

1 We thank the reviewers for their very insightful and helpful comments and address their questions and remarks below.

2 **Relationship to Mean Teacher/Temporal Ensembling:** We are grateful to reviewers for pointing out this relevant  
3 related work that we missed in our original literature review. We will add the two citations in the related work section  
4 along with a paragraph on semi-supervised learning and connect MT and our current ablation study. Indeed, Table 5b  
5 (line 7: no predictor,  $\beta = 0$ , 0.2%) corresponds to using an MT-like approach in unsupervised learning (i.e., removing  
6 MT’s classification loss) and we show that this approach does collapse. BYOL’s novelty over MT is to perform well even  
7 without labels or classification loss, thanks to the addition of a predictor.

8 **Code and reproducibility:** We will release an open-source version of our full pretraining pipeline, the pretrained check-  
9 points, and the linear evaluation pipeline on ImageNet within the next two weeks. To further improve reproducibility  
10 and accessibility, we will also provide a single-GPU setup for pretraining on the smaller Imagenette dataset.

11 **Importance of a near-optimal predictor:** As the predictor is only applied to the online branch, its role is to make  
12 the architecture asymmetric rather than just making the network deeper. Table 5b already shows the importance of  
13 combining a predictor and a target network: the representation does collapse when either is removed. Furthermore,  
14 new experiments show that keeping the predictor near-optimal at all times is key to preventing collapse, which may be  
15 one of the roles of BYOL’s target network. We further found that we can remove the target network without collapse by  
16 making the predictor near-optimal, either by (i) using an optimal *linear* predictor (obtained by linear regression on the  
17 current batch) before back-propagating the error through the network (52.5% top-1 accuracy at 300 epochs), or (ii)  
18 increasing the learning rate of the predictor (66.5% top-1). By contrast, increasing the learning rates of both projector  
19 and predictor (without target network) yields poor results ( $\approx 25\%$  top-1).

20 **Explaining BYOL’s non-collapse:** Similarly to GANs, BYOL uses two sets of parameters that are not minimizing  
21 the same objective. Thus, there is no *a priori* reason for BYOL’s dynamics to converge to a global minimum of  
22  $\|\bar{q}_\theta(z_\theta) - \bar{z}'_\theta\|^2$ , as they are not following the gradient of this loss ( $\mathcal{L}_\theta^{\text{BYOL}}$  uses  $\bar{z}'_\theta$ ). While these dynamics still admit  
23 undesirable equilibria where all images are mapped to the same constant projection (e.g., all zeros), BYOL’s empirical  
24 performance seems to indicate that such equilibria may be unstable. We hypothesize that maintaining a near optimal  
25 predictor at all times is key to avoid collapsed solutions. When using an optimal predictor, BYOL minimizes the  
26 (expected) conditional variance of the target projection given the online projection. With a fixed target network, adding  
27 more information to the online projection can reduce this conditional variance, but cannot increase it. For example,  
28 training dynamics will always tend not to collapse features from the online network, as for any constant  $C$  and variables  
29  $X$  and  $Y$ ,  $\text{Var}(Y|X) \leq \text{Var}(Y|C)$ . More generally BYOL is encouraged to keep features from the online projection  
30 diverse by latching onto any source of variability  $Z$  (stemming, e.g., from noise in training dynamics) distinct from  
31 existing features, as  $\text{Var}(Y|X, Z) \leq \text{Var}(Y|X, h(X))$  for any variables  $X, Y, Z$  and any function  $h$ . We will add these  
32 additional discussions and experiments to our submission to clarify the role of the predictor.

33 **Note on Fig 3a/footnote 4 (batch size):** When dividing the batch size by  $N$ , we also average gradients for  $N$  steps.  
34 Without batch-norm (for BYOL), the two computations would be exactly equivalent. With batch-norm for BYOL, the two  
35 computations only differ in how the batch-norm statistics are computed.

36 **Note on robustness:** As described in Section 5, contrastive methods need to make the discrimination task challenging,  
37 which requires many negative examples (large batch size) and absence of uninformative features that are easy to  
38 discriminate (strong transformations). They stop learning once their prediction is sufficiently similar to the positives  
39 compared to the negatives. Instead BYOL does not rely on comparing positive and negative pairs and keeps latching on  
40 new information thanks to the predictor. It should therefore not be as sensitive either to batch-size or transformations.

41 **Answers to Reviewer 1:** We thank you for your positive comments, and we hope to have answered your questions in  
42 explaining BYOL’s non-collapse. As an optimal predictor seems sufficient to favor a diverse representation and stabilize  
43 the training, it removes the need for negative examples that are required to balance contrastive objectives.

44 **Answers to Reviewer 2:** We thank you for your in-depth questions and comments. Non-collapse and L.745-756: see  
45 discussion above. Table 1b vs 2b: Linear evaluation trains the linear layer on top of the representation with 100% of the  
46 labels, while semi-supervised learning only uses 1 or 10% of the dataset to finetune the full network (including the linear  
47 layer). This explains the difference in performance. Note that the same trend is observed in SimCLR and other related  
48 literature. L.224-228: these changes allow us to use the same hyperparameters as those optimized for SimCLR in their  
49 paper; these parameters are however not optimal at 1000 epochs, though the gap is low. L.241-251: the cosine similarity  
50 in BYOL is not between the same terms, and crucially, involves the output of the predictor. See also note on robustness  
51 above.  $\ell^2$  vs cosine: we wanted to emphasize that BYOL also works well with just an  $\ell^2$  loss (no normalization, as per  
52 table 19). Table 5: table 5 provides the trend of the target update, we further sweep on this hyperparameter in our  
53 main experiments. L.559: No, it is disjoint. Equation 5: eq. (5) is a generalization of the InfoNCE loss with added  
54 temperature parameter  $\alpha$ , in expanded form. We will add the derivation from the  $I_{NCE}$  equation (10) of Poole et al. in  
55 Appendix F.4., to clarify differences with our equation (5).

56 **Answers to Reviewer 3:** Thank you for your insights and your time. Reproducibility: see discussion above.  
57 Detection results: we are using the same setup as in MoCo’s Table 4 (fine-tune the representation on trainval2007  
58 only, not trainval2007+12), for which BYOL gains  $+2.6AP_{50}$  over MoCo. L34: see note on robustness above.