with a batch size of 6144.
14 N-pairs achieves only 57.4% Top-1 on ImageNet (compare to our loss that achieves 78.7%). We believe this is due
15 to multiple factors missing from N-pairs loss compared to supervised contrastive: the use of multiple views; lower
16 temperature (e.g., see [5], Table 5) and many more positives. We show some results of the impact of the number of
17 positives per anchor in the Appendix (Sec. 6), and the N-pairs result is inline with them. We also note that the original
18 N-pairs paper [2] has already shown the outperformance of N-pairs loss to triplet loss. We will add this experiment and
19 relevant ablations to the paper.

20 **R1/R2/R3: “Concerns over compute cost”:** We agree with the reviewers that reducing the expense of contrastive
21 methods in general would be desirable. To this end, we experimented with memory based alternatives [1]. On ImageNet,
22 with a memory size of 8192 (requiring only the storage of 128-dimensional vectors), a batch size of 256, and SGD
23 optimizer, running on 8 Nvidia V100 GPUs, the supervised contrastive loss is able to achieve 79% top-1 accuracy on
24 ResNet-50 architecture. This is in fact slightly better than the 78.7% accuracy with 6144 batch size (and no memory);
25 and with significantly reduced compute and memory footprint. We will add this experiment to the paper. Finally, the
26 drop in performance by training for fewer epochs (e.g. 350 vs 700 epochs) is very small (difference of 0.1%). We will
27 add a sweep over training epochs to the paper.

28 **R3: “Experiment reporting is confusing”:** We agree with R3 and will adjust the table reporting to be more clear and
29 will highlight further SimCLR’s performance for transfer learning.

30 **R1/R3: “Improved baselines for cross-entropy; additional experiments”:** R3 requested that we show comparisons
31 of cross-entropy and supervised contrastive losses using the same optimizers. In Table 4 of the Supplementary
32 we include a sweep over optimizers used for supervised contrastive. In Table 1, we provide an equivalent sweep
33 over optimizers used for cross-entropy using the same network, augmentation strategy, and number of train epochs;
34 cross-entropy performs best with the momentum optimizer (the results show no improvement over those reported in
35 the paper). We additionally ran experiments to see if the larger effective batch size used for supervised contrastive
36 could explain its superior performance, as R1 suggested. We trained with cross-entropy loss using a batch size
37 of 12,288, but this only achieved 77.5% Top-1 accuracy. We also tried training for twice as many epochs (1400)
38 with the original batch size of 6,144, but this was even lower, at 77.0%. We will add these results to the final
39 version. Please note in Figure 1 and Table 3 of the paper, baseline numbers are taken directly from previous papers.
40

41 **R2: “Table 4 cross-entropy outperformed Supervised Contrastive; consider
42 a joint loss”:** Note that Table 4 is about transfer learning where we see that
43 SimCLR (unsupervised contrastive learning) actually outperforms both cross-
44 entropy and supervised contrastive slightly. We agree with R2 that a joint cross-
45 entropy and supervised contrastive loss could probably combine the benefits of
46 both and certainly make it simpler to use the supervised contrastive loss in practice.
47 We have experimented with this and seen strong results, although none that so far
48 outperform the two stage approach. We will add a discussion of this in the final
49 paper, but we feel that it’s important for the sake of future research to isolate the impact of the supervised contrastive
50 loss, so we chose to make the two-stage approach the focus of the paper.

Optimizer	Top-1 Accuracy
LARS	76.2%
RMSProp	76.2%
Momentum	78.2%

Table 1: Cross-entropy (ResNet-50) with different optimizers.

51 References

- 52 [1] Momentum Contrast for Unsupervised Visual Representation Learning, He et. al., 2019.
53 [2] Improved Deep Metric Learning with Multi-class N-pair Loss Objective, K. Sohn, 2016.
54 [3] Discriminative Feature Representation for Person Re-identification by Batch-contrastive Loss, Zhang et. al., 2018.
55 [4] Dimensionality reduction by learning an invariant mapping, Hadsell et. al., 2006.
56 [5] A Simple Framework for Contrastive Learning of Visual Representations, Chen et. al., 2020.